



McGRAW-HILL
ENCYCLOPEDIA OF
SCIENCE &
TECHNOLOGY

www.MHEST.com

6 **EBE-EYE**



Ebenales — Eyepiece

Ebenales

An order of flowering plants, division Magnoliophyta (Angiospermae) in the subclass Dilleniidae of the class Magnoliopsida (dicotyledons). The order consists of 5 families and about 1750 species, the Sapotaceae (about 800 species) and Ebenaceae (about 450 species) being the largest and most familiar families. The Ebenales are woody, chiefly tropical, sympetalous plants (those with flowers have the petals joined by their margins, at least toward the base, forming a basal tube, cup, or saucer) with usually twice as many stamens (including staminodes) as corolla lobes. When there is a single set of stamens, these are generally opposite the corolla lobes. Ordinarily there are only one or two ovules in each locule and the placentation is axile (that is, the placenta forms a central axis in an ovary which is divided by longitudinal partitions into two or more chambers). In a cross section of such an ovary, the partitions radiate from the central axis (placenta) to the lateral wall of the ovary. Chicle (from *Manilkara sapota*, in the Sapotaceae) and ebony (*Diospyros ebenum*) are obtained from members of the Ebenales. *See* CHICLE; DILLENIIDAE; EBONY; MAGNOLIOPHYTA; PLANT KINGDOM.

Arthur Cronquist

Ebola virus

Ebola viruses are a group of exotic viral agents that cause a severe hemorrhagic fever disease in humans and other primates. The four known subtypes or species of Ebola viruses are Zaire, Sudan, Reston, and Côte d'Ivoire (Ivory Coast), named for the geographic locations where these viruses were first determined to cause outbreaks of disease. Ebola viruses are very closely related to, but distinct from, Marburg viruses. Collectively, these pathogenic agents make up a family of viruses known as the Filoviridae.

Infectious agent. Filoviruses have an unusual morphology, with the virus particle, or virion, appearing as long thin rods. Virions have a diameter of 0.08 micrometer and a minimal length of 0.7–0.9 μm ; this length is comparable to that of a small bacterium. However, when these viruses are grown in cell cultures, they can form long filamentous and branched forms that reach 14 μm or more. This is especially true for Ebola viruses.

A filovirus virion is composed of a single species of ribonucleic acid (RNA) molecule that is bound together with special viral proteins, and this RNA-protein complex is surrounded by a membrane derived from the outer membrane of infected cells. Infectious virions are formed when the virus buds from the surface of infected cells and is released. Spiked structures on the surface of virions are formed by three molecules of a single glycoprotein and are firmly embedded in the virion membrane. These structures project from the virion and serve to recognize and attach to specific receptor molecules on the surface of susceptible cells. This recognition and binding allows the virion to penetrate into the cell. Then, with the virion free to operate within the cytoplasm, the genetic information contained in the RNA molecule directs production of new virus particles by using the cellular machinery to drive synthesis of new viral proteins and RNA. *See* RIBONUCLEIC ACID (RNA); VIRUS.

Pathogenesis. Although much is known about the agents of Ebola hemorrhagic fever disease, the ecology of Ebola viruses remains a mystery. The natural hosts of filoviruses remain unknown, and there has been little progress at unraveling the events leading to outbreaks or identifying sources of filoviruses in the wild. Fortunately, the incidence of human disease is relatively rare and has been limited to persons living in equatorial Africa or working with the infectious viruses. However, in this time of rapid transportation over large distances, the threat of agents such as Ebola viruses spreading to remote areas is

taken very seriously by public health professionals. People infected with the virus are not as contagious as, say, persons suffering from a cold or measles, and the virus is spread primarily through close contact with the body of an infected individual, his or her body fluids, or some other source of infectious material.

Diagnosis. Ebola virus hemorrhagic fever disease in humans begins with an incubation period of 4–10 days, which is followed by abrupt onset of illness. Fever, headache, weakness, and other flulike symptoms lead to a rapid deterioration in the condition of the individual. In severe cases, bleeding and the appearance of small red spots or rashes over the body indicate that the disease has affected the integrity of the circulatory system. Contrary to popular belief, individuals with Ebola virus infections do not melt or have their organs dissolve, but die as a result of a shock syndrome that usually occurs 6–9 days after the onset of symptoms. This shock is due to the inability to control vascular functions and the massive injury to body tissues.

Ebola viruses can be found in high concentrations in tissues throughout the body and are especially evident in the liver, spleen, and skin. From studies of human cases and experimentally infected animals, it appears that the immune response is impaired and that a strong cellular immune response is key to surviving infections. This immunosuppression may also be a factor in death, especially if secondary infections by normal bacterial flora ensue. No severe human disease has been associated with Reston virus infections, and it appears that this virus may be much less pathogenic for humans than it is for monkeys. *See IMMUNOSUPPRESSION.*

Epidemiology. The first described cases of Ebola virus disease occurred in 1976, when simultaneous outbreaks of two distinct subtypes occurred in northern Zaire and southern Sudan. Many of the human infections occurred in local hospitals, where close contact with fatal cases, reuse of contaminated needles, and a low standard of medical care led to most fatalities. An Asian form of Ebola virus was discovered in a type of macaque, the cynomolgus, exported from the Philippines to the United States in late 1989. It was determined that a single monkey breeding facility in the Philippines was the source of the virus. Named Reston virus after the Virginia city in which infected monkeys were first identified, it represented a new form of Ebola virus. A repeat of this incident took place in early 1992, when monkeys from the same breeding facility were shipped to Siena, Italy.

In 1994, another new species of Ebola virus was identified, associated with chimpanzee deaths in the west African country of Côte d'Ivoire. This episode signaled the reemergence of Ebola virus in Africa, which had not been seen since 1979 when an outbreak of the Sudan subtype occurred. A large outbreak of disease was identified in early May 1995 in and around the city of Kikwit, Zaire. This outbreak of a Zaire subtype was a repeat of the 1976 episode, but this time 600 mi (1000 km) to the south. Another

Ebola virus outbreak took place in 1995 in Gabon and resulted in human fatalities, but the details were slow to be revealed. This virus has been isolated, and preliminary reports indicate that it may be a variant of the Zaire subtype. This virus may have also been responsible for human fatalities in Gabon reported to have occurred in early 1996. Soon after this outbreak, yet another Reston virus outbreak occurred when infected monkeys were sent to a United States holding facility in Alice, Texas. That introduction was quickly detected and stopped by quarantine and testing measures implemented after the first Reston virus episode in 1989.

Control and prevention. Outbreaks of Ebola virus disease in humans are controlled by the identification and isolation of infected individuals, implementation of barrier nursing techniques, and rapid disinfection of contaminated material. Diagnosis of Ebola virus cases is made by detecting virus proteins or RNA in blood or tissue specimens, or by detecting antibodies to the virus in the blood. Such testing is important in tracking and controlling the movement of the virus during an outbreak.

Dilute hypochlorite solutions (bleach), 3% phenolic solutions, or simple detergents (laundry or dish soap) can be used to destroy infectious virions. No known drugs have been shown to be effective in treating Ebola virus (or Marburg virus) infections, and protective vaccines against filoviruses have not been developed. However, research is being directed at developing such treatment and should provide insights into the disease mechanisms used by filoviruses.

Evolution. Filoviruses have been shown to be genetically related to two other virus families, the Paramyxoviridae (including measles virus and mumps virus) and the Rhabdoviridae (including rabies virus and vesicular stomatitis virus). These viruses evolved from a common progenitor in the very distant past. These viruses evolved into very distinct lineages, yet have maintained a similar approach to reproducing.

The evolutionary profile, or phylogeny, for Ebola viruses and Marburg viruses has been determined. Ebola viruses have evolved into a separate and very distinct lineage within the filovirus family, and the individual species of Ebola viruses have also evolved into their own sublineages. Filoviruses show a great deal of diversity, indicating that they have likely co-evolved with their natural hosts over a long period of time. Ebola viruses within a given subtype did not appear to show a great deal of change over periods of many years. This genetic stability or stasis in the wild suggests that these viruses have evolved to occupy very specific ecological niches. Anthony Sanchez

Bibliography. Centers for Disease Control, Ebola-Reston virus infection among quarantined nonhuman primates—Texas, 1996, *MMWR*, 45:314–316, 1996; B. N. Fields et al. (eds.), *Fields Virology*, 3d ed., 1996; L. G. Horowitz, *Emerging Viruses: AIDS and Ebola—Nature, Accident or Intentional*, 1996; A. Sanchez et al., Reemergence of Ebola virus in Africa, *Emerg. Infect. Dis.*, 1:96–97, 1995;

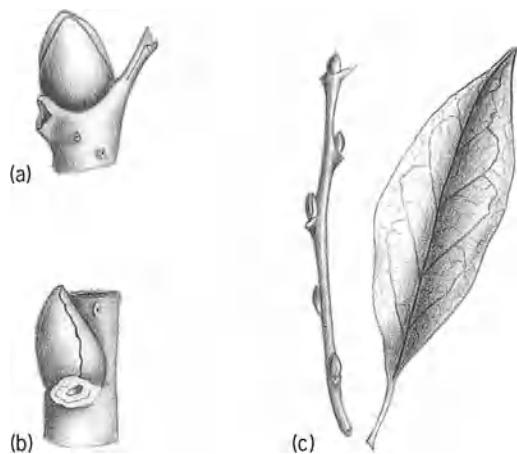
A. Sanchez et al., The virion glycoproteins of Ebola viruses are encoded in two reading frames and are expressed through transcriptional editing, *Proc. Nat. Acad. Sci. USA*, 93:3602-3607, 1996.

Ebony

A genus, *Diospyros*, of the ebony family, containing more than 250 species. Some species are important for succulent fruits, such as date plum, kaki plum, and persimmon, and several for timber, particularly the heartwood, the true ebony of commerce.

Although it is popularly supposed to be a black wood, most species have a heartwood only streaked and mottled with black. The heartwood is brittle, breaks with a conchoidal fracture, and is difficult to work, but it has long been in demand. The sapwood is white, becoming bluish or reddish when cut.

The use of ebony can be traced to the early Egyptians, who probably obtained it in Ethiopia. Ebony from India was known to the Greeks before 350 B.C. At present black ebony is used for knife handles, piano keys, finger boards of violins, hairbrush backs, inlays, and marquetry. Some of the woods called ebony, however, belong to different families, especially the pulse family, Leguminosae.



Persimmon (*Diospyros virginiana*). (a) Terminal bud. (b) Lateral bud. (c) Twig and leaf.

Persimmon (*D. virginiana*), of the southeastern United States, is one of numerous tropical or subtropical species (see **illus.**). Usually a medium-sized tree with black chunky bark, it attains 100-125 ft (30-37 m) in height with a trunk 20-30 in. (50-75 cm) in diameter. The sapwood is in demand for the manufacture of weaving shuttles and heads of golf clubs. The fruit is sweet and edible when slightly overripe, but when immature it is extremely pungent. The species in tropical America are too small or rare to be of economic value, although several of them have black heartwood used locally for making walking sticks, inlays, and miscellaneous articles of turnery and carving. See FOREST AND FORESTRY; TREE. Arthur H. Graves; Kenneth P. Davis

Ebriida

An order of flagellate Protozoa, subphylum Sarcocystophora, class Phytomastigophorea. Two genera, *Ebria* and *Hermesinum*, remain of the once numerous order, as determined by fossil remains. They possess an internal solid siliceous skeleton forming a shallow flattened or slightly arched structure that is enclosed by clear cytoplasm. Skeleton form and structure are distinctive. These phytoflagellates do not have chromatophores. A conspicuous nucleus lies anteriorly and two long and fine flagella emerge near it. Reserve materials apparently are fat. See CILIA AND FLAGELLA.

The organisms are common in the inshore waters of the Atlantic Ocean and the Gulf of Mexico. However, numbers sufficient to influence the general ecology have not been observed, nor has any ingestion of particulate matter. Their only occurrence seems to be marine, although somewhat lowered salinities are certainly tolerated. See PHYTOMASTIGOPHOREA; PROTOZOA. James B. Lackey

Ecdysone

The molting hormone of insects. It was isolated in crystalline form from the common silkworm (*Bombyx mori*) by A. Butenandt and P. Karlson in 1954, and its chemical structure was elucidated in 1965. It is a derivative of cholesterol and its structure can be described as $2\beta,3\beta,14\alpha,22R,25$ -penta-hydroxy- Δ^7 -5 β -cholestene-6-one. Shortly thereafter the molting hormone of crustaceans (crustecdysone) was isolated and identified as 20-hydroxy-ecdysone (**Fig. 1**). The sample compound was identified as one of the active compounds in *B. mori* extracts by P. Hocks and coworkers, and H. Hoffmeister and coworkers proposed the name ecdysterone. A group of Japanese chemists, working on the isolation of steroids from plants, came across some similarities of the spectra of 20-hydroxy-ecdysone and their compounds inokosterone and isoinokosterone; the latter was identical with 20-hydroxy-ecdysone. Inokosterone carries the last hydroxy group in position 26 instead of 25, which reduces the biological activity by a factor of 10. Several other structurally related and biologically active steroids have been isolated from plant sources, including ecdysone, but nothing is known about a possible biological function of these steroids in plants. See STEROID.

Ecdysone is synthesized in insects from cholesterol; 7-dehydro-cholesterol may be an intermediate. It has been shown that ecdysone is also readily metabolized and inactivated by an enzyme system; the inactivation occurs in the living insect, as well as in homogenates. The activity of the inactivating system varies markedly during development. The most striking physiological activity of ecdysone is the induction of puffs (zones of gene activity) in giant chromosomes of the salivary glands and other organs of the midge *Chironomus*. It has been shown that 20-hydroxy-ecdysone also induces puffs. The

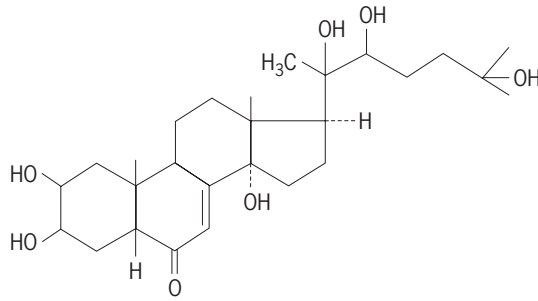


Fig. 1. Structure of 20-hydroxy-ecdysone.

induction of puffs has been visualized as the primary action of the hormone, indicating that ecdysone controls the activity of specific genes. The postulated mechanism for the action of ecdysone is shown in Fig 2; the primary effect is gene activation. (In the figure DNA is deoxyribonucleic acid, ATP is adenosine triphosphate, GTP is guanosine triphosphate, CTP is cytosine triphosphate, and UTP is uridine triphosphate.) The active gene produces messenger ribonucleic acid (mRNA), which in turn combines with ribosomes and becomes involved with protein

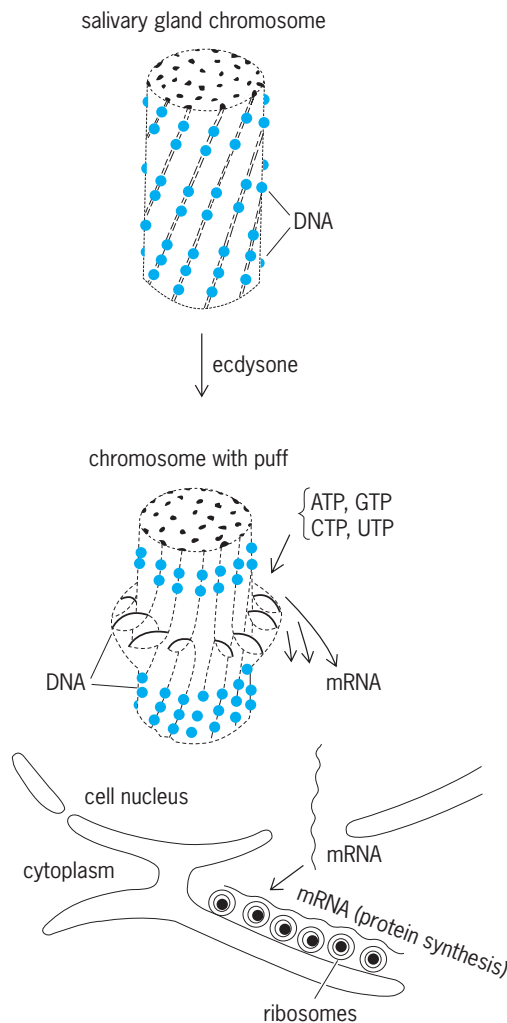


Fig. 2. Drawing of postulated mechanism of action of ecdysone in protein production in the midgut *Chironomus*.

synthesis. These proteins may be very specific and may be enzyme molecules. It has been shown that ecdysone stimulates the synthesis of mRNA, among which is the messenger for dopa decarboxylase. This enzyme is involved in the biosynthesis of the sclerotizing agent *N*-acetyl-dopamine. See CHROMOSOME; INSECT PHYSIOLOGY. Peter Karlson

Bibliography. F. Schwalm, *Insect Morphogenesis*, 1987; F. Taylor and R. Karban (eds.), *The Evolution of Insect Life Cycles*, 1986; A. Zaslavsky, *Insect Development*, 1988.

Ecdysozoa

A major division within the animal kingdom containing the majority of animal species and comprising the phyla Priapulida, Kinorhyncha, and Loricifera (no common names), Arthropoda (insects, crustaceans, spiders, etc.), Onychophora (velvet worms), Tardigrada (water bears), and Nematoda (roundworms). These forms include animals with disparate morphologies, but their close relationship is suggested by molecular DNA sequence comparisons—especially of both large and small subunit RNA molecules—and is consistent with myosin heavy-chain sequence comparisons and with a variety of uniquely shared genomic and morphological features. Ecdysozoans were named after their habit of molting (ecdysis), unique among invertebrates, and presumably all are descended from a common molting ancestor. See ANIMAL KINGDOM; MOLTING (ARTHROPODA).

General characteristics. Ecdysozoa is the major subdivision of the Protostomia (see illustration). In most protostomes the adult mouth forms on the site of the blastopore, and the secondary body cavity (the coelom) forms by splitting between tissue layers of the mesoderm. None of the ecdysozoan phyla have a eucoelom—a secondary body cavity that contains the major body organs and acts as a fluid skeleton—but may have smaller coelomic spaces, particularly during development, that function as ducts or contain specialized organs. Fluid skeletons are furnished by primary body cavities that either are pseudocoels (in the absence of a blood circulatory system) or, in arthropods and onychophorans, are hemo-coels (contain blood). In this feature they differ from most deuterostomes and lophotrochozoans (see illustration). The cuticle is largely flexible in some forms (priapulids, onychophorans) but is moderately to strongly sclerotized in others and even mineralized in some arthropods, and acts to stiffen the body wall and to serve as a site for muscle attachments, forming an exoskeleton in arthropods. Growth is accommodated by periodic molting and then secretion of a larger cuticle. None of the marine ecdysozoans has primary planktonic (free-floating) larvae, although in some arthropods several early molting stages are planktonic and serve for dispersal. See COELOM.

Priapulida, Kinorhyncha, and Loricifera. These phyla are pseudocoelomate worms, and all have minute

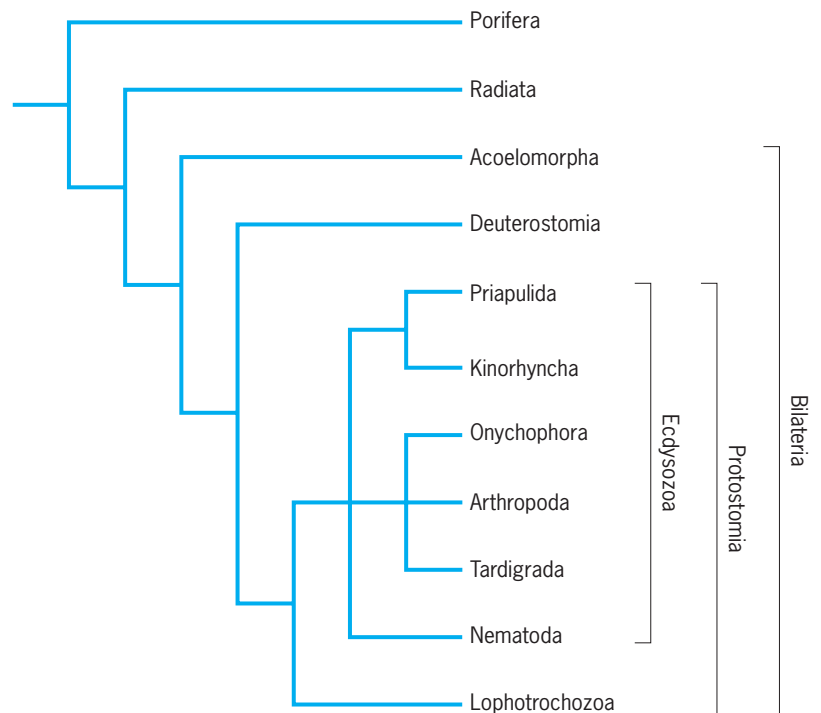
species that live interstitially in marine sediments; loriciferans are attached to sand grains, while the others are free-living. The cuticles bear secreted spines or plates that are sensory or aid in burrowing. Some priapulids also grow to much larger sizes [to 20 cm (8 in.)] and form prominent burrow structures. Kinorhynchs have a segmented cuticular exoskeleton and longitudinal muscle system, but the other two phyla are unsegmented. *See* KINORHYNCHA; LORICIFERA; PRIAPULIDA.

Arthropoda, Onychophora, and Tardigrada. These phyla are more closely related to each other than to other ecdysozoans on molecular evidence. All three are segmented to some extent and possess limbs. In onychophorans and tardigrades, the cuticle is rather flexible, and the limbs (lobopods) contain pseudocoelomic cavities that provide turgor. In arthropods, however, the cuticle is stiffened to form an exoskeleton, and the limbs act as levers operated by intrinsic muscles. All these phyla were originally marine but have their chief diversity in terrestrial habitats today; the onychophora are now entirely nonmarine. *See* ARTHROPODA; ONYCHOPHORA; TARDIGRADA.

Nematoda. This phylum consists of small, roundish worms that are unsegmented and pseudocoelomate, though minute species may essentially lack any body cavity. Many nematodes are parasitic, and one highly derived parasitic group, the Nematomorpha (horsehair worms), is sometimes considered as a phylum of its own. *See* NEMATODA (NEMATODA); NEMATOMORPHA.

Evolution: diversity. The earliest fossils that are very likely to be ecdysozoans are Early Cambrian trails and feeding marks left by scratching limbs, probably dating 525–520 Ma (million years ago), although burrows from as early as about 540 Ma could have been formed by priapulids. Body fossils of ecdysozoans appear around 530–525 Ma, represented then by stem groups (extinct branches) of priapulids, arthropods, and onychophorans, and by extinct arthropod allies that are segmented but evidently lack jointed appendages; some of these are quite large by Cambrian standards, up to 50 cm (20 in.) long. There is also an extinct Early Cambrian vermiform group, the paleoscoleoids, which may be pseudocoelomate ecdysozoans. Tardigrades are first known from about 510–500 Ma, while the other phyla do not have body fossil records.

Ecdysozoans dominate the early animal fossil record, with a great disparity of form (mostly extinct variants of arthropod architectures) and in terms of species richness. During the Cambrian (543–500 Ma), over 70% of known animal species are trilobites, an extinct arthropod class. Ecdysozoans continued to be prominent in the seas until today, where they are now challenged for dominance only by Mollusca among invertebrate phyla. The earliest indications of animal fossils from terrestrial environments are traces, burrows in soils deposited near 410 Ma; body fossil assemblages appearing near 400 Ma and slightly later are almost entirely of arthropods, including millipedes, springtails, and rarely insects; they are dominated by detritus feeders and preda-



Relationships among phylum-level taxa within the Ecdysozoa, and the relation of the Ecdysozoa to other major animal groups.

tors. Today, arthropod species richness, especially of insects, eclipses that of all other animal phyla in terrestrial habitats. For insects, the small body sizes and modular body plans within a secreted, segmented exoskeleton and with jointed appendages that can be specialized individually—together with the habitat heterogeneity afforded by the rise of diverse global plant communities—have produced unique opportunities for evolutionary radiations, which seem to have been fully exploited. *See* FOSSIL; TRILOBITA.

Evolution: genomics. The two best-known model organisms for evolutionary genetic and developmental studies, the fly *Drosophila* (Arthropoda) and the roundworm *Caenorhabditis* (Nematoda), are ecdysozoans. Understanding of the genetic processes underpinning body plan development has been worked out largely in these forms, which have both segmented and unsegmented architectures, respectively. As the key genes regulating these processes are widely distributed among other animal groups as well, including humans, comparative studies using the ecdysozoan models are helping to clarify the evolution of animal genomes during the branchings of the tree of life. *See* ANIMAL EVOLUTION; ORGANIC EVOLUTION.

James W. Valentine
Bibliography. R. C. Brusca and G. J. Brusca, *Invertebrates*, Sinauer Associates, Sunderland, MA, 2003; J. Cracraft and M. J. Donoghue (eds.), *Assembling the Tree of Life*, Oxford University Press, 2004; R. A. Fortey and R. H. Thomas (eds.), *Arthropod Relationships*, Chapman and Hall, 1998; M. W. Kirschner and J. C. Gerhart, *The Plausibility of Life*, Yale University Press, 2005; J. W. Valentine, *On the Origin of Phyla*, University of Chicago Press, 2004.

Echinacea

A superorder of Euechinoidea, having a rigid test, the periproct within the apical system, keeled teeth, a complete perignathic girdle, and branchial slits. J. Durham and R. Melville (1957) include five orders in this group. There were formerly distributed among the Stirodonta and Camarodonta in the classification of R. Jackson (1912). See ECHINODERMATA; ECHINOIDA; EUECHINOIDEA; HEMICIDAROIDA; PHYMOSOMATOIDA; SALENIOIDA; TEMNOPLUROIDA.

Howard B. Fell

Echinodermata

A phylum of exclusively marine animals with a peculiar body architecture dominated by a five-part radial symmetry. Echinodermata [from the Latin *echinus* (spine) + *dermis* (skin: "spiny skins")] include the sea stars, sea urchins, and related animals. The body wall contains an endoskeleton of numerous plates (ossicles) composed of calcium carbonate in the form of calcite and frequently supporting spines. The plates may be tightly interlocked or loosely associated. The spines may protrude through the outer epithelium and are often used for defense. The skeletal plates of the body wall, together with their closely associated connective tissues and muscles, form a tough and sometimes rigid test (hard shell) which encloses the large coelom. A unique water-vascular system is involved in locomotion, respiration, food gathering, and sensory perception. This system is evident outside the body as five rows of fluid-filled tube feet. Within the body wall lie the ducts and fluid reservoirs necessary to protract and retract the tube feet by hydrostatic pressure. The nervous system of these headless animals arises from the embryonic ectoderm and consists of a ring around the mouth with connecting nerve cords associated with the rows of tube feet. There may also be diffuse nerve plexuses, with light-sensing organs, lying below the outer epithelium. The coelom houses the alimentary canal and associated organs and, in most groups, the reproductive organs. The body may be essentially star-shaped or globoid. The five rows of tube feet define areas known as ambulacra, ambcs, or radii; areas of the body between the rows of tube feet are interambulacra, interambcs, or interradii.

The larvae are usually planktonic (free-floating) with a bilateral symmetry, but the adults are usually sedentary and benthic (bottom-dwelling in a marine environment). They inhabit all oceans, ranging from the shores to the greatest ocean depths.

The phylum comprises about 7000 existing species. The Echinodermata have a good fossil record with about 13,000 fossil species. They first appeared in the Early Cambrian and have been evolving, since pre-Cambrian time, for well over 600 million years. During this time several differing body plans have arisen. The surviving groups show few resemblances to the original stock. The existing representatives fall into five, possibly six, classes: Crinoidea (sea

lilies and feather stars); Asteroidea (sea stars); Ophiuroidea (brittle stars); Echinoidea (sea urchins, sand dollars, and heart urchins); and Holothuroidea (sea cucumbers). The recently described class Concentricycloidea (sea daisies) has been referred to the Asteroidea by some experts and retained as a distinct class by others.

The outline of classification for the phylum is shown below. Classes with living representatives are elaborated to the level of subclass or order.

- Phylum Echinodermata
- “Homalozoans”
- Class: Ctenocystoidea
 - Stylophora
 - Homostelea
 - Homoiostelea
- “Crinozoans”
- Class: Eocrinoidea
 - Paracrinoidea
 - Rhombifera
 - Diploporita
 - Crinoidea
 - Blastoidea
 - Coronoidea
 - Parablastoidea
 - Edrioblastoidea
- Subclass: Inadunata
 - Camerata
 - Flexibilia
 - Articulata
- Order: Millericrinida
 - Cyrtocrinida
 - Bourgetocrinida
 - Isocrinida
 - Uintacrinida
 - Roveacrinida
 - Comatulida
- “Asterozoans”
- Class: Somasteroidea
 - Asteroidea
 - Order: Platyasterida
 - Trichasteropsida
 - Paxillosida
 - Notomyotida
 - Spinulosida
 - Valvatida
 - Velatida
 - Forcipulatida
- Class: Concentricycloidea
 - Order: Peripodida
- Class: Ophiuroidea
 - Order: Stenurida
 - Oegophiurida
 - Phrynophiurida
 - Ophiurida
- “Echinozoans”
- Class: Helicoplacoidea
 - Camptostromatoidea
 - Edrioasteroidea
 - Cyclocystoidea
 - Ophiocystoidea
 - Echinoidea

- Subclass: Perischoechinoidea
 - Order: Bothriocidaroida
 - Echinocystitoida
- Subclass: Cidaroida
 - Order: Cidaroida
- Subclass: Euechinoidea
 - Order: Diadematoidea
 - Echinothurioida
 - Pedinoida
 - Arbacioida
 - Echinoida
 - Phymosomatoidea
 - Salenioida
 - Temnopleuroidea
 - Clypeasteroida
 - Holactypoida
 - Cassiduloida
 - Spatangoida
- Class: Holothuroidea
 - Order: Dendrochirotida
 - Dactylochirotida
 - Aspidochirotida
 - Elasipodiida
 - Molpadiida
 - Apodida

Morphology

Despite displaying an astonishing variety of body forms, existing echinoderms share common anatomical features: symmetry, body wall, skeleton, nervous system, coelom, alimentary system, reproductive organs, and water-vascular system.

Symmetry. Adult echinoderms usually show pentamerism, a five-part radial symmetry in which there are normally five axes radiating out from a central point in the body. The mouth may lie on the underside, and the anus may lie on top of the animal. There are numerous exceptions: in the crinoids the mouth and anus lie on the same side, and in the holothuroids the mouth and anus lie at opposite ends of a cylindrical body. The side of the animal bearing the mouth is termed the oral side, and the side away from the mouth is the aboral side. The relationship of the mouth with the substratum was at one time considered important so that those echinoderms which were essentially sessile, with a mouth facing away from the substratum, were classified as subphylum *Pelmatozoa*, while those mobile forms in which the mouth faced the substrate were classified as *Eleutherozoa*. These differing life habits have little bearing upon formal echinoderm classification, but the terms “pelmatozoan” and “eleutherozoan” are useful descriptors of lifestyles in the phylum.

It has been suggested that five-part symmetry in echinoderms confers strength upon the skeleton supporting the body wall. Pentamerism was not an original feature of the echinoderms. Certain early groups such as the carpoidea did not show it, and modern echinoderm larvae have bilateral symmetry.

Body wall. The body wall is cloaked in a thin epidermis, a layer of epithelial cells. The external layer is often ciliated, and the electron microscope has

shown that these cells have an outer membrane that bears a great number of minute tubular extensions like microvilli. This arrangement of the epidermis is similar in appearance to the absorptive epithelium of many intestinal tissues. Below the epidermis, the dermis contains the skeletal elements. Internal to the dermis are the body-wall muscles and connective tissue. In the crinoids, ophiuroids, some asteroids, and some echinoids, the skeletal elements make up a great part of the body-wall volume. In the more flexible asteroids, the spherical echinoids, and the holothuroids, the volume of these elements is reduced. The inner surface of the body wall is also lined with epithelium and borders the coelom.

Skeleton. The skeleton is unique and consistent in all groups. Almost every plate, spine, or ossicle is composed of a single calcite crystal. During development each plate grows as a result of calcite secretion by a group of cells. The plate grows like a three-dimensional lattice (stereom) to produce a reticulate (having or resembling a network of fibers or lines) ossicle, with the tissue that secretes and maintains it (stroma) occupying the spaces within (Fig. 1).

Sensory and neuromuscular system. Echinoderms have simple sensory systems. The nervous system is ectodermal in origin and comprises a radial nerve cord lying along each ambulacrum. Each radial nerve connects with the circumesophageal nerve ring, and it is believed that interradian coordination takes place via this ring. The absence of a head means that sensory receptors and large areas of integrative nervous elements are not aggregated in one place. The nervous system is essentially diffuse, although specialized sense organs, such as the recently discovered “eyes” of certain brittle stars, occur. It appears that the circumoral ring and radial nerve cords take

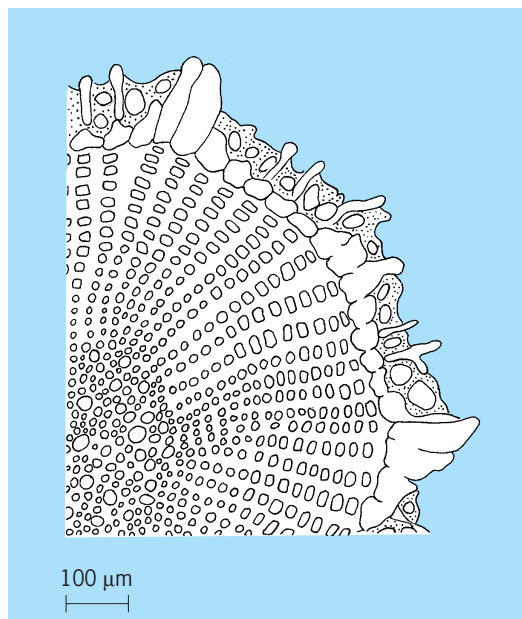


Fig. 1. Cross section through the spine of an echinoid showing the microscopic structure of the skeleton. The stereom forms a continuous mesh, and the stroma that secretes it lies in the interspaces which, in this case, are arranged in radial rows.

responsibility for gross activities such as posture, coordinated locomotion, and food collecting, but in the asteroids and the echinoids there are well-developed basiepithelial nerve networks that coordinate the numerous test appendages. Echinoderm muscles are principally smooth, but striated muscles have been detected in a few instances. See NERVOUS SYSTEM (INVERTEBRATE).

Coelom. The coelom is extensive and usually arises by enterocoely, or pouching, from the gut; but when there is direct development, with larval stages greatly reduced or absent, it may be schizocoelic, involving splitting in the mesoderm. The perivisceral coelom surrounds the viscera. During development, part of the coelom forms the water-vascular system, so the spaces within the tube feet and vessels are coelomic. See COELOM.

Alimentary system. In general the mouth is situated on the same side as the tube feet. From it arises the alimentary canal, which runs through the body to the anus (absent in the ophiuroids and some asteroids). The gut is lined by endoderm, and in some groups digestive caeca (blind pouches at the beginning of the large intestine) increase the area for absorption.

Particulate food is collected by the tube feet in many crinoids, ophiuroids, and holothuroids. Most asteroids are carnivores and engulf the prey, or evert part of the stomach over it, while some echinoids have well-developed chewing teeth for rasping at encrusting plants and animals. Many holothuroids, some echinoids, and a few asteroids ingest mud or sand, obtaining sustenance from associated organic material such as bacteria and diatoms. See DIGESTION (INVERTEBRATE).

Reproductive system. The reproductive organs lie within the coelom in the interradial position, except in the crinoids in which they arise on the arms. In most echinoderms there are five compact gonads, but the holothuroids have just one, and certain echinoids may have two, three, or four. The sexes are usually separate. The gonads discharge gametes by short ducts into the surrounding seawater. In some echinoderms, especially in polar regions, fertilized gametes are retained, and the young develop in special sacs in the body wall.

Water-vascular system. This is a multifunctional fluid-filled coelomic system that probably evolved first as a respiratory system but later took on increasingly important roles in food collection and locomotion (Fig. 2). It comprises a number of protrusible, hollow, tentacle-like tube feet, arranged externally in rows along each radial area of the body. These areas, composed of single to multiple rows of tube feet, are termed ambulacra. A small modified test ossicle, the madreporite, perforated by many irregularly shaped pores, connects with the water-vascular system via a calcified stone canal. It was thought that the water-vascular system opened directly to the surrounding seawater via the madreporite so that extra fluid could be drawn in to replenish the water-vascular fluid. The fluid volume of the water-vascular system may be more constant than was previously believed, and rel-

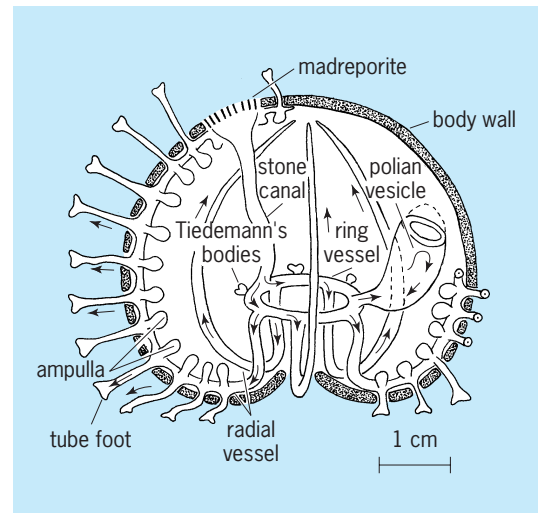


Fig. 2. Diagram of the water-vascular system in an echinoid. Arrows show direction of fluid flow.

atively little fluid passes in or out of the madreporite. Each tube foot is associated with a small bulbous reservoir inside the body. The feet are extended hydrostatically, and retracted by longitudinal muscles in their walls.

The form and the activity of tube feet vary greatly between the groups of echinoderms. In the crinoids they are fine, tapering structures suited for food collection and having no locomotory role. In the brittle stars they lack terminal suckers and are more important for food gathering than locomotion. In some asteroids, particularly the burrowers, they are tapering and pointed, while in others they are suckered and capable of adhering to the substratum. This is also the case in regular echinoids (which have no front or back end and can move in any direction), but in irregular echinoids (which have a definite front and back and do move in a particular direction) the tube feet are modified to suit the burrowing way of life. In holothuroids the tube feet surrounding the mouth are modified as tentacles of varying form and complexity and are used for food collecting. The rest may be present as locomotor or respiratory organs. In a few groups, such as the apodous and molpadiid holothurians, tube feet are missing from the body wall. When tube feet serve an important locomotor function, postural muscles govern the stepping movements.

Each tube foot is under the control of the radial nerve cord and is connected by a short lateral branch to the radial water-vascular canal as well as to the individual ampullae. All the radial canals are linked by the circumoral ring canal so that water-vascular fluid may be withdrawn from one part of the system to be supplied to another.

Saclike structures are associated with the ring canal. These are polian vesicles, which act as reservoirs for water-vascular fluid, and Tiedemann's bodies, which act as glands in which wandering coelomocytes (coelom corpuscles) are formed.

Embryology

Most echinoderms have an indirect development with a prolonged, food-gathering larval stage (Fig. 3). The larva feeds, usually on microscopic phytoplankton, using a ciliary mechanism. Two well-marked larval types occur: the pluteus group, with long-armed, bilaterally symmetrical, easel-shaped forms, common to ophiuroids and echinoids; and the auricularia group, barrel-shaped forms with a winding ciliated band which may be produced into lobes. The latter group is common to asteroids and holothuroids. In most asteroids the auricularia stage is followed by a similar but more complex larva, the bipinnaria, or sometimes also by an anchored final larval stage, the brachiolaria. Surviving crinoids have essentially a direct development, sometimes with a simple yolky larva, a vitellaria, which does not feed. This is also

the case for several members of other extant echinoderm groups. See INVERTEBRATE EMBRYOLOGY.

Phylogeny

The auricularia larva presents close and striking resemblances to the tornaria larva of some enteropneusts, and the enterocoelous development parallels that in primitive chordates. Hence echinoderms and chordates have long been regarded as related. However, the significance of similarities in the larvae of echinoderms and protochordates must be viewed in the context of molecular, morphological, and paleontological research. It seems likely that the similarities between the larvae of ophiuroids and echinoids, and asteroids and holothuroids, are due to convergent evolution and not to common evolutionary origins. The results of paleontology, molecular

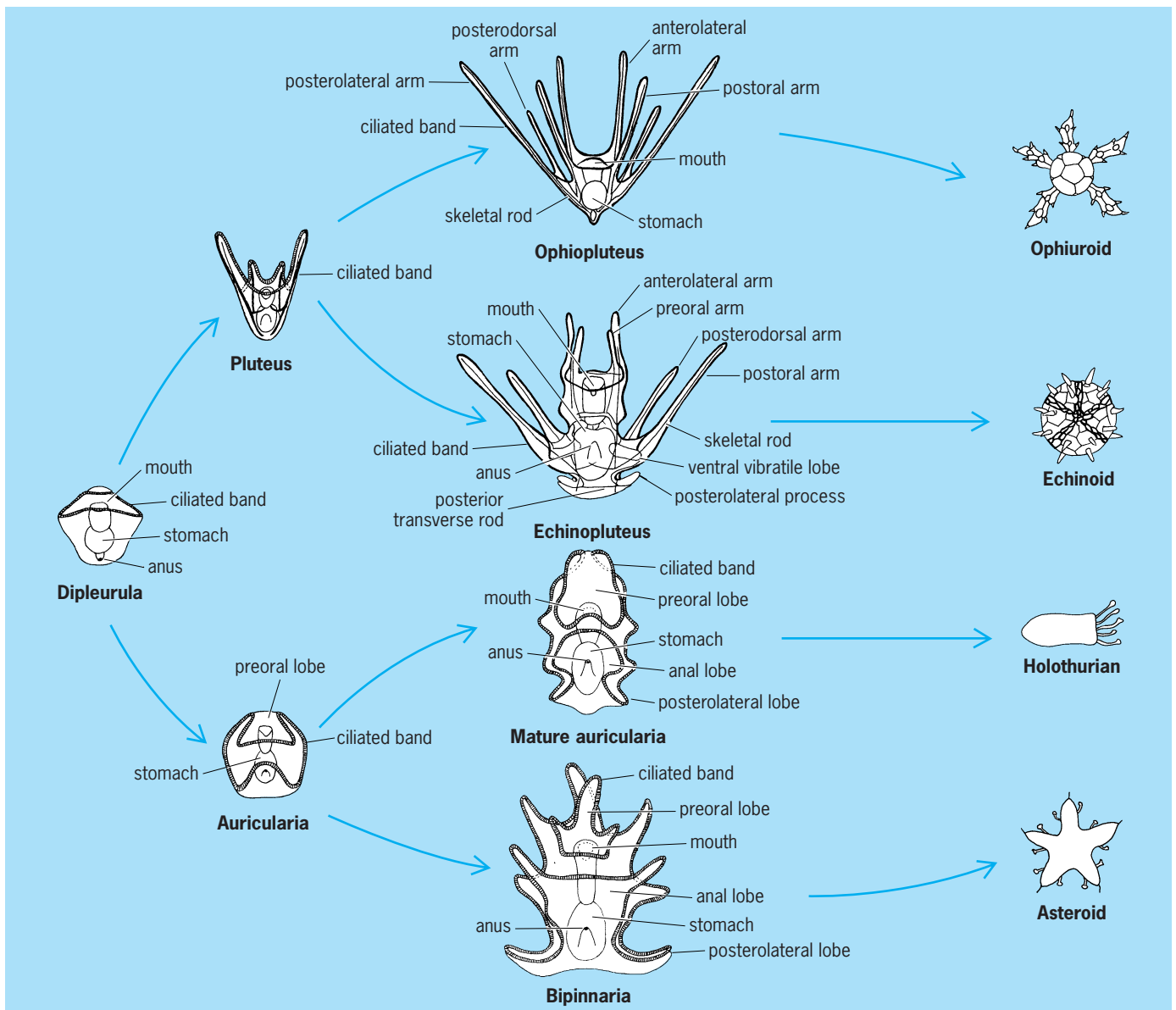


Fig. 3. General scheme indicating relationships of bilaterally symmetrical echinoderm larvae. (After H. B. Fell)

studies, and morphology indicate that ophiuroids and asteroids are closely related. Therefore it follows that within the phylum larval similarities do not indicate phylogenetic affinities. It is inadvisable to try to extrapolate beyond the phylum, so as to infer phylogenetic affinity between hemichordates and echinoderms solely on the ground that the auricularia resembles the tornaria. E. Marcus expressed the opinion that indirect development, with possession of pelagic larvae, must be prototypical for echinoderms and protochordates, and the asteroids and ophiuroids must be closely related, and therefore broad phylogenetic conclusions cannot be drawn on the basis of their larvae. L. Hyman grouped the extant echinoderms as their larval similarities suggest, and concluded that the arrangement adopted by paleontologists must be wrong. H. B. Fell and later authors have reaffirmed that the paleontological evidence overrules nebulous embryological considerations.

David L. Pawson; Andrew C. Campbell

Ecology and Feeding

Echinoderms occur everywhere in the world's oceans, but they are usually rare in areas where the salinity of the water is greatly reduced or where pollution levels are significant. In many areas, echinoderms are dominant invertebrates in terms of numbers or biomass, and locally sea stars can be the top predators. On a rocky shore, predatory sea stars and vegetarian sea urchins may be common on and under rocks. Brittle stars are typically concealed under rocks; they may emerge at night to forage for small organisms. In deeper water, where conditions are suitable for suspension feeding, huge aggregations of brittle stars or feather stars may occur, with their arms extended into the water for feeding. Echinoderms flourish on and near coral reefs, and it is in these tropical areas that they can achieve their greatest diversity. In polar regions, especially around Antarctica, echinoderms, along with sponges, dominate in most habitats (Fig. 4). Sandy to muddy bottoms may be populated with large numbers of burrowing detritus-eating sea urchins, long-armed brittle stars, and sea cucumbers. In such habitats, burrowing sea cucumbers may extend their sticky tentacles into the surrounding water to capture small drifting organisms. In deeper water, below 100 m (330 ft), stalked sea lilies can be common on hard substrates, efficiently suspension-feeding with their arms deployed in the form of a parabolic bowl. In the deep sea, echinoderms can also be dominant; on abyssal (relating to great ocean depths) plains, sediment-swallowing sea cucumbers can make up more than 95% of the total biomass on the sediment surface. Most echinoderm groups occur in the deep sea to depths in excess of 9000 m (5.6 mi).

David L. Pawson

Fossils

Echinodermata are an especially important group of fossil invertebrates. Except for reworked fragments found in some nonmarine deposits, echinoderm remains occur exclusively in strata laid down on sea



Fig. 4. Echinoderms in the deep sea near Antarctica at a depth of 595 m (1950 ft). Three sea cucumbers (*Scotoplanes globosa*) are feeding on seafloor sediments. At bottom right is a brittle star, probably *Ophiomusium* species. Scattered on the seafloor but barely visible are several feather stars. (Courtesy of the U.S. National Science Foundation)

bottoms—chiefly those of shallow seas. Many of these deposits, ranging in age from Cambrian to late Tertiary, contain abundant echinoderm fossils. In addition, fossil echinoderms undoubtedly are widely distributed beneath all oceans, although remains of these organisms in deep-water sediments, even of recent origin, are virtually unknown because of their inaccessibility.

The echinoderms are well adapted to preservation as fossils owing to their abundant calcareous skeletal elements. Paleontological importance of the group is explained partly by this fitness but more by the diversity of their kinds, the generally short-lived existence and wide geographic distribution of most recognized taxonomic units, and the clearness with which evolutionary trends can be defined.

Classification and Description

For some years, fossil and living echinoderms have been arranged into four or five subphyla. As formal groupings, the subphyla have been largely abandoned, for it is recognized that the subphylum definitions and distinctions are artificial rather than truly reflecting phylogeny. However, the subphylum names are useful for grouping the various echinoderm classes, and they are used here in an informal sense: homalozoans, crinozoans, echinozoans, and asterozoans. Similarly, the terms eleutherozoan (free-moving) and pelmatozoan (sedentary or sessile), previously used in formal classifications, are useful today in categorizing the lifestyles of echinoderms. See ELEUTHEROZOA; PELMATOZOA.

Homalozoans. The distinctive features that separate the extinct homalozoans (or carpoids) from all other Echinodermata are a complete lack of radial

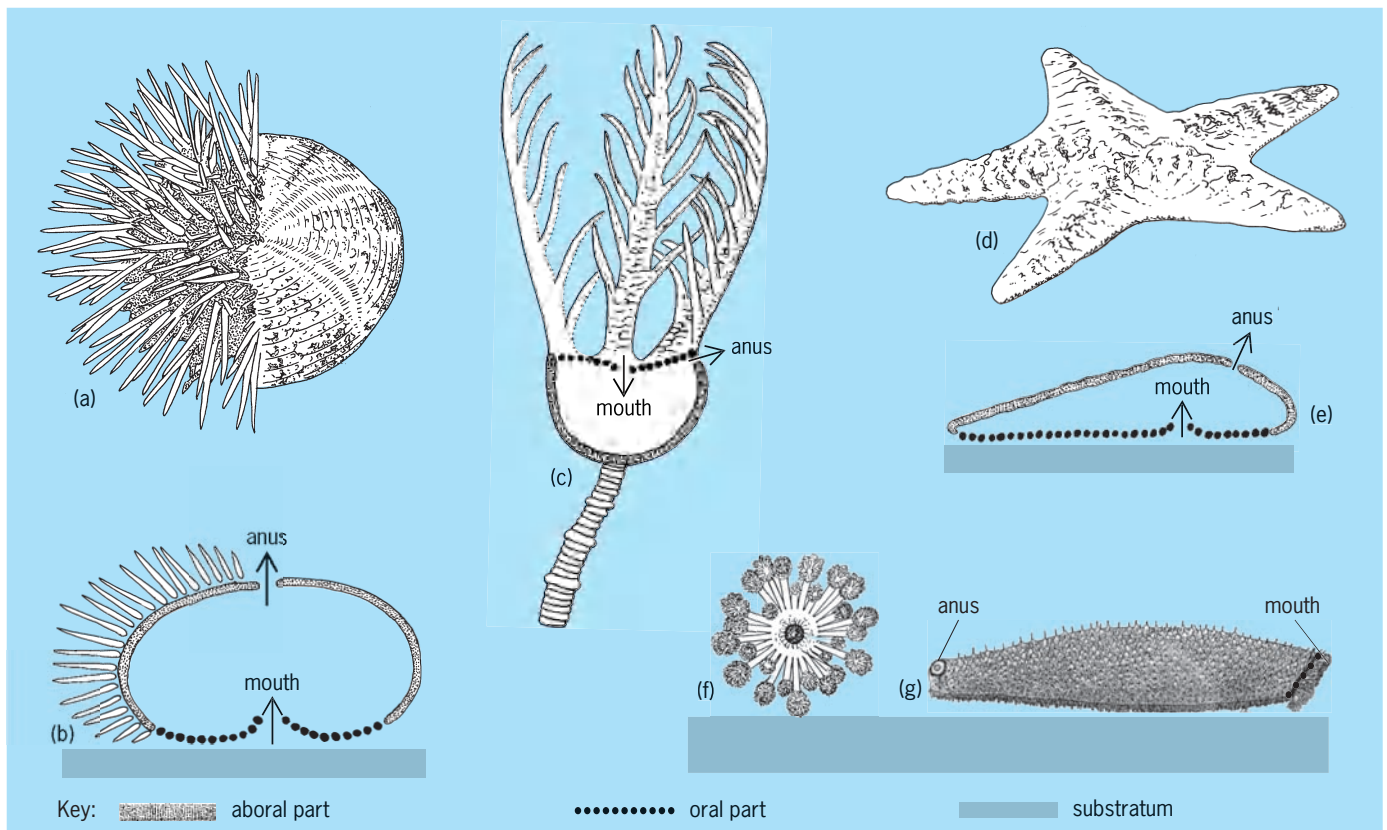


Fig. 5. Representative types of echinoderms. (a) Regular echinoid (*Lytechinus*, Recent), oblique aboral view, right half with spines removed, showing two ambulacra and three interambulacra. (b) Diagrammatic section through a, showing the stoutly built test. (c) Crinoid showing the stem and three arms, calyx sectioned. (d) Asteroid (*Dermasterias*, Recent), aboral oblique view. (e) Diagrammatic section through the ray at left in d and opposite the interambulacrum. (f) Holothuroid, oral view showing tentacles around the mouth. (g) Holothurian, lateral view.

symmetry in arrangement of skeletal parts and the flattened body form, from which one or more specialized slender appendages may extend. Plates enclosing the body are irregular in shape and size and vary in number between forms. There are four quite dissimilar classes of carpooids: Homostelea, Ctenocystoidea, Stylophora, and Homoiostela. Some experts, especially R. Jefferies, believe that some of the homalozoans, the Stylophora in particular, are not echinoderms at all, but are "calcichordates," groups forming the stem of the chordates. There has been considerable debate in the scientific literature about this matter, with persuasive evidence presented on both sides. Most paleontologists tend toward the idea that all homalozoans are true echinoderms. The homalozoans range in age from earliest Cambrian to Carboniferous. See CARPOIDS; HOMALOZOAN.

Crinozoans. The crinozoans evidently include types of echinoderms least modified from the ancient progenitors of the phylum. Crinozoans range from Cambrian to Recent, but only the Crinoidea, abundant in Paleozoic formations, have survived to the present. In some classifications, the Crinoidea and Paracrinoidea are treated as the subphylum Crinozoa in the strict sense, and the other classes are grouped together in the subphylum Blastozoa. A characteristic crinozoan feature is orientation of the body with the oral (ventral) side directed upward and ab-

oral (dorsal) side downward; thus, although commonly used, the designations ventral and dorsal are not suited to these particular echinoderms (Figs. 5c and 6). There are nine classes of crinozoans:

(1) Eocrinoidea (the crinozoan stem group) and (2) Paracrinoidea are crinozoan groups attached by a stalk. They have ovoid bodies enclosed by irregularly arranged plates that commonly lack well-developed radial symmetry and possess a combination of characters typical of the two following classes. (3) Diploporita and (4) Rhombifera (previously grouped in the polyphyletic Cystoidea, with saclike bodies, and each possibly polyphyletic) are a diverse group of primitive crinozoans distinguished especially by pores or tubular canals that penetrate irregularly arranged plates enclosing the body. Attachment to the substrate was by means of a stem composed of superposed discoid plates. (5) Crinoidea (sea lilies) include the most abundant crinozoans, most of the fossil forms anchored to the sea bottom by means of a short to very long stem and distinguished by upwardly directed arms which functioned for gathering food. Some ancient and most modern crinoids are stemless, known as feather stars, and adapted for a free-living existence. True radial symmetry and regularity of skeletal construction are attributes that distinguish crinoids. (6) Blastoidea (budlike forms) are small to

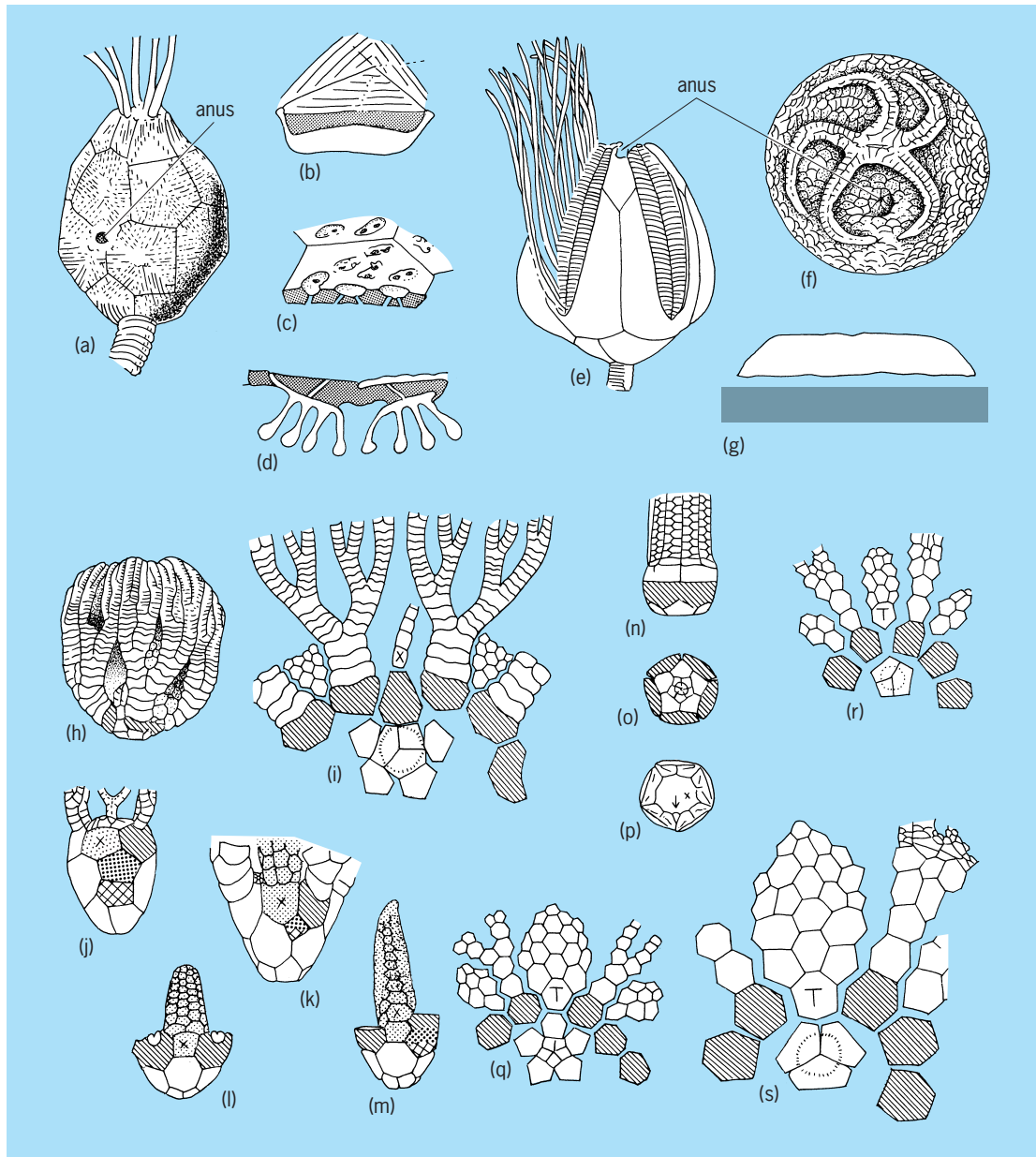


Fig. 6. Fossil pelmatozoan echinoderms. (a) Rhombiferan cystoid (*Echinoencrinus*, Ordovician). (b) Diagrammatic oblique view and section through a cystoid pore rhomb (half of one rhomb toward front bisected by a suture between plates). (c) Oblique view and section of a diploporitan cystoid plate. (d) Section through the ambulacrum of a blastoid (*Pentremites*, Mississippian) showing hydrospires beneath the lancet and side plates. (e) Blastoid (*Pentremites*) with brachioles restored along the side of one ambulacrum. (f) Edrioasteroid (*Carneyella*, Ordovician), oral view. (g) Section through a stalkless edrioasteroid. (h, i) Flexible crinoid (*Taxocrinus*, Mississippian), posterior side and partial plate diagram showing the right posterior plane of bilateral symmetry defined by the infrabasal circlet. (j–p) Inadunate crinoids, some showing anal sacs and parts of arms: (j) *Carabocrinus*, Ordovician; (k) *Botryocrinus*, Devonian; (l) *Cyathocrinites*, Mississippian; (m) *Cupulocrinus*, Ordovician; (n–p) *Delocrinus*, Pennsylvanian. (q) Diplobathrid camerate crinoid (*Ptychocrinus*, Ordovician), plate diagram. (r, s) Monobathrid camerates: (r) *Macrostylocrinus*, Devonian; (s) *Periechocrinites*, Mississippian, plate diagrams showing noteworthy distinctions in basal and radial circlets. (Parts f, g after L. H. Hyman, *The Invertebrates*, vol. 4: *Echinodermata*, McGraw-Hill, 1955)

medium-sized crinozoans having a highly developed radial symmetry and very regular arrangement of the few plates surrounding the body. The body was attached to the substrate by a very slender stem, which had a circular cross section. (7) *Coronoidea* is a small group of middle Paleozoic crinozoans, similar to blastoids and probably ancestral to them. Coronoids differ in the structure of the ambulacral plate and possess erect, as opposed to recumbent, am-

bulacra. (8) *Parablastoidea* is another small, principally Ordovician class of crinozoans with cystoid-like thecae but with uniserial ambulacra. (9) *Edrioblastoidea* comprises a single Ordovician genus with a theca that is blastoid-like in shape and symmetry but lacks brachioles or hydrospires (or pore rhombs). See CRINOZOA.

Eocrinoids. This small group of Cambrian to Silurian crinozoans combines characteristics of

cystoids and crinoids, yet differs significantly from both. They resemble cystoids in their mode of branching of the ambulacral grooves and ventrolateral location of the anus but lack thecal pores or distinct pore rhombs; they are like crinoids in plate structure and similarity of plate arrangement in the calyx (plated body). The eocrinoids have stems and unbranched or bifurcating arms. They were probably ancestral to all other crinozoans. *See* EOCCRINOIDEA.

Paracrinoids. Paracrinoida are stem-bearing echinoderms that also combine features of cystoids and crinoids, but in a manner quite unlike that of the eocrinoids. The paracrinoids, now known only from Middle Ordovician to Lower Silurian deposits, have numerous plates of the calyx that are not arranged in series and that lack a ventrally differentiated area corresponding to the tegmen of crinoids. The plates have a cystoid-like pore structure, but the arms are comparable to those of crinoids.

Cystoids. The Ordovician to Devonian cystoids, comprising the classes Diploporita and Rhombifera, are extinct crinozoans of globose, subcylindrical, or flattened ellipsoidal form that are characterized mostly by irregularity of the body plates (Fig. 6a) and a short, weak stem. A variable number of slender appendages (brachioles) on the upper side of the calyx, supplemented by ambulacral grooves, served for gathering food. The mouth was located centrally at the summit of the calyx, and an anus was fairly well down on one of the sides, which accordingly is defined as posterior. Other small orifices, interpreted as hydropore and gonopore, occur near the anus. In one class, Diploporita, the numerous calyx plates are perforated by numerous pairs of minute tubular openings that allow seawater to circulate through them (Fig. 6c). Remaining cystoids are named Rhombifera (rhomb-bearing forms) because some or all of their relatively few calyx plates are penetrated by rhomb-shaped groups of tubes running parallel to each other and to the plate surfaces (Fig. 6a). The two halves of any rhomb lie on adjoining plates, with the tubes crossing the plate boundaries at right angles (Fig. 6b). The pattern of calyx plates suggests that the blastoids and crinoids may be descendants of rhombiferan cystoids. *See* DIPLOPORITA; RHOMBIFERA.

Crinoids. In terms of abundance of fossil remains, the crinoids outrank all other echinoderms combined (Fig. 5c). There are about 5000 extinct species and at least 700 living kinds. Some Paleozoic deposits hundreds of feet thick in areas measured in hundreds of square miles are largely composed of fossil remains of trillions of crinoids. *See* CRINOIDEA.

Skeletal features. Although many modern and some ancient crinoids are stemless as adults, this group of pelmatozoans typically is attached to the sea bottom by a more or less elongate stem composed of superposed calcareous discs (columnals) perforated centrally by a circular, pentagonal, or pentastellate canal. At the opposite extremity of the stem, which exceptionally may be 15 m (50 ft) tall, is the crinoid body, encased in regularly arranged plates and surmounted by branched or unbranched free-moving

arms. The conjoined plates below the free arms make up the so-called dorsal cup; this cup, along with the plates of the ventral surface composing the tegmen, makes up the crinoid calyx (Fig. 6b-s). The calyx and its attached arms are termed the crown. The mouth is located on the ventral surface. The anus lies on the tegmen, is raised above it on an anal sac (Fig. 6k-m), or may lie on the side of the dorsal cup. The posterior side of the crinoid is defined by the position of the anus in one of the interrays or by extra plates introduced in such a position on one side of the dorsal cup; the ray opposite to the posterior interray is defined as anterior. This establishes a plane of bilateral symmetry that more or less modifies the fundamental radial symmetry of the crinoid (Fig. 7a).

The plates of the crinoid dorsal cup are arranged in a regular pattern of successive circlets, but with differences in various groups that furnish a basis for classification. Each circlet normally contains five plates. At the base of each ray is a radial plate. Beneath the circlet of radials are five basals that are interradial in position; some crinoids have a still lower circlet of infrabasals that alternate with the basals and hence occur in radial position. Crinoids with a single circlet of plates below the radials are termed monocyclic (Fig. 6r and s), and those with two circlets are dicyclic (Fig. 6b-q).

In some crinoids the tegmen is stoutly constructed of small, irregularly arranged plates, while in others it consists of a flexible, leathery integument that may be studded with calcareous ossicles. Five subequal oral plates larger than others of the tegmen may occur interradially around the mouth. The arms of crinoids are extremely variable in plan and construction. They may be unbranched or moderately to highly branched, with or without very numerous tiny branchlets called pinnules, and composed of a single or double series of arm plates (brachials). These characters, along with the mode of articulation among plates of the rays, are important also in classification.

Main types. Four subclasses of crinoids are recognized, three of them distributed from Ordovician to Permian; one of the three (Inadunata) persists to Middle Triassic. The fourth subclass (Articulata) ranges from the Mesozoic to present day.

The Inadunata ("not united" forms, referring to lack of incorporation of lower arm plates in dorsal cup) are crinoids with a relatively small dorsal cup containing one or two circlets of plates below the radials and having the arms entirely free above the cup (Fig. 6j-p). They include 1750 or more species that exhibit great variety in form and evolutionary trends. In a majority the anteroposterior plane of bilateral symmetry is well marked, without other deviation from a regular pentamerous plan (Fig. 7a), but in one group (superfamily Homocrinacea) a surprising degree of bilateral symmetry was developed in the plane of the left anterior ray (Fig. 7c). Some of these crinoids have a strongly downbent crown that was hinged on the summit of the stem. In another group (superfamily Heterocrinacea) a subordinate plane of

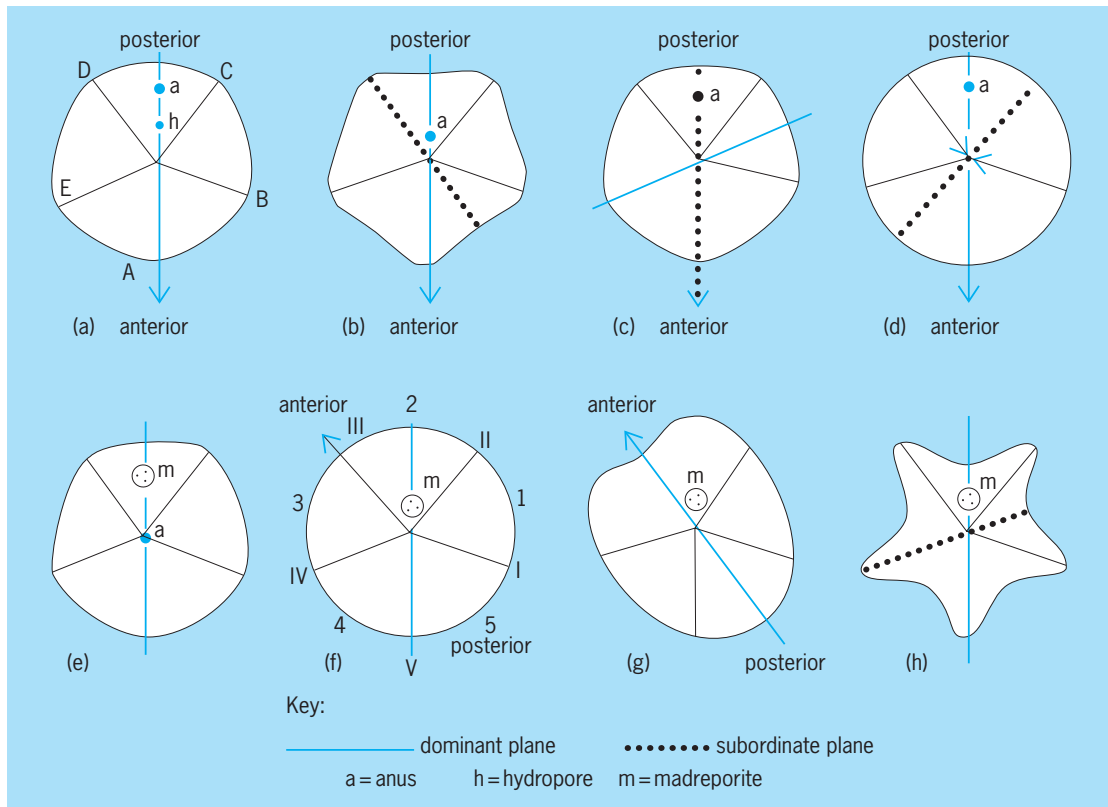


Fig. 7. Bilateral symmetry developed in various echinoderms. Anterior and posterior directions are indicated where distinguished. All diagrams represent aboral views. (a) Crinozoans (in general, most crinoids, cystoids, edrioasteroids); letters denote ray designations according to the Carpenter (P. H. Carpenter, 1884) system: A, anterior; B, right anterior; C, right posterior; D, left posterior; E, left anterior. (b) Blastoids, heterocrinoids, with subordinate bilateral symmetry in the left posterior plane. (c) Homocrinoids, with primary bilateral symmetry in the left anterior plane. (d) Glyptocrinoids, flexible crinoids, and rhombiferan cystoids, with subordinate bilateral symmetry in the right posterior plane. (e) Holothuroids. (f) Regular echinoids, rays marked according to the Lovén (S. Lovén, 1874) system, III being considered anterior. (g) Irregular echinoids, with prominent bilateral symmetry in the left posterior plane of crinozoans. (h) Asteroid, with subordinate bilateral symmetry in the left anterior plane.

bilateral symmetry is oriented in the left posterior plane (Fig. 7*b*). See INADUNATA.

The Camerata (chamber or “box” forms) are most numerous among ancient crinoids in both variety and quantity of individuals (Fig. 6*q-s*). More than 2500 species of these fossils have been described, two-thirds of which come from Mississippian rocks alone. The camerates are distinguished by the stout construction of their calyx, which incorporates lower ray plates and interradials in the dorsal cup, and subtegmina location of the mouth. Types with both one and two circlets of plates beneath the radials are common. Anteroposterior bilateral symmetry is developed almost universally in these crinoids as a modification of the dominant pentamer pattern (Fig. 7*a*); a secondary plane of bilateral symmetry directed through the right posterior ray prevails in the monobathrid suborder Glyptocrinina, as in all flexible crinoids (Fig. 7*d*). See CAMERATA.

The Flexibilia (flexibles) are a distinctive assemblage of exclusively dicyclic crinoids characterized by movable ligamentous union between most of the plates and several constant features in the organization of the calyx (Fig. 6*b* and *i*). They include approximately 300 described species, all of which exhibit a secondary bilateral symmetry in the right posterior

plane in addition to their generally well-marked primary bilateral symmetry directed anteroposteriorly (Fig. 7*d*). See FLEXIBILIA.

The Articulata (forms divided into joints) comprise Mesozoic and Cenozoic crinoids, represented by about 500 fossil and about 700 living species. The stalked articulate crinoids have been classified into four orders: Millericrinida, Cyrtocrinida, Bourgetocrinida, and Isocrinida. All of these orders have fossil and extant representatives. Some experts have abandoned the order categories for the 100 species of living stalked crinoids and prefer to arrange them into 11 families. There are three orders of unstalked crinoids: Uintacrinida and Roveacrinida are extinct, while the Comatulida (feather stars) has about 600 living species. Extant sea lilies are confined to depths below about 100 m (330 ft). The stalk is usually attached to a hard substratum—a rock, or a piece of shell embedded in soft mud—by means of a terminal disc cemented to the substratum or by means of claw-like whorls of cirri. In a few species, the end of the stalk is branched and rootlike for anchoring the stalk in soft sediments. When detached from the substratum and lying on the seafloor, sea lilies are capable of slowly moving to a new site. In contrast, the feather stars are

highly mobile, often capable of active swimming by thrashing of the arms in an up-and-down motion. Feather stars range from shallow waters to great ocean depths. They are common on coral reefs, and in certain suitable habitats can reach population densities of hundreds of individuals per square meter. Sea lilies and feather stars are passive suspension feeders; they extend their arms into the water and capture small drifting organisms and particulate matter on their sticky tube feet. Captured items are conveyed to the mouth along food grooves by means of ciliary action. In extant crinoids, the gonads develop inside the pinnules near the bases of the arms, causing conspicuous swelling. At the appropriate time for breeding, the pinnules burst to release eggs or sperm into the surrounding seawater. See ARTICULATA (ECHINODERMATA).

Blastoids. Blastoidea, Ordovician to Permian, display a regular, fivefold radial symmetry (Fig. 6e), on which are superposed a first order of bilateral symmetry in the anteroposterior plane and a second order directed through the left posterior ray (Fig. 7b). The calyx of average-size specimens is small, with diameter of about 15 mm (0.6 in.) and height of 20 mm (0.8 in.). It is composed of 23 or 24 plates, of which five (lancet plates) are not visible externally on unweathered specimens. Within forks of the radials are a multitude of very diminutive so-called side plates that conceal the lancets; they form ambulacral areas bordered laterally by rows of threadlike free armlets (brachioles) that function as food gatherers. Particles of food are conducted to midlines of the ambulacra and thence upward to the mouth, which is at the center of the ventral surface. Beneath the ambulacra are extremely delicate, longitudinally folded, calcareous lamellae that enclose narrow troughs for circulation of water admitted through pores at the base of the brachioles; these structures, termed hydrospires, correspond to the slitlike parts of pore rhombs in the rhombiferan cystoids (Fig. 6d). More than 500 species have been described. See BLASTOIDEA.

Asterozoans. The asterozoans are star-shaped, free-moving echinoderms with well-marked radial symmetry in arrangement of their skeletal parts. They are easily recognized by their strongly developed arms and central disk (Fig. 5d and e). The arms of most asterozoans, with their tube feet functioning for locomotion and predation, are highly mobile. They range in age from Early Ordovician to Recent but are more abundant and varied in modern seas than in the fossil record. They are divided into four classes: Somasteroidea, Asteroidea, Concentricycloidea, and Ophiuroidea. The extinct Somasteroidea ranged from the Ordovician to the Devonian. Somasteroids resemble true sea stars in general shape, with broad arms, but the arms are unique in having a pinnate arrangement of rodlike plates. Asteroidea (sea stars) are distinguished by prominent flexible rays that join the central disk without clearly shown demarcation (Fig. 5d and e). Concentricycloidea are flattened and discoidal with no evident arms. The Ophiuroidea (brittle stars) have long, slender, snakelike arms which are

sharply set off from the central discoid body. They are very active echinoderms, able to crawl rapidly in any direction.

Somasteroids. These rare, broad-armed sea stars are especially characterized by the featherlike arrangement of parallel, rodlike plates in their wide petaloid rays. The rods extend outward from the medially placed ambulacral plates. A food-carrying groove is contained in each ray. The aboral side of the skeleton has a rough, coarse meshwork.

Asteroids. The sea stars have a long fossil history, ranging back to the Ordovician. Sea stars are comparable to the ophiuroids in their general stellate form (Fig. 5d and e, Fig. 7b). They differ from the brittle stars in two main respects: lack of strongly marked separation between the body and its radially disposed hollow arms, and the open ambulacral grooves along the ventral side of the arms. The skeletal elements tend to be loosely joined; consequently, asteroids are not well suited for fossilization in a manner showing the skeletal arrangement of the entire animal. Sea stars may range in size from a few millimeters to over a meter (3.3 ft) in diameter. Ossicles along the ambulacra occur in two or four series differentiated as ambulacrals and adambulacrals. Some asteroids lack obvious arms and are essentially pentagonal with ambulacral grooves on their ventral (oral) side. Others, such as the Heliasteridae, develop as many as 44 arms as they grow, losing their fundamental pentamerous symmetry but structurally resembling common types of asteroids. A single madreporite lies on the upper (aboral) surface, placed between the arms (interradially). An anus is present in the center of the aboral surface in most groups and absent in others (notably the Paxillosida).

Behavior. Asteroids are universal symbols of the ocean. They are conspicuous in a great variety of habitats, from mud to sand to rock to coral reefs. In some areas they can be present in great numbers, causing local habitat damage. The multiarmed, crown-of-thorns sea star, reaching 1 m (3.3 ft) in diameter, is common on Indo-Pacific coral reefs. It can move across reefs in large swarms, eating the soft coral tissues and causing short- to long-term devastation of reefs. Some groups of sea stars, notably the Paxillosida, have pointed rather than suckered tube feet, and they tend to swallow prey whole (such as small clams and sand dollars), ejecting the shells from the mouth when digestion is complete. The Forcipulatida are notorious predators; a forcipulatid may attach its suckered tube feet to the two shells of a clam or a mussel and then, for several hours if necessary, gently pull on the two shells, until the bivalve's muscles tire, the shells open a fraction of an inch, and the sea star drops its stomach into the bivalve, digesting the animal inside its own shell. Some deep-sea asteroids are mud-swallowers, digesting from the mud whatever organic material may be present.

Classification. The classification of the Asteroidea has been revised by many authors over the past 50 years or so. In recent years, the most commonly used classification of both fossil and extant forms is that proposed by D. B. Blake. Extinct

Paleozoic orders are *Platyasterida*, *Pustulosida*, *Hemizonida*, and *Uractinida*. Post-Paleozoic sea stars include the extinct order *Trichasteropsida*, and seven extant orders, all of which also have fossil representatives: *Paxillosida*, *Notomyotida*, *Velatida*, *Valvatida*, *Spinulosida*, *Forcipulatida*, and *Brisingida*. See ASTEROIDEA; FORCIPULATIDA; NOTOMYOTIDA; PAXILLOSIDA; VALVATIDA.

Concentricycloids. The sea daisies comprise three extant deep-sea species, two known from the Pacific Ocean and one from the Atlantic. Sea daisies are flattened, discoidal echinoderms up to 15 mm (0.6 in.) in diameter. The upper surface is covered with delicate overlapping plates; the lower surface has a mouth frame and a single peripheral ring of tube feet. These animals have been found only in association with pieces of wood on the deep-sea floor. Presumably, they attach to the wood substratum using their suckered tube feet. This class was discovered and diagnosed in 1986, and its formal status has been debated ever since. Several scientists have argued persuasively that the concentricycloids should be placed in the Asteroidea, while others maintain that they represent a separate class of echinoderms.

Ophiuroids. The brittle stars or serpent stars composing the Ophiuroidea are highly mobile echinoderms closely related to the sea stars. They are distinguished readily by their external form, since the small, rounded to pentagonal or scalloped disk is sharply set off from the symmetrically placed long, slender arms that extend radially outward from it. The arms may be smooth or spiny. Almost invariably they are five in number and unbranched, but in a few kinds, such as the Recent basket stars (*Gorgonocephalidae*), the arms are repeatedly bifurcated. All known fossil ophiuroids have simple arms. The skeleton of the central disk is composed of many regularly arranged plates, without any deviation on either the aboral or oral surfaces from perfect pentamerous symmetry. The mouth is at the center of the oral surface; an anus is lacking. A single interradial madreporite is present on a plate near the mouth; in some basket stars there are five interradsial madreporites around the outer margin of the mouth skeleton. The arms are solid, not hollow as in asteroids. The arm skeleton is distinctive in being internal, consisting of "vertebrae" formed from fused ambulacral ossicles. Brittle stars are common in a great variety of habitats—living under rocks, concealed in coral reefs, or scattered in enormous numbers across great expanses of mud on the deep-sea floor. Many brittle stars are selective detritus feeders; others feed by extending their arms into the seawater and capturing small drifting organisms or particles on their sticky tube feet. Recent investigations have shown that certain brittle stars have well-developed eyes in their upper-arm skeletons, and they can respond rapidly to changing light regimes, seeking shelter in crevices as necessary.

Fossil ophiuroids are known from Lower Ordovician to Pleistocene (Recent), but they are not abundant. Fewer than 100 species have been described from Paleozoic, Mesozoic, and pre-Recent Cenozoic formations, as compared with approximately 2000

known living species. Ophiuroids are arranged in four orders: the extinct *Stenurida* (Ordovician to Devonian) and *Oegophiurida* (Ordovician to Carboniferous), and the extant *Phrynophiurida* and *Ophiurida* (Ordovician to Recent). See OPHIUROIDEA.

Echinozoans. These echinoderms generally have an ovoid or globose body, and they lack armlike appendages. Some echinozoans may have a low discoid body or an elongate cylindrical to almost wormlike form. All have conspicuous radial symmetry of their skeletal parts, and in almost all the symmetry is pentamerous. The skeleton of some echinozoans consists of a rigidly constructed test to which movable jaw parts and external spines are attached, whereas in others all skeletal parts either are joined together flexibly or are reduced in size and separated from one another by leathery tissue. The echinozoans comprise seven classes: *Helicoplacoidea*, *Campptostromatoidea*, *Edrioasteroidea*, *Cyclocystoidea*, *Ophiocystioidea*, *Echinoidea*, and *Holothuroidea*, the first five of which are long-extinct groups, while the others are represented by both fossil and living forms. (1) The *Helicoplacoidea* are Lower Cambrian forms with a flexible, expansible test composed of spirally arranged plates. (2) The Lower Cambrian *Campptostromatoidea* are conical or domal animals with plates of varying size overlapping on the lower theca. (3) *Edrioasteroidea* ("sessile sea stars") range from the Ordovician to the Permian. They have a many-plated test, ranging from rigid to somewhat flexible, of generally discoid shape with five distinct, straight or curved rays on the upper surface. Some were able to move about, but nearly all are found attached to some hard foreign object, such as a brachiopod or other invertebrate shell. (4) The *Cyclocystoidea* are Ordovician-to-Devonian echinoderms with a disk-shaped body, the upper surface of which is covered by plates arranged in concentric rings. (5) The *Ophiocystioidea* are Ordovician-to-Devonian echinoderms that differ from other classes in having large fine-plated tube feet, but share some skeletal features with the *Holothuroidea*, with which they have sometimes been aligned. (6) The *Echinoidea* (sea urchins) are characterized by a rigid to flexible skeleton of varied shapes, though mostly ovoid to subglobose, covered by numerous long or short movable spines (Fig. 5*a* and *b*). Pentamerous symmetry is strongly developed, modified in all orders by more or less distinct bilateral symmetry and in some by differentiation of anterior and posterior regions. Very abundant in present-day seas, sea urchins range from Early Ordovician onward. (7) The *Holothuroidea* (sea cucumbers) are soft-bodied, leathery-skinned echinozoans with an elongate cylindrical body and skeletal parts usually consisting of discrete, generally microscopic, ossicles ("little bones") embedded in the body wall and internal tissues (Fig. 5*f* and *g*). See ECHINOZOA.

Edrioasteroids. A unique assemblage of mainly early Paleozoic attached echinoderms (Fig. 6*f* and *g*), *Edrioasteroidea* were distributed from Lower Cambrian to Upper Carboniferous. They had no stem for anchorage but were fixed by their entire base. The upper (ventral) side was covered by many

flexibly joined plates, which are sharply differentiated into rows of paired ambulacrals with adjoining adambulacrals and irregular interambulacrals. The five ambulacral tracts curve outward from the centrally located mouth, with their extremities often deflected rather strongly in a consistent way. The anus lies in one of the interrays that accordingly is regarded as posterior. See EDRIOASTEROIDEA.

Helicoplacoids. Numerous representatives of Helicoplacoidea have been found in the oldest fossil-bearing strata (*Olenellus Zone*) of the Lower Cambrian at localities in western Nevada and southern California. The body is ovoid to pyriform, enclosed in a test of loosely joined, spirally arranged rows of elongate plates, with columns, interpreted as ambulacral, and intervening ones, as interambulacral, originating at an apical pole and extending to the opposite oral pole. The test could be distended and retracted so as to change its shape and volume. Other echinoderms associated with these oldest echinozoans include representatives of the pelmatozoan class Eocrinoidea. See HELICOPLAÇOIDEA.

Echinoids. The class Echinoidea is the most important group of fossil eleutherozoan echinoderms (Fig. 5*a* and *b*; Fig. 8). They are distinguished from crinozoans by freedom from a sessile mode of life, and they share with the holothuroids a downward orientation of their oral surface. Modern species of echinoids total approximately 800, as shown by T. Mortensen's several-volume monograph (1928–1952) covering all known living kinds. In comparison, 2000 species of sea stars, 2000 species of brittle stars, and 1200 species of sea cucumbers are alive today. All eleutherozoan groups are known as far back as Ordovician time, and yet fossil echinoids are more numerous than all others of the assemblage combined. Such disparity reflects inequality in fitness of the different sorts of skeletons to be preserved. Not all parts of echinoid skeletons may be found intact. Multitudes of movable

spines attached to the test during life tend to be separated and scattered, as does the echinoid masticatory apparatus, which is known as the Aristotle's lantern. Thus, knowledge of the entire skeletal structure of most fossil echinoids is incomplete. See ECHINOIDEA.

Morphology. The echinoid test is typically globular, unevenly ovoid, or discoid in shape (Figs. 5 and 8). Mainly it is formed by 10 meridionally disposed bands of plates, of which five are ambulacral, distinguished by porelike openings between or through the plates for passage of the tube feet; the five bands of plates that alternate with the ambulacrals are interambulacral and lack perforations for tube feet. In most echinoids, both living and fossil, the plates are joined together rigidly, but in a minority the union is somewhat flexible and imbricating (having overlapping edges). Virtually all post-Paleozoic echinoids have a constant arrangement of the ambulacral and interambulacral groups of plates, each of the 10 bands comprising a double column, making 20 columns of plates in the test as a whole. The plates are widest at the equator (ambitus) of the test and increasingly narrow toward the oral and aboral poles.

The mouth, located on the underside, is surrounded by a peristome, an area of naked integument or flexibly united small plates. The anus generally is placed on the aboral side of the test opposite the mouth, but it may be shifted backward to a posterior position or even to the oral surface. The anus is surrounded by a bare or small-plated periproct. Another part of the test important for orientation and classification is the so-called apical system of 10 special plates at the aboral pole, five so-called ocular plates alternating with five genital plates. One of these genital plates is an interambulacrally placed sieve plate (madreporite) that is a landmark for identifying homologous parts, not only among diverse kinds of fossil and living echinoids but among echinoderms generally. Fossil remains of echinoids include numerous spines not associated with the test that bore them. Even so, the spines are commonly distinctive and useful as paleontological tools (Fig. 8*c* and *d*).

Orientation. Seemingly, there should be little reason for discussing the orientation of echinoids because the downward-directed oral surface clearly is ventral and, at least among irregular echinoids which are characterized by obvious bilateral symmetry, locomotion consistently is in a single direction along the path defined by the plane of this symmetry. If such direction of movement is considered to be forward, the ambulacral ray on this side of the test must be anterior and the opposite interambulacrum is posterior (Fig. 7*g*). In irregular echinoids, such as the spatangoids, the spines nearly all point backward. In many irregular genera, the mouth is found to have shifted forward and the anus rearward. These observations have led to adoption by specialists of a scheme of numerical designation of the ambulacral and interambulacral rays, introduced by S. L. Lovén (1874) and known as the Lovén system (Fig. 7*f*). It is applied to fossils as well as living echinoids and to both irregular and regular types. Orientation of

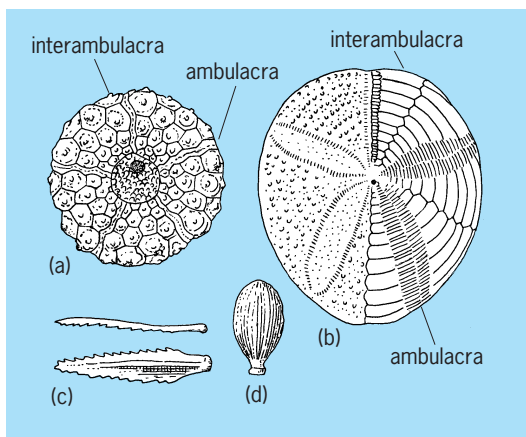


Fig. 8. Fossil echinoids. (a) Regular form (*Echinocrinus*, Mississippian), aboral view showing the anus and periproct, wide interambulacra and narrow ambulacra. (b) Irregular form (*Eupatagus*, Paleogene), aboral view, left side showing an unweathered appearance with many spine bases, right side showing the plate structure with ornamentation omitted; anterior ray directed upward. (c, d) Types of echinoid spines: (c) *Porocidaris*, Paleogene; (d) *Balanocidaris*, Jurassic.

the regular echinoids seemingly should not be easy, since they move in any direction with equal ease and without detectable preference; also, except for the madreporite, the test has perfect radial symmetry. Inclusion of the regulars is possible only by use of the off-center location of the madreporite on the aboral surface as a reference point, assuming that relative position is the same in all echinoids. The interambulacrum with the madreporite is identified as number 2 of the Lovén system.

Among crinozoans, which almost certainly embrace the ancestors of echinoids, the posterior interradius is readily and positively identified. If this orientation is applied to echinoids, a discrepancy becomes evident at once, for the adopted antero-posterior plane of echinoids clearly is that coinciding with the left posterior ray of the crinozoans (Fig. 7*e* and *f*). Interambulacrum 2 (Lovén), considered as right anterior by echinoid students, is equivalent to the posterior interambulacrum of crinoids, for example. This does not mean that echinoid orientation is wrongly conceived but merely that evolution of these echinoderms has pursued a divergent path of its own, which incidentally duplicates the subordinate plane of bilateral symmetry in some crinozoans (blastoids, heterocrinoids).

Habitats and behavior. Echinoids can be extremely common from the rocky intertidal down to abyssal ocean depths. The irregular urchins, such as sand dollars, sea biscuits, and heart urchins, are always associated with soft substrates—sand, sandy mud, or mud. Sand dollars and sea biscuits live at the surface of sand or burrow just beneath the surface. Many of these echinoids ingest sand, obtaining sustenance from small organisms, such as bacteria, associated with the sand grains. Heart urchins burrow into the substratum and remain buried, feeding on particulate matter in the mud. The regular urchins are usually associated with hard substrates, typically rock, coral, or rubble, but numerous species may occur on sand or mud in quiet waters. Most of the regular urchins are vegetarians, feeding primarily on algae. Typically, echinoids reproduce by liberating eggs or sperm into the seawater from genital pores on top of the body. There are usually five pores, but in irregular urchins there may be four, three, or two pores. Fertilization of eggs takes place in the seawater, and the developing embryo becomes a planktonic larva. Within a few weeks, the urchin develops on the left side of the larva, breaks away, and sinks to the bottom as a juvenile. Some urchins have a nonfeeding yolky larva, and in others the larval stage is omitted altogether, the young being brooded by the female parent in special chambers.

Classification. Former taxonomic arrangement of the echinoids recognized two main subclasses, the regular echinoids (Regularia) and irregulars (Irregularia). Present knowledge, based on fossils as well as living forms, shows that this arrangement is artificial. The so-called regulars are a composite assemblage that on the one hand contains several extinct primitive kinds, exclusively Paleozoic, and on the other, advanced forms of modern type, all Mesozoic and Cenozoic, that include the stocks (Diademat-

acea) from which irregular echinoids were derived. Accordingly, classification was considerably revised (J. W. Durham and R. V. Melville, 1957; A. B. Smith, 1984; and others) so as to take account both of Mortensen's comprehensive work on living echinoids and of evidence from paleontology.

Ophiocistioids. Ophiocistioida is an aberrant echinozoan group containing five known genera distributed from Lower Ordovician to Middle Devonian. The chief peculiarity of the group is the presence of numerous armlike appendages attached to the oral surface in an ambulacral position. They are interpreted as much-enlarged podia that are covered by abundant minute plates, a character seemingly shared with at least one early holothurian. Small ossicles in the ophiocistioid body wall closely resemble ossicles of holothurians. An early relationship with the holothuroids has been suggested.

Holothuroids. The sea cucumbers, or Holothuroidea, are a group of echinozoans that must have diverged very early from other echinoderm stocks. They clearly exhibit a fivefold radial symmetry characteristic of the phylum but otherwise differ radically from any other echinoderm assemblage. They are greatly elongated along the oral-aboral axis so as to have a subcylindrical form, and they lie on one of their sides identified as ventral because it is in contact with the substratum (Fig. 5*f* and *g*, Fig. 7*e*). The mouth, surrounded by feeding tentacles, lies at one extremity and the anus at the opposite end. The tentacles can assume a variety of forms, from richly branching (dendritic) to shield-shaped, to featherlike, to digitiform (fingerlike). Sea cucumbers are unique among living echinoderms in possessing an internal madreporite and a single gonad. In a small percentage of species, the body is enclosed in a test of conspicuous plates, but in most the body is flexible, the skeleton consisting of microscopic calcareous plates loosely distributed in the dermis. When the animal dies, decay of the uncalcified tissues liberates the skeletal parts, which almost invariably become scattered about. Less than a dozen whole-animal fossil holothuroids have been described. The fossil remains of the holothuroids usually consist only of discrete microscopic ossicles. Even so, about 200 species of these ossicles have been described from rocks ranging from Devonian to Pleistocene. Most of the fossil forms can be referred to one or the other of the orders of extant holothuroids. Sea cucumbers range from rocky shores to the greatest ocean depths. On the deep-sea floor they can be present in enormous numbers and may comprise 95% of the total weight of animals on the seafloor. Typically sea cucumbers ingest the soft sandy to muddy substrata on which they occur. Such feeding may be nonselective mud-swallowing or selective detritus ingestion. The sea cucumbers with branching tentacles are suspension feeders, extending their sticky tentacles into the water and capturing small drifting organisms or organic particles.

About 20 tropical species are prized as food, especially in Asia, and as a consequence these species are more or less fished out throughout the tropics. Six orders of living sea cucumbers are recognized,

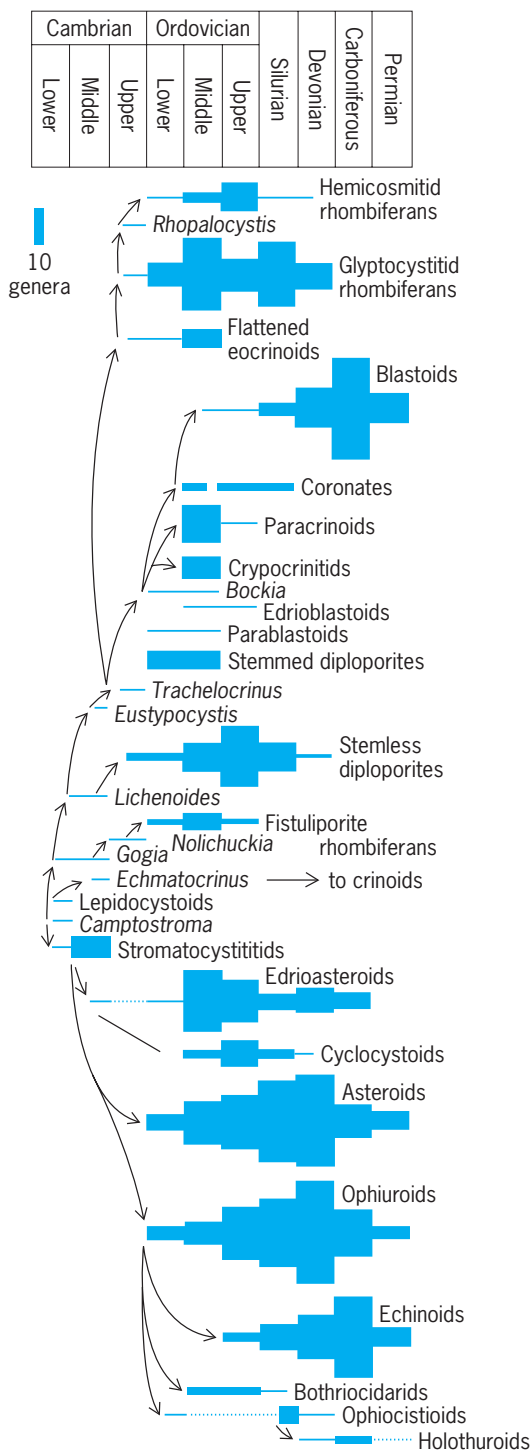


Fig. 9. Phylogeny and diversity range chart for the Paleozoic echinoderms. (After C. R. C. Paul and A. B. Smith, *The early radiation and phylogeny of echinoderms*, *Biol. Rev.*, 59:443-481, 1984)

based upon characters of the feeding tentacles and skeleton. See HOLOTHUROIDEA.

Origin and Phylogeny

The evolution and relationships of the major groups of echinoderms are under active discussion, and several major issues have yet to be resolved. The extraxial/axial theory (EAT) of B. David and R. Mooi offers a new approach to discussion of skeletal homolo-

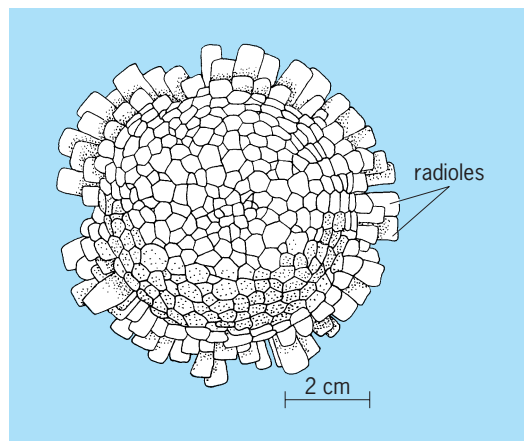
gies. Echinoderms evolved very rapidly near the beginning of the Paleozoic Era, and Lower Cambrian deposits contain such divergent branches of the phylum as homalozoans Helicoplacoidea, Edriasteroidea, and Eocrinoidea (Fig. 9). These are primitive types of echinoderms. Cystoids, crinoids, and blastoids, as well as all recognized main groups of asterozoans and echinozoans (except holothurians), appear in Ordovician strata. During the Paleozoic, numerous well-marked evolutionary trends are discernible in nearly all echinoderm groups, including free-moving forms (especially echinoids) as well as crinozoans. Many small classes of echinoderms became extinct during the Paleozoic, and the surviving groups, especially the crinoids, lost many members at the Late Permian mass extinction. All groups of modern echinoderms have their origin in early Paleozoic stocks, and the lines of their phylogeny are mostly indicated by the fossil record. Echinoids predominate in Mesozoic and Cenozoic echinoderms.

David L. Pawson; Raymond C. Moore;
J. John Sepkoski, Jr.

Bibliography. B. David and R. Mooi, Embryology supports a new theory of skeletal homologies for the phylum Echinodermata, *C.R. Acad. Sci. Paris Ser. 3*, 319:577-584, 1996; L. H. Hyman, *The Invertebrates: Echinodermata*, vol. 4, 1955; D. T. J. Littlewood, Echinoderm class relationships revisited, in *Echinoderm Research*, ed. by H. Roland et al., A. A. Balkema, Rotterdam, 1995; C. G. Messing, Living comatulids, *Paleontol. Soc. Pap.*, 3:3-30, 1997; R. Mooi, Not all written in stone: Interdisciplinary syntheses in echinoderm paleontology, *Can. J. Zool.*, 79:1209-1231, 2001; C. R. C. Paul and A. B. Smith (eds.), *Echinoderm Phylogeny and Evolutionary Biology*, 1988; A. B. Smith, Echinoderm larvae and phylogeny, *Annu. Rev. Ecol. Sys.*, 28:219-241, 1997; A. B. Smith, *Echinoid Palaeobiology*, 1984.

Echinoida

An order of Echinacea with a camarodont lantern, smooth test, imperforate noncrenulate tubercles, ambulacral plates of echinoid type, and shallow



Colobocentrotus atratus, aboral aspect, a Pacific species adapted for life on wave-exposed coral reefs.

branchial slits. There are numerous tropical and temperate species, some of them remarkably adapted to living on coral reefs (see *illus.*). See ECHINOIDEA.

The five included families are principally distinguished by characters of the pedicellariae. The Paraseleniidae are oblong forms with trigeminate ambulacral plates. They range from the Eocene to the present day. The Echinidae possess trigeminate or polyporous plates with the pores in a narrow vertical zone. Strongylocentrotidae are polyporous, with the pores in 2–4 vertical series. The Echinometridae show a variety of forms, which include polyporous types with an oblong test, and trigeminate or polyporous types with a spherical test. See ECHINACEA; ECHINODERMATA.

Howard B. Fell

Bibliography. A. Smith, *Echinoid Palaeobiology*, 1984.

Echinoidea

A class of Echinodermata known as the sea urchins, also including sand dollars, sea biscuits, and heart urchins. In echinoids the body is enclosed in a hard shell, or test, formed from regularly arranged plates that bear movable spines (**Fig. 1**). There are no arms, but radii are represented by five double rows of tube feet arranged as meridians between the upper and lower poles of the body.

There are about 800 living species, and some 5000 fossil species have been recorded, included in 225 genera. Echinoids are classified in 15 orders, which are grouped in three subclasses. Sea urchins range in size from approximately 5 mm (0.2 in.) across the test to 20 cm (8 in.). Although many species are dull or dark in color, some are brilliant shades of purple, red, green, or orange. Others species have particolored striped spines, and deep-sea forms may be white. The earliest echi-

noids occur in the Late Ordovician. Morphological and molecular analyses suggest that the echinoids are most closely related to the holothuroids (sea cucumbers) within the Echinodermata. See ECHINODERMATA; EUECHINOIDEA; SAND DOLLAR; SEA URCHIN.

Relation to humans. Most sea urchins are harmless, although some tropical species have hollow, brittle spines that cause septic wounds if they break off after penetrating the skin. The genus *Araeosoma* carries venomous spines that can inflict dangerous wounds, and a related genus, *Asthenosoma*, is said to be able to kill an adult human. The egg masses, or roe, of sea urchins are highly prized in Asian cuisine. As a consequence, active fisheries for sea urchins and sea urchin farms exist in the United States, Chile, and other countries. See ECHINOTHURIOIDA.

Ecology. Sea urchins occur in all seas from low-tide level downward. *Pourtalesia* reaches a depth of more than 4.5 mi (7250 m) in the Banda Trench. The rounded sea urchins feed mainly on algae, often being partially covered by day under stones, shells, or pieces of seaweed held over the test by tube feet, and emerging at night or at high tide. Heart urchins, sand dollars, and the like live buried in mud or sand, feeding on organic matter associated with ingested sand or selectively feeding on detritus.

Sea urchins move slowly, using muscles at the bases of the spines to swing the spines like stilts. The tube feet are used in concert with spines to aid in locomotion. The suctional tube feet can be used to ascend steep surfaces and act as anchors.

Numerous sea urchin parasites are recorded. Among these are protozoans, nematodes, and gastropods, of which several genera burrow into the test. In a few species, crabs may live in the rectum or on the test, feeding on spines and tube feet. Other animals shelter among the spines or attach themselves to exposed hard parts.

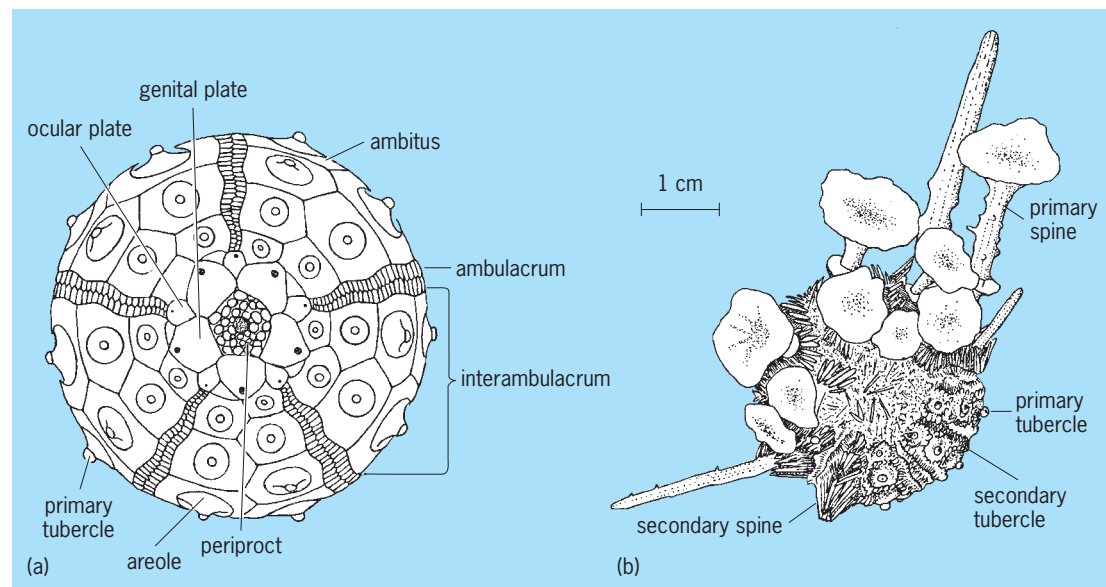


Fig. 1. Structure of the echinoid test of *Goniocidaris parasol*. (a) Naked test. (b) Test with spines.

Skeleton. The skeleton comprises the test and spines.

Test. The test is composed of regularly arranged plates that collectively form a rigid, or occasionally flexible, investing shell. In all echinoids since the Triassic, the test comprises 10 meridional areas, each composed of two vertical columns of plates. The meridians converge above and below at the upper and lower poles. The equatorial zone is termed the ambitus. Of the 10 meridional areas, five are ambulacra (or amb); each amb consists of a double row of plates carrying tube feet. Alternating with the amb are five meridional interambulacra (or interamb). The interamb is usually significantly wider than the amb (Fig. 2).

Paleontology reveals that the test evolved only after much trial and error. The earlier Paleozoic forms showed extreme instability in the number of columns of plates. The cidaroids were the first group in which stability was achieved (in the Permian and Triassic), but a rigid, spherical shape was not adopted until the Jurassic. All of the bilaterally symmetrical irregular echinoids evolved from cidaroid ancestors. See PERISCHOECHINOIDEA.

Spines. Most echinoids carry spines that articulate with tubercles on the test plates and are moved by muscles at the spine bases. Large primary spines articulate with large primary tubercles, smaller secondary spines with secondary tubercles. If a liga-

ment links the spine to the tubercle, the tubercle carries a small hole where the ligament is attached and is termed perforate. Tubercles without such a ligament are imperforate. The spine muscle may attach to an indented pattern on the edge of the boss (a protuberant part) on which the tubercle stands; if so, the tubercle is termed crenulate; if indentations are absent, the tubercle is noncrenulate. The muscle is attached to the test plate on a saucer-shaped depression around the tubercle, the areole. All these features are used in systematic diagnoses of the orders and families.

At the upper pole lies a cirlet of 10 plates, the apical system. Of these, five are at the top of the amb and the five that alternate with them are at the top of the interamb. Each plate atop the interamb carries a gonopore, and is termed a genital plate. One of these also carries a madreporite (a delicately perforated sieve plate that takes in water to the vascular system). The other five are termed ocular plates, an old and misleading name for each ocular plate carries a tentacle, not an eye. The sea urchin adds new plates to its skeleton at the site of the ocular plates. Thus, the youngest plates in the test are nearer to the oculars, and plates more distant from the oculars, that is, toward the mouth, are the oldest. Individual plates grow by adding calcium carbonate to their margins. The resulting growth rings on the plates are seldom correlated with elapsed growth time.

In the spherical (regular) echinoids (Fig. 2a), the anus lies within the apical system on a membrane termed the periproct. These forms are sometimes termed endocyclic. In the sand dollars and heart urchins, the anus becomes displaced outside the apical system, it enters an interamb termed the posterior interamb, and the echinoid is said to have become irregular (exocyclic) [Fig. 2b]. Irregular echinoids tend to be secondarily bilaterally symmetrical, although a five-part radial symmetry is always evident. These features were formerly used in classification and are still valuable in determining trends of development in evolution. Irregular forms may exhibit modification of the apical system. The anus, for example, tends to obliterate the posterior genital plate (and its gonad) as it migrates backward, and the lost structures may not be replaced, or considerable distortion may ensue. See IRREGULARIA; REGULARIA.

Lantern. The mouth lies on the lower (oral) side of the test, surrounded by soft or plated skin, the peristome. In regular forms the mouth is central; in irregular forms it is central or may be displaced to a radius opposite the interamb that contains the anus. The radius containing the mouth is then termed anterior. In regular forms and in some irregular forms, the mouth is furnished with a ring of five powerful jaws, each with one large tooth. The jaws and teeth collectively make up the so-called lantern, first described by Aristotle. The lantern is moved by muscles that are attached to an intumed flange of the test, around the peristome, the perignathic girdle. Variations in the structure of the girdle, lantern, and teeth are used in classification. The teeth may be

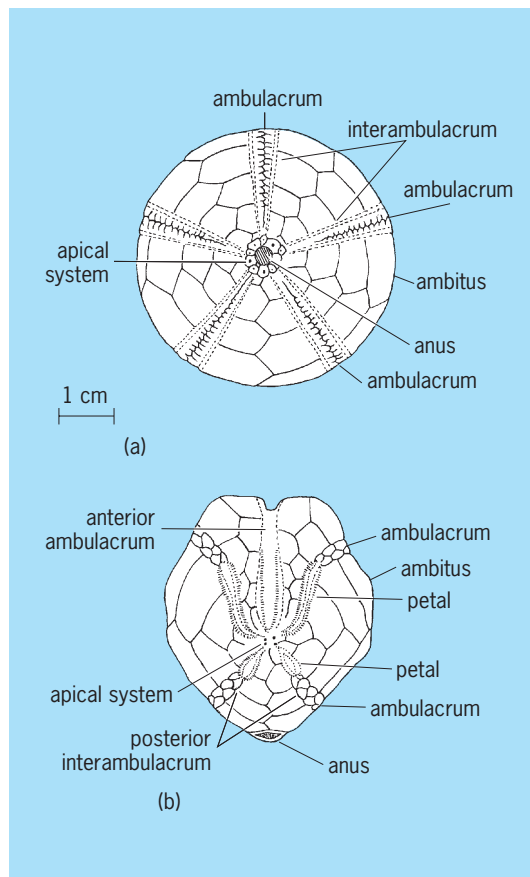


Fig. 2. Tests and related structures of (a) regular (endocyclic) echinoid and (b) irregular (exocyclic) echinoid.

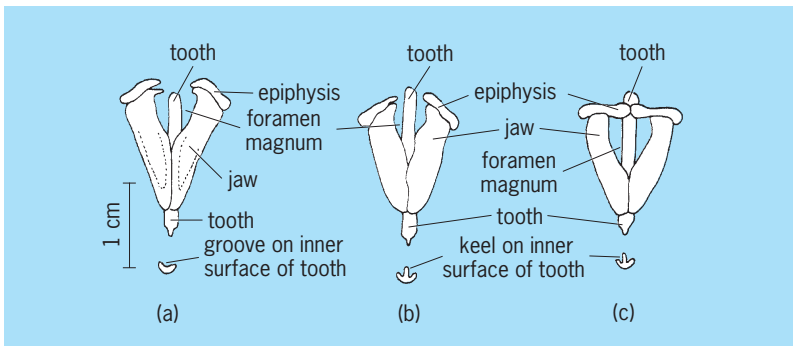


Fig. 3. Lantern structure and types of dentition for each main type. A single interradial jaw is illustrated, and below each is a cross section of the tooth. (a) Aulodont. (b) Stirodont. (c) Camarodont.

grooved or keeled, and the jaw may be partly or completely roofed over by epiphyses (Fig. 3). Three main types of dentition are distinguished: (1) aulodont, in which the teeth are grooved and epiphyses do not meet, so that there is an open foramen magnum in the jaw (Fig. 3a); (2) stirodont, in which the teeth are keeled within and the foramen magnum is open (Fig. 3b); and (3) camarodont, in which the teeth are keeled and the foramen magnum is closed by the epiphyses (Fig. 3c). These characters are useful in taxonomy. However, it was shown by T. Mortensen (during 1928–1951) that parallel dentitional evolution has occurred in various groups, and the three orders formerly based on dentition have been abandoned. As concise morphological descriptive terms, aulodont, stirodont, and camarodont remain valuable, and they are used in that sense in taxonomic diagnoses.

Water-vascular system. This system is highly developed and greatly influences the form of the ambulacral plates. The ring vessel rests upon the lantern and bears small reservoirs—polian vesicles. The stone canal passes upward through the coelom to the madreporite. The radial vessels and their ampullae lie on the inner surface of the amb, within the test. The tube feet alone emerge to the exterior by way of paired pores.

Each tube foot traverses the test wall by means of two pores, termed a pore pair. One pore serves for the outward flow of hydrocoel fluid, the other for inward flow. The pore pairs lie on the amb plates (Fig. 4). In young stages and in the adults of Cidaroida and the irregular echinoids, each amb plate bears only one pore pair. Such amb plates are termed simple (Fig. 4a). In the older stages of noncidaroid regular echinoids, the amb plates tend to fuse into compound plates, which therefore carry more than one pore pair. The most common arrangement is that in which an arc of three pore pairs occurs on a plate. Such a plate is termed trigeminate or oligoporous (Fig. 4b, d, and f). Plates with four or more pore pairs are termed polyporous (Fig. 4c, e, and g). One component of a plate is usually larger than the others and is termed the primary. The position of the primary and the arrangement of the pore pairs provide characters used in taxonomy.

Diadematooid plates may be trigeminate or polyporous with the primary immediately over the lowest element (Fig. 4b and c). Arbacioid plates are similar, but the pore pairs lie in a vertical series, and the secondary elements (demiplates) are rectangular (Fig. 4d and e). Echinoid plates differ from diadematooid plates in having the primary as the lowest element (Fig. 4f and g). The ambulacra in some irregular echinoids tend to change from simple meridians into petal-shaped areas around the upper pole and around the mouth. Such amb, are described as petaloid, and the parts that surround the apical system are called the petals, whereas those around the mouth are termed phyllodes. A further

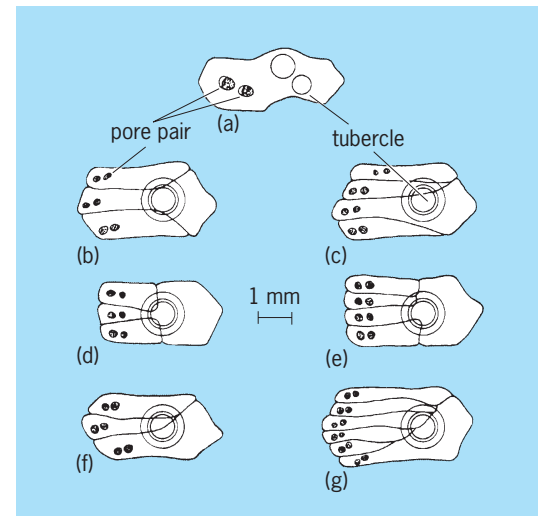


Fig. 4. Types of ambulacral plates as they are used in diagnoses of the orders. (a) Simple. (b) Trigeminate diadematooid. (c) Polyporous diadematooid. (d) Trigeminate arbacioid. (e) Polyporous arbacioid. (f) Trigeminate echinoid. (g) Polyporous echinoid.

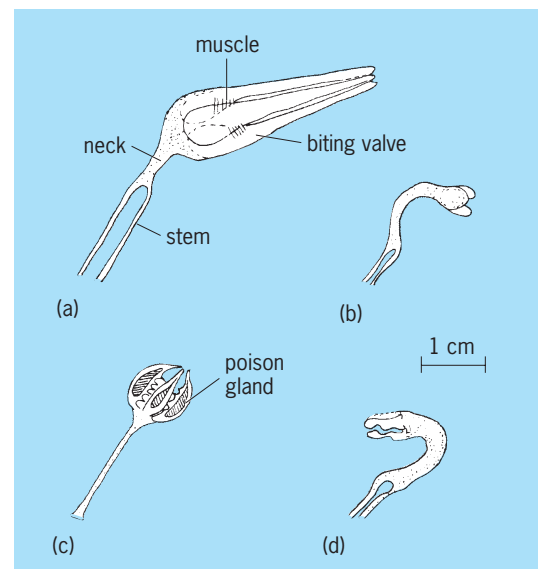


Fig. 5. Pedicellariae of echinoids. (a) Tridentate type. (b) Triphyllous or trifoliate. (c) Globiferous or gemmiform. (d) Ophicephalous.

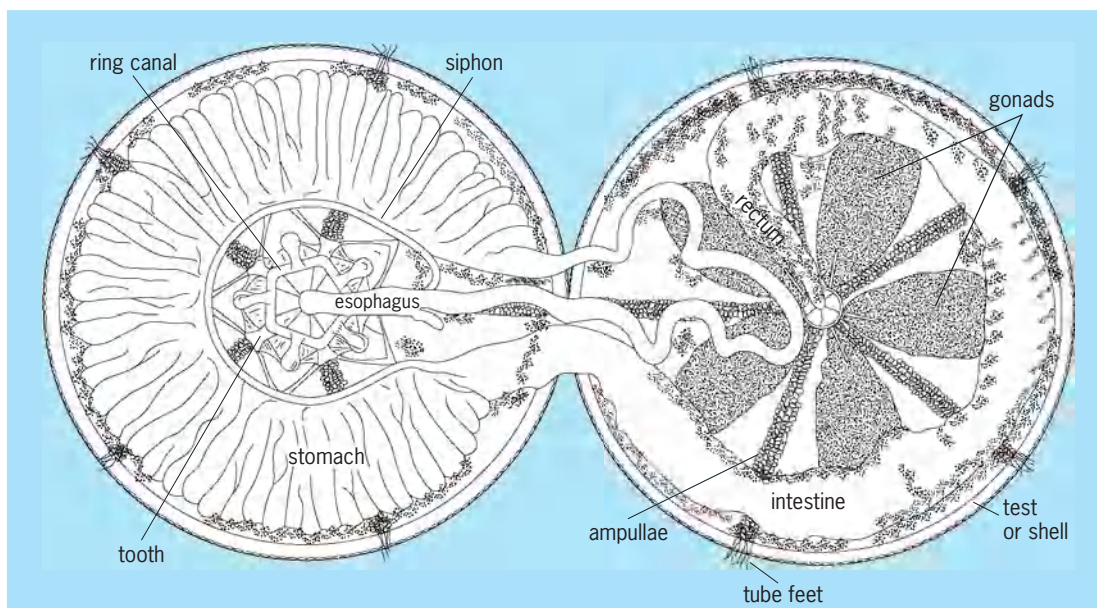


Fig. 6. Echinoid cut horizontally across the ambitus. Left is the adoral hemisphere, gonads removed. (Adapted from W. R. Coe, *Echinoderms of Connecticut, 1912*; republished in T. I. Storer and R. L. Usinger, *General Zoology, 3d ed.*, McGraw-Hill, 1957)

development occurs in the Cassiduloidea, where a flowerlike structure, the floscelle, is present.

The surface of the test plates is usually smooth but sometimes sculptured; that is, a raised pattern of ridges (epistoma) or of ridges and grooves ramifies among the tubercles. See TEMNOPLEROIDA.

Some echinoids have special respiratory organs or gills attached to the peristome. The gills, if present, usually notch the margin of the peristome. The notches, termed gill slits, may vary in size and shape or may be absent. They provide characters used in taxonomy. Many irregular forms use the tube feet as respiratory organs.

Small grasping organs, pedicellariae, are well developed in echinoids, in which they take the form of a beak carried on a stalk. The beak is made up of three (rarely two or four) movable jaws, operated by muscles and sometimes provided with venom glands. They respond to tactile stimuli and seize any small organisms or particles that may touch the skin. The intrusive material is passed from one to the other until one ambital pedicellaria drops it off the urchin. The chief types are shown in Fig. 5.

Nervous system. The nervous system follows the same pattern as the water-vascular system. The tube feet, especially those near the mouth, may serve as tactile and taste organs. A few sea urchins have photosensitive eyespots on the upper surface of the test, scattered in the ectoderm. The so-called ocular pores are not sensitive to light. Minute spherical stalked bodies attached to the skin in some echinoids are believed to be organs of balance.

Alimentary system. The alimentary canal is tubular. It lies in the coelom, attached to the wall of the test by mesenteries. The stomach runs from the esophagus in a counterclockwise coil (as viewed from above), and the intestine retraces the route in reverse. The

rectum passes upward to the anus (Fig. 6).

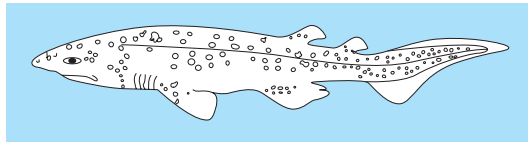
Reproduction. Like sea stars, sea urchins seem to become sexually mature after 1 year but continue to grow for several years. The life span is unknown but may average 5–6 years in medium-sized species, and some deep-sea and cold-water species may live much longer. The sexes are normally separate, although fertile hermaphrodites are known. The gonads may be reduced in number to only three or two. If a pelagic larva is present, it is an echinopluteus. In some species the young are brooded among the spines or in sunken petals, as in *Abatus*. Spines and some other organs are regenerated after injury. Autotomy (the process of self-amputation of appendages) and autoevisceration (expulsion of the digestive tract and associated organs through the anus) are unknown.

David L. Pawson; Howard B. Fell

Bibliography. J. W. Durham and R. V. Melville, A classification of echinoids, *J. Paleontol.*, 31(1):242–272, 1957; G. Hendler et al., *Sea Stars, Sea Urchins and Allies: Echinoderms of the Florida Keys and the Bahama Islands*, 1995; D. T. J. Littlewood, Echinoderm class relationships revisited, in *Echinoderm Research*, ed. by H. Roland et al., A. A. Balkema, Rotterdam, 1995; R. Mooi, Not all written in stone: Interdisciplinary syntheses in echinoderm paleontology, *Can. J. Zool.*, 79:1209–1231, 2001; T. Mortensen, *Monograph of the Echinoidea*, 5 vols., 1928–1951; D. Nichols, *Echinoderms*, 1962; S. P. Parker (ed.), *Synopsis and Classification of Living Organisms*, 2 vols., 1982; J. S. Pearse and R. A. Cameron, Echinodermata: Echinoidea, in A. C. Giese, J. S. Pearse, and V. B. Pearse (eds.), *Reproduction of Marine Invertebrates*, vol. 6, pp. 514–662, Boxwood Press, Pacific Grove, 1991; A. B. Smith, Echinoderm larvae and phylogeny, *Annu. Rev. Ecol. Sys.*, 28:219–241, 1997; A. Smith, *Echinoid Palaeobiology*, 1984.

Echinorhiniformes

An order of sharks (bramble sharks) represented by one family, one genus, and two species. The family (Echinorhinidae) differs from all other sharks in possessing the following combination of characters: body and fins covered with thornlike denticles; two small spineless dorsal fins placed posteriorly to the pelvic fin origin; anal fin lacking; minute spiracles placed well behind the eyes; teeth similar in both jaws; and caudal fin without a subterminal notch. Maximum length is 2 m (6.6 ft).



Bramble shark, *Echinorhinus brucus*. (Courtesy of J. S. Nelson, 2006)

Echinorhinus brucus (see **illustration**) occurs in the western North Atlantic, Indian, and western Pacific oceans, whereas *E. cookei* is limited to parts of the Pacific Ocean. *Echinorhinus brucus* has relatively few large denticles, differing from *E. cookei*, with relatively numerous small denticles. Both species inhabit temperate waters of continental and insular shelves and slopes, usually to depths greater than 200 m (660 ft) and less than 900 m (2950 ft). Their food consists of fishes, including small sharks, and crabs. Bramble sharks are ovoviparous (producing eggs that develop internally and hatch before or soon after extrusion).

The family was previously placed in the Squaliformes. See CHONDRICHTHYES; ELASMOBRANCHII; SELACHII.

Herbert Boschung

Bibliography. M. R. de Carvalho, Higher-level elasmobranch phylogeny, basal squalians, and paraphyly, pp. 35–62, in M. L. J. Stiassny, L. R. Parenti, and G. D. Johnson (eds.), *Interrelationships of Fishes*, Academic Press, San Diego, 1996; J. S. Nelson, *Fishes of the World*, 4th ed., Wiley, 2006.

Echinothurioida

An order of Echinothuriacea with solid and hollow primary radioles, diademoid ambulacral plates, noncrenulate tubercles, and the anus within the apical system. The extant members of the family Echinothuriidae are all deep-water forms. The Echinothuriidae have a large flexible test which collapses into a disk at atmospheric pressure, and the middle element of the diademoid plates is much larger than the other two elements. Some species carry venomous spines. Echinothuriids range from the Late Jurassic to present day. See ECHINODERMATA; ECHINOIDEA; PYGASTEROIDA.

Howard B. Fell

Echinozoa

A subphylum of free-living echinoderms in which the body is essentially globoid with meridional symmetry. They lack arms, brachioles, or other appendages, and do not at any time exhibit pinnate structure. The Echinozoa range from the Early Cambrian to the present day. There are four classes that can definitely be placed here: (1) Edrioasteroidea, Lower Cambrian to lower Carboniferous echinozoans in which the mouth and anus were both directed upward and ambulacra (three to five in number) served as food-collecting areas; (2) Echinoidea, the existing and fossil sea urchins, originating in the Middle Ordovician; (3) Ophiocystioidea of the Lower Ordovician through Middle Devonian, with a domed aboral surface cover with large polygonal plates and a flat adoral side with a mouth and five radiating ambulacra; and (4) Holothuroidea, the existing and fossil sea cucumbers, which apparently first appeared in the Devonian. Two other extinct echinoderm classes can be placed here for convenience: (5) Lower Cambrian Camptostromatoidea, conical or domal animals with plates of varying size that overlapped on the lower theca; and (6) Lower Cambrian Helicoplacoidea, cylindrical animals with a spirally plated test and three ambulacra on the surface. The latter class may be the sister group of both echinozoans and crinozoans. See CAMPTOSTROMATOIDEA; ECHINOIDEA; HELICOPLAOIDEA; HOLOTHUROIDEA.

The oldest definite echinozoans are stromatocystitids of the Lower and Middle Cambrian. This group, which may have camptostromatids as its sister group, may have been ancestral to other edrioasteroids and, perhaps, other echinozoans. The asterozoans (asteroids and ophiuroids) may have been derived from the echinozoans, making the subphylum paraphyletic. See ECHINODERMATA.

Howard B. Fell; J. John Sepkoski, Jr.

Bibliography. R. C. Moore and C. Teichert (eds.), *Treatise on Invertebrate Paleontology*, pt. T, 1978; C. R. C. Paul and A. B. Smith, The early radiation and phylogeny of echinoderms, *Biol. Rev.*, 59:443–481, 1984; C. R. C. Paul and A. B. Smith (eds.), *Echinoderm Phylogeny and Evolutionary Biology*, 1988.

Echiura

A small group of unsegmented, wormlike, marine animals once linked with the Sipuncula and Priapulida under the taxon Gephyrea. Since 1940, echiurans have been regarded as a separate phylum of the animal kingdom with affinities to the segmented annelid worms. They range from the tropical to the polar seas, burrowing in sea-floor sediments from the intertidal area to depths of 30,000 ft (9140 m). See ANNELIDA; PRIAPULIDA; SIPUNCULA.

Morphology. The echiuran body consists of a round or elongated trunk with an anterior prostomium (see **illustration**). The mouth is located at the anterior end of the trunk where the tongue-like prostomium

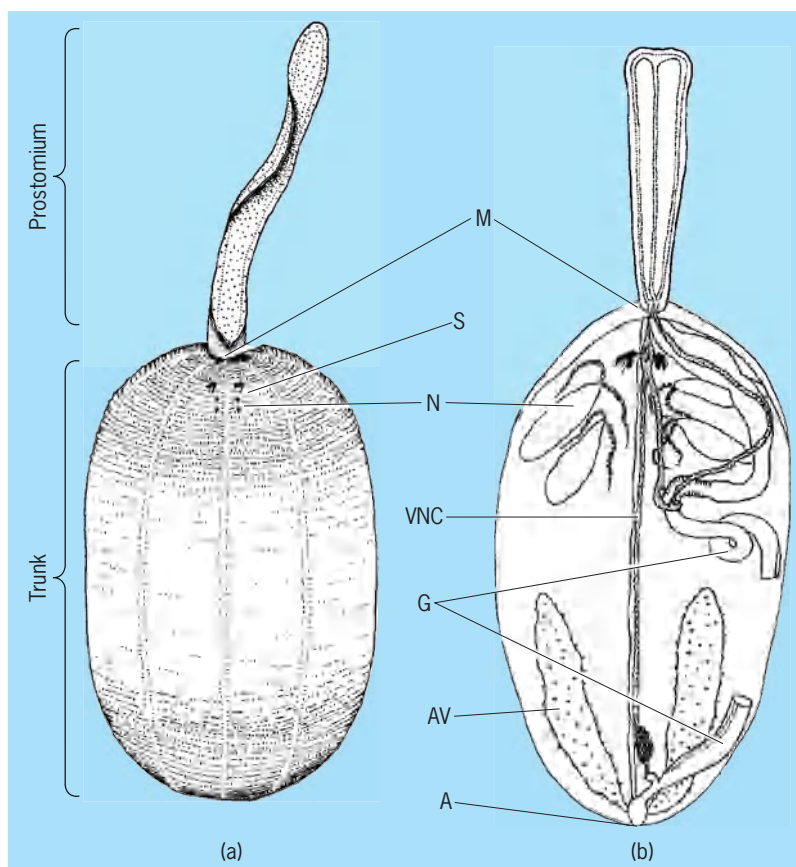
is attached. The anus is at the posterior end. The body wall of the trunk may be thin enough to see red blood and internal organs within it, or it may be green and opaque as in the Bonelliidae. The muscles of the body wall may be a uniform sheet or visibly gathered into varying numbers of bundles that run from anterior to posterior (front to back). There is usually a pair of setae (hooked bristles) situated ventrally a short distance behind the mouth, and in the genera *Echiurus* and *Urechis* there are one or more rings of setae that surround the anus. Within the burrow the trunk undergoes rhythmic contractions that move water through the burrow.

The prostomium is muscular and capable of considerable extension. It is important in the identification of species but often becomes accidentally detached and lost during collection. This can make identification difficult. In most echiurans the prostomium is elongated with a rounded or blunt anterior end. In the bonellid echiurans the tip of the prostomium is forked. Living in burrows in marine sediments, echiurans extend the prostomium out onto the surface where it collects and transports nutrient-rich sediments to the mouth to be consumed. Species in the genus *Urechis* have a short stubby prostomium. These species secrete a mucus net within their burrow and then pump water through it to filter out food.

The mouth, located at the base of the prostomium, leads into the gut, which is divided into distinct regions designed to systematically process the food. Two characteristic organs, the anal vesicles, are present at the posterior end of the gut near the anus. These organs may be tubular or branched and can extend for varying distances toward the anterior end of the body cavity. Water is intermittently drawn in through the anus to fill these sacs and then is expelled. Both respiratory and excretory functions have been proposed for these organs. Nephridia (organs typically known as having an excretory function in most invertebrates) are present on the ventral body wall behind the anterior setae. In echiurans the nephridia have been modified to gather ripe gametes (eggs or sperm) from the body cavity and store them prior to spawning. A ventral nerve cord runs from the anterior to the posterior end of the trunk. *See DIGESTION (INVERTEBRATE); EXCRETION.*

The sexes are separate and indistinguishable in most echiurans. In the Bonelliidae, however, the sexes are separate and very dissimilar, the male being minute and living within the female's egg storage organ. There they fertilize the eggs as they are spawned. The sex of bonellid offspring is determined largely by environmental factors.

Classification. Traditional morphological and embryological evidence has suggested a phylogenetic (evolutionary) relationship between these animals and the annelids, but the echiurans have remained a separate phylum. Embracing this taxonomic status, A. C. Stephen and S. J. Edmonds recognized three orders—Echiuroinea, Xenopneusta, and Heteromyota—the last order with one problematic species attached. Most species are included in the



Echiuran anatomy. (a) External anatomy, ventral view. (b) Internal anatomy, dorsal view. A, anus; AV, anal vesicles; G, gut (middle section removed for clarity); M, mouth; N, nephridia (egg and sperm storage organs); S, setae; VNC, ventral nerve cord (modified from A. C. Stephen and S. J. Edmonds, 1972).

Echiuroinea, which is divided into families Bonelliidae and Echiuridae.

Contemporary evidence used to classify animals and help determine their relationships (systematics and phylogeny) now comes from molecular data (for example, DNA and RNA nucleotide sequences) and sophisticated microscopic techniques. The DNA sequence of one important gene in an echiuran has been found to be more similar to the same gene in annelid worms than it is to other closely related but nonannelid animals. This finding suggests that the echiurans may indeed be very closely related to the segmented annelid worms; however, their unsegmented body remained problematic. Recent fluorescent microscopic techniques focused on the structure of the ventral nerve cord have revealed structural repetition in the anterior to posterior axis of the organ that suggests the presence of a more subtle segmental body plan in echiurans. Taken together, these two findings have revived a hypothesis that considers echiurans to be members of a distinct class within the phylum Annelida. Classification below this taxonomic level has not been formally proposed. However, if one were, it would likely retain the orders currently recognized. John F. Pilger; Mary E. Rice

Bibliography. R. Hessling, The use of confocal laser-scanning microscopy (cLSM) and 3d reconstruction

for systematic zoology—on the phylogeny of the Echiura, *Microsc. Today*, 1(5):8-11, 2001; R. Hessling and W. Westheide, Are Echiura derived from a segmented ancestor?—Immunohistochemical analysis of the nervous system in developmental stages of *Bonellia viridis*, *J. Morphol.*, 252:100-113, 2002; D. McHugh, Molecular evidence that echiurans and pogonophorans are derived annelids, *Proc. Nat. Acad. Sci.*, 94:8006-8009, 1997; J. F. Pilger, Echiura, in F. W. Harrison and M. E. Rice (eds.), *Microscopic Anatomy of Invertebrates*, vol. 12, Wiley-Liss, New York, 1993; J. F. Pilger, Sipunculans and echiurans, in S. F. Gilbert and A. M. Raunio (eds.), *Embryology: Constructing the Organism*, Sinauer Associates, Sunderland, MA, 1997; A. C. Stephen and S. J. Edmonds, *The Phyla Sipuncula and Echiura*, British Museum (Natural History), 1972.

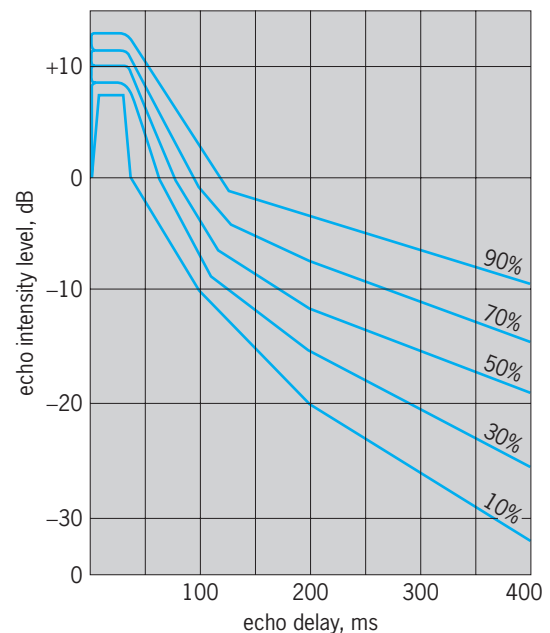
Echo

A sound wave which has been reflected or otherwise returned with sufficient magnitude and time delay to be perceived in some manner as a sound wave distinct from that directly transmitted. Multiple echo describes a succession of separately indistinguishable echoes arising from a single source. When the reflected waves occur in rapid succession, the phenomenon is often termed a flutter echo.

Echoes and flutter echoes are generally detrimental to the quality of the acoustics of rooms. They may be minimized through the proper selection of room dimensions, room shape, and distribution of sound-absorbing materials. Flutter echoes, for example, may be minimized by making the opposite walls of a room nonparallel or by making one of the walls highly sound absorptive. For a more complete discussion of the effect of room shape on echoes see ARCHITECTURAL ACOUSTICS

Individual echoes will generally be heard only if the time delay between the direct sound wave and a reflected sound wave is in excess of approximately 50 milliseconds. In a large auditorium where the reverberation time is of the order of seconds, many reflected waves will be present. These reflections will not be troublesome if their intensity is sufficiently below that of the initial sound. A relation between the approximate percent of listeners detecting some alteration in the acoustics of a room and the sound pressure level of the echo as a function of the delay time following the arrival of the direct sound is shown in the **illustration**. See REVERBERATION.

Echoes have been put to a variety of uses in measurement problems. For example, the distance between two points can be measured by timing the duration required for a direct sound originating at one location to strike an object at the other point and to return an echo to the location of the initial source. For the application of this principle to the detection of submerged objects see SONAR



Estimated percent of listeners disturbed by the echo related to the intensity of the echo (relative to initial sound intensity) and the time delay between the initial sound and perception of the echo. (After R. H. Bolt and P. E. Doak, *J. Acoust. Soc. Amer.*, 22:507-509, 1950)

Ultrasonic echo techniques have achieved considerable success in nondestructive testing of materials. When an ultrasonic wave is propagated through a metal, the presence of a crack or other flaw will cause a sound wave, or echo, to be reflected back to the initial source location. Observing the time delay between the original sound and the perception of the echo permits the location of the flaw to be determined. This technique has been found particularly useful in such problems as examining metal casting for internal defects and determining the location of cracks in pipes or welded structures. It has also been employed to locate brain tumors in humans. For further discussion of the application of echo reflection techniques see REFLECTION OF SOUND; SOUND; ULTRASONICS.

William J. Galloway

Echo sounder

A marine instrument used primarily for determining the depth of water by means of an acoustic echo. A pulse of sound sent from the ship is reflected from the sea bottom back to the ship, the interval of time between transmission and reception being proportional to the depth of the water.

Echo sounders, sometimes called fathometers, are used by vessels for navigational purposes, not only to avoid shoal water, but as an aid in fixing position when a good bathymetric chart of the area is available. Some sensitive instruments are used by commercial fishers or marine biologists to detect schools of fish or scattering layers of minute marine life. Oceanographic survey ships use echo sounders for charting the ocean bottom. **Figure 1** shows an

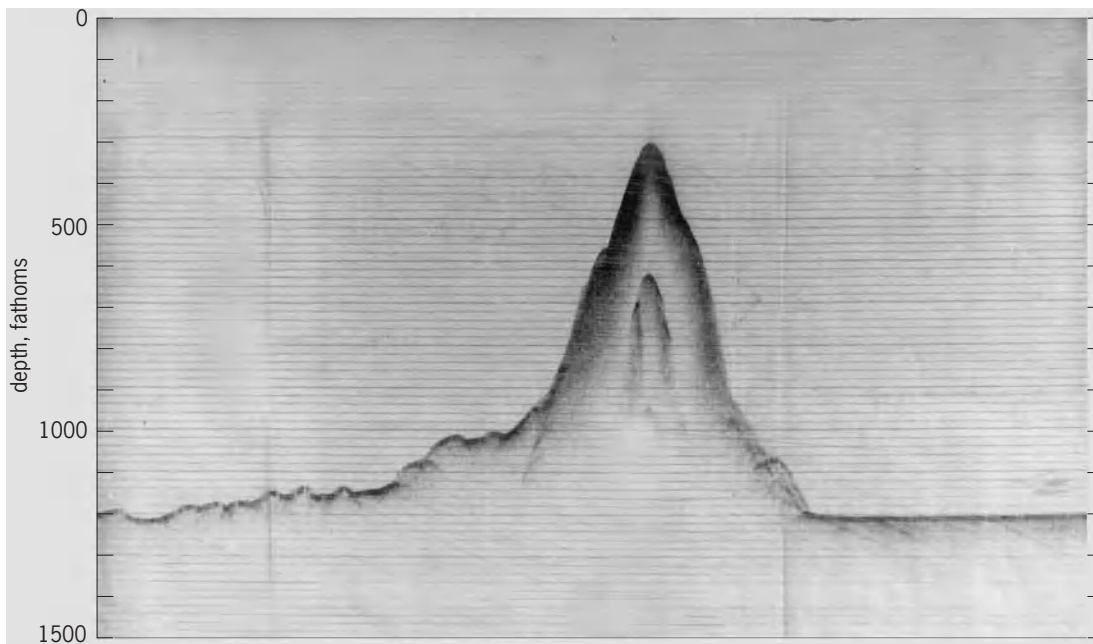


Fig. 1. Echo-sounder record of a seamount in the deep ocean. 1 fathom = 1.8 m. (Woods Hole Oceanographic Institute)

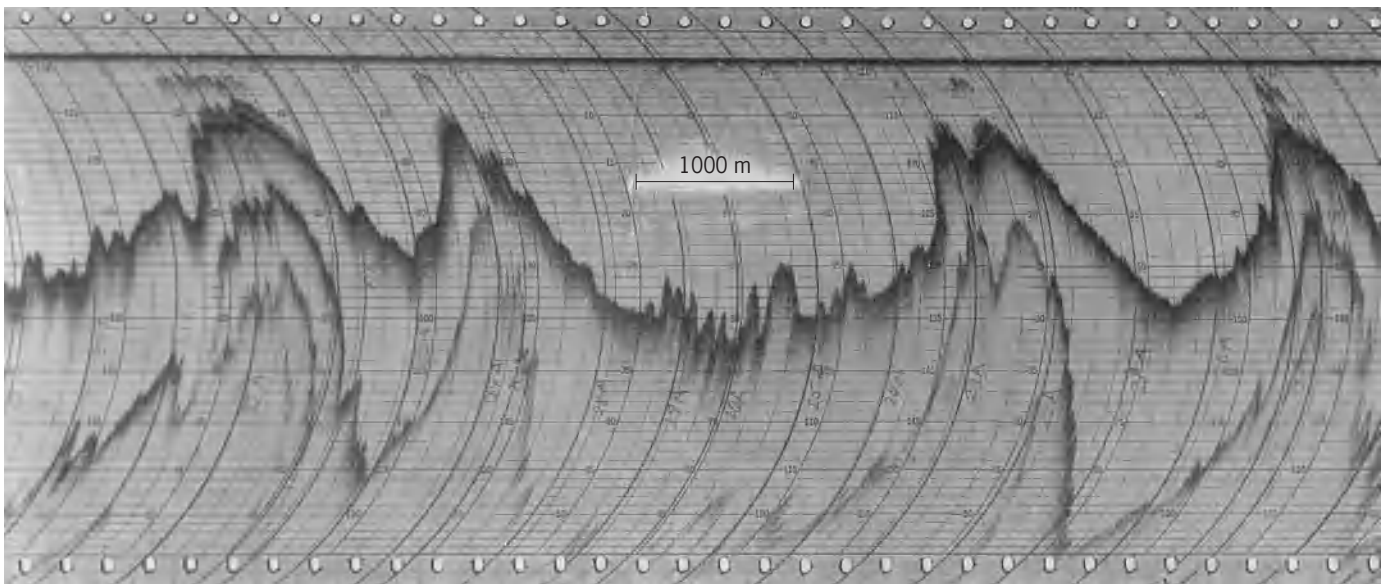


Fig. 2. Typical record of bottom profile obtained by an echo sounder. 1 m = 3.3 ft. (USCGS)

echo-sounder record obtained by oceanographers of a seamount (undersea mountain). See SCATTERING LAYER.

An echo sounder is really a type of active sonar. It consists of a transducer located near the keel of the ship which serves (in most models) as both the transmitter and receiver of the acoustic signal; the necessary oscillator, receiver, and amplifier which generate and receive the electrical impulses to and from the transducer; and a recorder or other indicator which is calibrated in terms of the depth of water. An echo sounder actually measures time differences, so some average velocity of sound must be assumed in order to determine the depth. The frequency generally employed is in the low ultrasonic

range (20,000–30,000 Hz). The depth display may be given by a trace-type recorder which supplies a continuous permanent record (Fig. 2) or, as is the case with less expensive commercial instruments, it may be a dial-type or digital display, giving instantaneous depth of water. See MARINE GEOLOGY; SONAR; UNDERWATER SOUND.

Robert W. Morse

Echocardiography

A diagnostic procedure that uses ultrasound at a frequency of 2.5–10 MHz to provide an image of the heart. It is based on the principle that the interface between tissues of different acoustical impedance

causes the ultrasound to be reflected to the transducer, which spends a fraction of each second receiving these echoes. There are many interfaces between blood and the various structures in the heart that contact blood, such as the heart walls, valves, and great vessels. Also, the surface of the heart reflects ultrasound because it is surrounded by the lungs, which are filled with air.

The ultrasound apparatus can estimate the distance from the transducer to the tissue interface that is reflecting the sound, because the velocity of sound in tissue is relatively constant and the time required to send and receive sound can be measured (distance = velocity \times time). A graphic representation of the locale of each interface can then be displayed on a video monitor or digital screen, and as the beam moves through the heart a composite two-dimensional image of cardiac structures can be constructed. The sampling rate of the ultrasound device is adequate to display heart motion in real time. See BIOMEDICAL ULTRASONICS; ULTRASONICS.

Procedure. An ultrasound image of the heart is generated on a monitor (see **illustration**). The image is generated by electronically moving the ultrasound beam repeatedly through an arc. The transducer is usually applied to the anterior chest by using a coupling gel devoid of air. This procedure is referred to as transthoracic echocardiography. Small transducers can also be attached to probes placed in the esophagus behind the heart—a procedure known as transesophageal echocardiography; during heart surgery the transducer can be placed directly on the heart—a procedure known as epicardial echocardiography; and a small transducer can be placed on a catheter to visualize the inside of vessels—a procedure known as intravascular ultrasound. The closer the transducer is to the heart, the higher the resolution of the resulting images. Axial resolution (parallel to the path of the sound beam) for all echocardiographic systems is quite high [approximately 1 mm

(0.04 in.)], but lateral resolution (perpendicular to the path of the sound beam) is reduced and deteriorates the farther away the structure is from the transducer. The frame rate of the image is about 30 per second.

Two-dimensional echocardiography. Two-dimensional echocardiography is completely harmless and provides excellent real-time images of the heart. It is very useful for determining the anatomy and function of heart valves, detecting abnormal amounts of pericardial fluid, and defining the complex anatomy of congenital heart defects. It is effective for estimating size and function of heart chambers, thickness and mass of the heart wall, and size of the great vessels. The equipment is portable, and imaging can be accomplished during stresses, such as exercise. Thus, two-dimensional echocardiography is frequently employed to evaluate suspected or overt heart disease.

Blood flow. Echocardiography can also be used to assess blood flow in the heart by employing the Doppler principle. The movement of a structure during an examination performed by ultrasound results in a shift in the frequency of the returning sound wave (echo). The shift in frequency is proportional to the velocity of the moving object. Since the ultrasound frequency transmitted (f_0), the speed of sound in tissue (C), and the angle of the sound beam to the direction of movement (θ) are known, velocity (V) can be calculated by measuring the frequency shift (Fd), using Eq. (1).

$$V = \frac{Fd \times C}{2f_0 \times \cos \theta} \quad (1)$$

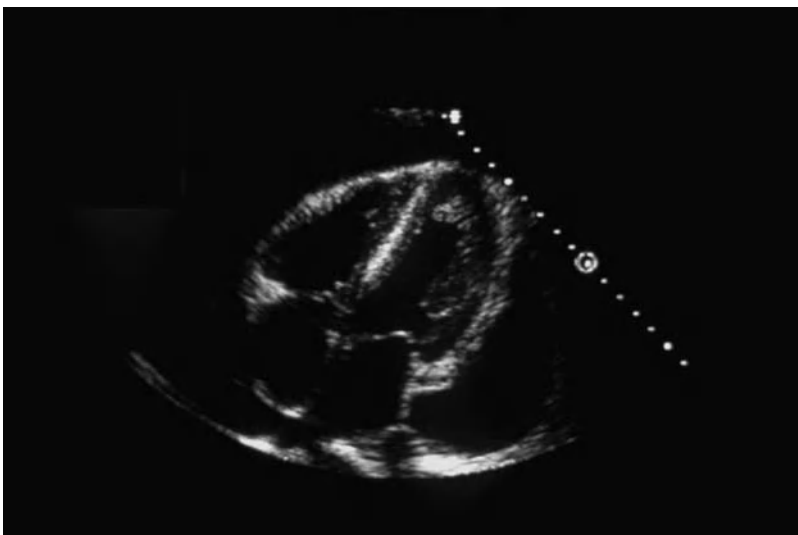
Blood, the object usually interrogated by Doppler ultrasound of the heart, moves because of a pressure gradient (ΔP) in the heart. The pressure gradient can be calculated from the simplified Bernoulli equation (2) if velocity (V) is known.

$$\Delta P = 4V^2 \quad (2)$$

Doppler systems. Three basic types of Doppler echocardiographic systems are available on most ultrasound machines: continuous-wave, pulsed, and color-flow.

In a continuous-wave (CW) Doppler system, sound is continuously transmitted and received. The advantage of CW Doppler is that the very high velocities of severe structural heart defects are accurately quantitated. The disadvantage is range ambiguity, that is, the inability to locate the origin of the high-velocity blood flow along the sound beam.

Pulsed Doppler avoids range ambiguity by being able to range-gate the returning echoes so that only certain areas in the sound beam are sampled. However, pulsed Doppler cannot resolve very high velocities because the pulse repetition frequency limits Doppler frequency-shift sampling. Pulsed and CW Doppler are useful for detecting and localizing increased blood velocities associated with obstructions to flow. See DOPPLER EFFECT.



Two-dimensional echocardiogram from a position near the apex of the ventricles. The left ventricle is on the right, separated from the left atrium below by the mitral valve. The right ventricle is on the left, and the right atrium is below it with the tricuspid valve in between.

Pulsed Doppler can also be directed at the myocardium to measure the velocity of muscle contraction or relaxation. Currently, these tissue Doppler velocities are used to assess relaxation properties of the left ventricle.

The third system is color-flow Doppler. Small sample volumes [1–2 mm³ (0.04–0.08 in.³)] throughout the two-dimensional echocardiographic image are sampled for velocity. Only velocities fast enough to result from blood movement are displayed, superimposed on the two-dimensional echocardiographic image and color-encoded. Movement (velocity) toward the transducer is usually red, and movement away is usually blue. The relative velocity is indicated by the intensity of the color. Turbulent flow is a mosaic color (for example, green). The resultant images display blood moving through the heart in real time as a color on the black-and-white anatomy of the heart and its internal structures. Color-flow Doppler is excellent for displaying abnormal flow due to valve leakage or other structural defects.

M mode. The complete echocardiographic examination also includes selective sampling along certain radians in the arced two-dimensional image, and displaying of the movement of the structures crossed by the radian over time in an analog fashion. In this format, sampling rates can be increased to 1000 per second so that complex motion patterns of a structure can be elucidated. This format is called time-motion echocardiography or M-mode echocardiography. It was actually the first (1950s) commercially available form of ultrasound for examining the heart. It is adjunctive to two-dimensional and color-flow Doppler echocardiography for answering specific questions about heart structural dimensions and function. See HEART DISORDERS.

Michael H. Crawford

Bibliography. M. D. Cheitlin et al., ACC/AHA/ASE 2003 guideline update for the clinical application of echocardiography, *J. Amer. Coll. Cardiol.*, 42:954–970, 2003.

Echolocation

The biological sonar that bats, porpoises, and certain other animals use to navigate without the visual system. Several different groups of animals have evolved the ability to perceive objects by emitting sounds and hearing the echoes that the objects reflect to their ears. The locations and characteristics of the objects are represented by acoustic properties of the echoes, and the ears and auditory systems of these animals act as the sonar receiver. The sense of hearing is specialized for converting echo information into displays of objects, which are perceived as acoustic images that guide the animal's behavior. The best-known examples of echolocating animals are bats (Microchiroptera) and porpoises and toothed whales (Cetacea). However, several other kinds of mammals (some flying foxes, shrews, and rats) and birds (oilbirds and cave swiftlets) also can echolocate. See SONAR.

Echolocation sounds. Bats produce their ultrasonic sonar sounds from the larynx and broadcast them through the open mouth or through a specialized transmitting antenna formed around the nostrils. Many kinds of bats (leaf-nosed bats, horseshoe bats) have oddly shaped nasal structures for sound emission that appear bizarre unless it is realized that they guide the beam of the sonar transmissions. Porpoises produce their sonar sounds from structures located beneath the blowhole on the top of the head (through which they breathe) and project them into the water through the rounded, protuberant forehead, which contains acoustically specialized tissue. The sonar sounds of porpoises and whales are very brief impulses, or clicks, which contain a wide range of ultrasonic frequencies, all occurring at the same instant. The rate of emission of these sonar clicks depends in part upon the distance to objects that interest the animal, and can vary from several sounds per second to several hundred. Porpoises use their sonar to find fish and presumably to perceive objects beyond the relatively restricted range of vision under water.

The sonar sounds of bats vary from one family and species to another and consist of short (0.5 to about 10 ms) frequency-modulated (FM) signals that are often coupled with short or long (0.5 to over 100 ms) constant-frequency signals (Fig. 1). Their frequency span, or bandwidth, is important because the quality or sharpness of the acoustic images produced from FM echoes improves as the bandwidth increases. Under favorable conditions, bats that use wide-band FM signals can perceive the distance to a target with an acuity of a fraction of a millimeter, for example. Constant-frequency sonar signals are well suited for initial detection of targets because they concentrate their energy near a known frequency to which the sonar receiver can be turned, and long-duration constant-frequency signals are useful for determining the velocity of a target from the Doppler shift of echoes. Bats that emit constant-frequency signals perceive small changes in the frequency of echoes and use this ability to determine how fast they are moving relative to objects. Bats that emit only FM signals observe changes in a target's position from one echo to the next to determine relative velocity. See DOPPLER EFFECT; FREQUENCY MODULATION.

Most bats are insectivorous and hunt at night. They emit sonar sounds in a stereotyped pattern correlated with stages in the interception of a flying insect (Fig. 1). The entire interception, from initial detection to capture, lasts less than a second. Bats produce sonar sounds at a low emission rate while searching for prey; the rate increases progressively while approaching the target after it has been detected. During the approach, the bat tracks the target and decides whether to attempt interception. During the final stage of interception, which lasts only about a tenth of a second, the sonar sounds appear at a rate as high as 100 sounds per second. The constant-frequency or FM structure of the sounds (Fig. 1) represents the acoustic strategy that different species of

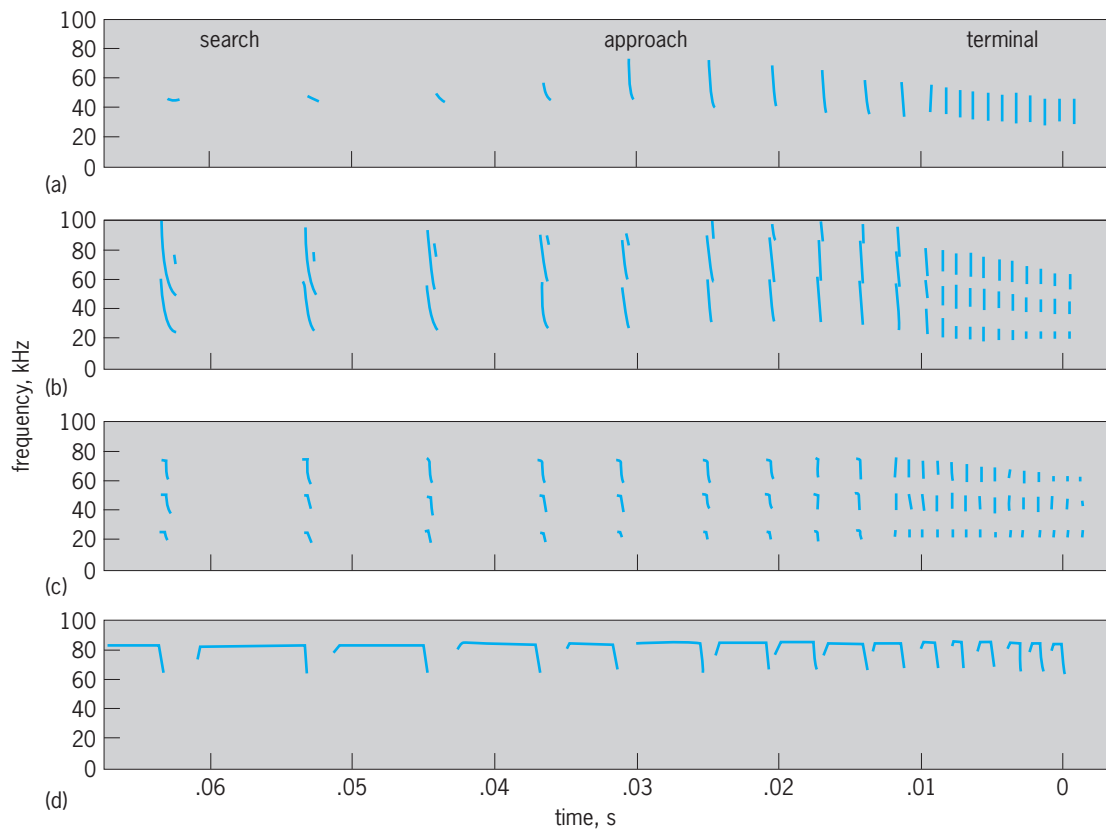


Fig. 1. Sound spectrograms representing sonar signals of (a) *Tadarid brasillensis*, (b) *Eptesicus fuscus*, (c) *Pteronotus personatus*, and (d) *Rhinolophus ferrumequinum* during search, approach, and terminal stages of pursuit of prey. Detection, identification, and tracking usually take place in less than 1 s. Spectrograms b and c use only FM signals, whereas a and d use constant-frequency signals during the search and change to FM signals during the approach and terminal buzz. (After J. A. Simmons, F. B. Fenton, and M. J. O'Farrell, *Echolocation and pursuit of prey of bats*, *Science*, 203:16–21, 1979)

bats have evolved to solve the problems of finding prey.

Auditory images. The images in echolocation are derived from either frequency-modulated or constant-frequency echoes.

Frequency-modulated images. The distance to a target, or target range, is represented acoustically by the delay of FM echoes. Bats first encode echo delay as the timing of nerve discharges but eventually represent delay by the location of neural responses within a discrete patch of the auditory cortex that functions as a biological map of target range. Some FM echo-locating bats combine their perception of target range with their perception of the shape of targets by perceiving the distance to different parts of the target along the same scale as the distance to the target as a whole.

The bat in **Fig. 2** emits an FM sonar sound and receives an echo from the target, here shown to be a flying moth. The target reflects echo components from parts located at different ranges (for simplicity, one echo component from its head and one from its wings). The target as a whole has a range, r , and the target's structure consists of the range separation, Δr . These target features are represented in the echo by the echo's delay, t , and the time separation of the two components within the echo, Δt , because the echo contains a separate component from each part

of the target, at delay t and delay $t + \Delta t$. The bat's inner ear represents the FM emission and the echo as spectrograms analogous to those shown in Fig. 1. The spectrograms in Fig. 2 are patterns of mechanical excitation produced by the inner ear itself. Mechanical excitation then evokes neural discharges that convey the spectrograms along parallel neural channels into the auditory nervous system. The time separation of the echo components, Δt , creates notches in the spectrum of the echo as a whole separated by the frequency interval, Δf , which equals $1/\Delta t$. These spectral notches appear in the neural spectrogram of the echo as a "scalloping" pattern of the array of discharges (solid circles) caused by the lengthening of discharge latency at the lower amplitudes of the echo within the notch as compared to outside the notch. Some neural discharges near the center of the notch may even fail to occur at all (open circles) if the echo is very weak in the notch.

The delay of echoes is presented in the bat's brain by spectrogram delays, t_f , at each frequency in the emission and the echo (Fig. 2). Spectrogram delays are further processed to be encoded, frequency by frequency, in terms of the location of neural activity on the target-range map in the auditory cortex, which is the bat's auditory display of targets. Spectrogram delays are somehow integrated across frequency to determine the distance to the

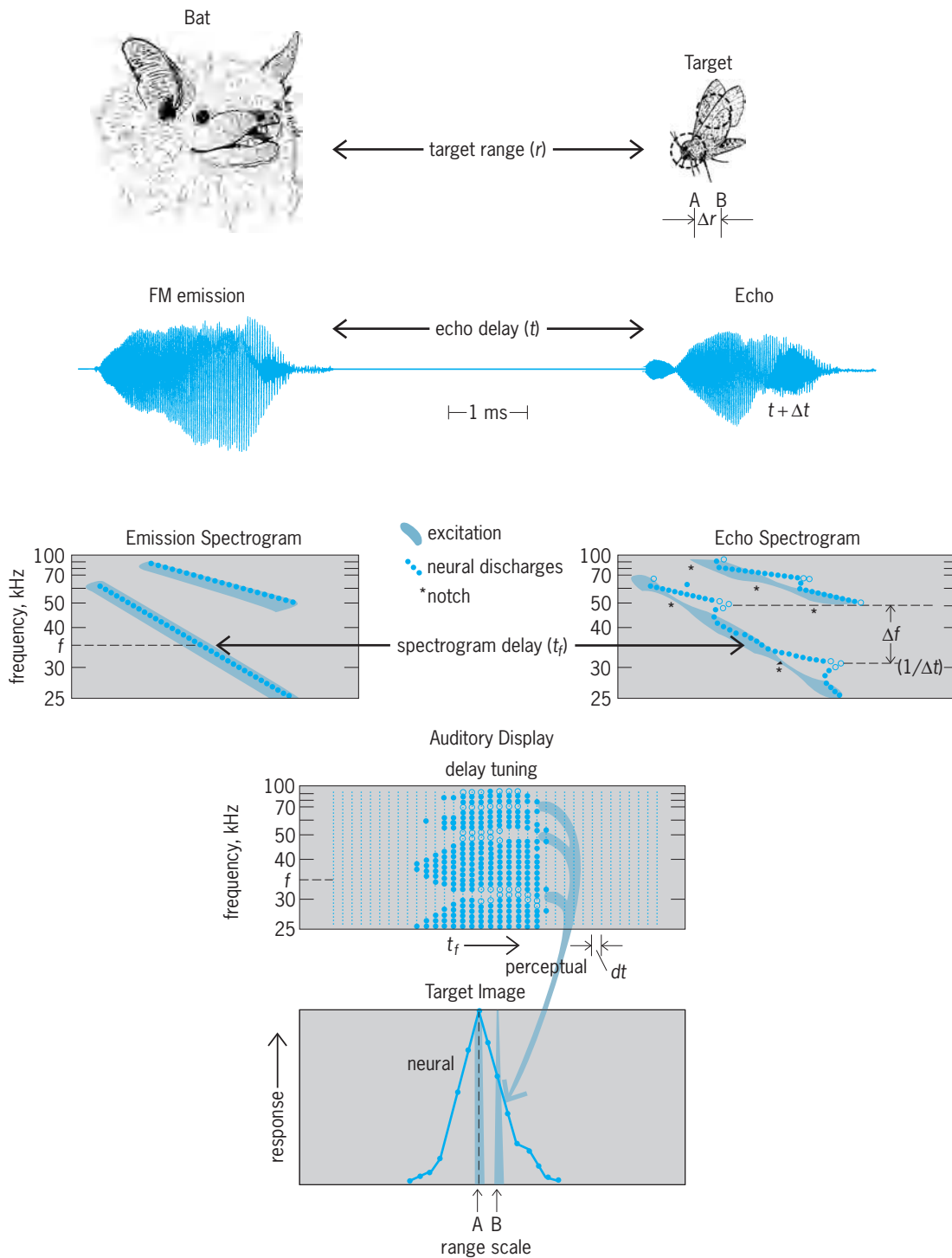


Fig. 2. Stages by which bats convert echoes of FM sonar sounds into images that incorporate the distance and the shape of the target. Details are discussed in the text. (After J. A. Simons, C. F. Moss, and M. Ferragamo, *Convergence of temporal and spectral information into acoustic images of complex sonar targets perceived by the echolocating bat, *Eptesicus fuscus**, *J. Compar. Physiol. A*, 1989)

leading edge of the target from the earliest component of the echo, at delay t . The auditory display in Fig. 2 shows the distribution of neural activity along the range axis of the auditory cortex. The arrow leading from the auditory display to the target image shows the effect of the locations of the notches in the spectrum that produce a wavy pattern of activity in the cortex across frequencies. This wave pattern is converted into an estimate of the

location of the echo's second component, at delay $t + \Delta t$. Somehow the single-peaked distribution of neural activity across the range map is converted into a two-peaked pattern to represent the target's image as two discrete components (shown as peaks labeled perceptual, A and B) located at ranges r and $r + \Delta r$. This conversion is an example of the brain creating a single, unitary psychological dimension out of neural representations that treat different features of stimuli

in different ways. It illustrates an important function of the brain that is not understood.

Constant-frequency images. The use by bats of constant-frequency echoes to determine target velocity from echo Doppler shifts is another example of neural information processing leading to perception. The bat's inner ear represents different frequencies of sound by the location of neural discharges within the population of neurons in the auditory nerve. The frequency of a constant-frequency echo activates only those neurons tuned to that particular frequency. This activity is passed along from the auditory nerve to higher centers of hearing in the brain. When the bat receives a particular constant-frequency echo frequency, neurons in a corresponding location on frequency maps in the auditory cortex respond, whereas other neurons do not. One of these cortical frequency maps becomes a map of target velocity for constant-frequency echoes. In this case, no neural computations are needed to create the frequency map because it already exists at the inner ear. However, neural computations are needed to adjust the relationship between each frequency and the size of the frequency change corresponding to each velocity, because the Doppler shifts are proportional to the frequency being shifted. Other frequency maps segregate neural responses to constant-frequency echoes to facilitate detection of the target's fluttering motions, such as the wing beats of insects, for identification of prey. See CHIROPTERA; PHONORECEPTION; ULTRASONICS.

James A. Simmons

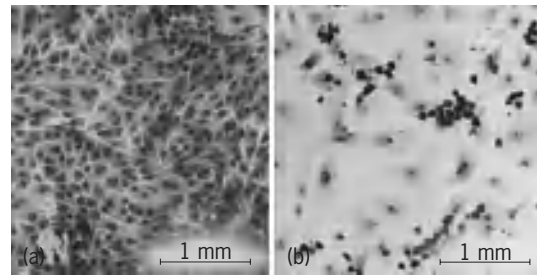
Bibliography. D. R. Griffin, *Listening in the Dark*, 1986; P. Nachtigall and P. W. B. Moore, *Animal Sonar Systems: Processes and Performance*, 1988; J. A. Simmons, M. B. Fenton, and M. J. O'Farrell, Echolocation and pursuit of prey by bats, *Science*, 203:16–21, 1979; J. A. Simmons, C. F. Moss, and M. Ferragamo, Convergence of temporal and spectral information into acoustic images of complex sonar targets perceived by the echolocating bat, *Eptesicus fuscus*, *J. Compar. Physiol. A*, 1989.

Echovirus

Members of the Picornaviridae family, genus *Enterovirus*. The name is derived from the term enteric cytopathogenic human orphan virus. More than 34 antigenic types exist. Only certain types have been associated with human illnesses, particularly with aseptic meningitis and febrile disease, with or without rash. Their epidemiology is similar to that of other enteroviruses. See ENTEROVIRUS; PICORNAVIRIDAE.

Echoviruses resemble polioviruses and coxsackieviruses in size (about 28 nanometers) and in many other properties. They are nonpathogenic for newborn mice, rabbits, or monkeys, but cytopathogenic for monkey kidney and other tissue cultures (see **illus.**). See COXSACKIEVIRUS; TISSUE CULTURE.

Diagnosis is made by isolation and typing of the viruses in tissue culture. Antibodies form during convalescence. See ANIMAL VIRUS. Joseph L. Melnick



Monkey kidney tissue culture. (a) Normal culture. (b) Three days after infection with echovirus type 1. (Courtesy of J. L. Melnick, Baylor College of Medicine)

Bibliography. N. J. Schmidt and W. Emmons (eds.), *Diagnostic Procedures for Viral, Rickettsial, and Chlamydial Infections*, 7th ed., 1995; Z. Woldehiwet and R. Miodrag (eds.), *Rickettsial and Chlamydial Diseases of Domestic Animals*, 1993.

Eclipse

The occultation (obscuring) of one celestial body by another. Solar and lunar eclipses take place at syzygies of the Sun, Earth, and Moon, when the three bodies are in a line. At a solar eclipse, the Moon blocks the view of the Sun as seen from the Earth. At a lunar eclipse, the Earth's shadow falls on the Moon, darkening it, and can be seen from wherever on Earth the Moon is above the horizon. See SYZYGY.

Eclipses of the Sun could be seen from other planets as their moons are interposed between the planets and the Sun, though their superposition is not as coincidental in angular size as it is for the Earth-Moon system. Eclipses of the moons of Jupiter are well known, occurring whenever the moons pass into Jupiter's shadow. An eclipse of the Sun visible from Uranus that occurred in 2006 was imaged by the *Hubble Space Telescope*. Certain binary stars are known to eclipse each other, and the eclipses can be followed by measuring the total light from the system. See BINARY STAR; ECLIPSING VARIABLE STARS; JUPITER.

Related phenomena are transits, such as those of Mercury and Venus, which occur when these planets cross the face of the Sun as seen from Earth. They are much too small to hide the solar surface. Transits of the Earth will be seen from spacecraft. Occultations of stars by the Moon are commonly seen from Earth, and are studied to monitor the shape and path of the Moon; a solar eclipse is a special case of such an occultation. Occultations of stars by planets and by asteroids are now increasingly studied; the rings of Uranus were discovered from observations of such an occultation, and the atmospheres of Titan and Pluto are best known from occultation studies. See OCCULTATION; TRANSIT (ASTRONOMY); URANUS.

Solar Eclipses

A solar eclipse can be understood as an occultation of the Sun by the Moon or, equivalently, the Moon's shadow crossing the Earth's surface. The darkest

part of the shadow, from which the Sun is entirely hidden, is the umbra (Fig. 1). The outer part of the shadow, from which part of the Sun can be seen, is the penumbra.

Solar eclipses can be central, in which the Moon passes entirely onto the solar disk as seen from Earth, or partial, in which one part of the Sun always remains visible. Central eclipses can be total, in which case the Moon entirely covers the solar photosphere, making the corona visible for the period of totality, or annular, in which case the Moon's angular diameter is smaller than that of the Sun because of the positions of the Earth and Moon in their elliptical orbits. At an annular eclipse, a bright annulus of photospheric sunlight remains visible; it is normally thousands of times brighter than the corona, leaving the sky too blue for the corona to be seen.

The plane of the Moon's orbit is inclined by 5° to the plane of the Earth's orbit (the ecliptic), so the Moon's shadow commonly passes above or below the Earth each month at new moon. But two to five times each year, the Moon's shadow reaches the Earth, and a partial, annular, or total eclipse occurs. The Moon is approximately 400 times smaller than the Sun but is also approximately 400 times closer, so its angular diameter in the sky is about the same as the Sun's. Thus the Moon fits approximately exactly over the photosphere, making the phenomenon of a total eclipse especially beautiful.

Phenomena. The partial phases of a total eclipse visible from the path of totality last over an hour. In the minute or two before totality, shadow bands—low-contrast bands of light and dark caused by irregularities in the Earth's upper atmosphere—may be seen to race across the landscape. As the Moon barely covers the Sun, photospheric light shines through valleys on the edge of the Moon, making dots of light—Baily's beads—that are very bright in contrast to the background. The last Baily's bead gleams so brightly that it appears as a jewel on a ring, with the band made of the corona; this appearance is known as the diamond-ring effect (Fig. 2). It usually lasts for 5–10 s, and in the clearest skies for as long as 40 s.

During the diamond-ring effect, the solar chromosphere becomes visible around the limb of the Moon, glowing pinkish because most of its radiation is in the form of emission lines of hydrogen, mostly the red hydrogen-alpha line. Its emission-line spectrum apparently flashes into view for a few seconds, and is called the flash spectrum. As the advancing limb of the Moon covers the chromosphere, the corona becomes fully visible (Colorplate a). Its shape is governed by the solar magnetic field; common are equatorial streamers and polar tufts. At the maximum of the solar activity cycle, so many streamers exist that the corona appears round when it is seen in projection, as viewed from Earth. At the minimum of the solar activity cycle, only a few streamers exist so that the corona appears more elongated in projection. See SOLAR CORONA; SUN.

Totality (Fig. 3) lasts from an instant up through somewhat over 7 min, 30 s. At its end, the phenom-

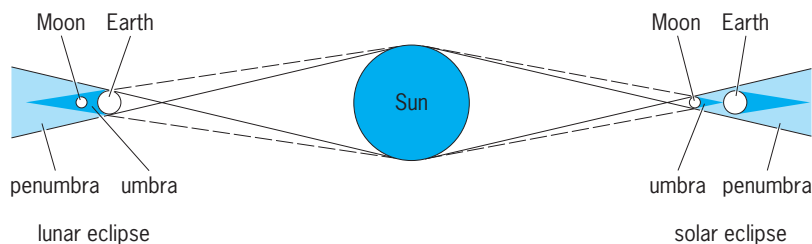


Fig. 1. Circumstances of solar and lunar eclipses (not to scale).

ena repeat in reverse order, including chromosphere, diamond ring, Baily's beads, shadow bands, and the partial phases.

The corona is now monitored continuously in x-rays and in ultraviolet and visible light by orbiting spacecraft such as the *Solar and Heliospheric Observatory (SOHO)*; ultraviolet and visible light), the *GOES Solar X-ray Images*, and the *Transition Region and Coronal Explorer (TRACE)*; ultraviolet). The spacecraft *Solar-B*, a successor to the *Yohkoh* spacecraft, was launched in late September 2006 to continue such studies. Eclipse observations can provide ground truth for comparison with these missions on the days of eclipses. See SATELLITE (ASTRONOMY).

Positions and timing. The paths of the Sun and Moon in the sky intersect at two points, the ascending node and the descending node. Only when both the Sun and the Moon are near a node can an eclipse occur. Thus eclipse seasons take place each year, whenever the Sun is near enough to the node so that an eclipse is possible. Each eclipse season is 38 days long. Because the Sun's gravity causes the orientation of the Moon's elliptical orbit to change with an 18.6-year cycle, the nodes slide along the ecliptic and a cycle of two eclipse seasons—an eclipse year—has a period of 346.6 days, shorter than a solar year. See MOON.



Fig. 2. Diamond-ring effect at the beginning of totality during the solar eclipse of March 29, 2006, observed from Kastellorizo, Greece. (Jay M. Pasachoff/Science Faction, Williams College Eclipse Expedition)



Fig. 3. Solar eclipse of April 8; 2005, observed during totality from a ship in the mid-Pacific Ocean, west of the Galápagos Islands. (Jay M. Pasachoff and Dava Sobel/Science Faction)

There must be at least one solar eclipse each eclipse season, so there are at least two each year. There may be as many as five solar eclipses in a calendar year (technically, a year), though most of these will be partial. Some of the partials will occur near the ends of the eclipse seasons, and it is also possible for there to be three eclipse seasons during a solar year. Adding lunar eclipses (including penumbral lunar eclipses, which may not be noticeable), there may be seven eclipses in a year.

Saros. An important coincidence relates lunar months and eclipse years. A total of 223 lunar months (technically, synodic months, the period of the phases) takes 6585.32 days. A total of 19 eclipse years—passages of the Sun through the same node of the Moon's orbit—takes 6585.78 days, and 242 nodical months—passages of the Moon through the node—take 6585.36 days. (Nodical months are also called draconic months, after the ancient Chinese dragon once thought to have been devouring the Sun at a solar eclipse.) Thus eclipses appear with this period of 18 years $11\frac{1}{3}$ days (plus or minus a day, depending on leap years), a period known as the saros. Further, 239 periods of the variation of distance of the Moon from the Earth, the anomalistic month, is 6585.54 days, so the relative angular sizes of the Sun and Moon are about the same at this interval. (The anomalistic month differs from the nodical and sideral months because the orientation of the Moon's elliptical orbit drifts around in its orbital plane.)

As a result of the saros, almost identical eclipses recur every 18 years $11\frac{1}{3}$ days. The significance of the $\frac{1}{3}$ day is that the Earth rotates one-third of

the way around, and the eclipse path is shifted on the Earth's surface. Thus the June 30, 1973, 7-min eclipse in Africa was succeeded in a saros series by the July 11, 1991, eclipse in Hawaii, Mexico, and Central and South America, which reached maximum duration of 6 min 54 s in Mexico, and will lead to the 6 min 40 s eclipse of July 22, 2009. After a saros, the Sun is slightly farther west than its original position, and the Moon is slightly north or south, depending on whether it is near an ascending or a descending node, so the eclipses in a saros drift from north to south or from south to north, starting near one pole and departing from the other. A complete series takes 1244 to 1514 years.

Motion over Earth's surface. The Moon's shadow travels at approximately 2100 mi/h (3400 km/h) through space. The Earth rotates in the same direction that the shadow is traveling; at the Equator, the resulting motion of the Earth's surface makes up about 1040 mi/h (1670 km/h), making the speed of the eclipse across the Earth's surface 1060 mi/h (1730 km/h). When eclipses cross higher latitudes, the speed of motion associated with the Earth's rotation is not as high, so the eclipse speed is even higher. The supersonic Concorde took advantage of an equatorial eclipse in 1973 to keep up with totality for 74 min. Ordinary jet aircraft cannot keep up with eclipses, so the term eclipse chasing is almost never accurate.

Scientific value. Even with advances in space technology, total solar eclipses are the best way of seeing the lower corona. Coronagraphs, telescopes for which special shielding and internal occulting allow observation of the corona from a few sites in the world on many of the days of the year, are limited to the innermost corona or to use of special filters or polarization. See CORONAGRAPH; SOLAR CORONA.

Coronagraphs have been sent into orbit, notably aboard *Skylab*, the *Solar Maximum Mission*, and *SOHO*, but limitations in spacecraft control lead to the necessity of overoccluding the photosphere, cutting out the inner corona. For example, the coronagraph on *Solar Maximum Mission* occulted 1.75 times the solar diameter. The innermost coronagraph, no longer in operation, on *SOHO*'s Large Angle Spectrographic Coronagraph (LASCO) overoccluded to 1.1 solar radii, and LASCO's two remaining coronagraphs image the Sun from 1.5 to over 30 solar radii. Other types of observations from space also apply to the corona, such as imaging x-ray observations. Further, any space observation is extremely expensive, over 100 times more for *SOHO* than for a given extensive eclipse expedition. So it is useful and necessary to carry on ground-based eclipse observations, ground-based coronagraph observations, and space-based observations to get the most complete picture of the Sun. Further, independent derivations of certain basic quantities must be carried out to provide trustworthy results; there have been cases in which data reduced from space observations had to be restudied because the need for different calibrations had shown up in eclipse work.

Colorplate *b* is a merger of eclipse and *SOHO* spacecraft images, and shows how ground-based images during a total solar eclipse fill a doughnut-shaped region between the inner and outer spacecraft images that is not visible from Earth at other times. The central image from *SOHO* shows the Sun's disk at temperatures around 60,000–80,000 K through a filter showing helium gas. The outer *SOHO* image shows gas at the millions of degrees typical of the Sun's corona, with LASCO's coronagraphs hiding the bright inner corona. Filling in the doughnut ring allows coronal streamers to be traced from their roots on the Sun's surface, through the eclipse corona, and into the solar wind, the expanding outer corona. See SOLAR WIND.

Related eclipse studies include use of the advancing edge of the Moon to provide high spatial resolution for radio observations of the Sun and, historically, of celestial radio sources. The atmospheric effects of the removal of incident sunlight from the Earth's atmosphere have also been studied at eclipses.

Historically, the test of the deflection of starlight carried out at the eclipse of 1919 and repeated in 1922 was the first verification of Einstein's general theory of relativity. These results were restudied for the 1979 Einstein centennial, and their accuracy improved. The experiment is a very difficult one, and has been attempted at some eclipses, notably 1970 in Mexico and 1973 in Africa, without improving on the early results. But the effect has been verified to higher accuracy by studies in the radio part of the spectrum and by the observation of gravitational lenses, so optical eclipse tests are no longer necessary. See GRAVITATIONAL LENS; RELATIVITY.

Annular eclipses. Central eclipses in which the Moon is sufficiently far from the Earth that it does not cover the solar photosphere are annular. The Moon's umbra is about $232,000 \pm 4000$ mi ($374,000 \pm 6400$ km) in length while the Moon's distance is $237,000 \pm 16,000$ mi ($382,000 \pm 25,000$ km), so the umbra sometimes falls short of reaching the Earth's surface. Annular eclipses, like total eclipses, occur about every 18 months on average.

Since the corona is 1,000,000 times fainter than the photosphere, if even 1% of photosphere is showing the corona is overwhelmed by the blue sky and cannot be seen. So most annular eclipses are of limited scientific use.

The annular eclipse of 1984, visible from the southeastern United States, provided 99.8% coverage. The major scientific work carried out involved detailed timing of the Baily's beads in order to assess the size of the Sun. Such work has also been carried out at total eclipses. Some of the results, compared with solar diameters deduced from historical eclipse paths, seem to show a possible shrinking of the Sun by a measurable amount in a time of decades or centuries, which would lead to impossibly large effects on geological time scales; this has led to the suggestion that the Sun could be oscillating in size. But the question of whether any real effect is present has not

been settled. The 1984 annular eclipse was so close to total that the corona could even be briefly seen and photographed, although no scientific studies of the corona were made. In the May 10, 1994, annular eclipse, widely viewed across the United States, only 94% of the Sun's diameter was covered, more typical of annular eclipses. The 1999 annular eclipse, whose path crossed Australia, provided 99% coverage, but since the Sun is about 1,000,000 times brighter than the full moon, that still left about 10,000 times more light than on a moonlit night.

TABLE 1. Total solar eclipses, 2005–2024[†]

Calendar date	Maximum duration of totality (min: s)	Geographic region of visibility
Apr. 8, 2005	0:42	Starts annular in western Pacific Ocean, becomes total in eastern Pacific, then becomes annular again for Costa Rica, Panama, Colombia, and Venezuela
Mar. 29, 2006	4:07	Eastern Brazil, Atlantic Ocean, Ghana, Togo, Benin, Nigeria, Niger, Chad, Libya, Egypt, Turkey, Russia
Aug. 1, 2008	2:27	Northern Canada, Greenland, Arctic Ocean, Russia, Mongolia, China
July 22, 2009	6:39	India, Nepal, Bhutan, Burma, China, Pacific Ocean
July 11, 2010	5:20	South Pacific Ocean, Easter Island, Chile, Argentina
Nov. 13, 2012	4:02	Northern Australia, south Pacific Ocean
Nov. 3, 2013*	1:40	Annular only at beginning and end of path: Atlantic Ocean, Gabon, Congo, Zaire, Uganda, Kenya, then annular in Ethiopia
Mar. 20, 2015	2:47	North Atlantic Ocean, Faeroe Islands, Arctic Ocean, Svalbard
Mar. 9, 2016	4:10	Indonesia (Sumatra, Borneo, Sulawesi, Halmahera), Pacific Ocean
Aug. 21, 2017	2:40	Pacific Ocean, United States (Oregon, Idaho, Wyoming, Nebraska, Missouri, Illinois, Kentucky, Tennessee, North Carolina, South Carolina), Atlantic Ocean
July 2, 2019	4:33	South Pacific Ocean, Chile, Argentina
Dec. 14, 2020	2:10	Pacific Ocean, Chile, Argentina, south Atlantic Ocean
Dec. 4, 2021	1:55	Antarctica
Apr. 20, 2023*	1:16	Total except at beginning and end of path: South Indian Ocean, western Australia, Indonesia, Pacific Ocean
Apr. 8, 2024	4:28	Pacific Ocean, Mexico, United States (Texas, Oklahoma, Arkansas, Missouri, Kentucky, Illinois, Indiana, Ohio, Pennsylvania, New York, Vermont, New Hampshire, Maine), southeastern Canada, Atlantic Ocean

*These are hybrid or annular-total eclipses. They are total along the middle portions of their paths but annular at the ends.

[†]Data courtesy of M. Littmann, K. Willcox, and K. Espenak, *Totally—Eclipses of the Sun*, Oxford University Press, 1999.

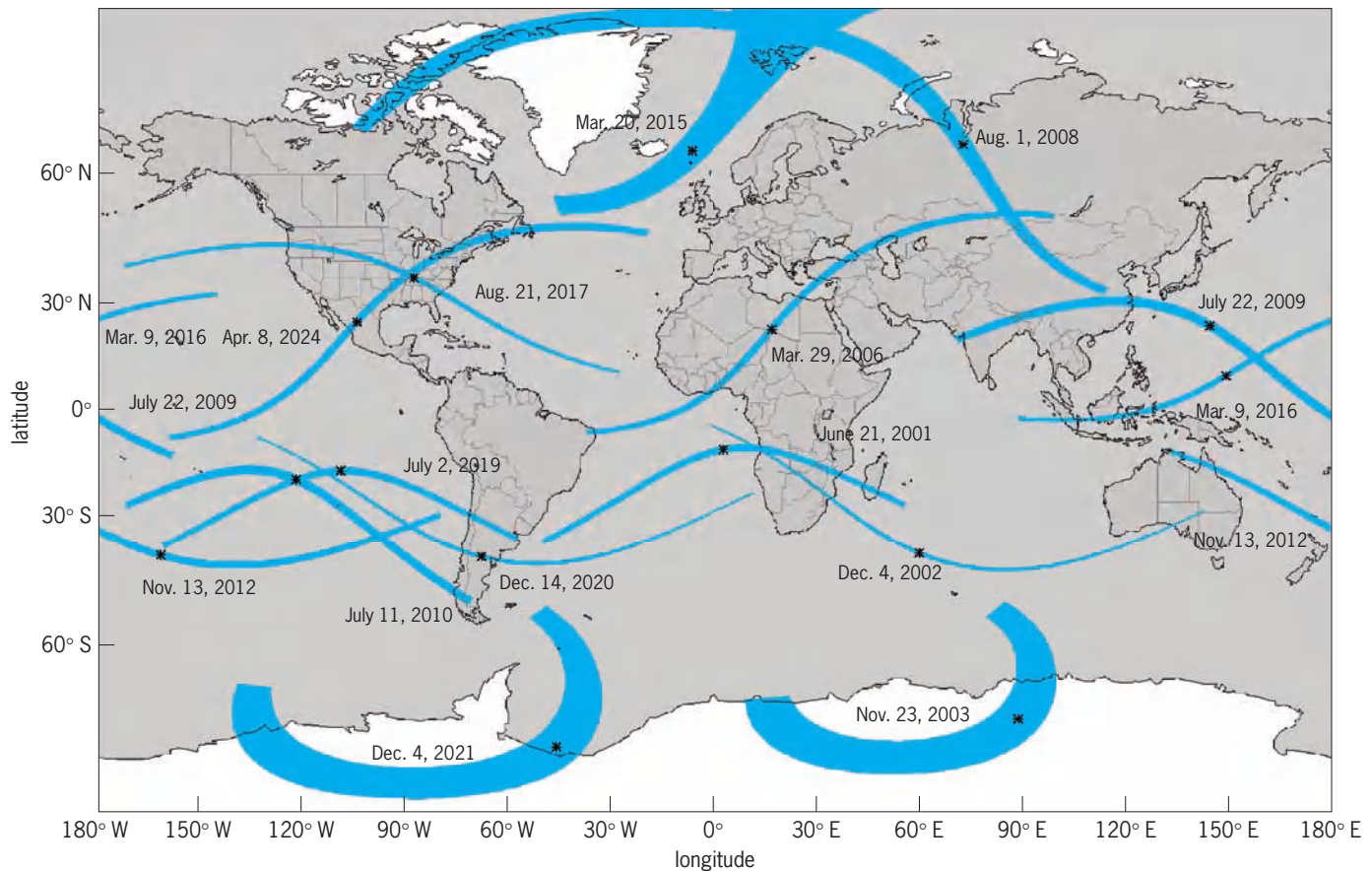


Fig. 4. Paths of all total solar eclipses from 2001 through 2025. (Paths of hybrid annular/total eclipses are not shown). The point of greatest eclipse is indicated by an asterisk on each eclipse path. (F. Espenak, NASA's Goddard Space Flight Center, <http://sunearth.gsfc.nasa.gov/eclipse/eclipse.html>)

The April 8, 2005, hybrid eclipse was annular at the ends of its path and total in the middle, over the Pacific Ocean, where the umbral cone reached the Earth's surface (Fig. 3)

Recent and future eclipses. Notable total eclipses in terms of duration and favorable weather were on June 21, 2001, in southern Africa; on December 4, 2002, in southern Africa and south Australia; and on March 29, 2006, in Africa, Turkey, and into Asia (Table 1; Fig. 4). The next total solar eclipses are on August 1, 2008, in Mongolia and China; on July 22, 2009, in China; and on July 11, 2010, on Easter Island and in southernmost South America. Major scientific expeditions carry out research on such occasions.

The next total eclipse to cross the continental United States will be in 2017, and the next eclipse to cross Canada will be in 2024.

Observing a solar eclipse. A solar eclipse is nature's most magnificent phenomenon. The false impression has grown that it is often hazardous to observe a solar eclipse, whereas, because of educational efforts on eye safety, actually few reports of eye injury exist and the incidence of lasting eye damage is even lower.

The total phase of an eclipse is completely safe to watch with the naked eye. The total brightness of the corona is only that of the full moon, so is equally safe to watch. The darkness at the disappearance of

the diamond ring effect comes so abruptly that people have no trouble telling when the safe time begins, and the diamond ring at the end of totality is so relatively bright that it is clear when it is time to look away. A glance at the Sun just before or after totality is not harmful; it is only staring for an extended time (more than a few seconds with the naked eye) or looking at the not totally eclipsed Sun through binoculars or a telescope that can cause harm.

During the partial phases, it is possible to follow what is going on without any special aid by watching the ground under a tree. The spaces between the leaves make a pinhole camera, and project the solar crescent myriad times onto the ground. A pinhole camera can be made individually by punching a small hole (approximately 0.1–0.2 in. or 2–5 mm in diameter) in a piece of cardboard, and holding it up to the Sun. A crescent image is projected onto a second piece of cardboard held 8 to 40 in. (20 cm to 1 m) closer to the ground, or onto the ground itself. An observer looks at the second cardboard, facing away from the Sun.

For direct observation of the partial phases, a special solar filter must be used. Fogged and exposed black-and-white (not color) film that contains silver and is developed to full density provides suitable diminution of the solar intensity across the entire spectrum. Inexpensive commercial solar filters

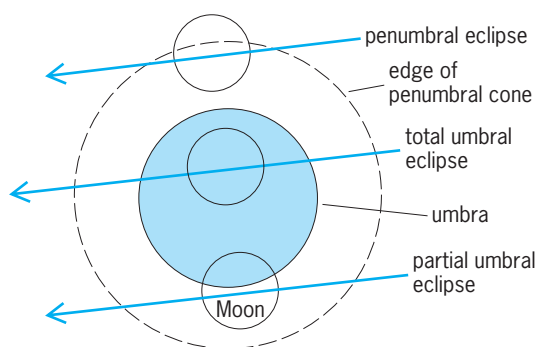


Fig. 5. Cross section through the Earth's shadow cone (umbra) at the Moon's distance, showing the motion of the Moon at three types of lunar eclipses: penumbral eclipse, partial umbral eclipse, and total umbral eclipse.

made of aluminized Mylar can also be used. Gelatin "neutral-density" filters are actually not neutral in the infrared, and so should not be used. Neutral-density filters made by depositing chromium or other metals on glass are safe if they are ND4 or ND5, and are commercially available, as is #14 welder's glass.

Jay M. Pasachoff

Lunar Eclipses

An eclipse of the Moon occurs when the Moon passes through some part of the Earth's two shadows. The inner shadow or umbra is that region where the Earth blocks all direct rays from the Sun. The outer shadow or penumbra is a zone where the Earth blocks part but not all of the Sun's rays. The two shadows result in three types of lunar eclipses (**Fig. 5**):

1. *Penumbral eclipse.* A part or all of the Moon passes through the Earth's penumbral shadow.
2. *Partial eclipse.* A part but not all of the Moon passes through the Earth's umbral shadow.
3. *Total eclipse.* The entire Moon passes through the Earth's umbral shadow.

A lunar eclipse can occur only during the full moon phase. Furthermore, the Moon must also be near one of the two nodes of its orbit. The nodes are the points where the plane of the Moon's orbit intersects the plane of the Earth's orbit. Since the Moon's orbit is tipped 5° to the Earth's orbit (the ecliptic), the Moon passes above or below the shadows during most full moons. But several times each year, full moon occurs

near one of the nodes and some portion of the Moon passes through one or both of the Earth's shadows. See PHASE (ASTRONOMY).

All partial lunar eclipses begin and end with penumbral phases. Similarly, all total eclipses begin with a penumbral phase followed by a partial phase. Totality is then followed by a partial and a penumbral phase in that order (**Fig. 6**).

Lunar eclipses are visible from much larger geographic areas than are solar eclipses. This is due to the fact that each lunar eclipse can be viewed from the entire night hemisphere of the Earth. Unlike solar eclipses, which require special filters for eye protection, lunar eclipses are completely safe to view using the unaided eye or binoculars.

Visibility. For each type of lunar eclipse, the Moon's maximum depth into either the penumbra or umbra is characterized by a value called the eclipse magnitude. This is the fraction of the Moon's diameter immersed in the shadow. For total eclipses, the umbral eclipse magnitude is always equal to or greater than 1.0.

Penumbral eclipses can last up to 5 h. However, these events are of only academic interest since they are quite faint and difficult to see. Only the deepest penumbral eclipses (penumbral eclipse magnitude greater than 0.7) are visible to the unaided eye as a subtle shading along one edge of the Moon.

Partial eclipses are far easier to see and appear as an obvious bite or missing section along one edge of an otherwise full moon. They can last as long as $3\frac{1}{2}$ hours but are usually shorter. Partial eclipses are relatively colorless although the edge of the Earth's umbra may exhibit a reddish tinge when viewed through a telescope. The umbral eclipse magnitude of a partial eclipse is always greater than 0.0 and less than 1.0.

Total lunar eclipses have a maximum duration of nearly $2\frac{1}{2}$ hours. Their appearance is dramatically different from both penumbral and partial eclipses. When the Moon is completely immersed within the Earth's umbral shadow, it is illuminated only by indirect sunlight that has been refracted through and filtered by the Earth's atmosphere. In the process, most of the blue and green colors present in sunlight are removed. The remaining light is deep red or orange and is greatly dimmed. This gives total lunar eclipses their



Fig. 6. Stages of the total lunar eclipse of January 21, 2000. (F. Espenak)



Fig. 7. Total lunar eclipse of December 9, 1992. (F. Espenak)

characteristically blood red color, which frightened many early civilizations. If the Earth had no atmosphere, the Moon would be completely black during a total eclipse. See REFRACTION OF WAVES; SCATTERING OF ELECTROMAGNETIC RADIATION.

The actual appearance of the Moon during totality can vary from one eclipse to the next. The colors range from dark gray or brown, through various shades of red, to bright orange and yellow. These visual differences are directly related to the quality and quantity of dust suspended in Earth's atmosphere. Smoke from forest fires, airborne particles from dust storms, and ash from volcanic eruptions contribute to the color and darkness of the Moon during totality. When large amounts of these particulates are present, they have a strong filtering effect and produce a correspondingly darker eclipse. For example, the 1991 eruption of Mount Pinatubo in the Philippines

pumped large amounts of volcanic ash 115,000 ft (35 km) high into the atmosphere where it was transported globally via the jet stream. Red sunsets were seen around the world, and the total lunar eclipse of December 1992 was so dark it rendered the Moon nearly invisible during totality (Fig. 7).

Frequency. The geometry and recurrence of lunar eclipses have much in common with solar eclipses. Lunar eclipses occur during eclipse seasons lasting about 39 days with a cycle of two seasons every 346.6 days. Lunar eclipses with similar characteristics repeat every 6585.32 days (about 18 years 11 days), which is the same Saros cycle described above for solar eclipses.

Every lunar eclipse is either preceded or followed by a solar eclipse within a time span of 15 days. Of course, the lunar eclipse will occur at full moon, while the solar eclipse will occur 2 weeks earlier or later at new moon. The link between lunar and solar eclipses means that there are a maximum of seven eclipses within a calendar year. They can occur in combinations of five and two, or four and three, with either solar or lunar eclipses in the majority.

In an average century, there are 243 lunar eclipses consisting of 89 penumbral, 84 partial, and 70 total eclipses. Every year has a minimum of two and a maximum of five lunar eclipses.

Future lunar eclipses. The twenty-first century will have 230 eclipses of the Moon, consisting of 87 penumbral, 58 partial, and 85 total eclipses. During this epoch, 76 years will have two lunar eclipses each, 18 years will have three eclipses, and 6 years will have four eclipses. No year contains more than four eclipses. The last time five lunar eclipses occurred during one calendar year was 1879, and the next will be 2132.

During the 22-year period 2004 through 2025, there will be 20 total eclipses of the Moon (Table 2). Deep eclipses have large magnitudes and long

TABLE 2. Total lunar eclipses, 2004–2025*

Date [†]	Umbral magnitude	Duration of totality	Geographic region of eclipse visibility
May 4, 2004	1.309	1 h 16 min	South America, Europe, Africa, Asia, Australia
Oct. 28, 2004	1.313	1 h 21 min	Americas, Europe, Africa, Central Asia
Mar. 3, 2007	1.238	1 h 14 min	Americas, Europe, Africa, Asia
Aug. 28, 2007	1.481	1 h 31 min	Eastern Asia, Australia, Pacific, Americas
Feb. 21, 2008	1.111	0 h 51 min	Central Pacific, Americas, Europe, Africa
Dec. 21, 2010	1.262	1 h 13 min	Eastern Asia, Australia, Pacific, Americas, Europe
June 15, 2011	1.705	1 h 41 min	South America, Europe, Africa, Asia, Australia
Dec. 10, 2011	1.110	0 h 52 min	Europe, Eastern Africa, Asia, Australia, Pacific, North America
Apr. 15, 2014	1.296	1 h 19 min	Australia, Pacific, Americas
Oct. 8, 2014	1.172	1 h 00 min	Asia, Australia, Pacific, Americas
Apr. 4, 2015	1.006	0 h 12 min	Asia, Australia, Pacific, Americas
Sept. 28, 2015	1.282	1 h 13 min	Eastern Pacific, Americas, Europe, Africa, western Asia
Jan. 31, 2018	1.321	1 h 17 min	Asia, Australia, Pacific, western North America
July 27, 2018	1.614	1 h 44 min	South America, Europe, Africa, Asia, Australia
Jan. 21, 2019	1.201	1 h 03 min	Central Pacific, Americas, Europe, Africa
May 26, 2021	1.016	0 h 19 min	Eastern Asia, Australia, Pacific, Americas
May 16, 2022	1.419	1 h 26 min	Americas, Europe, Africa
Nov. 8, 2022	1.364	1 h 26 min	Asia, Australia, Pacific, Americas
Mar. 14, 2025	1.183	1 h 06 min	Pacific, Americas, western Europe, western Africa
Sept. 7, 2025	1.367	1 h 23 min	Europe, Africa, Asia, Australia

*Data courtesy of Fred Espenak, <http://sunearth.gsfc.nasa.gov/eclipse/eclipse.html>.

[†]Date of mid-eclipse (Greenwich Mean Time).

durations (for example, June 15, 2011, and July 27, 2018).
Fred Espenak

Bibliography. F. Espenak, *Fifty Year Canon of Lunar Eclipses: 1986–2035*, NASA RP-1178, Sky Publishing, 1989; F. Espenak, *Fifty Year Canon of Solar Eclipses: 1986–2035*, NASA RP-1216, Sky Publishing, 1987; L. Golub and J. M. Pasachoff, *Nearest Star: The Surprising Science of Our Sun*, Harvard University Press, 2001; P. Guillermier and S. Koutchmy, *Total Eclipses: Science, Observations, Myths and Legends*, Springer-Praxis, 1999; M. Littman, K. Willcox, and F. Espenak, *Totality: Eclipses of the Sun*, 2d ed., Oxford University Press, 1999; B.-L. Liu and A. D. Fiala, *Canon of Lunar Eclipses 1500 B.C.—A.D. 3000*, Willmann-Bell, 1992; J. Meeus and H. Mucke, *Canon of Lunar Eclipses: –2002 to +2526*, 3d ed., Astronomisches Büro, Vienna, 1992; H. Mucke and J. Meeus, *Canon of Solar Eclipses: –2003 to +2526*, 2d ed., Astronomisches Büro, Vienna, 1992; J. M. Pasachoff, *A Field Guide to the Stars and Planets*, 4th ed., updated, Houghton Mifflin, 2006; J. M. Pasachoff, *The Complete Idiot's Guide to the Sun*, Alpha Books, 2003; J. B. Zirker, *Total Eclipses of the Sun*, expanded ed., Princeton University Press, 1995.

Eclipsing variable stars

Double star systems in which the two components are too close to be seen separately but which reveal their duplicity by periodic changes in brightness as each star successively passes between the other and the Earth, that is, eclipses the other. Studies of the light changes and the radial velocity changes of each component permit the computation of the radii, masses, and densities of the components—important quantities that cannot be measured directly in single stars. In addition, these close double stars are useful in studies of mass loss and of stellar evolution. Since eclipsing stars are variable in light, they are included in general variable star catalogs under the same system of nomenclature. See BINARY STAR; VARIABLE STAR.

Periods. The periods of light variation range from less than 3 h for very close systems to over 27 years for the peculiar system Epsilon Aurigae. However, the majority of the periods lie between 0.5 and 10 days. In many cases the periods are not constant but change with time. In a few cases the variation is caused by a slow change in the orientation of the major axis of an elliptical orbit; in such cases the rate of change combined with other quantities gives information concerning the manner in which the density of the star increases from the outer layers to the center. In most cases, however, the changes are unpredictable and are probably connected with ejections of matter from one of the stars.

Velocity curves. The radial velocity (velocity of approach or recession) of each star can be determined at any time by the displacement of the spectral lines. A plot of velocities against time over one period of orbital revolution is known as a velocity curve. The

maximum radial velocity of approach or recession depends on the true orbital velocity of the star and the “inclination,” or the amount by which the plane of the orbit is tilted relative to the line of sight to the Earth. (Technically, the inclination is the angle between a perpendicular to the orbit plane and the line of sight; when the inclination is 90° , the eclipses are central.) If the inclination is known, the orbital velocity of each star can be calculated from the radial velocity. The orbital velocity multiplied by the period will give the circumference of the orbit and from this the radius of each star about the center of mass. From the size of the orbit and the period, by using the law of gravitation, the mass of each star can be calculated in terms of the Sun's mass. However, the inclination cannot be determined from the velocity curve alone; the quantities finally determined are $m_1 \sin^3 i$, $m_2 \sin^3 i$, and $a \sin i$, where i is the inclination, m_1 and m_2 the stellar masses, and a the radius of a circular orbit or half the major axis of an elliptical one.

Light curve. The light curve shows the changes in brightness of the system throughout one orbital revolution. The manner in which the light changes during the eclipse of each star by the other depends very strongly on three factors: the size of each star relative to the radius of the orbit (r_1/a and r_2/a) and the inclination. The determination of precisely what relative sizes and inclination give a computed curve which will approximate satisfactorily the observations is one of the more difficult problems of modern astronomy. However, by use of complex sets of tables the problem can be solved in many cases. Numerical values for these three quantities are thus found. Several alternate methods have been developed which make use of the capabilities of electronic computers. These take into account simultaneously the various interaction effects such as distortion of the shape of the components in the earlier methods. These had to be removed from the light curve before a solution could be carried out. See ECLIPSE; LIGHT CURVES.

Absolute dimensions. The masses, radii, and densities of the stars, usually expressed in comparison with those of the Sun, are called the absolute dimensions of the system. Combining the results from the light and velocity curves yields this fundamental information which cannot be obtained from either approach alone. The inclination, determined from the shape of the light curve, can be substituted in the quantities $m_1 \sin^3 i$, $m_2 \sin^3 i$, and $a \sin i$ to find m_1 , m_2 , and a . Thus the masses of the stars in terms of the Sun's mass and the size of the orbit in terms of the Sun's radius can be found. Since the radii of the stars in fractions of the size of the orbit have been determined from the light curve, the sizes of the stars relative to the size of the Sun (either in miles or kilometers) can now be computed. **Figure 1** shows the relative sizes and separation of the stars in a typical eclipsing system.

Complications. Before the light curve can be “solved” to give the above quantities, corrections must be made for other factors which influence

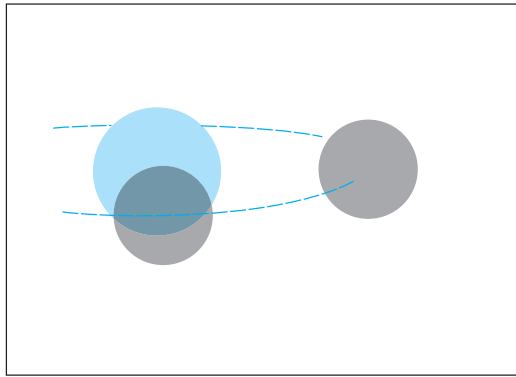


Fig. 1. Relative sizes and separation of components of R Canis Majoris. Shaded star represents cooler component. It is shown in two positions: at the middle of an eclipse and when separation would be greatest as viewed from Earth. Even at the greatest separation, all such systems appear as single stars.

the light changes. This is done by studying the light between eclipses where, were it not for these complications, the brightness of the system would not change. One of these effects is called ellipticity. The tidal attraction of each star for the other has caused distortions until the stars in extreme cases resemble footballs more than baseballs in shape. The technical term is prolate ellipsoids, although when the stars differ in size and mass they will also differ in shape. Further, the radiation of each star falling on the side of the other nearest it will cause this side to be brighter than the side turned away. The difference will be most marked for the cooler component where the effect of the intense radiation from the nearby hotter star is most strongly evident.

There are other effects, some very poorly understood, which cause light changes between eclipses, and all of these must be carefully studied before the analysis of the eclipse begins.

Evolutionary changes. Studies of single stars indicate that, when the hydrogen in the center of the star has been converted into helium, the star undergoes a relatively rapid expansion in size. The presence of a nearby companion complicates the picture considerably, but it does seem clear that much of the mass of the star must be lost to the system or possibly transferred to the other component.

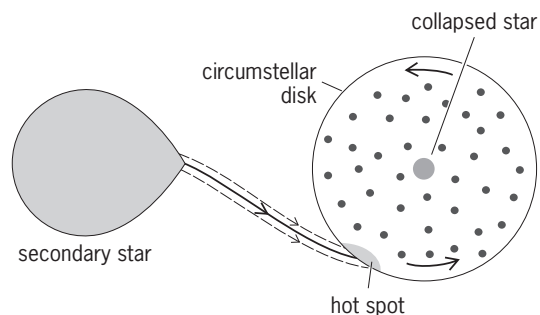


Fig. 2. Model generally accepted for certain types of eclipsing binary stars, in particular those known as dwarf novae.

The mass-losing star eventually becomes a collapsed object—a white dwarf, a neutron star, or a black hole. Each of these types has been identified in at least one binary system. Then when the secondary, originally less massive star begins its expansion, matter from it is transferred to the collapsed object, often with dramatic results. See BLACK HOLE; NEUTRON STAR; WHITE DWARF STAR.

It is now believed that all explosive variables (novae, recurrent novae, and so forth), with the exception of supernovae, are members of close binary systems. At least some of the x-ray sources are close binaries in this state, although the detection (by instruments carried aboard satellites) of x-radiation from Algol and other systems which are not yet in this state indicates that other physical mechanisms may also be responsible. Some systems show intermittent bursts of radiation at radio frequencies. Evidence indicates the presence of clouds of circumstellar material.

In a few of the eruptive variables, particularly those known as dwarf novae, rapid scintillation is found, presumably from a hot spot where the transferring mass collides violently with a circumstellar disk of relatively low-density material revolving around the collapsed star (Fig. 2); the scintillation stops periodically when the spot is eclipsed by the other component. Instruments on satellites have extended observations to the far ultraviolet, as well as the x-ray, regions of the spectrum. Thus, in addition to the classical reasons for studying eclipsing variable stars, observation of them leads into many branches of astrophysics. See ASTROPHYSICS, HIGH-ENERGY; CATAclysmic VARIABLE; NOVA; STAR; STELLAR EVOLUTION; X-RAY ASTRONOMY. Frank Bradshaw Wood

Bibliography. A. H. Batten (ed.), *Algols*, *Space Sci. Rev.*, 50 (1-2):384, 1989; A. H. Batten, *Binary and Multiple Systems of Stars*, 1973; E. F. Milone, *Light Curve Modeling of Eclipsing Binary Stars*, 1995; J. Sahade and F. B. Wood, *Interacting Binary Stars*, 1978.

Ecliptic

The apparent path of the Sun across the sky over the course of a year. The Earth's mean orbit defines the ecliptic plane, and an observer on Earth sees the Sun, traveling in this plane over the course of a year, follow the ecliptic. The orbits of the Moon and planets are at slight angles to the ecliptic, but these objects never appear too far in the sky from the ecliptic, which is often drawn on star maps. The ecliptic was identified by Greek and Babylonian scientists in about the fifth century B.C.

The ecliptic intersects the celestial equator at an angle of approximately $23\frac{1}{2}^\circ$, an angle known as the obliquity of the ecliptic. The points of intersection are the vernal equinox and the autumnal equinox. As a result of this obliquity, the Sun can be overhead at different times of year for locations on Earth between latitudes $+23\frac{1}{2}^\circ$ and $-23\frac{1}{2}^\circ$, that is, between the Tropic of Cancer and the Tropic of

Capricorn, respectively. See ASTRONOMICAL COORDINATE SYSTEMS; EARTH ROTATION AND ORBITAL MOTION.

Jay M. Pasachoff

Bibliography. J. Mitton, *Cambridge Dictionary of Astronomy*, Cambridge University Press, 2001; J. M. Pasachoff and A. Filippenko, *The Cosmos: Astronomy in the New Millennium*, 3d ed., Brooks/Cole, 2007; P. K. Seidelmann (ed.), *Explanatory Supplement to the Astronomical Almanac*, University Science Books, 1992.

Eclogite

A very dense rock composed of red-brown garnet and the grape-green pyroxene omphacite. Eclogites possess basaltic bulk chemistry, and their garnets are rich in the components pyrope ($\text{Mg}_3\text{Al}_2\text{Si}_3\text{O}_{12}$), almandine ($\text{Fe}_3\text{Al}_2\text{Si}_3\text{O}_{12}$), and grossular ($\text{Ca}_3\text{Al}_2\text{Si}_3\text{O}_{12}$), while the pyroxenes are rich in jadeite ($\text{NaAlSi}_2\text{O}_6$) and diopside ($\text{CaMgSi}_2\text{O}_6$). Accessory phases depend on the environment of crystallization, but may include quartz, kyanite, amphibole, olivine, orthopyroxene, zoisite, mica, and rutile, but never plagioclase.

Occurrences. Eclogite occurrences may be subdivided into three broad categories:

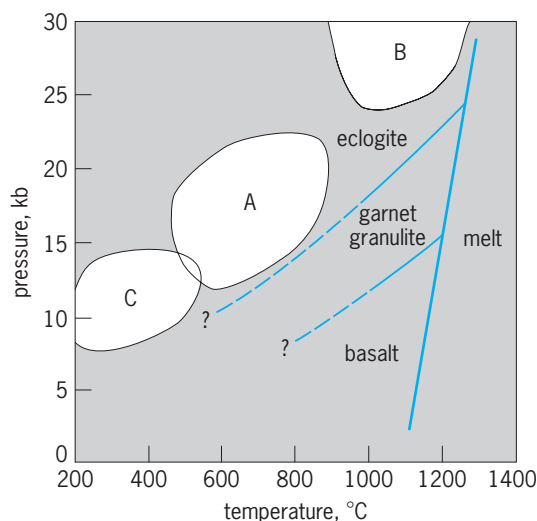
Group A eclogites are found as layers, lenses, or boudins in schists and gneisses seemingly of the amphibolite facies. Quartz, together with zoisite or kyanite, commonly occurs in these rocks. Amphibole of barroisitic composition may also be present. These rocks originate as metamorphosed basic bodies (dykes, sills, flows) or layers of dolomitic marl intercalated with metasedimentary schists. See QUARTZ.

Group B eclogites are found as inclusions in kimberlites and basalts. They are frequently accompanied by xenoliths of garnet peridotite; both eclogites and peridotites can be diamond-bearing, suggesting that these rocks come from deep in the mantle. See BASALT.

Group C eclogites are found as blocks and lenses in schists of the glaucophane-schist facies. Such eclogites do not contain kyanite, rarely contain quartz, but bear amphibole, epidote, rutile, or sphene. The presence of pillow lavas converted to eclogite indicates that these rocks had a crustal origin.

Conditions of formation. Eclogites can form over a wide range of pressure and temperature, but are clear indicators of rather high pressure, as their very high density might suggest. Experimental work has shown that rocks of basaltic composition are converted to eclogite at high pressures (see *illus.*).

The distribution of Fe^{2+} and Mg^{2+} between garnet and omphacite has been determined as a function of temperature, and allows approximate temperature regimes to be assigned to the three categories of eclogite. Eclogites of group C formed at relatively low temperatures and probably crystallized isofacially with rocks of the glaucophaneschist facies. Group A eclogites probably formed over a wider range of temperature and pressure, but at somewhat



Possible pressure-temperature fields for eclogites from gneiss areas: A from kimberlites; B from blueschist terrains; C basalt-eclogite transition. 1 kilobar = 10^5 kilopascals. $^{\circ}\text{F} = (^{\circ}\text{C} \times 1.8) + 32$.

higher temperatures than group C. The very high temperatures estimated for group B eclogites are consistent with their origin in the upper mantle. See GARNET; OMPHACITE.

Eclogite facies. Eclogite is the name given to the highest-pressure facies of metamorphism; the critical mineral assemblage defining this facies is garnet + omphacite, together with kyanite or quartz in rocks of basaltic composition. Where sedimentary and granitic rocks have been metamorphosed under eclogite facies conditions, they result in spectacular omphacite + garnet + quartz-bearing mica schists and metagranitic gneisses such as are found in the Sezia-Lanzo zone of the western Alps. See METAMORPHIC ROCKS.

The relationship between the eclogite and other facies is often unclear. Eclogites, for instance, can form in relatively anhydrous parts of blueschist and amphibolite facies terrains where normally their hydrated equivalents would crystallize. Some such eclogites are inferred, on textural grounds, to have formed in a dry environment. In other cases the coarse textures, the presence of veins of diagnostic eclogite facies minerals such as omphacite + kyanite + quartz, and fluid inclusions suggest that their formation was accompanied by copious fluid. Although very high pressures are required to stabilize eclogite under hydrous conditions, the presence of diluent species (for example, CO_2 or halides) in the fluid lowers the required pressures.

Eclogites and the mantle. The high density of eclogite, together with its elastic properties, makes it a candidate for upper mantle material. Large quantities of basaltic oceanic crust are returned to the mantle through the process of subduction, where prevailing high pressures convert it to eclogite. The quantity and distribution of eclogite within the mantle is not known; that it occurs is known from the nodules brought up in kimberlite pipes and in basalts. See HIGH-PRESSURE PROCESSES. Timothy J. B. Holland

Bibliography. M. G. Best, *Igneous and Metamorphic Petrology*, 1982; B. W. Evans and E. H. Brown (eds.), *Blueschists and Eclogites*, 1986; L. A. Raymond, *Petrology: The Study of Igneous, Sedimentary, and Metamorphic Rocks*, 1994; D. C. Smith (ed.), *Eclogites and Eclogite-Facies Rocks*, 1989; M. Suk, *Petrology of Metamorphic Rocks*, 1983.

Ecological communities

Assemblages of living organisms that occur together in an area. The nature of the forces that knit these assemblages into organized systems and those properties of assemblages that manifest this organization have been topics of intense debate among ecologists since the early years of this century. On the one hand, there are those who view a community as simply consisting of species with similar physical requirements, such as temperature, soil type, or light regime. The similarity of requirements dictates that these species be found together, but interactions between the species are of secondary importance and the level of organization is low. On the other hand, there are those who conceive of the community as a highly organized, holistic entity, with species inextricably and complexly linked to one another and to the physical environment, so that characteristic patterns recur, and properties arise that one can neither understand nor predict from a knowledge of the component species. In this view, the ecosystem (physical environment plus its community) is as well organized as a living organism, and constitutes a superorganism. Between these extremes are those who perceive some community organization but not nearly enough to invoke images of holistic superorganisms. See ECOSYSTEM.

Emergent and collective properties. The crux of this debate is whether communities have emergent properties, a concept that itself is sufficiently confused, so that part of the ecological debate may be semantic. Some denote by emergent property any property of a group that is not a property of the component individuals. Others call any group trait a collective property, and reserve as emergent only those collective properties that do not derive trivially from properties of component individuals or from the very definition of the group. An emergent property, in this conception, represents a new level of organization, and cannot be predicted even with complete knowledge of all properties of individuals in the group. Thus, the number of species in a community is a property that is not defined for any particular species, but once one knows all the species in a community one can easily measure the property "number of species" for that community. In the restricted view of emergent property, the number of species in a community is a collective property, but not an emergent one.

Whether any properties of communities are emergent is debatable, but all ecologists agree that communities have collective properties that are signifi-

cant both biologically and practically. For example, it is not a trivial matter to understand why the number of species in a given community is S and not $2S$, and what the consequences are of having S species. Also, human exploitation of ecological communities often rests on a thorough knowledge of certain collective properties, such as biomass produced per unit time, or rate of cycling of certain nutrients. See BIOMASS.

Size. Every community comprises a given group of species, and their number and identities are distinguishing traits. Most communities are so large that it is not possible to enumerate all species; microorganisms and small invertebrates are especially difficult to census. However, particularly in small, well-bounded sites such as lakes or islands, one can find all the most common species and estimate their relative abundances. The number of species is known as species richness, while species diversity refers to various statistics based on the relative numbers of individuals of each species in addition to the number of species. The rationale for such a diversity measure is that some communities have many species, but most species are rare and almost all the individuals (or biomass) in such a community can be attributed to just a few species. Such a community is not diverse in the usual sense of the word. Patterns of species diversity abound in the ecological literature; for example, pollution often effects a decrease in species diversity. The most popular species diversity statistic is given in the equation below, where

$$H' = \sum_{i=1}^S p_i \log(p_i)$$

S is the number of species and p_i is the fraction of all the individuals (or all the biomass) in the community contributed by species i . For many sets of data on diversities of groups of communities (or parts of communities, such as all birds or all plants), values of H' are highly correlated with species richness, suggesting that H' may not be the ideal expression of biotic diversity.

The main patterns of species richness that have been detected are area and isolation effects, successional gradients, and latitudinal gradients. Larger sites tend to have more species than do small ones, and isolated communities (such as those on oceanic islands) tend to have fewer species than do less-isolated ones of equal size. Later communities in a temporal succession tend to have more species than do earlier ones, except that the last (climax) community often has fewer species than the immediately preceding one. Tropical communities tend to be very species-rich, while those in arctic climates tend to be species-poor. This observation conforms to a larger but less precise rule that communities in particularly stressful environments tend to have few species.

Composition. Species composition is important, largely in assisting in the classification of communities. If the complete species list for each community and the proportional representation of each species

in it were known, no two communities would be found to be identical. The criteria are arbitrary for how similar these lists must be before two communities are viewed as the same.

Communities are usually denoted by the presence of species, known as dominants, that contain a large fraction of the community's biomass, or account for a large fraction of a community's productivity. Dominants are usually plants. Determining whether communities at two sites are truly representatives of the "same" community requires knowledge of more than just the dominants, however. "Characteristic" species, which are always found in combination with certain other species, are useful in deciding whether two communities are of the same type, though the designation of "same" is arbitrary, just as is the designation of "dominant" or "characteristic."

Although there is no objective answer to the question of whether two communities are the same, the question bears heavily on the matter of whether communities are in fact integrated superorganisms with emergent properties. One would not expect a super-organismic entity to have an infinite variety of representations, for the same cohesive forces that bind it into a holistic unit should constrain it to a limited number of forms. Organisms for the most part do not form a continuous gradient of all conceivable phenotypes, but rather are constrained by a number of physiological and morphological relationships (such as allometry). Thus there is generally no trouble in saying that individual X is a human, individual Y is a longleaf pine, and so on. Although there is phenotypic variation in all species, it is limited and centered on recognizable modes.

An active area of research among animal ecologists has been the search for limits to how similar two species can be and still occur in the same community. That there exists some limiting similarity beyond which coexistence is impossible is a theoretical consequence of conceptions of the ecological niche; but clear evidence, in any instance, that such limiting similarity is the reason why a given species is absent from a given community is exceedingly difficult to amass. There are two main difficulties. First, it must be demonstrated that some resource is limiting so that species that are too similar will find it in short supply. Second, an appropriate index of a species' use of the resource (its niche) must be determined. Much attention has been focused on the size of an animal or of its trophic appendages as such an index, though there is much evidence that many other forces affect these sizes. For few, if any, animal communities has it been shown that size similarity restricts species composition or that patterns of size differences among species are highly predictable or regular.

Spatial distribution patterns. A related matter bearing on the aptness of the superorganismic community metaphor is that most organisms have clear spatial boundaries; communities often do not. There is no difficulty discerning where one human ends and another begins, but there are not always clear boundaries to communities. Occasionally, very sharp

limits to a physical environmental condition impose similarly sharp limits on a community. For example, serpentine soils are found sharply delimited from adjacent soils in many areas, and have mineral concentrations strikingly different from those of the neighboring soils. Thus they support plant species that are very different from those found in nearby non-serpentine areas, and these different plant species support animal species partially different from those of adjacent areas. Here two different communities are sharply bounded from each other.

Usually, however, communities grade into one another more gradually, through a broad intermediate region (an ecotone) that includes elements of both of the adjacent communities, and sometimes other species as well that are not found in either adjacent community. One has little difficulty telling when one is in the center of either of the adjacent communities A and B, but exactly when one passes from A to B is not easily discerned. The reason is that, though each species in a community is adapted to its physical environment and to some extent to other species in the community, the adaptations to the physical environment are usually not identical, and most of the adaptations to one another are not obligatory.

The environment created by the dominant species, by their effects on temperature, light, humidity, and other physical factors, and by their biotic effects, such as allelopathy and competition, may entrain some other species so that these other species' spatial boundaries coincide with those of the dominants. The mangrove skipper, *Phocides pygmalion*, can feed only on red mangrove, *Rhizophora mangle*, so whatever aspects of the physical environment limit the mangrove (especially temperature), the skipper will also be limited to the same sites. However, many other species that feed on the mangrove (for example, the io moth, *Automeris io*) also feed on other plants, and their spatial boundaries do not coincide with those of the mangrove. Nor do most of the species in a community share identical physical requirements. Black mangrove (*Avicennia germinans*) co-occurs with red mangrove in most sites, and the two are normal constituents of a community often termed mangrove swamp. But *Avicennia* can tolerate much colder weather than can *Rhizophora*, and so it is found much farther north in the Northern Hemisphere. Eventually it too ceases, and the mangrove community is replaced in more northerly areas by salt marsh communities. In some intermediate regions, salt marsh grasses and black mangrove are found together. There is no clear boundary between the two communities.

This continuous intergradation of most communities argues against the superorganism concept, but there are aspects of the spatial arrangement within communities that suggest that the component species are far from independent, and indicate, if not complete holism, at least that some community properties are not easily predicted from exhaustive knowledge of the component species. One example

is stratification, the vertical arrangement of canopy layers in most forests. Individuals of the different species do not have heights that are independently and continuously distributed from the ground to the top of the tallest tree. Instead, there are a few rather distinct strata, with each species at maturity characteristically occupying one of these. Tropical forests from all parts of the world, even though they may have completely different species compositions, usually contain five fairly clear strata: a top-most layer composed of the tallest tree species, two lower layers of smaller trees, a shrub layer, and a ground layer. There are doubtless good physical reasons why the diffusion of light can explain this characteristic structure given a knowledge of evolution and plant physiology (though no completely compelling explanation has yet surfaced), so it may be that this is an elaborate collective property rather than an emergent one. In either case, there is clearly a high degree of multispecies organization in this spatial arrangement. *See* PHYSIOLOGICAL ECOLOGY (PLANT).

In addition to vertical arrangement, horizontal locations of individuals of different species are usually not random. Usually, individuals of a given species are clumped; they are found on average closer to one another than one would have predicted. Probably the major reason for this is response to habitat heterogeneity: conspecific individuals tend to favor more similar habitats than do heterospecific individuals. Individuals of different species may also be nonrandomly arranged with respect to one another. Competitive interactions may cause two species typically to be found in different microsites, while mutualistic interactions or preference for a similar physical habitat may cause two species to be associated spatially. *See* POPULATION ECOLOGY.

Succession. More or less distinct communities tend to follow one another in rather stylized order. As with recognition of spatial boundaries, recognition of temporal boundaries of adjacent communities within a sere is partly a function of the expectations that an observer brings to the endeavor. Those who view communities as superorganisms are inclined to see sharp temporal and spatial boundaries, and the perception that one community does not gradually become another community over an extended period of time confirms the impression that communities are highly organized entities, not random collections of species that happen to share physical requirements. In this superorganismic view, ecological succession of communities in a sere is analogous to the life cycle of an organism, and follows a quite deterministic course. That secondary succession following a disturbance often leads to a community resembling the original one is the analog in the superorganism to wound repair in an organism. The driving force for succession, in this conception, is that the community or its dominant species modify the environment so that these dominant species are no longer favored, and when the dominant species are replaced, the bulk of the community, complexly linked to the dominant species and to one another,

disappears as well, to be replaced by the next community in the sere.

This superorganismic conception of succession has been replaced by an individualistic succession. Data on which species are present at different times during a succession show that there is not abrupt wholesale extinction of most members of a community and concurrent simultaneous colonization by most species of the next community. Rather, most species within a community colonize at different times, and as the community is replaced most species drop out at different times. Thus, though one can usually state with assurance that the extant community is of type A or type B, there are extended periods when it is difficult to assign the assemblage of species at a site to any recognizable community.

That succession is primarily an individualistic process does not mean that there are not characteristic changes in community properties as most successions proceed. Species richness usually increases through most of the succession, for example, and stratification becomes more highly organized and well defined. A number of patterns are manifest in aspects of energy flow and nutrient cycling. *See* ECOLOGICAL SUCCESSION.

Functional organization. Living organisms are characterized not only by spatial and temporal structure but by an apparent purpose or activity. In short, they are “doing something,” and this activity has been termed teleonomy. In humans, for example, various physiological functions are continuously under way until death intervenes. Communities have functions analogous to physiology.

In the first place, the various species within a community have different trophic relationships with one another. One species may eat another, or be eaten by another. A species may be a decomposer, living on dead tissue of one or more other species. Some species are omnivores, eating many kinds of food; others are more specialized, eating only plants or only animals, or even just one other species. These trophic relationships certainly unite the species in a community into a common endeavor, the transmission of energy through the community. This energy flow is patently analogous to an organism's mobilization and transmission of energy from the food it eats.

One aspect of energy flow is a candidate for an emergent property: the topology of the food web. Examination of this topology for a number of webs suggests that they are highly constrained in structure. For example, the maximum number of trophic levels in a web rarely exceeds five. One reason may be that the total amount of energy that has not already been degraded by the time the energy has passed through three or four levels may not be enough to sustain a viable population of a species that would feed at still higher levels. An alternative explanation is that the population dynamics of a web with so many levels would probably confer mathematical instability on the web, so that one or more species would be eliminated. Other properties of food webs that have been discerned include a low number of

omnivore species (those feeding on more than one level), and a tendency for the number of predator species and the number of prey species to be in the ratio of 4:3. No explanation is forthcoming for the latter observation. The former observation is held to reflect mathematical instability that arises from the presence of omnivores in a web. Other workers contend that neither pattern will be maintained when much more comprehensive data are available on what paths calories actually follow as they flow through a community. If these patterns do not turn out to be artifacts of incomplete knowledge, they would appear to be emergent properties reflecting a high degree of organization. *See* FOOD WEB.

Just as energy flows through the communities, so do nutrients move. A calorie of energy does not move in the abstract from organism to organism; rather, the calorie is bound up in a molecule that moves when one organism eats (or decomposes) another. Or the calorie may be respired away as a result of the metabolism of the organism that ingests it. A calorie of energy, once respired by some member of the community, is no longer available to the community. But the molecule associated with that calorie, or a new molecule produced from it, is still present and can go through the food web again. Thus nutrients can cycle within a community, while energy flow, once the energy is transformed to heat, is one-way. Nutrient cycling is analogous to circulation in an organism, and combines with energy flow to support the superorganism metaphor. Different nutrients cycle at different rates, and for several nutrients the cycle within the community is linked to cycles in other communities. Nutrients exist in abiotic entities, as well as biotic organisms, so nutrient cycling is as much an ecosystem trait as a community one. A number of properties of nutrient cycling (such as rate, turnover times, and sizes of different pools or compartments) have been measured. Whether any of these are emergent as opposed to collective properties has yet to be determined.

Productivity. By virtue of differing rates of photosynthesis by the dominant plants, different communities have different primary productivities. Tropical forests are generally most productive, while extreme environments such as desert or alpine conditions harbor rather unproductive communities. Agricultural communities are intermediate. Algal communities in estuaries are the most productive marine communities, while open ocean communities are usually far less productive. The efficiency with which various animals ingest and assimilate the plants and the structure of the trophic web determine the secondary productivity (production of organic matter by animals) of a community. Marine secondary productivity generally exceeds that of terrestrial communities. *See* AGROECOSYSTEM; BIOLOGICAL PRODUCTIVITY.

Reproduction. A final property that any organism must have is the ability to reproduce itself. Communities may be seen as possessing this property, though the sense in which they do so does not support the superorganism metaphor. A climax commu-

nity reproduces itself through time simply by virtue of the reproduction of its constituent species, and may also be seen as reproducing itself in space by virtue of the propagules that its species transmit to less mature communities. For example, when a climax forest abuts a cutover field, if no disturbance ensues, the field undergoes succession and eventually becomes a replica of the adjacent forest. Both temporally and spatially, then, community reproduction is a collective rather than an emergent property, deriving directly from the reproductive activities of the component species. *See* ALTITUDINAL VEGETATION ZONES; BOG; CHAPARRAL; DESERT; ECOLOGY; GRASSLAND ECOSYSTEM; MANGROVE; MUSKEG; PARAMO; PUNA. Daniel Simberloff

Bibliography. T. F. H. Allen and T. B. Starr, *Hierarchy: Perspectives for Ecological Complexity*, 1988; M. Begon, M. Mortimer and D. J. Thompson, *Population Ecology: A Unified Study of Animals and Plants*, 3d rev. ed., 1996; P. G. Digby and R. A. Kempton, *Multivariate Analysis of Ecological Communities*, 1987; J. R. Trabalka and D. E. Reichle (eds.), *The Changing Carbon Cycle*, 1986.

Ecological competition

The interaction of two (or more) organisms (or species) such that, for each, the birth or growth rates are depressed and the death rate increased by the presence of the other organisms (or species). Competition is recognized as one of the more important forces structuring ecological communities, and interest in competition led to one of the first axioms of modern ecology, the competitive exclusion principle. The principle suggests that in situations where the growth and reproduction of two species are resource-limited, only one species can survive per resource.

Lotka-Volterra equations. The competitive exclusion principle was originally derived by mathematicians using the Lotka-Volterra competition equations. This model of competition predicts that if species differ substantially in competitive ability, the weaker competitor will be eliminated by the stronger competitor. However, a competitive equilibrium can occur if the negative effect of each species on itself (intraspecific competition) is greater than the negative effect of each species on the other species (interspecific competition). Because the competitive exclusion principle implies that competing species cannot coexist, it follows that high species diversity depends upon mechanisms through which species avoid competition. *See* MATHEMATICAL ECOLOGY.

Species coexistence. One way that the problem of competitive coexistence has been addressed is to extend the competitive exclusion principle to a continuum of resources through the theory of limiting similarity. Limiting similarity is defined as the maximum level of similarity in the use of a set of resources that will allow competing species to coexist. However, this approach still assumes that several species

cannot coexist on the same resources unless they use those resources differently.

A more realistic approach to understanding the coexistence of competing species is to view community structure as dynamic, in which species densities do not remain constant over time at each spatial location. Fluctuations in densities may be a result of floods, fires, and other abiotic disturbances, or of biotic interactions, such as predation, herbivory, and disease. The key to coexistence within a homogeneous patch in such a dynamic world is the prevention of competitive equilibrium and hence competitive exclusion.

In general, competitive exclusion can be prevented if the relative competitive abilities of species vary through time and space. Such variation occurs in two ways. First, dispersal rates into particular patches may fluctuate, causing fluctuations in the numerical advantage of a species in a particular patch. This may occur, for example, through differences among species in dispersal ability. Second, competitive abilities of species may be environmentally dependent and, therefore, fluctuate with local environmental changes. Different species will be favored under different sets of environmental conditions, allowing each species to have periods of strong recruitment.

Competitive exclusion can also be avoided if fluctuations in environmental factors reduce the densities of potentially competing species to levels where competition is weak and population growth is for a time insensitive to density. Therefore, intensities of competition within and between species fluctuate with time, and competing species are able to coexist.

In patchy environments, local extinctions or population reductions resulting from predators or environmental fluctuations can allow species to colonize and grow for some time in patches in the absence of competition. In this way, even inferior competitors can persist if they are, for example, relatively less susceptible to predation or are relatively better dispersers.

Coexistence is not merely a result of environmental harshness or fluctuations but also involves the critical element of niche differentiation (that is, species must differ from one another if they are to coexist). However, the focus is not how species coexist by partitioning resources, but how species can coexist on the same resources by differing sufficiently in their responses to environmental conditions and fluctuations. See ECOLOGICAL COMMUNITIES; ECOLOGICAL SUCCESSION.

Applications. Competition theory has been applied to human-manipulated ecosystems used to produce food, fiber, and forage crops as well as in forestry and rangeland management. Although many characteristics of agricultural systems are similar to those of natural ecosystems, agricultural communities are unique because they are often managed for single-species (sometimes multispecies) production and they are usually characterized by frequent and intense disturbance. Studies of competition in agricul-

ture have primarily examined crop loss from weed abundance under current cropping practices, and have evaluated various weed control tactics and intercropping systems. Because both total density and species proportion influence the outcome and interpretation of competition experiments, the ability to quantify the intensity of competition and to separate the effects of competition between individuals of the same species versus competition between individuals of different species (such as weed versus crop) depends on experimental design. Experimental approaches that systematically vary total and relative plant densities provide better quantification of competition in agroecosystems. Factors that influence competition in agroecosystems include the timing of plant emergence, growth rates, spatial arrangements among neighbors, plant-plant-environment interactions, and herbivory. See ECOLOGY; ECOLOGY, APPLIED. Paul C. Marino

Bibliography. E. F. Connor and S. Simberloff, Competition, scientific method, and null models in ecology, *Amer. Sci.*, 74:155-162, 1986; S. R. Reice, Nonequilibrium determinants of biological community diversity, *Amer. Sci.*, 82:424-435, 1994; J. Weins, On competition and variable environments, *Amer. Sci.*, 65:590-597, 1977.

Ecological energetics

The study of the flow of energy within an ecological system from the time the energy enters the living system until it is ultimately degraded to heat and irretrievably lost from the system. It is also referred to as production ecology, because ecologists use the word production to describe the process of energy input and storage in ecosystems.

Ecological energetics provides information on the energetic interdependence of organisms within ecological systems and the efficiency of energy transfer within and between organisms and trophic levels. Nearly all energy enters the biota by green plants' transformation of light energy into chemical energy through photosynthesis; this is referred to as primary production. This accumulation of potential energy is used by plants, and by the animals which eat them, for growth, reproduction, and the work necessary to sustain life. The energy put into growth and reproduction is termed secondary production. As energy passes along the food chain to higher trophic levels (from plants to herbivores to carnivores), the potential energy is used to do work and in the process is degraded to heat. The laws of thermodynamics require the light energy fixed by plants to equal the energy degraded to heat, assuming the system is closed with respect to matter. An energy budget quantifies the energy pools, the directions of energy flow, and the rates of energy transformations within ecological systems. See BIOLOGICAL PRODUCTIVITY; FOOD WEB; PHOTOSYNTHESIS.

The peak of studies in ecological energetics occurred in the 1960s and early 1970s largely because a major concern of the International Biological

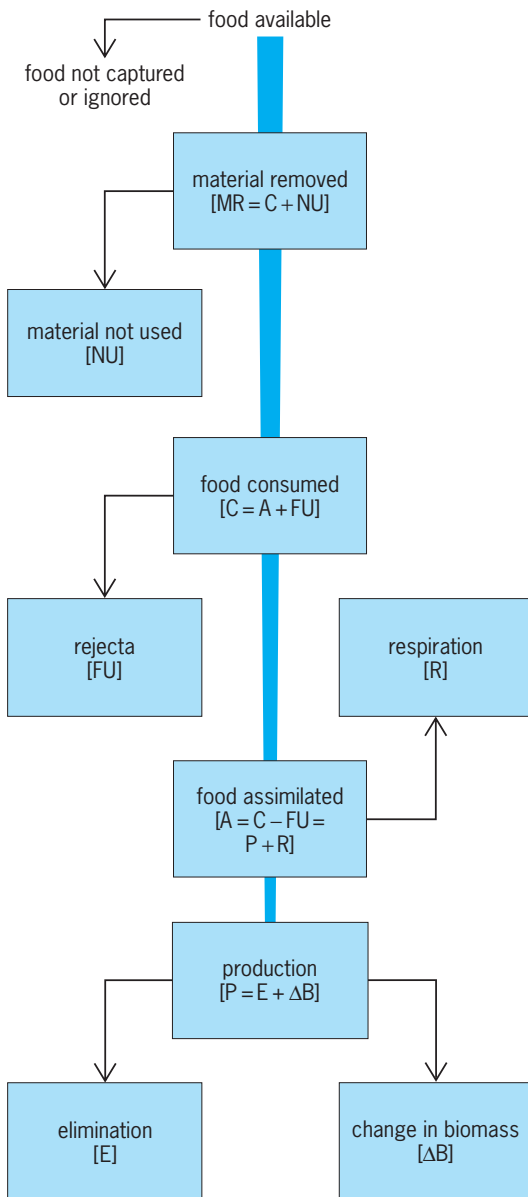


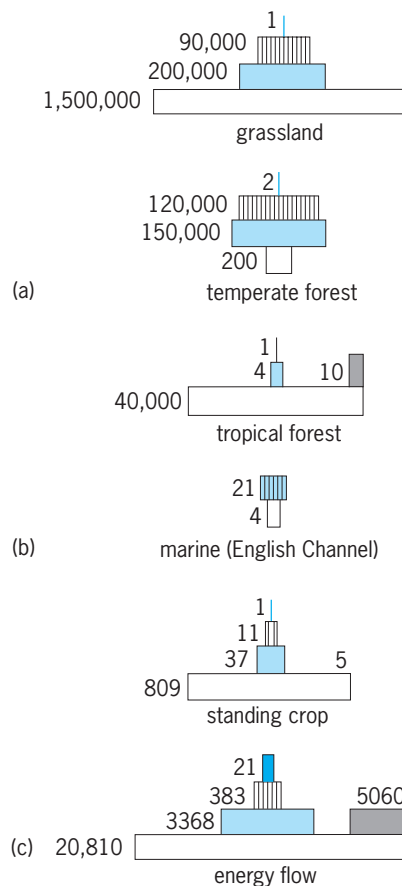
Fig. 1. Diagrammatic representation of energy flow through an ecological unit; abbreviations are explained in the text.

Program was an appraisal of the biological productivity of terrestrial and aquatic communities in relation to human welfare. Initially considered to have the potential of becoming a unifying language in ecology—an ecological Rosetta Stone—the subject has yielded little in the way of general theory.

The essentials of ecological energetics can be most readily appreciated by considering the schema (Fig. 1) of energy flowing through an individual; it is equally applicable to populations, communities, and ecosystems. Of the food energy available, only part is harvested (MR) in the process of foraging. Some is wasted (NU), for example, by messy eaters, and the rest consumed (C). Part of the consumed food is transformed but is not utilized by the body, leaving as fecal material (F) or as nitrogenous waste (U), the by-product of protein metabolism. The remaining energy is assimilated (A) into the body, part

of which is used to sustain the life functions and to do work—this is manifest as oxygen consumption. The remainder of the assimilated energy is used to produce new tissue, either as growth of the individual or as development of offspring. Hence production is also the potential energy (proteins, fats, and carbohydrates) on which other organisms feed. Production (P) leads to an increase in biomass (B) or is eliminated (E) through death, migration, predation, or the shedding of, for example, hair, skin, and antlers.

Pathways. Energy flows through the consumer food chain (from plants to herbivores to carnivores) or through the detritus food chain. The latter is fueled by the waste products of the consumer food chain, such as feces, shed skin, cadavers, and



Key:
 □ = producers
 □ = herbivores
 ▨ = carnivores
 ■ = top carnivores
 ■ = detritores

Fig. 2. Trophic levels of a number of ecosystems represented in different units. (a) As numbers of individuals per 1000 m² of grassland and temperate forest community in summer; microorganisms and soil animals excluded. (b) The standing crop or biomass (grams dry weight per meter squared) of terrestrial (Panamanian tropical rainforest) and marine (English Channel) communities; note the inversion of the marine pyramid. (c) The aquatic community of Silver Springs, Florida, represented as standing crop (kilocalorie per meter) and energy flow (kilocalories per meter per year). (After E. P. Odum, *Fundamentals of Ecology*, 3d ed., W. B. Saunders, 1971)

nitrogenous waste. Most detritus is consumed by microorganisms such as bacteria and fungi, although this food chain includes conspicuous carrion feeders like beetles and vultures. In terrestrial systems, more than 90% of all primary production may be consumed by detritus feeders. In aquatic systems, where the plants do not require tough supporting tissues, harvesting by herbivores may be efficient with little of the primary production passing to the detritivores.

Pyramids of biomass are used to depict the amount of living material, or its energetic equivalent, present at one time in the different trophic levels (Fig. 2). Although the energy flow cannot increase at higher trophic levels, pyramids of biomass may be inverted, especially in aquatic systems. This occurs because the index P/B is inversely related to the size of the organisms. Hence a small biomass may support a high level of production if the biomass is composed of small individuals (Fig. 3).

Units. Traditionally the calorie, a unit of heat energy, has been used, but this has been largely replaced by the joule. Confusion is possible, especially in the field of nutrition, because with an initial capital, Calorie may denote kilocalories. Biomass or standing crop is expressed as potential energy per unit area, but the other compartments in Fig. 1, for example P and R, are expressed in terms of energy flux or rates. The efficiency values such as P/A are dimensionless, but the ratio P/B is a rate—the inverse of the turnover time.

Measurement of energy flow. For illustrative purposes some general methods for assessing biological productivity are described here in the context of energy flow through a population. Production is measured from individual growth rates and the reproductive rate of the population to determine the turnover time. The energy equivalent of food con-

sumed, feces, and production can be determined by measuring the heat evolved on burning a sample in an oxygen bomb calorimeter, or by chemical analysis—determining the amount of carbon or of protein, carbohydrate, and lipid and applying empirically determined caloric equivalents to the values. The latter three contain, respectively, 16.3, 23.7, and 39.2 kilojoules per gram of dry weight. Maintenance costs are usually measured indirectly as respiration (normally the oxygen consumed) in the laboratory and extrapolated to the field conditions. Error is introduced by the fact that animals have different levels of activity in the field and are subject to different temperatures, and so uncertainty has surrounded these extrapolations. Oxygen consumption has been measured in animals living in the wild by using the turnover rates of doubly labeled water (D₂O).

Levels of inquiry. Ecological energetics is concerned with several levels of inquiry: the partitioning of energy between the compartments denoted in Fig. 1; the pathways traced by the energy as it passes through the trophic levels; and the efficiency of energy transfer between trophic levels. The ratio of energy flux through one compartment in Fig. 1 to any previous compartment is referred to as an efficiency. Numerous efficiencies can be calculated both within and between trophic levels. The most common are the assimilation efficiency (A/C), namely the proportion of energy assimilated by the body from the food consumed, and the production efficiency (P/A), which denotes the proportion of energy assimilated which ends up as new tissue. These various efficiencies combine to limit the energy available to the higher trophic levels. The ratio of food consumed or ingested at one trophic level to that ingested by the next lower level is termed ecological efficiency. A value of 10% for this efficiency is often cited; consideration of the A/C and P/A efficiencies

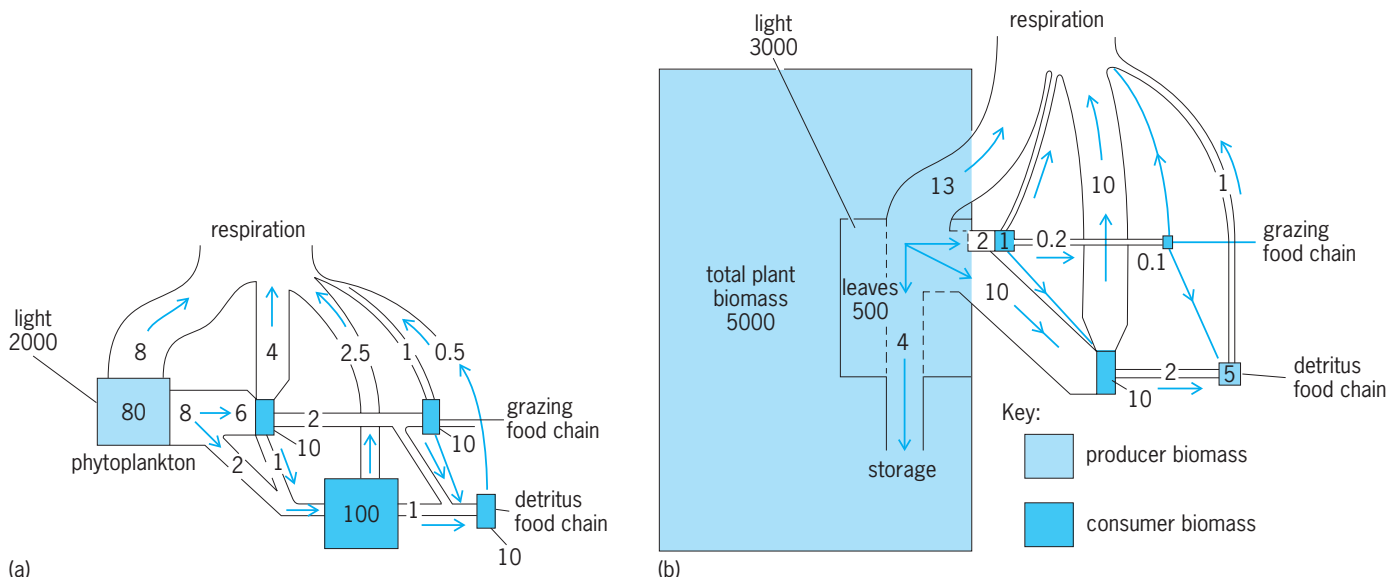


Fig. 3. Models of energy flow through two contrasting ecosystems: (a) a marine bay and (b) a forest. The boxes represent the biomass in kilocalories per meter and the flow lines show the energy flux (kilocalories per meter per day) along the various pathways. The boxes and the flux lines are scaled to indicate their relative magnitudes. (After E. P. Odum, *Relationships between structure and function in the ecosystem*, *Jap. J. Ecol.*, 12:108–118, 1962)

Production efficiency of populations of various classes of animals living in the wild*

Animal group	Production efficiency (P/A), %
Shrews	0.9
Birds	1.3
Other mammals	3.1
Fish, ants, and termites	9.8
Invertebrates other than insects	25.0
Herbivores	20.8
Carnivores	27.6
Detrivores	36.2
Insects except ants and termites	40.7
Herbivores	38.8
Detrivores	47.0
Carnivores	55.6

*After W. F. Humphreys, Production and respiration in animal populations, *J. Anim. Ecol.*, 48:427–454, 1979.

of most organisms shows that it could seldom exceed 15–20%. However, the effect of heat losses at each trophic level in limiting the length of food chains in nature remains controversial.

Factors affecting efficiency. Respiration rate of organisms is scaled as the three-quarters power of body weight. Hence larger organisms have proportionately slower rates of respiration. This scaling factor seems to affect many rate processes in the body so that size does not influence those efficiencies which are the focus of ecological energetics. However, different types of organisms of the same size have different metabolic rates. For example, warm-blooded animals have much higher weight-specific respiration rates than cold-blooded ones. Analysis of energy budgets derived for wild-living animals shows that a number of taxonomic and trophic groups can be distinguished according to characteristic production efficiencies (see **table**). Production efficiency appears to be related to the general level of metabolic activity—animals with high rates of metabolism generally having low production efficiency.

Due to the loss of usable energy with each transformation, in an area more energy can be diverted into production by plants than by consumer populations. For humans this means that utilizing plants for food directly is energetically much more efficient than converting them to eggs or meat. See BIOMASS; ECOLOGICAL COMMUNITIES; ECOSYSTEM.

W. F. Humphreys

Bibliography. S. Brody, *Bioenergetics and Growth*, 1945; R. M. May, Production and respiration in animal communities, *Nature*, 282:443–444, 1979; E. P. Odum, *Fundamentals of Ecology*, 3d ed., 1971; K. Petruszewicz and A. Macfadyen, *Productivity of Terrestrial Animals: Principles and Methods*, 1970; J. Phillipson, *Ecological Energetics*, 1966; S. L. Pimm, *Food Webs*, 1982.

Ecological methods

The methods required for the measurement of features of plant and animal populations or for quantifying energy flow or other aspects of ecological

communities. Ecologists are concerned with such a wide range of variables that they may use many other techniques, including in particular those for the measurement of weather conditions or for determining features and composition of air, water, and soil. See ECOLOGICAL COMMUNITIES.

Patterns of distribution. Ecologists often attempt to describe the way that individual plants or animals are distributed in the habitat. One pattern is the clumping of organisms. This may occur because only part of the habitat is suitable for them (as when rushes occur in the damp patches of a field) or because they represent the offspring of an original colonist and have not moved far from the parent. Another pattern is known as random distribution because each organism is independent of the other (as when this-tledown settles on a plowed field). A regular pattern describes the situation when organisms are virtually equidistant from each other. It is often a reflection of competition.

Such patterns are normally determined by sampling the organisms from unit areas—generally squares known as quadrats. For vegetation sampling, long pins resembling knitting needles are dropped into the ground and plants that touch them are recorded. These pins are regarded as infinitely small quadrats and are termed point quadrats. Various statistics, such as the mean (\bar{x}) and the variance (s^2), may be calculated from the results of such sampling, and their magnitudes and relationships to each other will show whether the pattern is clumped, random, or regular: if the ratio $s^2/\bar{x} \sim 1$ the individuals are randomly distributed; if the ratio is greater than 1 they are clumped; but if less than 1 they are in a regular pattern. See STATISTICS.

Estimation of populations. For plants, individuals (such as trees) may be counted directly, and the main problem is often (as with grasses) deciding what is a plant. However, the populations of only a few animals can be estimated by direct counting; in these cases, the organism is either sessile (like barnacles) or the individuals are very large, and all the individuals in an area can be seen (or at least recorded on a photograph) at the same time. Aerial photography and remote sensing techniques have provided powerful approaches for the estimation of the numbers of large animals (such as cattle or elephants) over wide areas or the extent and type of vegetation cover. However, most animals are active and small, and only a part (sample) of the population can be seen at any one time. Alternatively, it may be tedious and time-consuming to count the animals, so only those in small parts (samples) can be counted (for example, mites per leaf rather than per tree). See AERIAL PHOTOGRAPHY; REMOTE SENSING.

One approach to estimating population size is to mark a number of the animals (a), perhaps with a colored dye or paint, and release these into the population. If it can be assumed that they have mixed completely (and the marking has not affected their behavior), the proportion that they constitute of the next sample will be the same as the proportion that the original number marked constitutes of the

total population (N). This is shown in the equation below, where r is the number of marked individu-

$$\frac{r}{n} = \frac{a}{N} \quad \text{or} \quad N = \frac{an}{r}$$

als recaptured, n is the total number in the sample after marking, a is the number originally marked, and N is the population. This calculation also assumes that there are no births or deaths, or any migration in the period between the samples. Often these assumptions are not justified, and many more complex systems, of both sampling and marking (on several occasions) and of calculation, have been devised.

An alternative approach is to take a number of samples and determine the total number of animals in them. The average number per sample is then multiplied by the total number of sample units in the habitat (for example, leaves on a tree) to obtain the population for that habitat. Special techniques are often needed to determine the total animals in a sample. For example, mites and insects in samples of soil (or litter) are usually extracted by placing the soil in a special apparatus, the Berlese-Tullgren funnel, where a light bulb or other heat source slowly dries out the soil so that the animals move down and eventually drop through the funnel to a collecting tube. Other soil animals (for example, many types of worms) require moist conditions, and they can be driven out of wet samples by lowering the oxygen level (again by heat) or by chemicals. There is a wide variety of methods for different organisms and for various habitats. Special methods may also be necessary to take the samples: a common method of sampling insects from grassland is the D-vac sampler, basically a large, specially designed vacuum cleaner. See MATHEMATICAL ECOLOGY.

Population dynamics. The aim of population estimation is often to contribute toward an understanding of how the population is changing over time, such as to determine trends or the major mortalities. Methods are devised to measure the amount of mortality (for example, the remains of dead animals or plants can be found, or the level of parasitism assessed), the numbers born (natality), and the organisms moving in (immigration) and out (emigration) of the habitat. Special computational techniques are used to estimate the numbers of individuals in a population and to construct life tables. All these methods make certain assumptions, but since these are not always justified, there is great value in making population estimates by a number of different methods. See POPULATION ECOLOGY.

Energy flow. The contribution of organisms of different sizes to the functioning of an ecosystem may be best compared in terms of their total weight—the biomass. The biomass in a habitat is termed the standing crop. It is measured in wet weight or dry weight, the latter obtained by freeze drying. The dynamics of an ecosystem, that is, the exchanges between the different components, may be expressed in terms of energy transfers to form a picture of the energy flow. The energy content of a given biomass

is often determined by burning a sample in a bomb calorimeter; the heat rise is a measure of the calorific value of the sample. Then it is possible to determine the amount of animal or plant material produced (the production) and, for animals, the amount of food consumed. See BIOLOGICAL PRODUCTIVITY; BIOMASS.

Two equations are particularly useful for determining the energy relationships of organisms:

$$\begin{aligned} \text{Consumption} - (\text{fecal and urinary waste}) \\ &= \text{assimilation} \\ \text{Assimilation} &= \text{respiration and production} \end{aligned}$$

Consumption and production can usually be measured in terms of biomass and calorific value. Respiration is often measured with special respirometers, and assimilation can sometimes be measured directly by using an inert marker in the food.

Other methods are used to measure the flow of specific nutrients in ecosystems; special computational methods are used to calculate indices to describe the diversity of animal and plant communities; many types of traps or similar devices are used to capture animals; and radioactive isotopes can be used as markers for organisms or materials. Ecologists are always devising new methods for the many and varied problems they encounter. See ECOLOGICAL ENERGY; ECOLOGY; ECOSYSTEM. T. R. E. Southwood

Bibliography. R. Clarke (ed.), *The Handbook of Ecological Monitoring*, 1986; P. A. Erickson, *Environmental Impact Assessment: Principles and Applications*, 1987; J. A. Ludwig and J. F. Reynolds, *Statistical Ecology: Primer on Methods and Computing*, 1988; R. Southwood and P. A. Henderson, *Ecological Methods*, 3d ed., 2000; H. Walter and S. W. Breckle, *Ecological Systems of the Geobiosphere*, 1985.

Ecological modeling

The use of computer simulations or mathematical equations to address questions that cannot be answered solely by experiments or observations. Ecological models have two major aims: to provide general insight into how ecological systems or ecological interactions work; and to provide specific predictions about the likely futures of particular populations, communities, or ecosystems.

Models can be used to indicate general possibilities or to forecast the most likely outcomes of particular populations or ecosystems. Models differ in whether they are “basic” or are intended to address management decisions. As ecology has grown in its sophistication, models are increasingly used as decision support tools for policy-makers. See POPULATION ECOLOGY.

Areas of modeling. The earliest ecological models were produced by mathematicians attempting to explain cycles of predators and prey, disease epidemics, pest outbreaks, and human population

growth. These early models are the foundations of ecology and exemplify the value of mathematics as a tool for formulating biological questions. For example, the Italian mathematician Vito Volterra was asked to explain why the fish that are eaten by sharks decreased in relative abundance when fishing halted in the Mediterranean Sea during World War I. The models developed by Volterra not only accounted for this decline in shark prey with the war-related curtailment of harvest but also indicated that insect pests might actually increase in numbers if chemical insecticides killed both their predators and the insects themselves [a phenomenon observed 40 years later after the invention of dichlorodiphenyltrichloroethane (DDT) and other chemical insecticides]. More generally, Volterra's models predicted a tendency for predators and their prey to fluctuate even if the environment stayed the same from year to year—another prediction that is generally upheld by data.

Models of human population growth stimulated the development of mathematical tools for converting schedules of birth and death into predictions about annual rates of population change. Although these demographic models originated over 50 years ago, modifications of them provide the foundation for much modern conservation biology. For instance, by inputting rates of mortality and reproduction for loggerhead sea turtles into a simple model, it was found that extinction of the turtles could best be prevented if large young adults were protected from shrimp trawlers; legislation for this was passed by the U.S. Congress.

Models of virtually every possible type of ecological interaction have been developed (competition, parasitism, disease, mutualism, plant-herbivore interactions, and so forth). The models vary in their level of detail. Some models simply keep track of the density of organisms, treating all organisms of any species as identical (mass action models). At the other extreme, the movement and fate of each individual organism may be tracked in an elaborate computer simulation (individual behavior models).

Building a model. The best way to understand how ecological modeling works is to actually build a model. One of the most important general models in ecology describes the consequences of habitat loss for species that require specific types of habitat. The type of model described below is called a metapopulation model because it keeps track of many populations in an archipelago of patches. In order to model how such species are likely to be influenced by losses of habitat, ecologists have examined models that range from simple ("pencil-and-paper") exercises that can be expressed with a few equations to intricate computer models that involve hundreds of complex mathematical relationships and lengthy computer programs.

A useful pencil-and-paper model about the consequences of habitat destruction may be devised by carefully expressing ideas in mathematical language. For example, begin by simplifying the world

to patches of habitat that are either occupied by the species of interest (a fraction x) or not occupied by the species ($1 - x$). No species will occupy habitat forever, since there will inevitably be local death or even disappearance of an entire population within some patch of forest or alpine meadow. Conversely, empty or unoccupied habitat will not remain vacant forever, as long as there is some source of colonists that can find their way to the open habitat. A simple equation (1) summarizes the change

$$\text{change} = \left(\begin{array}{c} \text{gains in } x \\ \text{due to colonization} \\ \text{of empty patches} \end{array} \right) - \left(\begin{array}{c} \text{losses in } x \\ \text{due to disappearance} \\ \text{from occupied patches} \end{array} \right) \quad (1)$$

in species abundance expected in such a world of habitat patches (occupied or unoccupied) with perpetual turnover; here **change** represents the change in the fraction of occupied patches per unit time.

This expression may be made more precise by hypothesizing that there is a constant loss or disappearance of the species from occupied patches at a rate e ; this means the "minus term" of Eq. (1) becomes $-ex$. It is reasonable to assume that vacant patches of habitat ($1 - x$) get colonized in proportion to the fraction of occupied habitat from which colonists might arise (c times x , where c indicates the dispersal or colonization ability of the species, and x is the proportion of habitat occupied as mentioned above). This is a specific model that predicts the fraction of habitat patches that should be occupied by an organism, and how that fraction changes through time [Eqs. (2)–(4)]. The species will increase as long

$$\text{change} = cx(1 - x) - ex \quad (2)$$

$$\text{change} = cx(1 - x - D) - ex \quad (3)$$

$$(1 - D - e/c) \quad (4)$$

as **change** is greater than zero. **Change** is greater than zero if x is less than $(1 - e/c)$; and **change** is less than zero if x is greater than $(1 - e/c)$. Hence x will tend to move toward the fraction given by $(1 - e/c)$, and it should be clear why $(1 - e/c)$ can be called the equilibrium fraction of habitat occupied by the species of interest. As long as colonization rates are sufficiently high (that is, c is greater than e), the species will persist forever. Habitat destruction permanently removes a fraction D of the vacant habitat; thus, instead of $(1 - x)$ fraction of habitat being available to be colonized, only $(1 - x - D)$ habitat is available to be colonized. With habitat destruction, Eq. (2) is modified to Eq. (3). With habitat loss incorporated through the D term, the equilibrium fraction of habitat occupied by a species is now (4). Thus, if habitat loss D is too large, the species will become extinct because the equilibrium will be less than zero, which means the species will decline until it has totally disappeared. This model makes several important points. For example, even though there is empty habitat (which means $1 - x$ is greater than 0), the species may still relentlessly

decline simply because so much habitat has been destroyed that the colonization of empty habitat occurs too infrequently to sustain the species. This exposes the fallacy of arguments that have been made regarding old growth forests and the endangered spotted owl species in the Pacific Northwest: the notion that there is ample old growth forest for this endangered species because some of the forest is devoid of owls is clearly wrong. The model also shows that there is a threshold for habitat destruction (for example, when habitat loss D becomes greater than $1 - e/c$), which when crossed implies inevitable gradual decline of a species to extinction.

A key feature of ecological modeling is a parameter estimation. The parameters in the model related to Eq. (4) are e , c , and D . If the value of these parameters was known, it would be possible to predict whether or not the species of interest could persist, and if it could persist, what fraction of habitat would be expected to be occupied by that species. Thus, in order to apply the model even crudely to any situation, it is necessary to estimate the extinction rate or disappearance rate (fraction lost per year or week), the colonization rate, and the fraction of habitat that has been lost altogether. This might seem like an easy task, but in most circumstances the factor that most thwarts ecological modelers is estimating the parameters. Because parameters may be difficult to estimate, there is typically a great deal of uncertainty surrounding the application of any model. For example, suppose it is known that the disappearance rate for spotted owls in old growth forest patches is 3% loss per year and the colonization rate is 5%, but it is uncertain whether or not 35% or 45% of the habitat has been destroyed. If $D = .35$ the species could persist (because $1 - .35 - .03/.05$ is greater than 0), but if $D = .45$ the species would go extinct (because $1 - .45 - .03/.05$ is less than 0). Any uncertainty about habitat loss D would yield uncertainty about the fate of the species.

Simple algebraic models such as Eqs. (1)-(4) are very useful for indicating general principles and possibilities. They were part of the testimony considered in court cases regarding the halting of logging in Washington and Oregon so that spotted owls would not cross the threshold for habitat loss, at which point the species would be doomed. But something as simple as Eq. (4) could never be practically applied to the management of any specific forest and spotted owl population. In order to be a management tool, the model must be more complicated and detailed to reflect the specific situation under examination. For example, instead of a few equations, ecologists have modeled spotted owl populations and old growth forests in Washington using a detailed computer simulation that keeps track of habitat in real maps at the scale of hectares. In these simulation models, owls are moved as individuals from one hectare to another, and their fate (survival, death, or reproduction) is recorded in the computer's memory. By tracking hundreds or even thousands of owls moving around in this computer world, different forestry practices corresponding to differ-

ent logging scenarios can be examined. See ECOLOGY, APPLIED; MATHEMATICAL ECOLOGY; SYSTEMS ECOLOGY.

Advantages and limitations. A model is a formal way of examining the consequences of a series of assumptions about how nature works. Such models refine thinking and clarify what results are implied by any set of assumptions. As models become more complicated and specific, they can also be used as surrogate experimental systems; in other words, ecologists may use models to conduct experiments that are too expensive or impractical in the field. For example, to protect against extinction of northern spotted owls in the northwestern United States, investigators do not have the time or the resources to engage in thousand-acre experiments that entail different forestry practices (for example, logging in a way that leaves behind clusters of old growth forest versus logging randomly across the landscape). However, if enough is known about how owls move through forests, a simulation model can be devised that "does the experiment," which may suggest the best logging practice.

One danger of ecological modeling is the uncertainty of the models and the shortage of supporting data. Properly used, models allow exploration of a wide range of uncertainty, pointing out the limits of current knowledge and identifying critical information required prior to management decision making. Increasingly, the public will be asked to contribute to decisions about logging, sprawling urban development, removing dams from rivers, and so forth; critical information regarding these decisions will be provided by ecological models. It is therefore essential that citizens understand how models are formulated and applied. Conclusions based on models should not be quickly dismissed because the model may be the best line of reasoning available; however, it would not be prudent to rely solely on the output of any model. Ecological modeling is a branch of science that requires expertise in ecology, mathematics, statistics, and usually computer programming. Done well, ecological modeling can be tremendously valuable as a source of general scientific principles and as a decision-support tool in arenas of resource management and conservation. See ECOLOGY. Peter Kareiva

Bibliography. N. Gotelli, *A Primer of Ecology*, 2d ed., Sinauer Press, Sunderland, MA, 1998; J. Maynard Smith, *Models in Ecology*, Cambridge University Press, Cambridge, 1974; J. Roughgarden, *A Primer of Ecological Theory*, Prentice Hall, Upper Saddle River, NJ, 1998.

Ecological succession

A directional change in an ecological community. Populations of animals and plants are in a dynamic state. Through the continual turnover of individuals, a population may expand or decline depending on the success of its members in survival and reproduction. As a consequence, the species composition of communities typically does not remain static with

time. Apart from the regular fluctuations in species abundance related to seasonal changes, a community may develop progressively with time through a recognizable sequence known as the sere. Pioneer populations are replaced by successive colonists along a more or less predictable path toward a relatively stable community. This process of succession results from interactions between different species, and between species and their environment, which govern the sequence and the rate with which species replace each other. The rate at which succession proceeds depends on the time scale of species' life histories as well as on the effects species may have on each other and on the environment which supports them. See ECOLOGICAL COMMUNITIES; POPULATION ECOLOGY.

The course of ecological succession depends on initial environmental conditions. Primary succession occurs on novel areas such as volcanic ash, glacial deposits, or bare rock, areas which have not previously supported a community. In such harsh, unstable environments, pioneer colonizing organisms must have wide ranges of ecological tolerance to survive. In contrast, secondary succession is initiated by disturbance such as fire, which removes a previous community from an area. Pioneer species are here constrained not by the physical environment but by their ability to enter and exploit the vacant area rapidly.

As succession proceeds, many environmental factors may change through the influence of the community. Especially in primary succession, this leads to more stable, less severe environments. At the same time interactions between species of plant tend to intensify competition for basic resources such as water, light, space, and nutrients. Successional change results from the normal complex interactions between organism and environment which lead to changes in overall species composition. Whether succession is promoted by changing environmental factors or competitive interactions, species composition alters in response to availability of niches. Populations occurring in the community at a point in succession are those able to provide propagules (such as seeds) to invade the area, being sufficiently tolerant of current environmental condi-

tions, and able to withstand competition from members of other populations present at the same stage. Species lacking these qualities either become locally extinct or are unable to enter and survive in the community.

Primary succession. In some cases, seres may take hundreds of years to complete, and direct observation at a given site is not possible. Adjacent sites may be identified as successively older stages of the same sere, if it is assumed that conditions were similar when each seral stage was initiated.

Glacier Bay. In the Glacier Bay region of Alaska glaciers have retreated, in phases, some 61 mi (98 km) since 1750, leaving a series of moraines of known ages supporting a range of seral vegetational types. Soil factors become modified by vegetation, enabling new species to become established (Table 1). Acidic decomposition products of alder leaves sharply reduce soil pH, whereas virtually no change occurs on bare glacial till or under other vegetation. Pioneer species must tolerate low nitrogen levels, but alder is able to fix atmospheric nitrogen (N) by the presence of microbial symbionts in root nodules and is correlated with a rapid increase in soil nitrogen by way of leaf fall. After invasion by spruce, levels of accumulated nitrogen fall as nitrogen becomes incorporated in forest biomass, and the annual additions from alder are reduced. Soil organic matter increases progressively, and influences the structural development of the soil. The mature forest remains only on well-drained slopes. In areas of poorer drainage, invasion by *Sphagnum* moss leads to replacement of trees by wet acidic bog or muskeg.

Lake Michigan. Another example of primary succession has been recorded on a sequence of dune ridges bordering Lake Michigan, ranging in age to 12,000 years. Succession is initiated on a bare sand surface either following a fall in lake level, as occurred in phases since the last glaciation, or due to wind erosion of an existing dune redepositing sand (Table 2). Marram grass impedes the transport of sand across bare dune surfaces and so promotes accretion of sand. At the same time, it grows and branches vigorously, thus maintaining cover on the expanding dune (Fig. 1).

TABLE 1. Successional changes in vegetation and soils observed on an aged series of moraines at Glacier Bay, Alaska

Years	Vegetation	Soil environment
0	Initial colonizers: mosses, fireweed, horsetail, <i>Dryas</i> , dwarf willows; later, willows form dense scrub	Initially pH 8.0–8.4 due to CaCO ₃ ; soil N and organic matter lacking
50	Alder invades, forming dense thickets less than 33 ft (10 m) tall	pH falls to 5.0 in 30–50 years due to acidic products of alder leaf decomposition; marked increase in soil N due to fixation in alder root nodules; soil organic matter accumulates
170	Sitka spruce invades, forming dense forest	Reduction in soil N by incorporation in forest biomass; progressive increase in soil organic matter
250+	Western and mountain hemlock enter (climax forest on well-drained slopes) <i>Sphagnum</i> bog with occasional pines replaces forest in poorly drained areas	Soil becomes waterlogged, deoxygenated, acidified

TABLE 2. Succession on sand dunes of Lake Michigan

Years	Vegetation	Soil environment	Ground invertebrates
0	Marram grass (<i>Ammophila</i>) invades by rhizome migration	Initial sand pH 7.6	
6		Sand accumulates, dune builds up	White tiger beetle (<i>Cicindela lepida</i>) Sand spider (<i>Trochosa cinerea</i>)
20	Marram declines in vigor	Sand surface stabilized	White grasshopper (<i>Trimerotropis maritima</i>) Longhorn grasshopper (<i>Psinidia fenestralis</i>) Burrowing spider (<i>Geolycosa pikei</i>)
50	Jack pine and white pine become established, with light-demanding understory		Digger wasps (<i>Bembex</i> , <i>Microbembex</i>) Bronze tiger beetle (<i>Cicindela scutellaris</i>) Ant (<i>Lasius niger</i>)
100	Black oak becomes established, with shade-tolerant understory		Migratory locust (<i>Melanoplus</i>) Sand locusts (<i>Agoeotettix</i> , <i>Spharagemon</i>) Digger wasp (<i>Sphex</i>)
150			Ant lion (<i>Cryptoleon</i>) Flatbug (<i>Neuroctenus</i>)
•			Grasshoppers (six species)
•			Wireworms (<i>Elateridae</i>)
1000		Soil N builds up to 0.1%	Snail (<i>Mesodon thyroides</i>)
•			
10,000		Sand pH drops to 4.0 due to CaCO ₃ ; soil matures to a nutrient-poor brown humic sand	

Succession on Lake Michigan dunes demonstrates that local factors may modify the typical pattern (Fig. 2). In damp depressions with impeded drainage, a grassland community develops. Sheltered pockets on lee slopes have a moister microclimate and tend to accumulate leaf litter from more exposed parts of the dunes, as well as receiving protection from frequent fires. As a result, a more nutrient-rich soil can develop, and succession proceeds via basswood to more nutrient-demanding oak-hickory and finally maple-beech woodland, typical of normal soils of the region. The black oak association appears to be stable on dune soils, since it is tolerant of low nutrients and water limitation, and tends to maintain these conditions by returning few nutrients in the leaf litter.

Climax community. Early stages of succession tend to be relatively rapid, whereas the rates of species

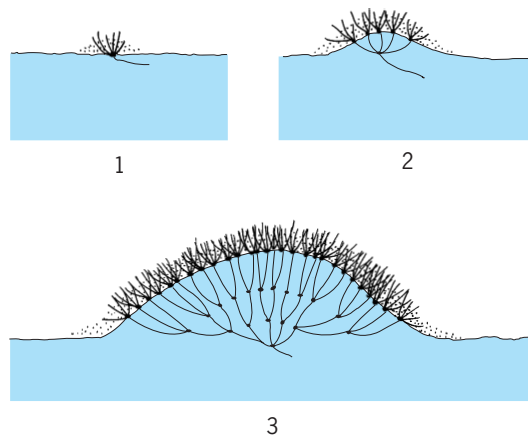


Fig. 1. Dune formation by the gradual deposition of wind-carried sand particles around aerial shoots of *Ammophila arenaria*. 1, 2, and 3 indicate successive periods of time. (After K. A. Kershaw, *Quantitative and Dynamic Ecology*, 2d ed., Arnold, 1973)

turnover and soil changes become slower as the community matures. Eventually an approximation to the steady state is established with a relatively stable community, the nature of which has aroused considerable debate. Earlier, the so-called climax vegetation was believed to be determined ultimately by regional climate and, given sufficient time, any community in a region would attain this universal condition. This unified concept of succession, the monoclimax hypothesis, implies the ability of organisms progressively to modify their environment until it can support the climatic climax community. Although plants and animals do sometimes ameliorate environmental conditions, evidence suggests overwhelmingly that succession has a variety of stable end points. This hypothesis, known as the polyclimax hypothesis, suggests that the end point of a succession depends on a complex of environmental factors that characterize the site, such as parent material, topography, local climate, and human influences.

In the Lake Michigan sand dunes, the course of succession and its climax appear to be determined by physiographic conditions at the start (Fig. 2). Similarly, the transformation of glacial moraine forest to muskeg depends on local drainage. In the tropical rainforest of Moraballi Creek, Guyana, five apparently stable vegetation types have been distinguished on different soil types under the same climate. A mixed forest is present on red loam, whereas the Wallaba forest occurs on bleached sand, with the Mora forest type in areas liable to flooding.

Autogenic vs. allogenic factors. In the examples of succession discussed above, the chief agent in modifying the environment is the community itself: thus marram stabilizes the sand dune surface, and alder increases the soil nutrient status of moraine soil. These actions of the community on the environment, termed autogenic, provide an important driving force promoting successional change, and

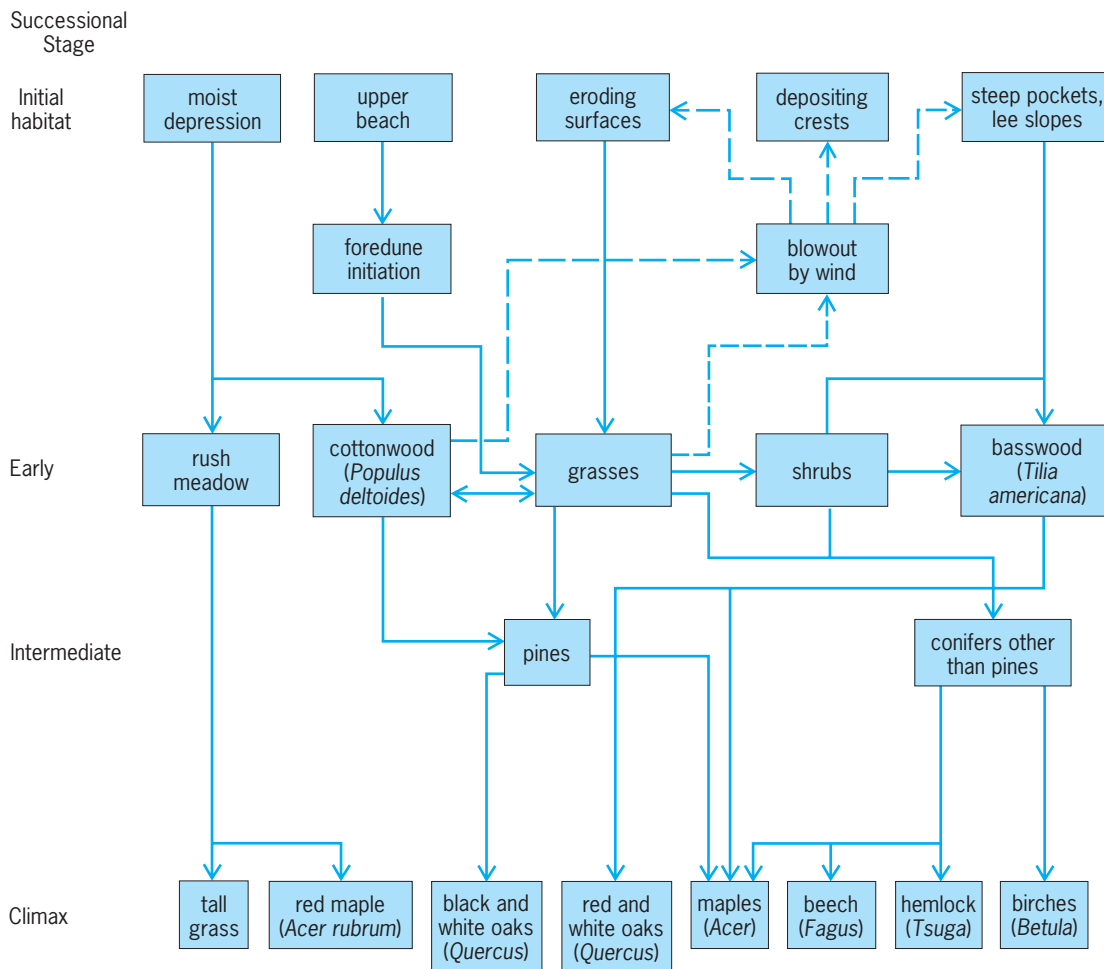


Fig. 2. Alternative patterns of primary plant succession on Lake Michigan dunes, depending on initial conditions, extrinsic variables, and colonization-invasion patterns. (After S. J. McNaughton and L. L. Wolf, *General Ecology*, 2d ed., Holt, Reinhart, and Winston, 1979)

are typical of primary succession where initial environments are inhospitable. Alternatively, changes in species composition of a community may result from influences external to the community called allogenic. For example, in aquatic ecosystems (the hydrosere) the community commonly develops from open water with submerged and floating aquatic plant species toward a swamp community in which rooted emergent plants dominate in shallower water, until finally the marsh is colonized by land plants from the surrounding area as sediment deposition continues and the soil dries out. Reduction in water depth, enabling colonization by marsh species and finally terrestrial species, occurs with input of waterborne and airborne sediments—thus the aquatic phase of the hydrosere is controlled by input of materials from outside the system. Similarly, lakes are typically subject to enrichment of nutrients from surrounding areas, resulting in increased productivity. Extremely high production occurs in culturally eutrophic lakes, which receive nutrient inputs from human activities. In aquatic systems where the influence of allogenic factors such as siltation are apparently minimal, vegetation tends to develop via a series of productive intermediate steps toward an

oligotrophic community, that is, one with low productivity, dominated by *Sphagnum* moss (Fig. 3). See EUTROPHICATION.

Whereas intrinsic factors often result in progressive successional changes, that is, changes leading from simple to more complex communities, external (allogenic) forces may induce retrogressive succession, that is, toward a less mature community. For example, if a grassland is severely overgrazed by cattle, the most palatable species will disappear. As grazing continues, the grass cover is reduced, and in the open areas weeds characteristic of initial stages of succession may become established.

Heterotrophic succession. In the preceding examples of succession, the food web was based on photosynthetic organisms and there was a slow accumulation of organic matter, both living and dead. This is termed autotrophic succession. In some instances, however, addition of organic matter to an ecosystem initiates a succession of decomposer organisms which invade and degrade it. Such a succession is called heterotrophic. In an Illinois pasture, cow dung is degraded by a seral community of some 40 to 60 invertebrate species over a period of 30 days. The newly deposited cow pat is immediately visited by

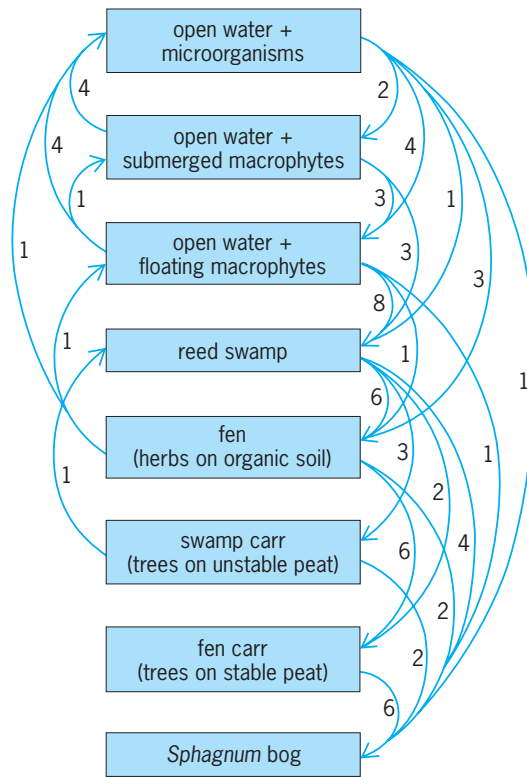


Fig. 3. Transitions between vegetation “stages” free from obvious allogenic influences derived from 20 pollen diagrams for lake sediments and peat. Figures show numbers of observed transitions.

the horn fly (*Haematobia irritans*), which quickly lays eggs and returns to the cow. It is followed by several other dung flies whose larvae are eaten by beetles such as *Sphaeridium scaraboides* which burrow through the dung and lay their eggs. A parasitic wasp (*Xyalophora quinquelinata*) deposits eggs inside maggots of *Sarcophaga* flies. As the dung ages and dries out, it is inhabited by a wider variety of species. Little is known in detail of the importance of the various saprovores in degrading the dung and in modifying the microhabitat, or of their dependence on the activities of saprophytic fungi.

Discharge of organic effluent into a river is detectable downstream by a progression in chemical factors and in the biota. Succession in time is here equivalent to the change in species composition resulting from the decline in organic effluent that corresponds to the distance from the discharge. Marked reduction in dissolved oxygen directly below an outfall results from respiration of microorganisms as they degrade organic matter. Detritivores tolerant of low oxygen concentrations, such as *Tubificidae* and *Chironomidae*, attain high population densities in the bottom sediments. Subsequently, a bloom of algae is typical, utilizing released nitrate and phosphate. As the river flows downstream, the aquatic food web progressively changes from a heterotrophic to an autotrophic basis and productivity declines to its normal level as the “clean water” community returns. See BIOLOGICALS; FOOD WEB; PRODUCTIVITY.

Secondary succession. Following the partial or complete destruction of an established community by disturbing events such as fire or clearfelling, and similarly on the cessation of grazing or tillage, a sequence of species invasion and replacement ensues. Compared to the slow initial progress of primary succession in which amelioration of the environment plays an important part, secondary succession is characterized initially by rapid turnover of typically opportunist species which invade relatively congenial habitats.

Piedmont. Abandoned fields in Piedmont area of North Carolina show a rapid sequence of replacement of herbaceous species, apparently related to the life histories of the plants (Table 3). Horseweed produces seeds in late summer which germinate immediately so that the plant overwinters as a juvenile, grows rapidly the following year, and dies after seeding in the summer. Aster seeds do not germinate until the following spring, and seedlings grow slowly due to shading by established horseweed plants. In addition, decaying horseweed roots inhibit its own growth and, to a greater extent, that of aster. Horseweed attains dominance in the first year by efficient seed dispersal and rapid establishment. Being a perennial, aster is able to outcompete horseweed in the second year despite the inhibitory effect of the latter. Seedlings of aster are present in abundance in third-year fields but are less drought-resistant than those of broomsedge, which outcompetes aster except in fields with more available water, where aster survives the competition longer. See ALLELOPATHY.

Broomsedge seeds are not available to colonize initially because seeds are not produced until the end of the plant’s second year and require a period of cold dormancy before germination. Seedling growth is apparently enhanced by decomposition products of the previous colonists. The late establishment of broomsedge in the succession is dependent not on changes brought about by earlier colonists but on the life history of the plant.

TABLE 3. Secondary succession on abandoned fields in the Piedmont area of North Carolina*

Years after last cultivation	Dominant plant	Other common species
0 (autumn)	Crabgrass (<i>Digitaria sanguinalis</i>)	
1	Horseweed (<i>Erigeron canadensis</i>)	Ragweed (<i>Ambrosia elatior</i>)
2	Aster (<i>Aster pilosus</i>)	Ragweed
3	Broomsedge (<i>Andropogon virginicus</i>)	
5-15	Shortleaf pine (<i>Pinus echinata</i>)	Loblolly pine
50-150	Hardwoods (oaks)	Hickory

*From C. Krebs, *Ecology*, 2d ed., Harper and Row, 1978.

After broomsedge, shortleaf pine invades the herb community. Pine seeds require mineral soil and minimal root competition to become established, and the seedlings are not shade-tolerant. Hence after about 20 years, under a dense pine canopy, reproduction of pines is almost lacking. Accumulation of litter and shade under pines causes the old-field herbs to die out.

Oak seedlings become established after about 20 years, when the depth of litter is adequate to prevent desiccation of acorns. Organic matter in the soil surface layer also increases, improving its water-holding capacity. After about 50 years, several oak species become established and gradually assume dominance as the pines fail to reproduce. Unlike pine, which is capable of germinating on bare soil, oaks and other hardwoods require changes in the soil resulting from pine litter before their seedlings can establish successfully.

While plant species are the main indicators of succession, it is important to note that animal species are also changing over time. In the Piedmont, the changes in bird species as succession proceeds have been well documented. Just a few species, such as meadowlarks and grasshopper sparrows, are found in the initial stages, while more complicated assemblages of species are common in the latter forest stages.

Nova Scotia forest. Following clearfelling in a Nova Scotia forest, the course of secondary succession involves invasion by shrubs (raspberry) followed by understory trees (pincherry, aspen), followed by shade-intolerant species (red maple, paper birch), and finally a shade-tolerant community (hard maple, yellow birch, white ash). Perhaps shade-tolerant species such as red maple, which are of low commercial value, could be inhibited if strip felling were practiced by the forestry industry as an alternative to clearfelling in large blocks. Commercially desirable shade-tolerant species such as white ash would be favored where there was greater local shading of the regenerating community.

Mechanisms of species replacement. Observed changes in the structure and function of seral communities result from natural selection of individuals within their current environment. Three mechanisms by which species may replace each other have been proposed; the relative importance of each apparently depends on the nature of the sere and stage of development.

1. The facilitation hypothesis states that invasion of later species depends on conditions created by earlier colonists. Earlier species modify the environment so as to increase the competitive ability of species which are then able to displace them. Succession thus proceeds because of the effects of species on their environment.

2. The tolerance hypothesis suggests that later successional species tolerate lower levels of resources than earlier occupants and can invade and replace them by reducing resource levels below those tolerated by earlier occupants. Succession proceeds despite the resistance of earlier colonists.

3. The inhibition hypothesis is that all species resist invasion of competitors and are displaced only by death or by damage from factors other than competition. Succession proceeds toward dominance by longer-lived species.

None of these models of succession is solely applicable in all instances; indeed most examples of succession appear to show elements of all three replacement mechanisms. In secondary succession on North Carolina croplands, stimulation of broomsedge growth by decomposition products of previous colonists, and the requirement of oak seedlings for a deep litter layer in which to germinate, exemplify facilitation. The ability of broomsedge to displace aster in competition for water suggest the tolerance mechanism, whereas the inhibition hypothesis is supported by the greater tolerance of horseweed seedlings than aster seedlings to horseweed decomposition products.

Deterministic vs. stochastic succession. Succession has traditionally been regarded as following an orderly progression of changes toward a predictable end point, the climax community, in equilibrium with the prevailing environment. This essentially deterministic view implies that succession will always follow the same course from a given starting point and will pass through a recognizable series of intermediate states (such as in Fig. 2). In contrast, a more recent view of succession is based on adaptations of independent species. It is argued that succession is disorderly and unpredictable, resulting from probabilistic processes such as invasion of propagules and survival of individuals which make up the community. Such a stochastic view reflects the inherent variability observed in nature and the uncertainty of environmental conditions. In particular, it allows for succession to take alternative pathways and end points dependent on the chance outcome of interactions among species and between species and their environment.

Consideration of community properties such as energy flow supports the view of succession as an orderly process. Early in autotrophic succession gross primary productivity (P_g) increases rapidly with community biomass (B), whereas community respiration (R) increases more slowly (Fig. 4). As a result, net primary productivity (P_n , where $P_n = P_g - R$) builds up early in succession, and the ratio P_g/B is at its

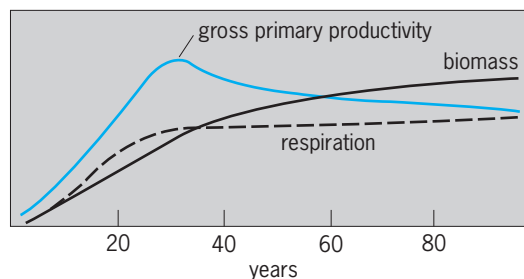


Fig. 4. The energetics of ecosystem development in a forest. The difference between gross primary productivity and respiration is the net primary productivity. (After E. P. Odum, *Ecology*, 2d ed., Holt, Reinhart, and Winston, 1975)

highest in the initial stages. As the community increases in biomass and complexity over time, more complete overall utilization of basic resources such as light limits further increase in primary productivity, whereas R continues to increase because of the increase in tissue to support. Hence zero and the biomass of a mature forest community no longer accumulates. The rate of gross primary productivity typically becomes limited also by the availability of nutrients, now incorporated within the community biomass, and declines to a level sustainable by release from decomposer organisms. Species diversity tends to rise rapidly at first as successive invasions occur, but declines again with the elimination of the pioneer species by the climax community.

Trends in community function, niche specialization, and life history strategy are summarized in **Table 4**. As the community acquires increasing maturity, P_n declines to zero, nutrients become incorporated in biotic pools, broad-niched species are replaced by those with more specific requirements,

and the structural organization of the community increases.

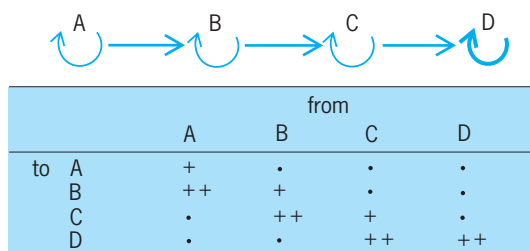
Regeneration of an area of subtropical rainforest in Queensland was observed after the vegetation and surface litter were removed with a bulldozer. Because of small environmental differences in the 10 quadrats observed, succession took four directions after the demise of the first ubiquitous colonizing species, resulting in four apparently stable plant associations. This divergence may result merely from small-scale variation in topography within the 65×130 ft (20×40 m) experimental site and from differing efficiencies of removal of surface litter between the quadrats. Hence the different plant associations detected could be interpreted as divergent products of succession or as phases within a larger-scale vegetation unit.

Stochastic aspects of succession can be represented in the form of models which allow for transitions between a series of different "states." Such models, termed Markovian models, can apply

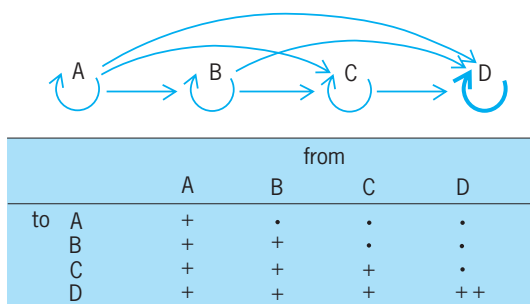
TABLE 4. Proposed successional trends in ecosystem structural and functional organization, species characteristics, evolutionary factors, and homeostasis*

Ecosystem property	Ecosystem stage	
	Successional (immature)	Climax (mature)
<i>Energy Flow</i>		
Gross productivity/respiration (P_g/R)	Autotrophic >1 Heterotrophic <1	
Biomass supported/unit energy flow (B/P_g)	Low	High
Net productivity (P_n)	High	Low
Type of food chains	Linear, grazing	Webs, detritus
<i>Nutrient Flow</i>		
Mineral cycles	Open	Closed
Flow rate: organism-environment	Rapid	Slow
Role of detritus	Unimportant	Important
<i>Community Structure</i>		
Total organic matter	Little	Much
Location of chemicals	Habitat pools	Biotic pools
Species richness	Low	High
Species evenness	Low	High
Biochemical diversity	Low	High
Spatial heterogeneity	Low	High
<i>Species Life History Characteristics</i>		
Niche breadth	Broad	Narrow
Organism size	Small	Large
Life cycles	Short, simple	Long, complex
<i>Selection Pressure</i>		
Growth form	Rapid growth (r -selection)	Feedback control (K selection)
Production	Quantity	Quality
<i>Overall Homeostasis</i>		
Internal symbiosis	Undeveloped	Developed
Nutrient conservation	Poor	Good
Resistance to perturbations	Poor	Good
Entropy	High	Low
Information	Low	High

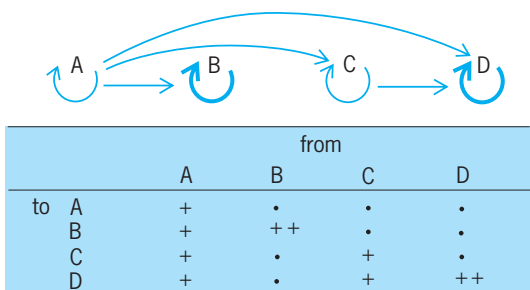
*After Odum, 1976.



(a)



(b)



(c)

Fig. 5. Schematic representation of three postulated mechanisms for species replacement in succession: (a) facilitation, (b) tolerance, and (c) inhibition. Arrows indicate the direction in which the systems tend to move between different states, A–D. Relative probabilities of movement are indicated by thickness of arrows and by symbols in the accompanying transition matrix, where • = close to zero, + = moderate, and ++ = high.

at various levels: plant-by-plant replacement, changes in tree size categories, or transitions between whole communities. A matrix of replacement probabilities defines the direction, pathway, and likelihood of change, and the model can be used to predict the future composition of the community from its initial state. With alternative transition matrices, this simple model could represent a linear progression toward a stable end state, or a cyclical, recursive sequence of communities. The three postulated mechanisms for species replacement discussed above can be illustrated by the topology of alternative Markovian models (Fig. 5). The facilitation model of succession is represented as a linear sequence with greater probabilities of progression toward the final state (D) than maintenance of intermediate states. In the tolerance model, later stages may develop from earlier stages depending on the availability of propagules of subsequent stages and their competitive ability. Again, only state

D has a high probability of self-replacement. In the inhibition model, there is a high probability that an intermediate state (in this case, B) will persist by strong self-replacement, thereby truncating the normal succession toward state D. A very low probability that B will change to C or D is assumed. Hence, a high degree of realism can be achieved with a simple model system, and alternative predictions of successional changes can be compared with observed data. See ECOLOGY; STOCHASTIC PROCESS.

Peter Randerson

Bibliography. D. H. Boucher (ed.), *The Biology of Mutualism: Ecology and Evolution*, 1985, reprint 1989; R. E. Ricklefs and G. L. Miller, *Ecology*, 4th ed., 1999; P. Yodis, *Introduction to Theoretical Ecology*, 1988.

Ecology

The subdiscipline of biology that concentrates on the relationships between organisms and their environments; it is also called environmental biology. Ecology is concerned with patterns of distribution (where organisms occur) and with patterns of abundance (how many organisms occur) in space and time. It seeks to explain the factors that determine the range of environments that organisms occupy and that determine how abundant organisms are within those ranges. It also emphasizes functional interactions between co-occurring organisms. In addition to its character as a unique component of the biological sciences, ecology is both a synthetic and an integrative science since it often draws upon information and concepts in other sciences, ranging from physiology to meteorology, to explain the complex organization of nature.

Environment is all of those factors external to an organism that affect its survival, growth, development, and reproduction. It can be subdivided into physical, or abiotic, factors, and biological, or biotic, factors. The physical components of the environment include all nonbiological constituents, such as temperature, wind, inorganic chemicals, and radiation. The biological components of the environment include the organisms. A somewhat more general term is habitat, which refers in a general way to where an organism occurs and the environmental factors present there. See ENVIRONMENT.

A recognition of the unitary coupling of an organism and its environment is fundamental to ecology; in fact, the definitions of organism and environment are not separate. Environment is organism-centered since the environmental properties of a habitat are determined by the requirements of the organisms that occupy that habitat. For example, the amount of inorganic nitrogen dissolved in lake water is of little immediate significance to zooplankton in the lake because they are incapable of utilizing inorganic nitrogen directly. However, because phytoplankton are capable of utilizing inorganic nitrogen directly, it is a component of their environment. Any effect of inorganic nitrogen upon the zooplankton,

then, will occur indirectly through its effect on the abundance of the phytoplankton that the zooplankton feed upon. See PHYTOPLANKTON; ZOOPLANKTON.

Just as the environment affects the organism, so the organism affects its environment. Growth of phytoplankton may be nitrogen-limited if the number of individuals has become so great that there is no more nitrogen available in the environment. Zooplankton, not limited by inorganic nitrogen themselves, can promote the growth of additional phytoplankton by consuming some individuals, digesting them, and returning part of the nitrogen to the environment.

Ecology is concerned with the processes involved in the interactions between organisms and their environments, with the mechanisms responsible for those processes, and with the origin, through evolution, of those mechanisms. It is distinguished from such closely related biological subdisciplines as physiology and morphology because it is not intrinsically concerned with the operation of a physiological process or the function of a structure, but with how a process or structure interacts with the environment to influence survival, growth, development, and reproduction.

Scope. There are a wide variety of approaches to ecology because of its broad, comprehensive character. Ecological studies can be characterized by the type of organisms studied, the habitat where studies take place, the level of organization that is of interest, and the methodology used. These are nonexclusive categories, and ecologists combine them in various ways while doing research. Major subdivisions by organism include plant ecology, animal ecology, and microbial ecology. Subdivisions by habitat include terrestrial ecology, the study of organisms on land; limnology, the study of freshwater organisms and habitats; and oceanography, the study of marine organisms and habitats.

The levels of organization studied range from the individual organism to the whole complex of organisms in a large area. Autecology is the study of individuals, population ecology is the study of groups of individuals of a single species or a limited number

of species, synecology is the study of communities of several populations, and ecosystem, or simply systems, ecology is the study of communities of organisms and their environments in a specific time and place.

Higher levels of organization include biomes and the biosphere. Biomes are collections of ecosystems with similar organisms and environments and, therefore, similar ecological properties. All of Earth's coniferous forests are elements in the coniferous forest biome. Although united by similar dynamic relationships and structural properties, the biome itself is more abstract than a specific ecosystem. The biosphere is the most inclusive category possible, including all regions of Earth inhabited by living things. It extends from the lower reaches of the atmosphere to the depths of the oceans. See BIOME; BIOSPHERE.

The principal methodological approaches to ecology are descriptive, experimental, and theoretical. Much of ecology through the first half of the twentieth century was descriptive, concentrating on describing the variety of populations, communities, and habitats throughout Earth. Experimental ecology, which involves manipulating organisms or their environments to discover the underlying mechanisms governing distribution and abundance, has become of increasing importance since the mid-1960s. Theoretical ecology uses mathematical equations based on assumptions about the properties of organisms and environments to make predictions about patterns of distribution and abundance. It also has become of increasing importance. All of these approaches, however, are evident in the origins of ecology. The science was more descriptive in its early phases and has become more experimental and theoretical in recent decades.

Ecosystem. An ecosystem is the organisms and physical factors in a specific location that are interrelated through the flow of energy and chemicals to form a characteristic trophic structure (Fig. 1). The concept of the ecosystem is fundamental to ecology and may be applied at various levels of organization, although it commonly encompasses several different species. The trophic structure of an ecosystem characterizes organisms according to their feeding level and how those feeding relationships of species result in specific patterns of energy flow and chemical cycling. The living mass of a given population or trophic level at any given time is called biomass. A change in mass with time is referred to as net productivity. See BIOLOGICAL PRODUCTIVITY; BIOMASS.

Primary producers, largely green photosynthetic plants, utilize the energy of the Sun and inorganic molecules from the environment to synthesize organic molecules. Those organic molecules serve as food for higher trophic levels. Primary consumers, or herbivores, feed on the producers. Secondary consumers, or carnivores, feed on other consumers, while omnivores feed at several trophic levels. Decomposers feed on the dead tissues of other organisms.

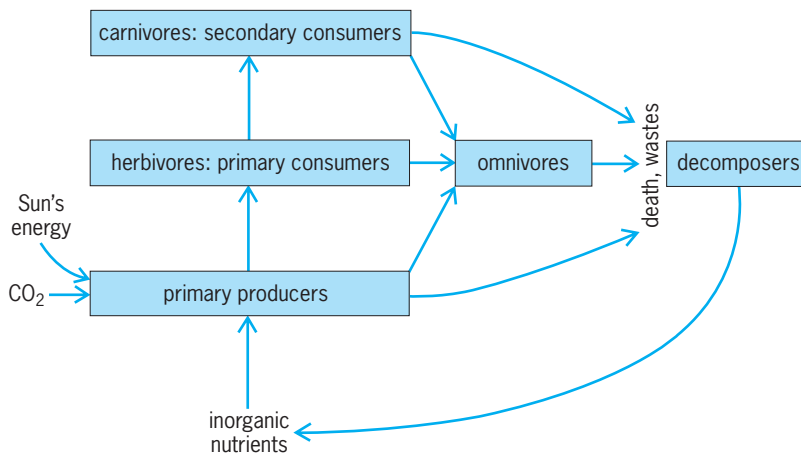


Fig. 1. Generalized trophic of ecosystems. Arrows indicate flow of energy and chemicals.

Each single pathway of energy and chemical flow is referred to as a food chain; the entire collection of food chains in an ecosystem is a food web, or trophic web. Since energy is expended to do work and chemicals are released back into the environment at each step of the trophic web, energy flow in ecosystems is unidirectional, constantly diminishing up the trophic web, while chemicals are recycled to the environment where they can be reutilized in organic syntheses. *See* FOOD WEB.

The ultimate energy source of all ecosystems is the radiant energy of photons from the Sun. Two major types of ecosystems can be distinguished, however, by their proximate energy sources. Autotrophic ecosystems have primary producers as a principal component and sunlight as the major initial energy source. Heterotrophic ecosystems depend upon preformed organic matter that is imported from autotrophic ecosystems elsewhere. The organisms in a stream, for example, may be highly dependent upon organic matter which has eroded from surrounding terrestrial ecosystems.

The ecosystem concept is flexible in application, but the way that it is applied must be clearly specified. An ecosystem study might concentrate on all the organisms and their environments in a specified forest. Alternatively, it might confine itself to the forest floor and soil, concentrating on trophic relationships within that more restricted habitat. *See* ECOSYSTEM.

Autecology

Autecology is the study of particular organisms, typically directed toward determining the traits they possess, and the traits they lack, that allow them to occupy, or be excluded from, certain environments. Early in the twentieth century, autecology placed considerable emphasis upon the morphology and anatomy of organisms, but the field has become diversified in the last few decades into a variety of different approaches. Autecology now also includes physiological ecology, which emphasizes the role of physiological processes, and chemical ecology, which emphasizes the role of biochemical traits. *See* PHYSIOLOGICAL ECOLOGY (PLANT).

One of the earliest generalizations of ecology was the principle of limiting factors, which states that organisms are limited by the factor or combination of factors that are farthest from the requirements of the organism. It was first stated by Liebig based on his studies of soil fertility. He recognized that there usually was a single chemical that would promote yield when it was added to the soil. When that initial limiting factor had been enriched, he found that addition of another nutrient might promote yield even more. The British physiologist Blackman restated the principle in its modern form based on his studies of photosynthesis revealing that the process sometimes could be limited simultaneously by more than one factor.

Though autecology is concerned with individual organisms, that information is often used to contrast the traits of organisms that occupy different environ-

ments. Thus there is a strong comparative character to autecological studies; two major habitat comparisons, discussed below, are marine vs. freshwater, and terrestrial vs. aquatic.

Marine vs. freshwater habitats. Life first evolved in the sea; the colonization of freshwater habitats that followed was accompanied by the evolution of traits associated with exploiting the more dilute solutions of freshwater habitats. Seawater is over 3.6% dissolved solids, largely sodium chloride, while freshwater contains 0.01–0.05% dissolved solids. Colonization of freshwater was coupled with the evolution of relatively impervious body surfaces, mechanisms allowing organisms to accumulate ions against a concentration gradient, and the production of large quantities of very dilute urine or other excretory waste. For example, many marine invertebrates are isotonic (equal in osmotic pressure) to seawater so they have no tendency to either lose or gain water. Elasmobranchs have high concentrations of urea in the blood, and although the concentrations of ions is lower than in seawater, the concentrations of urea plus ions is isotonic. Marine teleosts have blood osmotic concentrations well below that of seawater, and they therefore tend to lose water to the sea. They counteract this loss by drinking large quantities of seawater, absorbing it through the gut, and secreting the salt from special cells in the gills. Most freshwater animals are hypotonic (lower osmotic concentration) to the solution that bathes them so they tend to gain water. *See* OSMOREGULATORY MECHANISMS.

Terrestrial vs. aquatic habitats. The colonization of terrestrial habitats was accompanied by the evolution of a large number of traits that allowed organisms to occupy habitats where water was a comparatively rare component of the environment and the atmosphere and soils became principal media surrounding organisms.

Since the body was no longer bathed in fluid, plants and animals both evolved traits allowing them to conserve water and obtain it when it was scarce. Plants developed extensive absorptive surfaces which allowed them to grow into moist regions of the soil and obtain both water and nutrients from the soil solution. Animals often drank water directly and developed outer tissue layers that were extremely impervious to water. The evolution of vascular systems in both animals and plants allowed them to transport water from limited sites of acquisition to the rest of the body. Since the atmosphere is a much less dense and buoyant medium than water, terrestrial plants and animals also evolved traits providing rigid support, such as the skeletons of vertebrates and the strengthened woody stems of vascular plants.

Terrestrial habitats are also characterized by much greater temperature fluctuations than occur in aquatic habitats. Animals therefore evolved physiological traits that allowed them to maintain their body temperatures constant, as well as behavioral traits that involved the ability to move into habitats where favorable temperatures occurred. Physiologically,

evolution culminated in homeothermy in birds and mammals; these organisms expend energy to maintain body temperature within narrow limits. Other animals, called heterotherms, which are incapable of maintaining their body temperature over a wide range of environmental temperatures, often move into localized environmental areas providing appropriate temperatures, and spend unfavorable seasons in quiescent metabolic states. *See* THERMOREGULATION.

Finally, the colonization of terrestrial habitats was associated with many changes in reproductive methods due to the inability of organisms to effectively disperse gametes through the surrounding medium. Such adaptations include the behavioral mating systems of birds and mammals that allow mates to find one another, and the elaborate pollination mechanisms of flowering plants that may involve specific transfer of pollen between individuals of the same species by animal vectors. *See* POLLINATION; REPRODUCTIVE BEHAVIOR.

Physiological and chemical ecology. The postwar development of analytical instrumentation that was much more sensitive to many chemicals than previous analytical techniques led directly to the development of physiological and chemical ecology, which applies that instrumentation to studying organisms in their natural environments. The development of infrared gas analyzers that detect carbon dioxide and water vapor at concentrations that occur in nature led to many studies of photosynthesis, respiration, and transpiration by plants. Those studies have revealed two different types of photosynthesis that occur in plants occupying different types of environments. C_3 plants, so named because the initial stable product of carbon fixation is a three-carbon molecule, are characteristic of almost all aquatic habitats and of terrestrial environments with lower temperatures, less arid climates, and lower solar radiation intensities. C_4 plants, so named because the initial stable product of photosynthesis is a four-carbon molecule, are more common in arid, hotter environments with more intense solar radiation. Those differences are due to the fact that C_3 photosynthesis is intrinsically more efficient but is also more temperature-sensitive and generally more likely to be light-saturated well below the intensity of full sunlight. Thus, C_4 plants are common in tropical and subtropical, arid terrestrial ecosystems, while C_3 plants are more common in most other ecosystems. *See* PHOTOSYNTHESIS.

Chemical ecology studies have been largely devoted to studying rare molecules that have major biological effects. Many animals, for example, produce pheromones that are extremely dilute in the environment but can attract mates, repel predators, or kill pathogens. Many plants also produce organic molecules that are present in their tissues in very small concentrations but act to deter herbivores and pathogens. The evolution of highly toxic molecules by plants sometimes has led to the evolution of specific mechanisms of avoidance or detoxification by

a limited number of animals. That type of coupled evolution in which populations change genetically in response to the properties of each other is called coevolution. For example, some animals, such as the monarch butterfly whose larvae feed on milkweeds, have even appropriated aspects of their food plant's chemistry to protect them from predators. Milkweeds produce cardiac glycosides, chemicals that poison the heart muscles of vertebrates. The larvae of the monarch sequester those chemicals in their body, which protects them from such potential predators as birds and small mammals. *See* CHEMICAL ECOLOGY.

Population Ecology

Population ecology is the study of the vital statistics of populations, and the interactions within and between populations that influence survival and reproduction. Population ecologists are concerned with the balances between births and deaths that determine the rate of change of population size through time. When births exceed deaths, of course, population size increases, and when deaths exceed births, population size declines. Vital statistics are often expressed on a per capita basis or, in human populations, on a per-10,000-persons basis. Density, the number of organisms per unit area, is a fundamental property of populations.

One of the principal concerns of population ecologists is identifying and determining the importance of factors that control population densities. Since A. van Leeuwenhoek first calculated the massive potential of flies to reproduce—that is, one pair of flies could produce 746,496 flies after only 3 months of breeding—population ecology has been concerned with the factors that limit the complete realization of this reproductive potential. There are two fundamentally different types of factors limiting population density. Density-independent factors are environmental factors that reduce reproduction or increase death rate independently of the number of organisms in the population. Weather is believed to be a major density-independent factor that often affects populations in a catastrophic way to reduce density. Severely cold winters, deep snowfalls, exceptionally wet growing seasons, and extreme drought occur with a frequency and intensity that are independent of population size. Although these factors can determine population density, they cannot regulate that density since regulation involves maintaining density within certain limits. Density-dependent factors influence the survival or reproduction of individuals in a way that is proportional to density. They can, therefore, regulate population density. Food supply, predators, disease, and such behavioral interactions as territoriality can limit survival and reproduction with increasing intensity as population density increases.

Two vital statistics of populations of principal interest to population ecology are survivorship and reproductive schedules (**Fig. 2**): It is the balance between the probability that an individual will live to reach a certain age, and the average number of

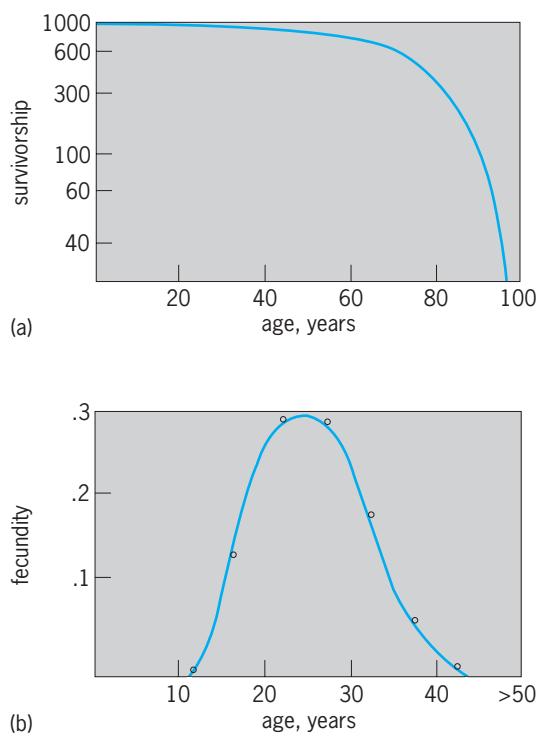


Fig. 2. Demographic statistics of females in the United States. (a) Survivorship, where the vertical axis is the number surviving, of 1000 born alive (on a logarithmic scale). (b) Reproductive schedule, where the vertical axis is the number of female offspring per female of given age, at 5-year intervals.

offspring produced by individuals of that age that determine population growth rate.

Population growth. One of the earliest quantitative expressions in ecology was a description of density-dependent population growth called the logistic equation. It was first described by P. F. Verhulst in 1838, but it was completely forgotten until its rediscovery by R. Pearl and L. Reed in 1920. Population growth rate is a consequence of population size (N) and per capita birth (b) and death (d) rates, as shown by Eq. (1), where dN/dt is rate of change

$$\frac{dN}{dt} = (b - d)N \quad (1)$$

of population size. If the symbol r is used to represent the realized per capita reproductive rate, that is, $b - d$, then Eq. (1) can be written in the form of Eq. (2). If r is a constant, a population described

$$\frac{dN}{dt} = rN \quad (2)$$

by Eq. (2) would grow exponentially with time. The logistic equation recognizes two factors lacking in Eq. (2) that can influence population growth: the maximum per capita reproductive potential of the population, r_m , and the carrying capacity of the environment, K . Maximum reproductive potential would only be realized in the most favorable environment when population density was low enough that individuals did not make simultaneous demands on the

same resources. Carrying capacity is a measure of both the amount of resources in the environment and the efficiency with which organisms use those resources. A population is at its carrying capacity when no additional individuals can be supported by an environment. Incorporating those constants into the growth equation above gives Eq. (3), the logistic

$$\frac{dN}{dt} = r_m N \left(1 - \frac{N}{K}\right) \quad (3)$$

equation. Maximum per capita reproductive potential is determined by the genetic ability to reproduce in an unlimited environment. A population growing according to this equation has a sigmoid, or S-shaped, pattern of population size through time (Fig. 3). The maximum population growth rate, which, again, is the product of population size and per capita reproductive rate, is reached when the population is halfway to K .

Dividing Eq. (3) through by N and rearranging gives Eq. (4). Plotting r against N gives a straight line

$$\frac{dN}{Ndt} = r_m - \left(\frac{r_m}{K}\right)N \quad (4)$$

with an intercept on the y axis that is equal to r_m and an intercept on the x axis that is equal to K .

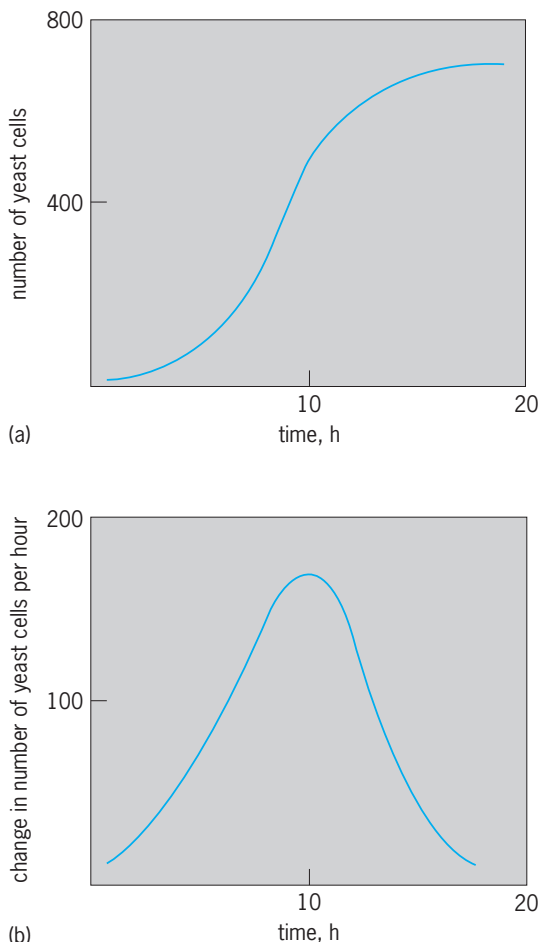


Fig. 3. A yeast population growing according to the logistic equation. (a) Population size. (b) Growth rate.

An equation as simple as the logistic equation provides merely a yardstick against which the growth of real populations can be compared since it contains a number of assumptions that will not be met in nature. Principal among those assumptions in the logistic equation are that carrying capacity is a constant; that each individual in the population is equivalent enough to all other individuals to be represented in the average population statistic, r_m ; that all individuals have an equal probability of mating; and that there are no time lags so that birth and death rates are adjusted immediately for each increment of N as population size increases to K . Much of the study of the ecology of real populations involves understanding how those assumptions are violated.

Competition is inherent in the logistic equation, and in population growth in nature. Competition occurs when the resources in the environment are in short supply relative to the demands that individuals make on those resources. In the logistic equation, that competition is expressed in the decline of r as N increases. This type of competition between members of the same species population is called intraspecific competition.

Population interactions. Competition also is a principal type of interaction between the members of different species, when it is termed interspecific competition. Another important type of interaction between populations is predation, when the individuals of one population feed on the individuals of another population. Both of these interactions were described mathematically in the 1920s by two ecologists working independently, A. Lotka and V. Volterra. The Lotka-Volterra competition equations are Eq. (5) and the exact equivalent, Eq. (6), where 1 and 2 refer

$$\frac{dN_1}{dt} = r_{m1}N_1 \frac{K_1 - N_1 - \alpha_{12}N_2}{K_1} \quad (5)$$

$$\frac{dN_2}{dt} = r_{m2}N_2 \frac{K_2 - N_2 - \alpha_{21}N_1}{K_2} \quad (6)$$

to the different species, with α_{12} being the competitive effect of an individual of species 2 on species 1, and α_{21} being the competitive effect of an individual of species 1 on species 2. These Lotka-Volterra

competition equations recognize that each species depresses the growth rate and carrying capacity of the other. Depending upon the respective competition coefficients, carrying capacities, and maximum per capita reproductive rates, there are four potential outcomes of competition: (1) species 1 may go extinct; (2) species 2 may go extinct; (3) the species may have an unstable equilibrium in which they coexist for a while until chance events allow one species to increase a bit above the equilibrium and the other species then goes extinct; or (4) there may be a stable equilibrium in which both species coexist at some population size below the carrying capacity for either species alone.

The Lotka-Volterra predation equations describe population ecology when members of a predator population, P , feed on members of a prey population, H . The prey equation is Eq. (7), where k_1 is

$$\frac{dH}{dt} = r_mH \left(1 - \frac{H}{K} \right) - k_1PH \quad (7)$$

a constant that is the predation rate per capita of predator. The equation for the predator is Eq. (8),

$$\frac{dP}{dt} = k_2PH - k_3P \quad (8)$$

where k_2 is the efficiency with which predators convert their food into more offspring and k_3 is the per capita predator death rate. Unlike the logistic and Lotka-Volterra competition equations, which predict that populations will reach some final constant size, the Lotka-Volterra predator-prey equations predict that predator and prey populations will tend to oscillate through time. At low densities of both predator and prey, the values of H/K and k_1PH in the prey population equation will both approach zero, and the prey will grow according to the unmodified logistic, approaching the maximum per capita reproductive rate. As the prey population increases, the predator population also can increase, and there is a tendency for predator density to overshoot prey density, leading to oscillations through time. Such oscillations are, in fact, often observed in nature, with a classic example being Canadian populations of snowshoe hare, the prey, and predators, such as lynx (Fig. 4). Whether such oscillations in real populations are due to the instability of predator-prey relations inherent in the Lotka-Volterra predation equations, or to the fluctuations of other factors such as weather or food supply available to the hares, has not yet been proven. See MATHEMATICAL ECOLOGY; POPULATION ECOLOGY.

Parasitism and disease are additional types of population interactions involving detrimental effects on one participant and beneficial effects on the other participant. A parasite is an organism that grows, feeds on, and is sheltered by another organism, the host. Disease is an abnormal condition of an organism, often due to infection, that impairs normal physiological activity. Parasitism and disease often are described mathematically by equations similar to the Lotka-Volterra predation equations, although they

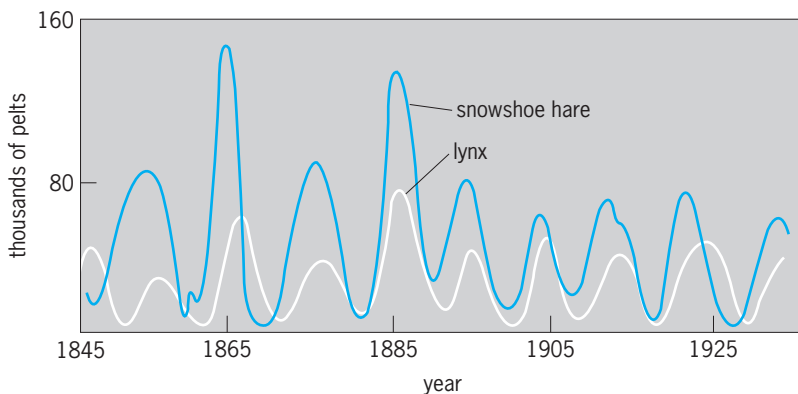


Fig. 4. Oscillations of predator and prey populations. Data are from pelt records of lynx and snowshoe hare of the Hudson's Bay Company, in Canada.

are commonly modified to account for the dispersal of parasites or pathogens between different hosts. These population interactions have been of growing interest in ecology, although they have been the subject of the entire field of epidemiology, an important aspect of medicine. The importance of parasitism and disease in natural ecosystems is evident from the extinctions of two important tree species, the American chestnut (*Castanea dentata*) and American elm (*Ulmus americana*), from the forests of eastern North America during this century. Those extinctions were caused, respectively, by chestnut blight and Dutch elm disease, both due to pathogens introduced into North America from Europe. *See* EPIDEMIOLOGY.

Symbiosis is a type of interaction between populations in which a participant is promoted by the interaction. There are two types of symbiosis: commensalism, which occurs when one participant is promoted by the interaction and the other is unaffected, and mutualism, which occurs when both participants are promoted by the interaction. Mutualism often involves organisms that are very different and have coevolved such a strong mutual interdependence that neither can exist in natural environments in the absence of the other. Symbiosis commonly is described mathematically by variations of the Lotka-Volterra competition equations in which the sign preceding the α becomes positive for one (commensalism) or both (symbiosis) populations. This transforms the competition coefficients into mutualism coefficients.

Mutualism is an extremely important type of ecological interaction since it often allows the symbionts to perform functions that have proven impossible for any single organism to accomplish with the same efficiency. Lichens are symbiotic associations between fungi and algae. They are capable of colonizing some of the most inhospitable environments on Earth, such as the bare faces of boulders and extremes of drought, cold, and heat. The fungal and algal species are so closely integrated that lichens have been classified by the binomial systems that are applied to other true species.

An example of the functional importance of symbiosis to ecology is the digestion of cellulose, which is accomplished largely by symbiotic associations. Cellulose is the most abundant organic molecule on Earth, but this vast energy source is largely inaccessible to living organisms except to symbionts. Among the most important symbiotic associations capable of digesting cellulose are ruminant animals and their gut flora, and termites and their gut protozoa. In the absence of these symbiotic associations, Earth would be covered with vast residues of undecomposed cellulose.

Synecology

Synecology is the study of groups of populations co-occurring in time and space; it is also called community ecology. Synecology is concerned with factors controlling the species composition of communities and why species composition changes in different

environments and with time in the same habitat. Species composition of a community may include density, biomass, or productivity of individuals of different species.

The study of communities is one of the oldest branches of ecology, arising strongly in the studies of naturalists in the sixteenth to nineteenth century, and one of the most descriptive. Only since the mid-1960s have ecologists begun to turn significantly to experimental studies that manipulate either the environment or the species composition of communities. Those experiments are designed to understand mechanisms responsible for the species compositions of communities and the interactions that occur between different species in the community. *See* ECOLOGICAL COMMUNITIES.

Community organization. A principal concern of community ecology is the life forms of organisms, the taxonomic composition of co-occurring species groups, and the spatial arrangement of species. Plant life forms are relatively easy to characterize, and communities in widely separated geographic areas but similar environments commonly have a similar life-form composition. Most of the species in tropical rainforests, for example, are trees and shrubs, while most of the species in grasslands are small shrubs, and perennials with buds belowground. Although the species compositions of communities in different geographic areas may be quite different, the life forms are often quite similar in similar environments.

Another aspect of community organization that is a concern of synecology is spatial patterns, or dispersion. Ecologists typically distinguish between dispersal, the process of movement, and dispersion, the arrangement of individuals in space. The three types of dispersion patterns are random, regular, and clumped (aggregated). In a random distribution, the probability of an individual occupying a given location in a community is not influenced by the presence of another individual at a nearby location, which implies that no important factor, such as environment or parental proximity, has influenced the dispersion of individuals. In a regular distribution, the probability of an individual occupying a location is reduced if another individual occupies a nearby location; this distribution may be due to competition for an evenly spread environmental resource. In a clumped distribution, which may be due to proximity to seed source, or environmental heterogeneity, the probability of an individual occupying a location increases if another individual occupies a nearby location. *See* POPULATION DISPERSAL; POPULATION DISPERSION.

Temporal distributions are also a concern of synecology. By substituting the phrase "at a given time" for the phrase "at a given location" in the above definitions, temporal distributions are judged by the same criteria. Spatial and temporal patterns of individuals are important aspects of community ecology because they can provide insight into the underlying distributions of limiting factors in the environments and how species are related to those factors. Species

with distributions in which individuals are clumped in space or time often utilize resources that also are clumped in space or time. Limited dispersal also can lead to a clumped distribution due to the clustering of offspring around their parents. Regular distributions can often result from competition between individuals for a resource that is evenly distributed in space or time. Desert plants, for example, frequently have a regular distribution pattern because they compete for soil moisture, and each individual must utilize the moisture available in a certain minimum soil volume if it is to survive.

Niche. The niche concept refers to the environmental factors that control a species, and its distribution in relation to those factors. A species' niche is a consequence of the genetic properties of the individuals of that species, the environmental factors that affect it, and interactions with other species. G. E. Hutchinson described the niche as an n -dimensional hypervolume occupied by a population, where n is the number of environmental factors affecting it. He also distinguished between a fundamental niche, which is the genetic potential of the species, and a realized niche, which is the range of environmental factors occupied in the presence of other organisms.

Species are often distributed along environmental gradients as a series of bell-shaped curves (Fig. 5). Those curves represent the realized niches of species. Removal of other species, a type of experiment that is becoming increasingly common in synecology, would be required to reveal a species' fundamental niche. A species that is narrowly distributed has a narrow niche and is often referred to as an ecological specialist, while a widely distributed species has a broad niche and is often referred to as an ecological generalist. It is often inferred that intense interspecific competition will lead to specialization, restricting a species' realized niche to only that portion of the fundamental niche where it is a strong competitor. Conversely, intense intraspecific competition will tend to lead to a broader niche that reduces the effects of different members of the same species upon each other. However, there are as yet insufficient experimental studies to strongly support either of these inferences.

Species diversity. Species diversity is a property of a community that encompasses both the number of

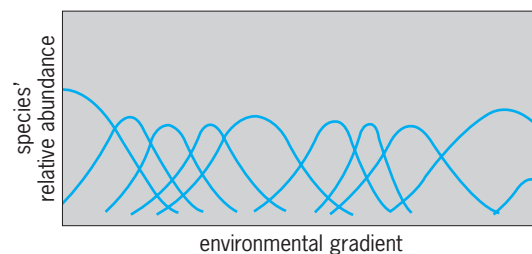


Fig. 5. A continuum of species distributions along an environmental gradient, where each line represents a different species. Note that peaks of species abundance do not overlap, and that breadth of curves differs between species.

species present (species richness) and the equality of relative abundances of the species present (equitability of evenness). One of the principal concerns of synecology since its origin has been why communities differ so dramatically in the number of species present and the similarities of their relative abundances. Naturalists attached to exploratory expeditions during the historical colonial period were particularly overwhelmed by the variety of plant and animal life in the wet tropics. Tropical forests seemed to consist of an almost countless number of plant and animal species compared with communities in the temperate climates with which the naturalists were familiar. Indeed, there is a general gradient of increasing diversity from the poles to the Equator.

Various explanations for the latitudinal differences in biotic diversity have been proposed, but none is wholly satisfactory and none has gained prominence. Chief among those explanations are: (1) Nontropical habitats have been subject to severe climatic fluctuations, including glaciation, that lead to repeated species extinctions. (2) The moderate climates of tropical locations allow species to become more specialized for narrow ranges of climate and other environmental factors. (3) Biological interactions, such as competition and predation, are stronger in tropical locations, leading to a proliferation of different adaptive types. (4) Productivity is higher in tropical locations, allowing more species to coexist due to greater available food. (5) Spatial variation of environments is greater in the tropics, leading to a proliferation of specialists in different areas. (6) There is greater environmental complexity in the tropics, partly because of more species, leading to greater specialization among species in a sort of positive feedback. See BIOGEOGRAPHY.

Community: reality vs. abstraction. One of the most contentious issues throughout the history of ecology has been the degree to which the species that co-occur in space and time represent an organized, meaningful biological unit and to what extent they are merely a fortuitous collection of individuals. Two of the principal ecologists in American ecology in the early decades of the twentieth century argued these points. E. Clements referred to communities as being like an organism, while H. A. Gleason said that communities were merely a happenstance collection of species with similar ecological tolerances. Clements's view prevailed until the middle of the twentieth century, and then Gleason's view gained ascendancy, largely because of quantitative studies of community composition revealing that species distributions along environmental gradients were commonly the type of bell-shaped curves in Fig. 5. Those bell-shaped curves suggested that each species population had its own, individualistic response to the environment, and that the species that occurred at a given point on the gradient did so fortuitously.

Still, there is no gainsaying that similar communities occur in similar habitats. Although the species compositions of communities may vary considerably

in different geographic areas, the occurrence of communities with very similar life forms in geographically separated locations with similar environments suggests that a community is not merely a random collection of species. Many experiments in marine intertidal communities indicate that the removal of one species has a cascading effect that alters the abundances of many other species, even leading to the extinction of some species and invasion of the community by other species that were previously absent. These results indicate that the species in a community are interconnected through the functional processes of competition, predation, and symbiosis in ways producing a community that is not just a random collection of species.

Succession. Succession is the process of change in community composition and environmental properties over time on a site. Primary succession occurs in newly formed habitats, such as glacial moraines, river levees, or volcanic ash. Secondary succession occurs in habitats where the previous community has been destroyed or severely disturbed, such as following forest fire, abandonment of agricultural fields, or epidemic disease or pest attack.

A general consequence of succession is an amelioration of physical factors and a reduction of their importance as controlling factors, and an increase in the complexity and importance of biological factors. In newly formed ash fields or newly abandoned fields, sunlight is intense, nutrients are often in poor supply, winds are strong, and evaporative stress is severe. As such sites are colonized by plants, all of these physical factors tend to be modified and trophic webs develop that depend upon the plants. Species diversity increases as more species colonize, and increasingly complex trophic patterns develop. There often is an increase in the stature of the vegetation, which reduces sunlight, wind, and evaporative stress below the canopy. That creates new conditions that allow species to invade that were incapable of withstanding the physically harsh conditions early in succession.

Early in succession, there is directional change in both environmental factors and species composition of the community. Ultimately a relatively constant environment is reached and species composition no longer changes in a directional fashion but, instead, fluctuates about some mean, or average, community composition. That stage is called the climax.

As in the nature of the community itself, there is disagreement among ecologists about the exact nature of succession, and of the climax. Nevertheless, there is a general tendency for communities to become more complex and for biological factors to increase in importance and the effects of physical factors to be reduced as time passes. *See* ECOLOGICAL SUCCESSION.

Systems Ecology

Systems ecology is the study of dynamic relationships between the units of the ecosystem, particularly those relationships that influence rates of energy flow and chemical cycling. The ecosystem

concept, as pointed out above, can be applied to many different levels of organization and varying degrees of complexity. Ecosystem science, or systems ecology, can be broadly divided into approaches concerned with relationships between trophic levels and those that employ with specific details of population dynamics of community organization within trophic levels. In both approaches the goal is to develop mechanistic explanations for the observed behavior of the system as a whole. Often this goal is accomplished through simulation modeling.

Simulation modeling involves the use of mathematical equations to characterize the state and dynamics of a system. Mathematical approaches to population dynamics, as seen earlier in this article, commonly use simplified equations, such as the logistic, to model the dynamics of one or a few populations. Simulation modeling differs principally in using a set of coupled equations to characterize the processes occurring in ecosystems. Those models often have much more detail than the simplified models of population dynamics, and require numerical solution using digital computers. A simulation model typically begins with a conceptual model or diagram consisting of boxes that represent components of system state and connecting arrows representing processes that interconnect the components. Simulation modeling uses rate equations to describe the processes and evaluate the mechanisms that are important to understanding both the dynamics and functional properties of ecosystems. Simulation models have as goals both isolating the general features of a broad range of similar ecosystems and providing sufficient detail to accurately predict how a specified system will respond to perturbations.

A system can be defined as a group of interacting and interdependent elements that form a collective entity. An ecosystem, therefore, is a collection of environmental factors and organisms that interact through specific, dynamic relationships defined by the food web of which they are members. Systems may be open or closed. A closed system is completely self-contained and does not interact with any other system; an example would be molecules of a gas enclosed in a container that was a perfect insulator. Those molecules would function independently of all other things. An open system interacts with other systems. All ecosystems are open systems. In particular, they are open to energy, whose ultimate source is the Sun. Proximally, however, energy can be imported into an ecosystem from some other system where the energy of photons from the Sun is converted into organic energy by photosynthesis.

Energy flow. Energy flow in ecosystems is approached in a way closely related to the laws of thermodynamics, that branch of physics that deals with the relationships among different forms of energy. The first law of thermodynamics states that the energy input to a system either is stored or is used to do work. It is generally given as Eq. (9), where Q

$$Q = \Delta E + W \quad (9)$$

is energy input, ΔE is a change in the energy content of the recipient system, and W is work done by that system. Ecologists restate the law as Eq. (10),

$$P_g = P_n + R \quad (10)$$

where P_g is gross productivity, P_n is net productivity, and R is respiration. The change in energy content, P_n , is therefore equal to the energy input minus the work done. That work has two components: one is maintenance respiration, the cost of maintaining the structure of the system at the time of energy input, and the other is growth respiration, the cost of synthesizing new biomass.

The second law of thermodynamics also is important to energy balance in ecosystems. That law states that the entropy of isolated systems always tends to increase. Entropy (S) is a measure of the randomness, disorder, or lack of organization of a system. In energy terms it is measured as the heat capacity of the system (joules \cdot degree⁻¹ \cdot mole⁻¹) at a particular temperature (T). In isolated systems, entropy always increases. Because ecosystems are open systems, they are able to maintain their organization and avoid decaying to a state of maximum entropy and zero free energy. The maintenance component of respiration is the cost of maintaining present structure, that is, replacing the free energy degraded to heat. The growth component of respiration is the cost, that is, the free energy intake necessary to build structure and organization. Those costs are expressed both as an increase in the entropy (thus an increase in heat content) of the environment to which organisms in an ecosystem are coupled as well as the loss of heat produced by "friction" representing the degree of thermodynamic inefficiency of the metabolic energy transformations. See OPEN-SYSTEMS THERMODYNAMICS (BIOLOGY); THERMODYNAMIC PRINCIPLES.

A fundamental generalization of systems ecology is that net productivity diminishes about one order of magnitude at each level in the trophic web. If primary P_n averaged 100 g/m² \cdot yr, secondary P_n at the herbivore level would be about 10 g/m² \cdot yr, and at the carnivore level would be about 1 g/m² \cdot yr. This so-called energy pyramid has two causes. First, the above-mentioned maintenance costs of respiration at each trophic level diminish the energy available to higher levels. Second, the ecological harvest of energy by any trophic level is not fully efficient. Therefore, the energy flow must diminish at each successive trophic level in an ecosystem.

In most ecosystems, certainly, the biomass of primary producers is greater than that of the herbivores which, in turn, is greater than carnivore biomass. Sometimes, however, particularly in aquatic habitats, biomass increases up the trophic web. That is possible only when lower trophic levels consist of organisms with very high rates of net productivity per unit of biomass. The ratio of energy flow to biomass is referred to as turnover time. If turnover time is short, as it is for small organisms with rapid life cycles, it is possible for the biomass pyramid to be inverted; that is probably most common in the open ocean where

large, long-lived organisms are at upper trophic levels and small, short-lived organisms are at lower trophic levels. See ECOLOGICAL ENERGETICS.

Chemical cycling. Another fundamental concern of systems ecology is the pattern of chemical flow that is coupled to the process of energy flow. Each time that energy is used to do work, some chemicals are released back into the environment. Death also releases chemicals back into the environment.

There are two different types of trophic webs in ecosystems. The grazing food web is based on the consumption of the tissues of living organisms. The detritus food web is based on the consumption of dead organic material, called detritus in aquatic systems and sometimes referred to as litter in terrestrial ecosystems. Chemicals are recycled to the environment from the grazing food web each time that work is done, and as the excretory wastes of living organisms. Chemicals are recycled to the environment from the detritus food web as organisms utilize excretory wastes and the dead tissues of organisms.

Decomposition is the process of degrading the energy content of dead tissues and simultaneously releasing chemicals back into the environment. When those chemicals are released in inorganic forms, the process is called mineralization. Decomposition and mineralization typically involve many steps and take place in a trophic web fully as complex as those of the grazing food web. There is a size-dependent hierarchy as particles are broken into progressively smaller particles; that subdivision increases surface-volume ratios, exposing progressively more of the substance to decay. There also typically is a chemical hierarchy. Fats, proteins, and simple carbohydrates are the most easily utilized organic chemicals, and they generally are rapidly attacked. Wood, chitin, and bones are only slowly decomposed. More resistant substances accumulate in terrestrial ecosystems as soil humus and in aquatic ecosystems as sediments, or ooze, beneath the water. In geological time, those sediments were transformed into the hydrocarbon deposits that serve as sources of coal, oil, and gas for modern industrial society.

Microorganisms are exceptionally important in the processes of decomposition and mineralization. The final steps of mineralization are almost invariably accomplished by microorganisms. For example, the cycling of nitrogen in ecosystems involves microorganisms at several critical steps (Fig. 6). The only significant pathway introducing nitrogen into trophic webs is the process of nitrogen fixation in which gaseous, molecular nitrogen is converted into organic nitrogenous compounds by nitrogen-fixing bacteria and blue-green algae. Some of these bacteria and algae are free-living, and some participate in symbiotic associations with other organisms. The organic molecules are degraded in the detritus food web by ammonifying bacteria that release ammonium ion or ammonia. Nitrite bacteria convert the ammonia into nitrite, and nitrate bacteria oxidize the nitrogen further into nitrate; those conversions are referred to collectively as nitrification. Both

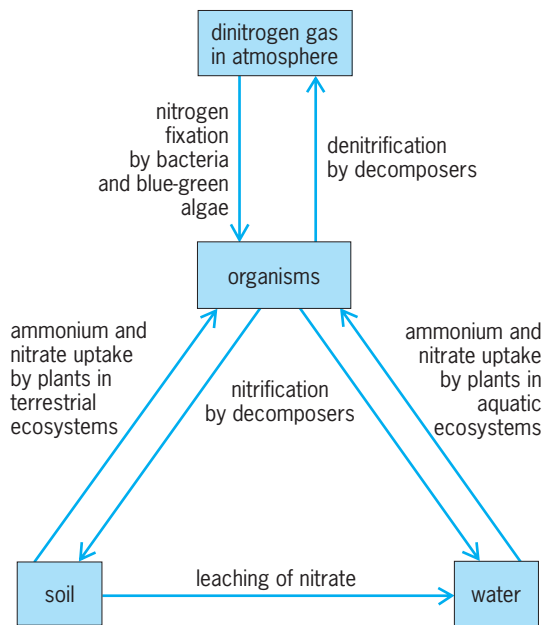


Fig. 6. Generalized nitrogen cycle. Arrows represent flow of nitrogen.

nitrate and ammonium ion can reenter the food web via uptake by plants. Denitrifying bacteria convert the nitrate into nitrogen gas, a process called denitrification that releases gaseous nitrogen back into the atmosphere. See BIOGEOCHEMISTRY; MATHEMATICAL ECOLOGY.

Applied Ecology

Applied ecology is that branch of ecology dealing with practical problems of immediate social importance. Ecology is often confused in the minds of laypersons with the environmental movement, a social interest group concerned with environmental degradation and problems of resource supply to human societies. Applied ecology makes a fundamental contribution to those social and political concerns by identifying environmental problems, gauging their significance, and suggesting potential solutions. Among the environmental issues to which ecology has made an important contribution have been problems of population growth and resource supply, acid rain, eutrophication, consequences of pollution, biological control of crop pests, range management, forestry, and the ecological consequences of nuclear war. From T. R. Malthus's essay on human population growth to recent scientific evaluations of nuclear war, ecology has always been concerned with problems that are important to the affairs of humans. See ECOLOGY, APPLIED. Samuel J. McNaughton

Bibliography. B. Glaeser (ed.), *The Green Revolution Revisited*, 1987; H. Howe and L. Westley, *Ecological Relationships of Plants and Animals*, 1988, reprint 1990; J. Keating, *Interdependence in the Natural World*, 1987; R. P. McIntosh, *The Background of Ecology: Concept and Theory*, 1985, reprint 1986; M. B. Rambler et al. (eds.), *Global Ecology: Towards a Science of the Biosphere*, 1988.

Ecology, applied

The application of ecological principles to the solution of human problems and the maintenance of a quality life. It is assumed that humans are an integral part of ecological systems and that they depend upon healthy, well-operating, and productive systems for their continued well-being. For these reasons, applied ecology is based on a knowledge of ecosystems and populations, and the principles and techniques of ecology are used to interpret and solve specific environmental problems and to plan new management systems in the biosphere. Although a variety of management fields, such as forestry, agriculture, wildlife management, environmental engineering, and environmental design, are concerned with specific parts of the environment, applied ecology is unique in taking a view of whole systems, and attempting to account for all inputs to and outputs from the systems—and all impacts. In the past, applied ecology has been considered as being synonymous with the above applied sciences. See SYSTEMS ECOLOGY.

Ecosystem ecology. Ecological systems, or ecosystems, are complexes of plants, animals, microorganisms, and humans, together with their environment. Environment includes all those factors, physical, chemical, biological, sociocultural, which affect the ecosystem. The complex of life and environment exists as an interacting system and is unique for each part of the Earth. The unique geological features, soils, climate, and availability of plants, animals, and microorganisms create a variety of different types of ecosystems, such as forests, fields, lakes, rivers, and oceans. Each ecological system may be composed of hundreds to thousands of biological species which interact with each other through the transfer of energy, chemical materials, and information. The interconnecting networks which characterize ecosystems are often called food webs (Fig. 1). It is obvious from this structural feature of interaction that a disturbance to one population within an ecosystem could potentially affect many other populations. From another point of view, ecosystems are composed of chemical elements, arranged in a variety of organic complexes. There is a continual process of loss and uptake of chemicals to and from the environment as populations are born, grow and die, and are decomposed. Ecosystems operate on energy derived from photosynthesis (called primary production) and from other energy exchanges. The functional attributes of ecosystems, such as productivity, energy flow, and cycling of chemical elements, depend upon the biological species in the ecosystem and the limiting conditions of the environment. See BIOLOGICAL PRODUCTIVITY; ECOLOGICAL ENERGETICS; FOOD WEB.

Ecological systems develop in accord with the regional environment. Although these systems have evolved to resist the normal expected perturbations encountered in the environment, unusual disturbances and catastrophic events can upset and even destroy the system. In this case, recovery can occur

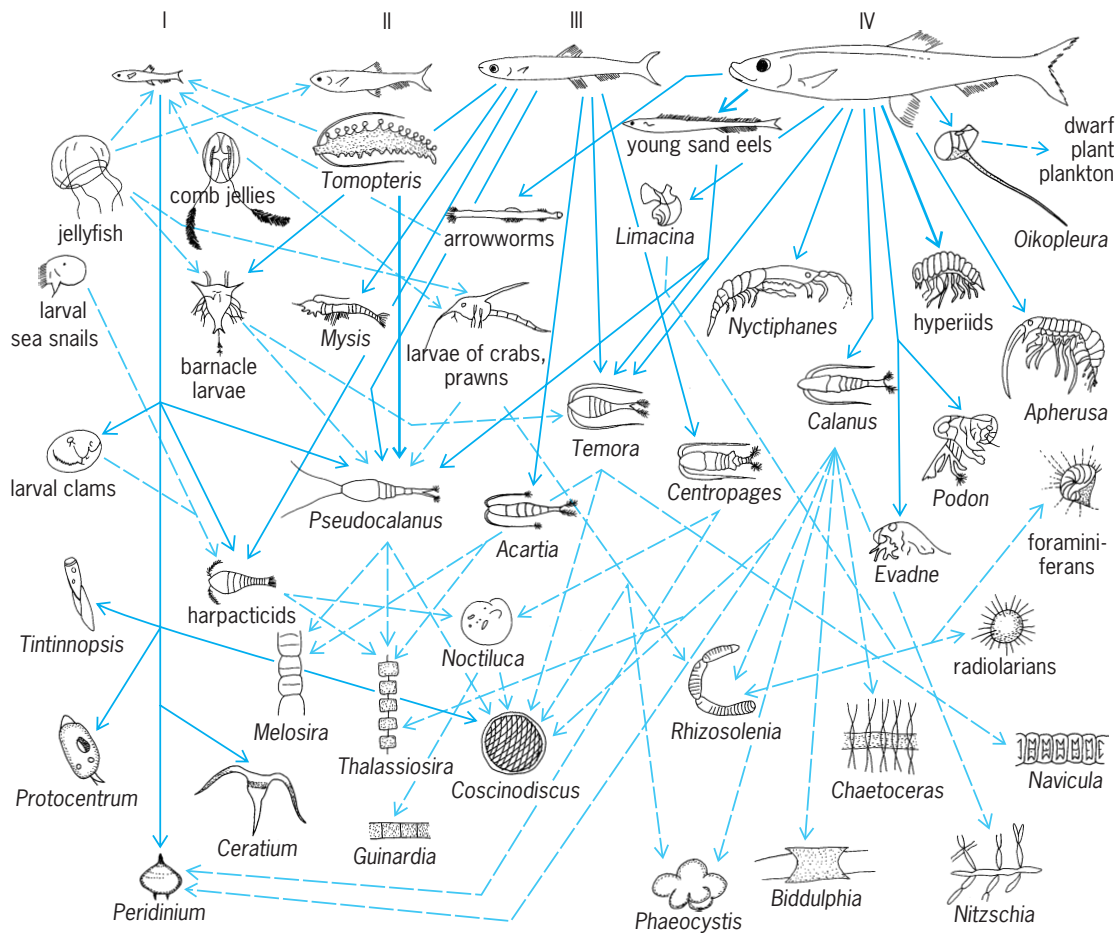


Fig. 1. Food relations of the herring at different stages in its life. Sizes of herring are (I) 0.6 to 1.3 cm, (II) 1.3 to 4.5 cm, (III) 4.5 to 12.5 cm, and (IV) over 12.5 cm. Solid lines indicate food eaten directly by herring. 1 cm = 0.4 in.

after the disturbance stops. Recovery is termed ecological succession since it comprises a sequence of communities which succeed each other until a dynamic steady state is reestablished. See ECOLOGICAL SUCCESSION.

Populations in these ecosystems fill a variety of structural and functional roles within the system. Often, groups of populations coevolve, so that they form a more or less isolated subunit. For example, the pollinators of a plant species, and their predators and parasites, form such a guild. Populations continually adapt and develop through natural selection, expanding to the limit of their resources. Population growth is, therefore, due to an increase in resources or a relaxation of limiting factors. See ECOSYSTEM; POPULATION ECOLOGY.

Ecosystem management theory. The objective of applied ecology management is to maintain the system while altering its inputs or outputs. Often, ecology management is designed to maximize a particular output or the quantity of a specific component. Since outputs and inputs are related (Fig. 2), maximization of an output may not be desirable; rather, the management objective may be the optimum level. Optimization of systems can be accomplished through the use of systems ecology methods which consider all parts of the system rather than a specific set of

components. In this way, a series of strategies or scenarios can be evaluated, and the strategy producing the largest gain for the least cost can be chosen for implementation.

The applied ecology management approach has been partially implemented through the U.S. National Environmental Policy Act. This act requires that an environmental-impact analysis be carried out by an interdisciplinary team of specialists

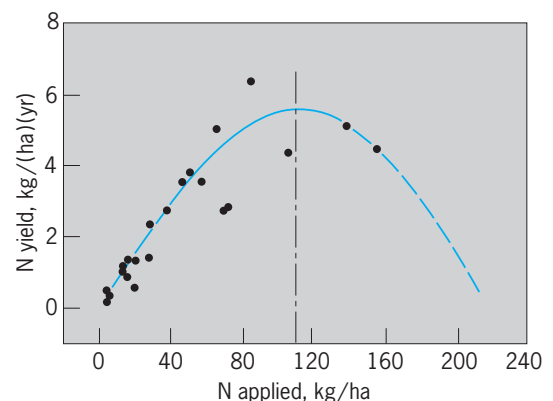


Fig. 2. Relationship between N yield and N consumption (1964-1969).

representing the subjects necessary for an ecosystem analysis. This team is required to evaluate a project on the basis of its environmental, social, and cultural features. One alternative that must be considered is that of no alteration of the system.

A variety of general environmental problems within the scope of applied ecology relate to the major components of the Earth: the atmosphere, water, land, and the biota. The ecological principles used in applied ecology are discussed elsewhere; a sequence of environmental problems of special importance to applied ecology is discussed below.

Atmospheric problems. The atmosphere is one of the most important components of the environment to consider from the viewpoint of applied ecology since it connects all portions of the Earth into one ecosystem. The atmosphere is composed of a variety of gases, of which oxygen and nitrogen make up the largest percentage. It is not uniform in its depth or its composition, but is divided into several layers or zones which differ in density and composition. Although most interaction with humans occurs in the zone nearest Earth, the most distant parts of the atmosphere are also important since they affect the heat balance of the Earth and the quality of radiant energy striking the surface. Disturbances to these portions of the atmosphere could affect the entire biosphere. *See* ATMOSPHERE.

The composition of the atmosphere varies according to location. The qualities of minor gases such as carbon dioxide, the amounts of various metallic elements, and the quantity of water vapor and dust all may differ, depending on the relative distance from land or sea. But, in addition, the atmospheric composition may change in time. For example, over the history of the planet, the percentage composition of oxygen has changed from a very oxygen-poor environment to the present atmosphere, with 20.95% oxygen by volume. *See* ATMOSPHERE, EVOLUTION OF; ATMOSPHERIC CHEMISTRY.

Human activities may introduce a variety of pollutants into the atmosphere. The principal pollutants are carbon dioxide, sulfur compounds, hydrocarbons, nitrogen oxides, solid particles (particulates), and heat. The amounts of pollutants that are produced may be quite large, especially in local areas, and have increased in amount as industrialization has become more widespread. Industrial and domestic activities have also been estimated to put 2.18×10^8 tons (1.96×10^8 metric tons) of sulfur into the atmosphere per year. In most cases, these pollutants have increased during the recent past, and in many areas have become a serious problem. *See* AIR POLLUTION.

Atmospheric problems can have a variety of impacts on humans. Numerous observers have attributed climatic change to atmospheric alteration, since any change in the gaseous envelope of the Earth could alter the heat balance and the climate. The Earth's climate is not constant, and it is difficult to establish an exact correlation between pollution and variation in temperature or solar radiation at the Earth's surface. The pollutants most likely to have an effect on climate are carbon dioxide and solid

particles. Pollution may also affect the chemical balance of regions of the Earth. These effects may be extremely complicated. For example, it has been reported that sulfur oxides produced in the industrial districts of northern Europe have moved north in the atmosphere over Scandinavia. The sulfur oxides react with water to form sulfuric acid, which rains out in precipitation. The acid rain changes the acidity of the soil in this region and may depress the activity of blue-green algae, which fix nitrogen from the atmosphere. Acid rain appears to cause reduction in the growth of trees. Other pollutants may act in a similar complicated fashion through the network of interactions in the Earth ecosystem. *See* ACID RAIN.

Finally, atmospheric pollution has direct effects on plants and animals and human activities. Pollutants, like other materials, can act as limiting influences on the growth, reproduction, and survival of plants and animals. A variety of plants, such as lichens and mosses, that are extremely sensitive to pollution can be used to indicate the degree of atmospheric deterioration. In some severe cases, all vegetation and animal life may be destroyed in the vicinity of the polluting industry. Gases and solids are taken into the lungs of humans, and cause disease or discomfort. In cities, such pollutants as asbestos and lead are exceedingly dangerous to the population. And finally, the impact of pollutants on buildings, clothes, artwork, and machines is costly.

Control of atmospheric pollution requires interception of pollutants at the point of discharge. Industrial control can be achieved by the use of special filters, precipitators, and other devices. Control of pollution for automobiles also may involve special equipment, as well as redesign of engines and fuel. Reduction in dust and similar general sources of air pollution may demand a change in the operation producing the problem.

Water problems. The aquatic environment is of equal importance to applied ecology since most of the Earth's surface is covered by the oceans, and the land is connected to the oceans by streams and rivers. Thus, like the atmosphere, the waters are a connection between distant parts of the biosphere and can carry a disturbance from its origin to another region.

The composition of waters varies widely, and it is essential in evaluating aquatic health to establish the base-line conditions which are stable under the normal or undisturbed conditions. Water pollution arises from a variety of sources—industry, domestic sewage, agricultural fertilizers and feedlots, construction activities, and forest practices. Principal pollutants are sediments, organic pollutants containing nitrites, and phosphates, lead, mercury, hydrocarbons, pesticides, and various synthetic chemicals. The impact of the pollutant depends both on its chemical nature and on the quantity released. All water bodies receive quantities of chemicals and solid materials; a variety of organisms which break down and utilize these inputs have evolved. Serious problems arise when the inputs to the water body become unusually large or contain synthetic materials

which cannot be decomposed by the extant organisms. Further problems may develop through concentrations of pollutants in the food web of the aquatic system. For example, if a chemical is not metabolized by organisms, but is concentrated in their tissues, as are some pesticides, then as each organism is eaten by another, the body burden of chemicals can increase. In this way, predators may obtain very large and dangerous amounts of pollutants. The decline in population of certain fish-eating birds has been attributed to this process of transfer and concentration.

Aquatic pollution has many consequences of significance to humans. An excess of chemical materials which enrich plant and animal growth can cause rapid increase in life. This process, termed eutrophication, may entail dramatic increases in the algal, planktonic, or rooted aquatic plant populations, with the result that the water body appears green in color or becomes clogged with vegetation. Toxic chemicals released to water bodies may directly kill aquatic life or, if present in sublethal amounts, may change the species of plants and animals present. Often aquatic pollution is not a dramatic either-or proposition, with all fish or other aquatic life killed; rather, more commonly a trend toward an increase in the more resistant species is seen, with the elimination of those forms which are especially susceptible to the pollutant. Aquatic pollution also involves heat, especially that derived from industrial activities, including nuclear power and fossil fuel plants. In these instances, water is used to cool the machines or reactors and is exhausted to the environment at elevated temperatures. Since all metabolic and chemical processes are influenced by heat, thermal pollution should have a significant effect on aquatic systems, but thus far it has been difficult to prove that such an impact occurs.

Other aquatic problems of interest to applied ecology concern alteration of water channels by impoundments or channelization and irrigation. In each instance the natural pattern of water movement is altered, and deterioration of the environment may result. Impoundments limit the natural movement of sediment and chemical elements; production patterns in the water and lands below the impoundment may be altered, and other changes may occur. But, on the plus side, impoundments often provide fisheries, electrical energy, recreation, and other advantages. Irrigation problems may involve the movement of salts from depths in the soil, with deposition near the surface. Disturbance of the chemical equilibrium of the soil, in turn, interferes with plant growth. *See* WATER POLLUTION.

Terrestrial and soil problems. Terrestrial environments constantly undergo a degrading and decomposing process owing to the action of water, frost, wind, and other environmental processes on the surface which involve the linkages joining land, water, and atmosphere. Human activities may accelerate these natural processes. In addition, the use of chemical materials on the land may have effects similar to those resulting from their addition to water.

Most terrestrial environmental problems are caused by agricultural, grazing, or forest practices. Probably the most serious effect concerns practices which increase the rate of surface erosion. Only a small percentage of the Earth's surface is suitable for agriculture, and the loss of soil from these areas is extremely serious. As a consequence, certain regions have been denuded and are no longer productive. Overgrazing may also remove the cover of vegetation and allow water and wind to erode the soil. Deserts have increased in extent almost everywhere because of overgrazing, and in India the increase in the Rajasthan Desert can be measured in feet per year. Dust from this desert blows as far east as Thailand. Overcutting trees and lack of reforestation programs also may increase soil erosion and nutrient losses in forest regions. These impacts are not solely the mark of modern civilization. Misuse of the land has been noted in many past civilizations and can even be a problem for present-day primitive societies. *See* DESERTIFICATION; EROSION; FOREST AND FORESTRY.

However, modern agriculture has added new problems to those of the primitive farmer. Various chemicals such as fertilizers, pesticides, and herbicides are used to increase agricultural production. These chemicals may be needed because of past misuse of the land or because of economic demands in a society that does not recognize the need to maintain and protect terrestrial resources. Organic and ecological agricultural practices seek to reestablish a pattern of land use without causing deterioration of the soil and biotic resources. Although reestablishment of the pattern may result in somewhat lower productivity, proponents of "ecoagriculture" argue that high productivity can be maintained without loss of soil through erosion and without a reduction in fertility. *See* AGRICULTURE.

Probably the most serious short-term impact arises from the use of chemicals on the land. In the most extreme cases, the health of the agriculturalist may be affected by the materials. But more commonly, the pesticide or chemical in the soil is taken up by the crop and then enters the human food chain. Many modern governments maintain agencies to advise farmers on the proper amounts of chemicals to apply so that buildup does not occur; other agencies periodically sample and analyze foodstuffs for residues. In this way, the consumer can be protected from misuse of chemicals or from excessive concentrations. Unfortunately, these agencies seldom consider the impacts of agricultural chemicals on other animal food chains. For example, the soil fauna and the natural nitrogen-fixing organisms present in the soil, as well as the terrestrial faunas living near the agricultural or forest plantations, can be significantly affected; and populations of animals, even beneficial species, may be reduced through misuse of chemical materials. However, extinction of plants and animals, which is also a serious applied ecology problem, usually is due to the destruction of their habitat. Pollution, disturbance of the land, and overhunting may provide the final cause of the destruction of a particular living species.

Nuclear energy. Industrialized societies require large quantities of energy. Energy production from nuclear reactors has been enthusiastically developed in many regions of the biosphere. However, nuclear energy also has environmental consequences that are of concern to applied ecology and that must be considered when these facilities are designed and operated, so that negative impacts on the environment do not occur.

Nuclear energy has three primary environmental consequences: the storage of radioactive products, release of radioactive material to the environment, and, as mentioned above, release of heat.

There are various kinds of radioactivity associated with the particular elements used in the reactor. In the process of the generation of energy, these fuel elements are changed into a suite of radioactive materials. Although these materials, in turn, form a new source of chemicals, the process of separation and concentration is very costly and dangerous. In either case, however, the processes result in radioactive waste that must be stored for periods of hundreds, even thousands of years. The fact that the potential danger of these wastes will require technical attention for periods of time longer than the histories of many modern societies is a prime argument against the widespread use of nuclear energy. However, proponents of the use of nuclear energy state that certain geological structures such as salt mines can be safely used for storage indefinitely.

A second environmental problem concerns the loss of relatively small quantities of radioactive materials to the environment during chemical processing in reactors or chemical plants. If these materials enter the body, they can cause disease and death. Like pesticides, radioactive chemical materials may be concentrated in food chains and can appear in relatively large concentrations. In certain fragile environments such as the arctic tundra, the food chains are very short. Thus, radioactive chemicals derived from testing atomic weapons pass through lichens or reindeer or caribou to humans. Concentrations in certain localities may be high enough to cause concern to public health authorities. *See* NUCLEAR REACTOR; RADIOACTIVITY.

Population problems. Applied ecology also is concerned with the size of the human population, since many of the impacts of human activities on the environment are a function of the number and concentration of people. The human population on Earth has increased exponentially, and in many countries this increase poses almost insurmountable problems. Control of environmental degradation, even a concern for environment, is nonexistent when the population is undernourished, starving, ill-housed, and underemployed. Social disorder, alienation, psychological disturbances, physical illness, and other problems have been correlated with overpopulation. Population also places demands on the resources of the Earth; as the standard of living rises, these demands increase.

The human population problem is exceedingly complicated because the growth of population is

controlled largely by the decisions of individual families. The family may visualize several different strategies—the number of children that is best for the family, best for their social group or tribe, or best for the human race—or have no strategy at all. Families may decide that a large family is best even under serious conditions of overpopulation. Considerable evidence indicates that the size of the family declines as the population becomes less rural and more industrial. Thus, some specialists urge economic development, regardless of environmental impact, as a means of solving the population problem. Others urge that the family be more directly influenced to reduce the number of children. Direct action might entail birth control advice, medical abortion and sterilization, and taxation. Yet others argue that no measures such as these can be significant and that the human population will be controlled by famine, war, or disease. Each of these positions leads toward a set of social policies, all of which have an environmental impact which, in turn, affects the human society. *See* HUMAN ECOLOGY.

Environmental planning and design. The foregoing discussion suggests that there is an optimum environment for the human race which is influenced by a variety of population densities and activities. Thus, although the population of the United States is relatively sparse, it has a large environmental impact. This, in turn, suggests that the human environment and society could be designed in such a way to minimize the negative impacts and provide a satisfactory productive life for the population. Environmental design considers economic and social policy, as well as the impact of designed rural, urban, transport, industrial, and other systems. It also considers the design of the individual environment of house, furniture, clothes, and so on. Considering the often violent impact modern society has had on the environment, a design revolution is required to reorganize the environment created by society so that these impacts can be reduced.

Environmental planning and design obviously have a deep political component, since the methods used to redesign society depend upon the control of individual demands. At one extreme, individual demand is allowed to express itself without limit, and education is used to create in the individual a realization that environmental constraints must be recognized. At the other extreme, the government or party controls demand through regulation. Most societies operate somewhere between these extremes.

Throughout human history, utopian designs have been developed for human societies and environments. Today these designs pay more attention to environmental features of society and are often labeled as ecological. Applied ecology, thus, considers not only the alteration of specific features of the modern industrial society to correct some environmental defect but also the fundamental reorientation of society to achieve a balance between humans and the natural world on which they depend. *See* ECOLOGY.

Frank B. Golley

Ecosystem

A functional system that includes an ecological community of organisms together with the physical environment, interacting as a unit. Ecosystems are characterized by flow of energy through food webs, production and degradation of organic matter, and transformation and cycling of nutrient elements. This production of organic molecules serves as the energy base for all biological activity within ecosystems. The consumption of plants by herbivores (organisms that consume living plants or algae) and detritivores (organisms that consume dead organic matter) serves to transfer energy stored in photosynthetically produced organic molecules to other organisms. Coupled to the production of organic matter and flow of energy is the cycling of elements. See ECOLOGICAL COMMUNITIES; ENVIRONMENT.

Autotrophic production. All biological activity within ecosystems is supported by the production of organic matter by autotrophs (organisms that can produce organic molecules such as glucose from inorganic carbon dioxide; Fig. 1). More than 99% of autotrophic production on Earth is through photosynthesis by plants, algae, and certain types of bacteria. Collectively these organisms are termed photoautotrophs (autotrophs that use energy from light to produce organic molecules). In addition to photosynthesis, some production is conducted by chemoautotrophic bacteria (autotrophs that use energy stored in the chemical bonds of inorganic molecules such as hydrogen sulfide to produce organic molecules). The organic molecules produced by autotrophs are used to support the organism's metabolism and reproduction, and to build new tissue. This new tissue is consumed by herbivores or detritivores, which in turn are ultimately consumed by predators or other detritivores.

Terrestrial ecosystems, which cover 30% of the Earth's surface, contribute a little over one-half of the total global photosynthetic production of organic matter (approximately 60×10^{15} grams of carbon per year; Fig. 2). The global rate of photosynthetic production of organic matter in terrestrial ecosystems is

regulated by temperature, water availability, and nutrient concentrations. Correspondingly, the rate of photosynthetic production in the tropics (average rate = $1800 \text{ g C m}^{-2} \text{ y}^{-1}$) is 25-fold greater than in deserts (average rate = $70 \text{ g C m}^{-2} \text{ y}^{-1}$) and 13-fold greater than in tundra and alpine ecosystems (average rate = $140 \text{ g C m}^{-2} \text{ y}^{-1}$).

Oceans, which cover 70% of the Earth's surface, produce approximately $51 \times 10^{15} \text{ g C y}^{-1}$ of organic matter (Fig. 2). Photosynthesis in oceans varies greatly across the Earth in response to nutrient availability (primarily nitrogen and iron). Photosynthesis is greatest along continental margins and in regions where deep, nutrient-rich water upwells into surface waters receiving light. Photosynthetic production is greatest in regions of deep upwelling (averaging $420 \text{ g C m}^{-2} \text{ y}^{-1}$), intermediate in coastal margins (averaging $250 \text{ g C m}^{-2} \text{ y}^{-1}$), and lowest in the open ocean (averaging $130 \text{ g C m}^{-2} \text{ y}^{-1}$).

Whereas most ecosystems receive sunlight to support photosynthesis, photosynthetic production does occur in the deep ocean and in caves. In a few unique ecosystems, such as deep-sea thermal vents and caves, the basis of organic matter production is by chemoautotrophic bacteria. These ecosystems harbor surprisingly diverse communities that function in the complete absence of sunlight.

Heterotrophic production. The organic matter produced by autotrophic organisms is ultimately consumed by heterotrophs (organisms that obtain energy and nutrients by consuming other organisms; Fig. 1). Globally, most organic matter produced by plants and algae is consumed by detritivores. Herbivores consume only about 5% of photosynthetic production, fire destroys another 5%, and a small fraction of autotrophic production is lost through burial within the Earth's crust (to potentially become fossil fuels). The proportion of photosynthetic production consumed by herbivores can vary considerably among ecosystems. In terrestrial ecosystems, much plant tissue is structural molecules such as cellulose and lignin that are difficult or impossible for most organisms to digest. In aquatic ecosystems, in which algae dominate, a much greater proportion of photosynthetic production is consumed by herbivores. Algae have little in the way of structural molecules and contain more molecules that can be consumed by grazing organisms. See BIOMASS; ECOLOGICAL SUCCESSION.

Food webs. Organisms are classified based upon the number of energy transfers through a food web (Fig. 1). Photoautotrophic production of organic matter represents the first energy transfer in ecosystems and is classified as primary production. Consumption of a plant by a herbivore is the second energy transfer, and thus herbivores occupy the second trophic level, also known as secondary production. Consumer organisms that are one, two, or three transfers from photoautotrophs are classified as primary, secondary, and tertiary consumers. Moving through a food web, energy is lost during each transfer as heat, as described by the second law of thermodynamics. Consequently, the total number of

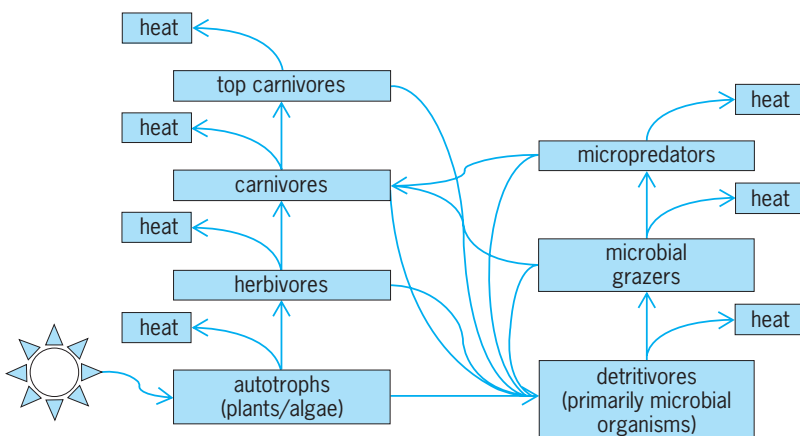


Fig. 1. General model of energy flow through ecosystems.

energy transfers rarely exceeds four or five; with energy loss during each transfer, little energy is available to support organisms at the highest levels of a food web. See ECOLOGICAL ENERGETICS; FOOD WEB.

Biogeochemical cycles. In contrast to energy, which is lost from ecosystems as heat, chemical elements (or nutrients) that compose molecules within organisms are not altered and may repeatedly cycle between organisms and their environment. Approximately 40 elements compose the bodies of organisms, with carbon, oxygen, hydrogen, nitrogen, and phosphorus being the most abundant. If one of these elements is in short supply in the environment, the growth of organisms can be limited, even if sufficient energy is available. In particular, nitrogen and phosphorus are the elements most commonly limiting organism growth. This limitation is illustrated by the widespread use of fertilizers, which are applied to agricultural fields to alleviate nutrient limitation. See BIOGEOCHEMISTRY.

Nitrogen cycle. Nitrogen commonly limits the rate of primary production in terrestrial, fresh-water, estuarine, and oceanic ecosystems (Fig. 3). In one turn of the biogeochemical cycle of nitrogen, (1) cyanobacteria and certain types of eubacteria transform atmospheric dinitrogen into organic molecules such as protein (nitrogen fixation); (2) organisms are consumed by other organisms and organic nitrogen is converted into ammonium as a waste product (decomposition); (3) ammonium either is assimilated by plants, algae, and bacteria into organic forms or is converted to nitrate by chemoautotrophic bacteria (nitrification); and (4) nitrate is converted back to dinitrogen by anaerobic bacteria (denitrification). Whereas nitrogen is very abundant in the atmosphere as dinitrogen (N_2), only a select group of organisms can transform dinitrogen into organic forms (amino acids), and the transformation requires a lot of metabolic energy. Humans have greatly impacted the global nitrogen cycle by doubling the amount of ammonium, nitrate, and organic nitrogen in the environment through the widespread use of fertilizers and combustion of fossil fuels. This massive change in nitrogen forms has led to problems with acid precipitation, alterations of global patterns of primary production rates, increased production of nitrous oxide (a greenhouse gas), and pollution of aquatic ecosystems.

Phosphorus cycle. Phosphorus also commonly limits primary production. In contrast to nitrogen, phosphorus lacks an atmospheric phase. The primary source of phosphorus to ecosystems is the weathering of the Earth's crust. This inorganic phosphorus is then assimilated by plants, algae, and bacteria and is incorporated into organic forms such as phospholipids and adenosine triphosphate (ATP). As organisms are consumed by other organisms, organic phosphorus is converted back to inorganic forms. Both organic and inorganic forms are slowly carried by water from plants and soils to streams and rivers, and ultimately to the world's oceans where phosphorus is buried in sediments. The global phosphorus cycle is completed over long, geologic time scales

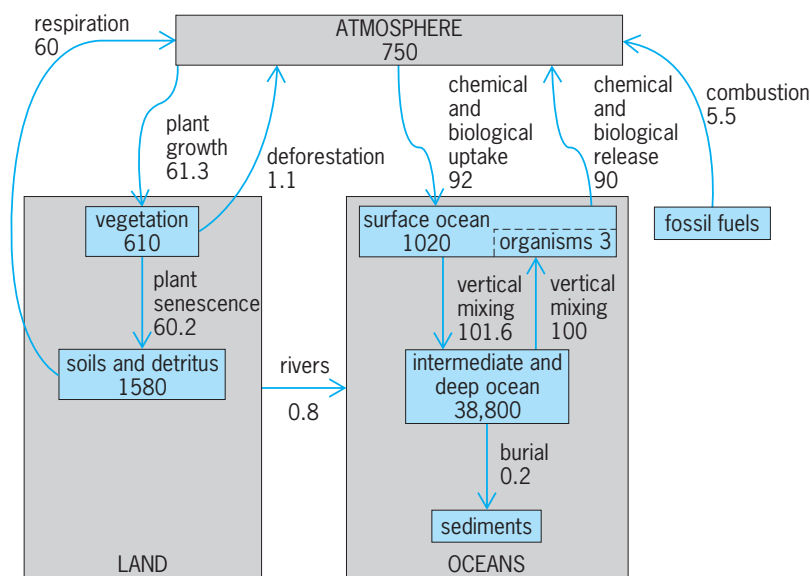


Fig. 2. Model of the global carbon cycle illustrating movement of carbon between the atmosphere and terrestrial and oceanic ecosystems. The storage of carbon in these ecosystems is in picograms ($1 \text{ pg} = 10^{15} \text{ g}$), and fluxes of carbon between boxes are in pg y^{-1} .

when ocean sediments are compacted into rock and this new rock is uplifted onto continents. As with the nitrogen cycle, humans have dramatically impacted the phosphorus cycle by application of fertilizers. The major source of phosphorus fertilizers is from the mining of guano. Current estimates project readily accessible sources of guano being depleted by the end of the twenty-first century.

Carbon cycle. The carbon cycle has also received considerable research focus, especially with the increase of carbon dioxide in the atmosphere due to the combustion of fossil fuels, and the impact on global climate. Carbon cycles between the atmosphere and terrestrial and oceanic ecosystems. This cycling results, in part, from primary production and decomposition of organic matter. Rates of primary production and decomposition, in turn, are regulated

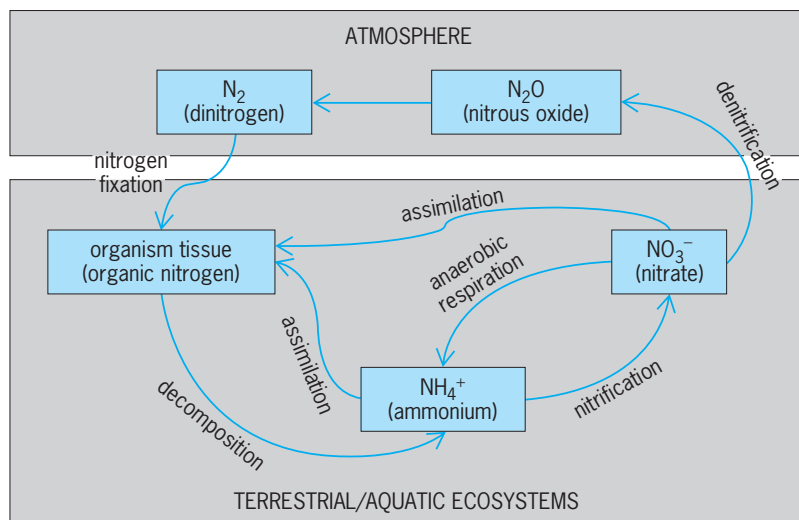


Fig. 3. General model of the nitrogen cycle in ecosystems.

by the supply of nitrogen, phosphorus, and iron. The combustion of fossil fuels is a recent change in the global cycle that releases carbon that has long been buried within the Earth's crust to the atmosphere. Carbon dioxide in the atmosphere traps heat on the Earth's surface and is a major factor regulating the climate. This alteration of the global carbon cycle along with the resulting impact on the climate is a major issue under investigation by ecosystem ecologists. See AIR POLLUTION; CONSERVATION OF RESOURCES; ECOLOGY, APPLIED; HUMAN ECOLOGY; WATER POLLUTION.

Jeremy B. Jones

Bibliography. F. B. Golley, *A History of the Ecosystem Concept in Ecology*, Yale University Press, 1993; J. B. Hagen, *An Entangled Bank: The Origins of Ecosystem Ecology*, Rutgers University Press, 1992; W. H. Schlesinger, *Biogeochemistry: An Analysis of Global Change*, 2d ed., Academic Press, 1991; R. H. Waring and S. W. Running, *Forest Ecosystems: An Analysis at Multiple Scales*, Academic Press, 1998.

Ecotone

A geographic boundary or transition zone between two different groups of plant or animal distributions. The term has been used to denote transitions at different spatial scales or levels of analysis, and may refer to any one of several attributes of the organisms involved. For example, an ecotone could refer to physiognomy (roughly, the morphology or appearance of the relevant organisms), such as between the boreal forest and grassland biomes; or it could refer to composition, such as between oak-hickory and maple-basswood forest associations; or it could refer to both. Ecotones are generally distinguished from other geographic transitions of biota by their relative sharpness. The ecotone between boreal forest and prairie in central Saskatchewan occurs over a hundred kilometers or so, in contrast to the transition from tropical forest to savanna in South America or Africa that is associated with increasing aridity and is dispersed over hundreds of kilometers. The "tension zone" between broadleaf deciduous forests in south-central Michigan and mixed forests to the north is similarly sharp. Ecotones are thought to reflect concentrated long-term gradients of one or more current environmental (rather than historical or human) factors. Though often climatic, these factors can also be due to substrate materials, such as glacial sediments or soils. Regardless of their specific environmental basis, most ecotones are thought to be relatively stable.

A good example of an ecotone is the abrupt prairie-forest transition in North America that extends southeastward from south-central Alberta, Canada, through central Minnesota, and ultimately (in more diffuse form) to approximately southern Lake Michigan. Separating the grassland biome in the western Midwest and Great Plains from the boreal and deciduous forest biomes to the northeast, this ecotone is thought to reflect long-term patterns of

precipitation and evaporation (warmer, drier to the southwest). These, in turn, are produced by modal (most frequent) atmospheric circulation patterns over North America as they repeatedly steer daily weather features across the region and ultimately generate environments, either directly or indirectly through patterns of disturbance, that are conducive to prairie or forest persistence in different parts of the region. Typically, specific sites within an ecotone will support communities that are more like one or the other broad groupings of biota the ecotone separates, depending on local factors including environmentally related patterns of disturbances such as fire.

Ecotones are often reflected in the distributions of many biota besides the biota used to define them. The prairie-forest ecotone, for example, is defined not only by the dominant vegetation components but also by many faunal members of the associated ecosystems, such as insects, reptiles and amphibians, mammals, and birds, that reach their geographic limits here. Across the glaciated midwestern United States, where landscapes of different geologic ages co-occur, the associated soil patterns are sometimes reflected in similarly bold ecological patterns. In central Illinois, for example, the transition from a Wisconsinan-aged glacial landscape (deposited approximately 20,000 years ago) in the north to an Illinoian-aged landscape (several hundred thousand years older) to the south corresponds to a transition to older (better-developed) soils, different dominant tree species in the upland wooded tracts, and even contrasts in the herpetofauna, with many species in one area historically absent from the other. See ALPINE VEGETATION; BIOME; ECOLOGICAL COMMUNITIES; ECOSYSTEM; FOREST ECOSYSTEM; GRASSLAND ECOSYSTEM; LIFE ZONES; SAVANNA; ZOOGEOGRAPHY.

Jay R. Harman

Bibliography. J. A. Harrington and J. R. Harman, Climate and vegetation in central North America: Natural patterns and human alterations, *Great Plains Quart.*, 11:103-112, 1991; J. B. Lachavanne and R. Juge (eds.), *Biodiversity in Land/Inland Water Ecotones*, CRC Press-Parthenon, 1997; K. E. Medley and J. R. Harman, Relationships between the vegetation tension zone and soils distribution across central Lower Michigan, *Mich. Bot.*, 26:78-87, 1987.

Eddy current

An electric current induced within the body of a conductor when that conductor either moves through a nonuniform magnetic field or is in a region where there is a change in magnetic flux. It is sometimes called Foucault current. Although eddy currents can be induced in any electrical conductor, the effect is most pronounced in solid metallic conductors. Eddy currents are utilized in induction heating and to damp out oscillations in various devices.

Causes. If a solid conductor is moving through a nonuniform magnetic field, electromotive forces (emfs) are set up that are greater in that part of the conductor that is moving through the strong part of

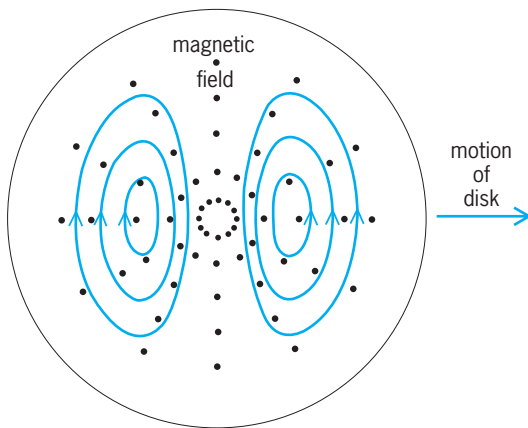


Fig. 1. Eddy currents which are induced in a disk moving through a nonuniform magnetic field.

the field than in the part moving through the weaker part of the field. Therefore, at any one time in the motion, there are many closed paths within the body of the conductor in which the net emf is not zero. There are thus induced circulatory currents that are called eddy currents (Fig. 1). In accordance with Lenz's law, these eddy currents circulate in such a manner as to oppose the motion of the conductor through the magnetic field. The motion is damped by the opposing force. For example, if a sheet of aluminum is dropped between the poles of an electromagnet, it does not fall freely, but is retarded by the force due to the eddy currents set up in the sheet. If an aluminum plate oscillates between the poles, it will be stopped quickly when the switch is closed and the field set up. The energy of motion of the aluminum plate is converted into heat energy in the plate. See ELECTROMAGNETIC INDUCTION; LENZ'S LAW.

Eddy currents are also set up within the body of material when it is in a region in which the magnetic flux is changing rapidly, as in the core of a transformer. As the alternating current changes rapidly, there is also an alternating flux that induces an emf in the secondary coil and at the same time induces emfs in the iron core. The emfs in the core cause eddy currents that are undesirable because of the heat developed in the core (which results in high energy losses) and because of an undesirable rise in temperature of the core. Another undesirable effect is the magnetic flux set up by the eddy currents. This flux is always in such a direction as to oppose the change that caused it, and thus it produces a demagnetizing effect in the core. The flux never reaches as high a value in the core as it would if there were no eddy currents present.

Laminations. Induced emfs are always present in conductors that move in magnetic fields or are present in fields that are changing. However, it is possible to reduce the eddy currents caused by these emfs by laminating the conductor, that is, by building the conductor of many thin sheets that are insulated from each other rather than making it of a single solid piece. In an iron core the thin iron

sheets are insulated by oxides on the surface or by thin coats of varnish. The laminations do not reduce the induced emfs, but if they are properly oriented to cut across the paths of the eddy currents, they confine the currents largely to single laminae, where the paths are long, making higher resistance; the resulting net emf in the possible closed path is small. Bundles of iron wires or powdered iron formed into a core by high pressure are also used to break up the current paths and reduce the eddy currents. See CORE LOSS.

Kenneth V. Manning

Testing. The testing of metallic materials through the detecting of eddy currents induced within them is widely used. A high-frequency alternating current in a coil placed close to a metallic material induces eddy currents within the material through electromagnetic induction. The flux caused by the eddy currents threads the coil producing the original flux and alters its apparent impedance. This alteration is a measure of the strength of the eddy currents. The flow of the eddy current is impeded by the material's electric resistance, which reflects characteristics such as hardness and chemical composition. Therefore, eddy current measurements reveal various properties of materials.

The hardness of aluminum alloys used in the aerospace industry can be related to their electrical conductivity. Eddy-current conductivity meters are widely used to nondestructively monitor these alloys at various stages of aircraft production, in particular after critical heat treatment processes. The rapid measurement of electrical conductivity is one of the many checks made by coin vending machines, so the conductivity of alloys used in the manufacture of coins is often determined using an eddy-current conductivity meter to assure close quality control.

Flaws such as open cracks within materials cause eddy currents to detour around them and effectively enhance electric resistance; accordingly, eddy current testing can be used to detect flaws. Differences in thickness of foils and plating cause differences in the electric resistance; hence, eddy current testing is also convenient for estimating thickness. Stronger

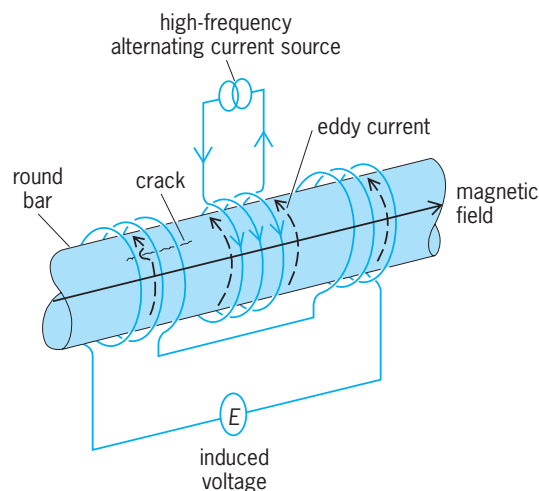


Fig. 2. The alignment of coils used to inspect a long bar.

magnetic fields are induced within ferromagnetic materials undergoing testing, so that the magnetic permeability also affects the eddy current. The evaluation of the hardness of ferritic steel through eddy current testing is based upon this effect of permeability, which is very susceptible to hardness. See FERROMAGNETISM; MAGNETIC MATERIALS; MAGNETISM.

Figure 2 shows the general alignment of coils to detect flaws within long bars. Annular coils encircle the round bar. The alternating current applied in the central coil generates a magnetic field inside the coil and within the bar in the direction of their axis, and thereby induces eddy currents along concentric circular paths within the bar. Simultaneously this eddy current produces a counter magnetic field which is superimposed on the primary one. The combined magnetic field gives rise to electromotive forces, that is, voltages between terminals of each coil. Two coils on either side of the central coil are wound in opposite directions, connected in series. The voltage between terminals is monitored continuously. When a flaw within a moving bar enters one of the two coils, the voltage reading corresponds to the difference between the voltages of the coils. This alignment permits the detection of flaws with high sensitivity, because it eliminates noises that originate from irregularities in diameter or physical properties of materials, whose variation is negligibly small between the closely spaced coils.

Since eddy current testing is of the noncontacting type, it can detect flaws within materials moving at high speed and even those raised to red heat in hot-rolling mills. Inspection of the latter is advantageous to exclude the influence of variation in permeability, because ferromagnetism vanishes at temperatures above the Curie point (1416°F or 769°C for iron). See CURIE TEMPERATURE; NONDESTRUCTIVE EVALUATION.

Kazuo Watanabe; A. E. Drake

Braking and damping. The principle outlined above, whereby a conductor which is moving through a nonuniform magnetic field experiences a retarding force, is used to brake or damp the oscillations of moving objects. In these applications, the field is that produced by a permanent magnet or an electromagnet supplied with direct current. The retarding force is proportional to the product of the flux produced by the magnet ϕ and the effective eddy current I produced by it. The effective eddy current, however, results from the rate of flux cutting, that is, the flux ϕ multiplied by the velocity of movement v . The retarding force is thus proportional to $\phi^2 v$. For a constant flux, which applies with a permanent magnet or electromagnet with constant current, the retarding force is proportional to the velocity of the moving object. With an electromagnet, the current can be varied, to adjust the ratio of retarding force to velocity. This ratio is also affected by the physical dimensions and materials used, for example, the resistivity of the moving object. Brakes are often fitted to rotating objects, and in these applications the retarding torque is proportional to the angular velocity ω .

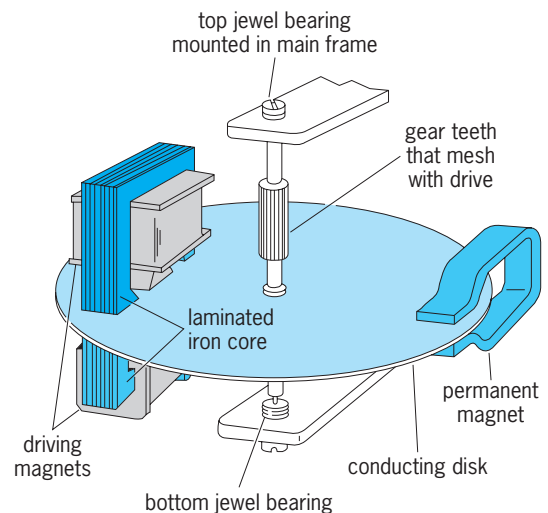


Fig. 3. Basic construction of electromechanical energy meter.

Eddy-current brakes are probably most widely used in the electromechanical energy meters which are installed in consumers' premises by electricity supply companies. In this application, the meters (**Fig. 3**) contain electromagnets which produce fluxes proportional to the system voltage and current. These fluxes induce emf's and eddy currents in a conducting disk, which is mounted in jewel bearings so that it may rotate freely. The eddy currents react with the fluxes to produce a torque proportional to the power being supplied. Eddy-current braking, produced by a flux from a permanent magnet cutting the disk, provides a retarding torque proportional to the disk speed. Because the frictional torques are negligible, the disk runs at a speed at which the driving and retarding torques are almost equal, and as a result the speed is proportional to the power being supplied to the consumer. A geared drive to indicating dials meshes with gear teeth on the disk's spindle. See WATT-HOUR METER.

The same form of braking is also employed widely in the induction-type relays to provide time-graded protection. Eddy-current brakes incorporating electromagnets are used to provide controllable load torques on test equipment for small machines. See RELAY.

Eddy-current damping is used to suppress undesirable mechanical oscillations and is used widely in deflection-type indicating instruments. See DAMPING.

Arthur Wright

Bibliography. Atomic Energy of Canada, Ltd, *Eddy Current Testing*, 2 vols., 1987; G. Birnbaum and G. Free (eds.), *Eddy-current Characterisation of Materials and Structures*, ASTM Spec. Publ., no. 722, American Society for Testing and Materials, 1981; B. I. Bleaney and B. Bleaney, *Electricity and Magnetism*, 2 vols., 3d ed., 1976, paper, 1989; D. Hagemeyer, *Fundamentals of Eddy Current Testing*, 1990; A. E. Drake and A. C. Lynch, AC conductivity standards for the calibration of eddy-current conductivity meters, *J. Phys. E*, 120:137-139, 1987;

H. L. Libby, *Introduction to Electromagnetic Non-destructive Test Methods*, 1971, reprint, 1979; E. M. Purcell, *Electricity and Magnetism*, 2d ed., 1985; S. S. Udpa and P. O. Moore (eds.), *Eddy Current Testing*, vol. 5 of *Nondestructive Testing Handbook*, 3d ed., American Society for Nondestructive Testing, 2004.

Edema

An abnormal accumulation of fluid in the cells, tissue spaces, or cavities of the body, also known as dropsy. An excess of fluid in the pleural spaces is referred to as hydrothorax, in the pericardial sac as hydropericardium, and in the peritoneal cavity as ascites. Anasarca is a generalized subcutaneous edema.

There are three main factors in the formation of generalized edema and a fourth which plays an important role in the formation of local edema. They are (1) permeability of the capillary wall, (2) colloid osmotic pressure of the plasma proteins, (3) hydrostatic pressure in the capillaries, and (4) lymphatic obstruction.

Permeability of capillary wall. Normally the capillary walls are freely permeable to water, salts, and dissolved gases but are almost impermeable to proteins. When the vessel wall is injured by toxins, anoxia, or paralytic dilatation, the capillary endothelium becomes permeable to proteins. With a diffusion of protein into the tissues, the plasma osmotic pressure is lowered and the osmotic pressure of the tissue is increased. Under these circumstances fluid collects in the tissue spaces.

Such a condition plays an important role in inflammatory edema. It is a factor in the edema of severe infections, metabolic intoxications, asphyxia, anaphylactic reactions, secondary shock, and acute nephritis. It also contributes to the edematous conditions when there is a fall in the level of plasma proteins.

Osmotic pressure of plasma proteins. A fall in plasma proteins tends to decrease the forces tending to reabsorb and hold fluid in the vascular compartment. Albumin is the protein of greatest importance in this regard. When the plasma protein level drops below 3 g/100 ml, the colloid osmotic pressure is no longer sufficient to maintain a balance with the hydrostatic pressure of the blood, which tends to drive fluid out into the tissue spaces. Hence more fluid goes out into the tissue spaces and remains there until a new equilibrium is reached.

This form of edema is seen in association with prolonged malnutrition (nutritional edema) as during a famine or with chronic nutritional or metabolic defects. With a marked loss of albumin in the urine as in nephrotic syndrome, there follows a lowering of the plasma albumin fraction and the development of edema.

In edema of kidney disease the protein content and specific gravity of the edema fluid are low. The blood reveals a markedly elevated cholesterol and a drop in total protein level with a relatively greater drop in the albumin fraction. As the plasma os-

motric pressure drops, water passes into the tissues and with it crystalloids. These substances, especially sodium chloride, are retained in the tissues. Thus as water is taken in, it passes rapidly into the tissues and is not eliminated into the urine. *See* OSMOREGULATORY MECHANISMS.

Capillary hydrostatic pressure. Under normal conditions the hydrostatic pressure in the arterial end of the capillary is sufficient to overcome the plasma osmotic pressure and drives fluid out of the capillary into the tissue spaces. During passage through the capillary the pressure drops to a level low enough to allow the osmotic pressure of the proteins to draw fluid back into the vascular compartment. An increase in hydrostatic pressure at the venous end of the capillary will upset the balance, resulting in a decreased absorption of tissue fluid by the osmotic pressure of the plasma proteins. Under these circumstances an increased amount of fluid will be returned via the lymphatics, but as the condition progresses edema will develop. Such a situation can follow venous congestion of long duration. *See* CIRCULATION.

Cardiac edema following the generalized venous congestion of cardiac failure is the commonest form of this type of edema. The fluid which collects in the tissue spaces is affected by changes in position, being more marked in the dependent portions of the body. Fluid also collects in the serous cavities. The lymphatics empty into the venous system, and therefore a rise in pressure in the venous system results in an increased pressure within the lymphatics, which also contributes to the edema formation. As the capillary walls are distended, they become more permeable. In addition, a state of chronic hypoxia may exist, causing a further insult to the capillary endothelium with a loss of protein into the tissue spaces.

Other factors seem to play a role in cardiac edema. With decreased cardiac output there is reduced renal blood flow and glomerular filtration rate, with a consequent reduced excretion of salt and water. This may be responsible for an increased volume of extracellular fluid and plasma which, in turn, is followed by a rise in venous pressure.

Pulmonary edema (edema within the lung) is usually a form of cardiac edema but may be secondary to other factors such as inflammation. *See* RESPIRATORY SYSTEM DISORDERS.

Another condition resulting from an increased hydrostatic pressure in the capillaries is postural edema. This occurs when an individual has been standing motionless for a long period of time; the fluid collects in the subcutaneous tissues of the feet and ankles.

Cirrhosis of the liver causes an impediment to the flow of blood through the portal circulation. There is a consequent rise in venous pressure and ascites forms. *See* CIRRHOSIS.

Lymphatic obstruction. A portion of the intercellular tissue fluids returns to the circulation via the lymphatics. Obstruction to this channel will contribute to local edema. Infestation by filaria is one cause of lymphatic obstruction, particularly in the

tropics. Lymphatic channels may be destroyed or obstructed by surgical procedures, resulting in localized edema. Milroy's disease is a chronic hereditary edema thought to be due to lymphatic obstruction. See LYMPHATIC SYSTEM.

Management of edema. The management of patients with edema is directed toward the treatment of the underlying medical condition. Diuretics are often used, most successfully in cardiac edema.

Romeo A. Vidone; Irwin Nash

Bibliography. A. C. Guyton and J. E. Hall, *Textbook of Medical Physiology*, 9th ed., 1996; J. M. Kissane (ed.), *Anderson's Pathology*, 9th ed., 1989; S. L. Robbins et al., *Pathologic Basis of Disease*, 6th ed., 1999; J. B. Wyngaarden and L. H. Smith (eds.), *Cecil Textbook of Medicine*, 17th ed., 1985.

Edentata

The former name of an order that included the Xenarthra (sloths, anteaters, armadillos), Pholidota (pangolins), and Tubulidentata (aardvarks). The order was so named because of the apparent lack of some or all teeth, although most of these species are not toothless but have at least vestigial teeth. It was subsequently realized that Edentata was polyphyletic—that it contained unrelated families and was thus invalid. It is now thought that any similarities between these groups is the result of similar adaptations to a common way of life and are not indicative of actual relationships. Thus, pangolins and aardvarks are now placed in separate orders, and the order Xenarthra was established to group the remaining families (which are all related). The term edentate is now used only as a term of convenience when referring to these unusual mammals; it has no taxonomic validity. See AARDVARK; ANTEATER; ARMADILLO; DENTITION; MAMMALIA; PHOLIDOTA; SLOTH; TOOTH; TUBULIDENTATA.

Donald W. Linzey

Bibliography. R. M. Nowak, *Walker's Mammals of the World*, 6th ed., Johns Hopkins University Press, 1999.

Ediacaran biota

A widely distributed group of soft-bodied marine organisms that are preserved as fossils in rocks of latest Proterozoic age (600–543 million years ago, or Ma). The biota characterizes a geological period, known as the Ediacarian or the Vendian, which precedes the widespread appearance of animals with mineralized skeletons. The name Ediacara refers to an abandoned mining area about 380 mi (600 km) north of Adelaide, South Australia.

Discovery and distribution. Although Ediacaran fossils were described in the 1930s from southwest Africa (Namibia), it was the discovery in 1946 of abundant fossil “jellyfish” at Ediacara that sparked international interest in this biota. Subsequently, sim-

ilar or identical fossils were found in central England (1957); southeastern Newfoundland (1967); northern Russia, the Ukraine, northern Siberia, and the Ural Mountains (1960s–1990s); North Carolina (1966); northwestern Canada (1979); southern Nevada (1997); and elsewhere. About 30 localities have been described, with the most diverse biotas found at Ediacara, Namibia, and on the coast of the White Sea in northern Russia. As some of these sites could not have been less than about 6000 mi (10,000 km) apart, no matter how the continents were arranged at the time, there is no doubt that the Ediacaran biota was a globally distributed marine biota. Ediacaran fossils have also been found in both shallow-water and deep-water facies. The biota appears to show a degree of biogeographic provinciality, but it has proved difficult to untangle the effects of geographic separation, paleoenvironmental differences, age differences, and sampling and preservation biases on the composition and preservation of the biota at each major site. A number of taxa are found in rocks formed under widely different conditions, and these, at least, may have been ecological generalists.

Age of assemblage. Recent radiometric dating has bracketed the Ediacaran biota between 600 and 543 Ma, with most sites between 565 and 543 Ma. Most known occurrences of Ediacaran organisms precede the earliest great radiation of skeletal fossils (archaeocyaths, trilobites, mollusks, brachiopods), or else they can be placed as latest Proterozoic on other evidence. Although unfossiliferous strata separate the Ediacaran biota from the earliest Cambrian fossils at many localities, Ediacaran biotas in Namibia and Nevada are found in rock sequences extending right up to the base of the Cambrian. Radiometric dates have confirmed that there was no significant time gap between the disappearance of the Ediacaran organisms and the Cambrian radiations. In fact, a few Ediacaran fossils have been found in Cambrian strata; the biota did not entirely die out before the Cambrian. See CAMBRIAN.

Fossil preservation. At almost all sites, the fossils are preserved as impressions in some kind of detrital sedimentary rock. Commonly, they are found on the bases of sandstone beds (South Australia), within sandstone beds (Namibia), or below volcanic ashes swept into deep water by wind and turbidity currents (Newfoundland). The organisms appear to have lived in continental shelf to slope environments and are normally preserved in sediments that were deposited under fairly quiet conditions below normal wave base. Organisms with resistant bodies are usually preserved as concave impressions in the casting medium, but other fossils appear to be the sediment fillings of either soft bodies or the cavities left by their decay. Most strata where Ediacaran fossils have been found also show evidence of extensive microbial growth that formed firm mats on the sediment surface. By stabilizing and firming the sediment, and possibly by speeding up its lithification, these mats are now thought to have been crucial to the fossilization of these soft-bodied organisms.

Taxonomic diversity. More than a hundred scientific names have been given to various kinds of Ediacaran fossils, but many of these are based on questionable material, or they are unnecessary synonyms for the same genus or species preserved in different ways. Nevertheless, at a conservative estimate there are probably 40–50 distinct genera or probable genera of Ediacaran organisms worldwide. Assigning these genera to higher taxa, however, has been controversial. Prior to 1985, many scientists who worked on the Ediacaran fossils attempted to relate them to the phyla of living invertebrates. A real problem is that few of these fossils can unequivocally be referred to living or extinct animal taxa. Because many of the fossils are simply circular structures with or without radial or concentric markings, they impart little information and are difficult to interpret. Some of these circular “medusoids” could be the remains of cnidarian jellyfish (Hydrozoa or Scyphozoa), but very few show diagnostic cnidarian characters such as fourfold symmetry, or even diagnostic metazoan characters such as definite guts. The recognition of many critical features is hampered by the nature of the preservation. In addition, many Ediacaran fossils, notably most forms from Namibia, are so unusual in shape that they cannot be placed firmly in any modern group. See CNIDARIA.

Many workers have placed the more unusual organisms in an extinct higher taxon of phylum grade, commonly named the Petalonamae. Others have used differences in symmetry and body organization to identify and characterize several major taxonomic groups regarded, in general, as extinct higher taxa of phylum or class grade. Beginning in the mid-1980s, much controversy and interest were sparked by proposals that all of the Ediacaran organisms were constructed on a single basic plan that was radically different from any other animal. According to this hypothesis, the Ediacarans were “quilted” organisms lacking heads, muscles, or digestive systems, and consisted of parallel sheetlike walls held together by regularly spaced internal partitions, and the whole organism was inflated by body fluids. Without digestive systems, these organisms would have lived by taking up dissolved nutrients or perhaps by photosynthesis. Proponents of this hypothesis classify the Ediacaran organisms in an extinct kingdom of multicellular life, the Vendobionta. Still others have proposed that the Ediacaran organisms belonged to extant kingdoms other than the animals; they were lichens, algae, or single-celled or colonial protists.

Controversy still reigns, but in all likelihood no hypothesis is entirely incorrect. Ediacaran fossils have been described which do show clear features diagnostic of living animal phyla. However, many of the Ediacaran organisms remain impossible to definitively link with living animals, and they may well represent extinct phyla or even kingdoms.

Distinctive genera such as the discoidal “worm” *Dickinsonia* (Fig. 1), the sac-shaped *Ernieetta*, the paddle-shaped “sea pen” *Charnia*, and the trifoliate *Pteridinium* seem to form a coherent group of organisms with a Vendobionta-like “quilted” or

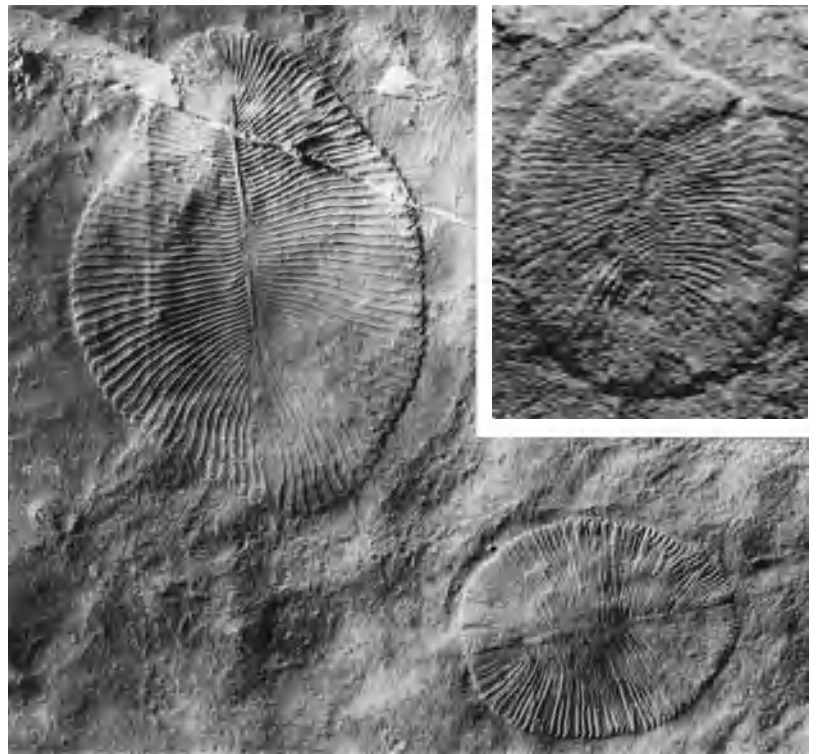


Fig. 1. Natural impressions of *Dickinsonia* from Australia and from Russia (inset). The Australian specimens, 5 and 3 in. (13 and 8 cm) in length, have about the same number of segments and may be regarded as the contracted and expanded states of a single individual.

“air-mattress” type of construction. Whether this common structural pattern is due to convergent evolution in response to particular environmental conditions remains uncertain, and so the Petalonamae or even the Vendobionta may be an unnatural group of unrelated taxa. There is, for example, evidence that “sea pens” such as *Charniodiscus* and *Rangaea* (Fig. 2) lived attached to the ocean floor by

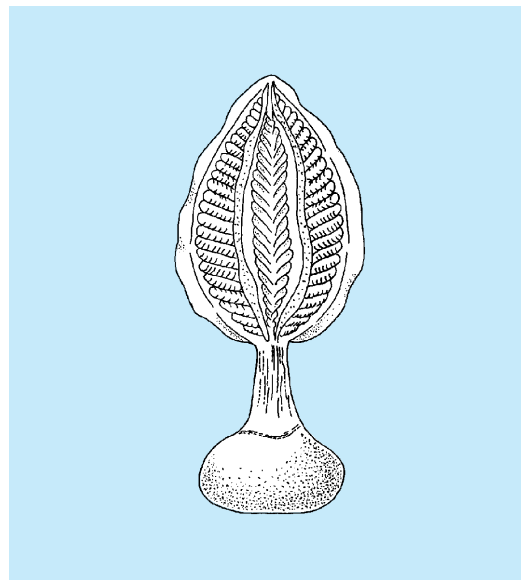


Fig. 2. Reconstruction of *Rangaea*, a complex, sessile genus of “sea pens” from Namibia.

globular holdfasts, whereas *Ernietta* lived partially buried in sand with no holdfast. Unnamed spindle-shaped “sea pens” from Newfoundland also show no trace of an attachment structure. *Dickinsonia* seems to have been free-living and mobile and may have had a gut.

The members of a second group of Ediacaran animals are unified by their threefold symmetry. All have a circular outline and low conical shape reminiscent of modern jellyfish, and they have been regarded as an extinct class of the Cnidaria by some researchers. However, the morphological complexity of genera such as *Tribrachidium* and *Albumares* argues against such a view. It is plausible that the tri-radial forms are members of a second extinct phylum, the Trilobozoa. True free-living cnidarians may be represented by “medusoids” such as *Ediacaria* and *Eoporpita*, which are now thought to be benthic polyplike organisms. A few forms, such as the tetradial *Conomedusites* and the chondrophorine-like *Ovatoscutum* (Fig. 3), may have been planktonic cnidarians. The Australian *Palaeophragmodictyon* shows impressions of a spicule network and is almost certainly a sponge.

Although several higher invertebrate phyla such as the Annelida (segmented worms), Arthropoda, and Echinodermata have been identified from the Ediacaran biota, only a few of the known genera could be used to make a plausible case for the presence of these phyla. *Spriggina floundersi*, known from tens of specimens from Ediacara and nearby sites, had a long, bilaterally symmetrical, segmented body with a boomerang-shaped feature at one end. *Spriggina* has been placed with the annelids, but it is equally likely that it was a primitive soft-shelled arthropod. In this interpretation, the boomerang-shaped structure is considered to be homologous with the arthropod head. A few *Spriggina*-like fossils show traces of branching caecae within this head, closely resembling the anatomy of trilobites and other arthropods from Cambrian sites such as the Burgess Shale. Yet an-



Fig. 3. Discoidal fossil *Ovatoscutum* from South Australia may be a link between the Ediacaran fauna and younger biotas. It is probably an impression of the float of a chondrophorine cnidarian.



Fig. 4. Mold of *Cloudina* or a related genus, an annulated tubular fossil that was mineralized in life. *Cloudina* has been found together with soft-bodied Ediacaran fossils at several localities; this one is from Nevada.

other fossil, *Kimberella quadrata*, appears to have borne a stiff shell on its dorsal side, which covered a soft muscular “foot” surrounded by numerous gill-like folds. *Kimberella* is likely to be closely related to the Mollusca, although it may not be closely related to any extant molluscan class.

Finally, a few organisms associated with the Ediacaran biota formed hard or even biomineralized parts or skeletons. Biomineralization did not absolutely originate in the Early Cambrian. The calcified tube *Cloudina* (Fig. 4) and the agglutinated tube *Archaeichnium* may have been formed by worms. Mineralized sponge spicules have been found in Ediacaran-age rocks in Mongolia. Other mineralized fossils are more difficult to decipher and include tubular, cup-shaped, conical, and hexagonal forms. Organic-walled tubes include the threadlike sabelliditids, possibly made by annelid or pogonophoran relatives; and the tetradial *Corumbella*, which resembles the tubes made by certain scyphozoan polyps. Small, unmineralized toothlike fossils (*Redkinia*), which somewhat resemble annelid jaws, are known from Ediacaran assemblages in Russia. See ANNELIDA; ARTHROPODA; BURGESS SHALE; ECHINODERMATA.

Associated trace fossils. The marks left in soft sediments by otherwise unknown animals provide another source of knowledge of late Precambrian animal life. The figure-8-shaped trail in Fig. 5 indicates that animals capable of directed, muscular, gliding motion (like that of a garden snail) coexisted with more typical members of the Ediacaran biota. In a similar fashion, strings of fecal pellets demonstrate the existence of animals with one-way guts, and closely meandering marks imply an ability for



Fig. 5. Ediacaran trace fossil *Gordia* found on the sole of a sandstone bed in South Australia gives clear evidence for the existence of mobile animals in the late Precambrian.

systematic grazing. There is little evidence for vertical burrowing in rocks of this age.

Significance. Although many Ediacaran fossils are enigmatic, there is sound evidence that sponges, cnidarians, bilaterian worms, and possibly arthropods and other phyla were present. This implies that the Animalia originated even farther back in time. The largest Ediacaran fossils, reaching up to 3 ft (1 m) in length, are flattened “fronds” that had large surfaces compared with their volumes. Proponents of the Vendobionta hypothesis have claimed that such large organisms, lacking guts or muscles, must have been photosynthetic or chemosynthetic and probably contained symbiotic, photosynthetic microorganisms. However, such a lifestyle is also found in modern reef corals and various other marine animals. Alternatively, the high ratio of surface to volume may represent adaptations to an atmosphere and hydrosphere relatively low in oxygen. See PALEONTOLOGY.

Bruce Runnegar; Ben Waggoner

Bibliography. J. G. Gehling, Microbial mats in terminal Proterozoic siliciclastics: Ediacaran death masks, *Palaos*, 14:40–57, 1999; J. P. Grotzinger et al., Bio-stratigraphic and geochronologic constraints on early animal evolution, *Science*, 270:598–604, 1995; G. M. Narbonne, The Ediacara biota: A terminal Neoproterozoic experiment in the evolution of life, *GSA Today*, 8(2):1–6, 1998; B. Runnegar, Vendobionta or Metazoa? Developments in understanding the Ediacara “fauna,” *N. Jb. Geol. Paläontol. Abh.*, 195:303–318, 1995; A. Seilacher, Vendobionta and Psammocorallia: Lost constructions of Precambrian evolution, *J. Geol. Soc. London*, 149:607–613, 1992; R. C. Sprigg, Early Cambrian (?) jellyfishes from the Flinders Ranges, South Australia, *Trans. Roy. Soc. S. Aust.*, 71:212–224, 1947.

Eel

A fish in the order Anguilliformes, one of several orders of the superorder Elopomorpha, which share a leptocephalus (a slender, transparent larva with a long narrow head) [Fig. 1] larval development. The name eel also applies to other fishes similar in appearance to true eels, that is, having a serpentine body,

lacking pectoral fins, and usually lacking scales. See ANGUILLIFORMES.

True Eels

The Anguilliformes consist of suborders Anguilloidei, Muraenoidei, and Congroidei.

Anguilloidei. This suborder comprises three families.

Freshwater eels (family Anguillidae). These eels (Fig. 2) are found in tropical and temperate seas except for the eastern Pacific and southern Atlantic; they consist of one genus and 15 species. They are catadromous (inhabiting freshwater and migrating to spawn in salt water), living practically all of their adult life in freshwater.

Mud eels (family Heterenchelyidae). These eels inhabit the tropical Atlantic, Mediterranean, and eastern Pacific; they consist of two genera and eight species. They typically burrow head-first into the substrate.

Spaghetti eels (family Moringuidae). These eels inhabit the tropical Indo-Pacific and western Atlantic; they consist of two genera and six species. Their bodies are extremely elongate and adapted for burrowing head-first into the substrate.

Muraenoidei. This suborder comprises three families.

False morays [family Chlopsidae (Xenocoagridae)]. These eels inhabit the tropical and subtropical Atlantic, Indian, and Pacific oceans; they consist of eight genera and 18 species. Their gill openings are restricted to small rounded, laterally placed apertures; pectoral fins are absent in two of the eight genera; and lateral-line pores are restricted to the head.

Myroconger eels (family Myrocongridae). These eels inhabit the eastern tropical Atlantic and Pacific oceans; they consist of one genus and four species. Their gill

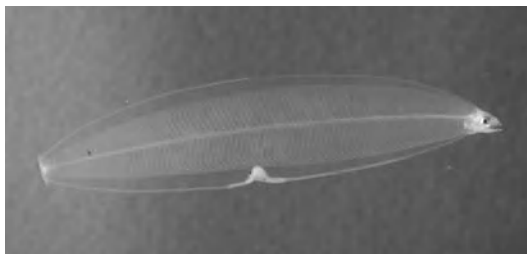


Fig. 1. *Leptocephalus* larva. (Photo by Antonio J. Ferreira, © California Academy of Sciences)



Fig. 2. *Anguilla anguilla*. (Drawing by Robbie Cada, www.FishBase.org)

openings are small but not greatly restricted; pectoral fins are present; their bodies are strongly compressed; and lateral-line pores are located in the branchial (pertaining to the gills) region above the pectoral fins.

Moray eels (family Muraenidae). These eels inhabit tropical and temperate seas of the world, with some species entering freshwater; they consist of 15 genera and about 185 species. Their gill openings are restricted to small rounded, laterally placed apertures; the pectoral fins are absent; and lateral-line pores are located on the head and branchial region. The giant moray (*Gymnotoborax javanicus*) of the Indo-Pacific attains a length of 3 m (10 ft), making it one of the largest eels.

Congroidei. This suborder comprises nine families.

Cutthroat eels (family Synbranchidae). These eels inhabit the Atlantic, Indian, and Pacific oceans; they consist of 10 genera and 32 species. The gill openings are low, at or below the insertion of the pectoral fins; hence, the name cutthroat eels.

Snake eels and worm eels (family Ophichthidae). These eels inhabit tropical and warm temperate seas of the world, with some species in freshwater; they consist of 52 genera and about 290 species. The family differs from other eels in the unusual form of the branchiostegal (pertaining to the membrane covering the gills) rays, which broadly overlap those of the other side, leaving the branchiostegal region inflated or slightly bulbous. *Ophichthus rex*, known only in the Gulf of Mexico from Florida to Texas, and *O. ophis*, from the tropical and subtropical western and eastern Atlantic, both with total lengths of about 210 cm (83 in.), are the largest ophichthids.

Shorttail eels (family Colococongridae). These eels inhabit the Atlantic, Indian, and western Pacific oceans; they consist of one genus and five species. Shorttail eels, with a blunt head and stubby body, are the least elongate of all the true eels.

Longneck eels (family Derichthyidae). These eels inhabit the Atlantic, Indian, and Pacific oceans; they consist of two genera and three species. These are mesopelagic [inhabiting the middle ocean region, from about 200 to 1000 m (650 to 3300 ft)] to bathypelagic [inhabiting the deep ocean region, from about 1000 to 4000 m (3300 to 13,000 ft)] eels, with a maximum length of about 60 cm (24 in.).

Pike congers (family Muraenesocidae). These eels inhabit the tropical Atlantic, Indian, and Pacific oceans; they consist of possibly six genera and 13 species. Pike congers have a large mouth equipped with large conspicuous teeth; their large eyes are covered with skin; and they have a conspicuous lateral line. Five species attain a total length of 200 cm (78 in.) or more, with the largest known being 250 cm (98 in.).

Snipe eels (family Nemichthyidae). These eels are mesopelagic in the tropical and temperate Atlantic, Indian, and Pacific oceans; they consist of three genera and nine species. Snipe eels are the most bizarre anguilliforms. Their bodies are extremely long and ribbonlike; one genus, *Nemichthys*, has over 750 vertebrae. The jaws are very long and narrow, and the teeth cannot be occluded (that is, the teeth do not come into contact with cusps of the opposing teeth

fitting together), except in sexually mature males, which undergo an extreme shortening of the jaws and loss of teeth. The two sexes are so different that at one time some species of snipe eels were thought to be different taxa, even at the suborder level.

Conger eels (family Congridae). These eels inhabit the tropical to temperate Atlantic, Indian, and Pacific oceans; they consist of 32 genera and 160 species. Most congrid eels are dull colored, usually tan or brown on the back and silvery on the sides and belly, with only a few having spots or stripes; the lateral line is complete; the dorsal and anal fins are well developed and continuous with the caudal fin; and the pectoral fins are usually well developed, except in garden eels, in which the pectorals are reduced to a tiny fleshy flap. Garden eels live in colonies, where they partly submerge themselves in tubelike burrows. With their head bent facing the current, they sway back and forth picking plankton as it drifts by. Congrid eels vary greatly in size—the smallest, *Bathycongrus odontostomus*, of the Indo-West Pacific, is only 4.5 cm (1.8 in.) in total length, whereas *Conger conger* of the eastern Atlantic attains a total length of 300 cm (118 in.). The longest true eel in American waters is *C. oceanicus*, with a total length of 230 cm (90 in.).

Duckbill eels (family Nettastomatidae). These eels inhabit the Atlantic, Indian, and Pacific oceans; they consist of six genera and about 38 species. They have elongate heads, large mouths, and flattened snouts—hence the name duckbill; their tails are attenuated; and pectoral fins are usually absent in adults. Their maximum length is about 100 cm (39 in.).

Sawtooth eels (family Serrivomeridae). These eels are pelagic (living in the open ocean) in the tropical and temperate Atlantic, Indian, and Pacific oceans; they consist of two genera and about 10 species. Their jaws are extremely elongate and slender, but less so than in snipe eels; their gill openings are connected ventrally; the tails are attenuated; and the vomerine teeth (that is, teeth on the vomer, which is the median unpaired bone on the roof of the mouth) are lancet-like and arranged in a sawlike row. Their maximum total length is about 78 cm (31 in.).

Other “Eels”

Other “eels” are classified in orders Saccopharyngiformes, Gymnotiformes, and Siluriformes.

Order Saccopharyngiformes. This order comprises two families.

Bobtail snipe eel (family Cyematidae). These fish are bathypelagic in the Atlantic, Indian, and Pacific oceans; they consist of two monotypic genera. Their bodies are relatively short and compressed; they lack lateral-line pores; their eyes are small to vestigial; the tip of the tail is blunt; the caudal fin is present; and maxillae are present. The maximum length is about 15 cm (6 in.).

Pelican eels (family Eurypharyngidae). These eels inhabit the tropical and temperate Atlantic, Indian, and Pacific oceans; they consist of one species. They have a huge head, a pelicanlike mouth, jaws with numerous minute teeth, and minute pectoral fins. This fish

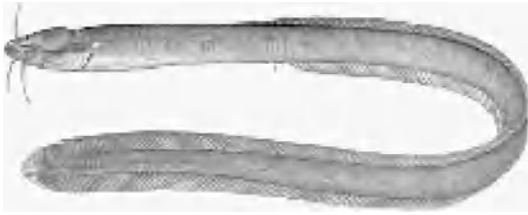


Fig. 3. Eel catfish (*Channallabes apus*). (Drawing by P. Mertens, www.FishBase.org)

is unique among teleosts in having five gill arches. The maximum length is about 74 cm (29 in.).

Order Gymnotiformes. This order comprises a single family.

Electric eel (family Gymnotidae). Another true-eel mimic is the electric eel (*Electrophorus electricus*), a member of the family Gymnotidae (nakedback knife-fishes). It occurs in the Orinoco and Amazon rivers of northern South America and is unique among the gymnotiforms in having electric organs capable of delivering discharges up to 600 volts. The electric eel is further identified by the lack of a dorsal fin, while displaying a very long anal fin terminating at the tip of the tail. It has an oral respiratory organ for breathing air and a rounded body that lacks scales, and can continuously add vertebrae throughout life. It attains a total length of 220 cm (87 in.). See ELECTRIC ORGAN (BIOLOGY).

Order Siluriformes. A catfish that could easily be mistaken for an eel inhabits the Congo River Basin of Africa. It is called the eel catfish (*Channallabes apus*) [Fig. 3] and is the only one among about 2870 species of catfishes that bears the name eel. The eel catfish, belonging to the siluriform family Clariidae (air-breathing catfishes), has reduced eyes, lacks pectoral and pelvic fins, has long dorsal and anal fins, and has the habit of burrowing in the substrate, all characteristics of many species of true eels. See SILURIFORMES.

Herbert Boschung

Bibliography. R. Froese and D. Pauly (eds.), FishBase, World Wide Web electronic publication, version 05/2005; J. S. Nelson, *Fishes of the World*, 4th ed., Wiley, New York, 2006.

Effective dose 50

This term is used chiefly to characterize the potency of a drug by the amount required to produce a response in 50% of the subjects to whom the drug is given. The term is also known as ED₅₀ or median effective dose. At one time it was usual to try to measure the effect of a drug by noting the amount which was just sufficient to produce a particular response; but when it was realized that this amount varied greatly from subject to subject, attention was turned to measuring the effect on a group of subjects. Suppose, for example, that a drug is being used to relieve a certain type of pain; then the median effective dose is of such a size that it controls the pain in 50% of the sufferers and is insufficient to control it in the remaining 50%.

The term median effective dose is most commonly applied in connection with drugs, but it may be used

when various other sources of stimuli, for example, x-rays, are under consideration. The response must be of the kind known as quantal, or all or nothing, where the investigator is simply able to report that the response either was or was not elicited; for example, convulsions did or did not occur, hemorrhage was or was not produced, pregnancy did or did not ensue, the animals did or did not survive.

The median effective dose is not a well-defined quantity until the test animal, the end response, and such factors as the route of injection of a drug and the state of nutrition of the animal are specified. In the work of an investigator who uses the appropriate controlled conditions, however, the ED₅₀ is a reproducible measurement.

Determination of the ED₅₀ of a drug requires the administration of at least two separate amounts of the drug, each one given to several subjects. Suppose that an animal physiologist wishes to measure the ED₅₀ of a hormone, estrone, that causes estrus, or heat. He has specified that he will use spayed female rats as experimental animals, and has also laid down various conditions that he regards as important in controlling the results of his investigation: the age of the animals, perhaps, and the route by which the hormone will be administered. It is unlikely that a reliable measurement could be made with fewer than 30 animals. Suppose that a total of 60 were used, 20 at each of 3 doses, with the following results:

Dose, international units	Animals showing estrus, %
2	20
3	45
4	80

Inspection of the results suggests that a dose somewhat greater than 3 units will cause estrus in 50% of the animals. Methods based on the theory of probability can be used to give a more objective analysis of the data and to provide a measure of the error of the estimate, but it is sufficient to note for the present purpose that the ED₅₀ is approximately 3. It is clear that in measuring the ED₅₀ some choices of the dosage to be used will work out better than others. Doses giving a response in a very small or very large percentage of the subjects contribute little information to the measurement of ED₅₀. The ideal is to have doses on either side of the ED₅₀, and the closer they are to ED₅₀, the better. However, in practice a pharmacologist sometimes determines an ED₅₀ as one part of a larger experiment to investigate the dose response curve of a drug, or to compare the potency of two preparations of a drug. In such a case, though he would still want to have doses on each side of the ED₅₀, he would not want them close to the ED₅₀. More information about the total curve would be given if the low dose had an effect in, say, 25% of the animals and the high dose in about 75%.
See BIOASSAY.

Colin White

Bibliography. C. Janney and J. Timpke, *Calculation of Drug Dosages*, 4th ed., 1993; R. Wiederhold, *Dosages and Calculations*, 1991.

Effector systems

Those organ systems of the animal body which mediate overt behavior. Injury to an effector system leads to loss or to subnormal execution of behavior patterns mediated by the system, conditions termed paralysis and paresis, respectively.

Overt behavior consists of either movement or secretion. Movement results from contraction of muscle. Secretion is a function of glands. Neither muscular contraction nor glandular secretion is autonomous but is regulated by an activating mechanism which may be either neural or humoral. In neurally activated systems the effector organ, whether muscle or gland, is supplied by nerve fibers originating from cell bodies situated in the central nervous system or in peripherally located aggregates of nerve cell bodies known as ganglia. The nerve fibers make intimate contact with, but are not protoplasmically continuous with, the cells of the effector organ. Activation of the effector organ occurs when the nerve cell body is excited and generates a nerve impulse, an electrochemical alteration which is conducted along the nerve fiber to its terminations on the effector organ cells. Here the nerve impulse releases from the nerve ending a chemical transmitter that generates a similar electrochemical alteration in the effector cells, and this alteration, in turn, leads to either contraction or secretion. Such systems composed of a muscle or a gland along with their regulating nerves are termed neuromuscular or neuroglandular effector systems, respectively. Examples are the skeletal muscles together with their motor nerve supplies, and the medullae of the adrenal glands along with their innervation (splanchnic nerves). In many neurally regulated effector systems (for example, skeletal muscle and adrenal medulla), function is totally dependent on intact innervation, and denervation leads to functional paralysis. In other organs (for example, the salivary glands), denervation causes only temporary paralysis. When recovery occurs, the gland may oversecrete continuously (paralytic secretion), apparently because the denervated gland cells become unusually sensitive to certain blood-borne chemical agents (denervation supersensitivity). *See* BIOPOTENTIALS AND IONIC CURRENTS; ENDOCRINE MECHANISMS; MOTOR SYSTEMS.

In other effector systems (humeromuscular and humeroglandular), the activating agent is normally a blood-borne chemical substance produced in an organ distant from the effector organ. Uterine smooth muscle is uninfluenced by the uterine nerve activity but contracts vigorously when the blood contains oxytocin, a chemical substance elaborated by the posterior lobe of the hypophysis; sensitivity to oxytocin increases progressively during pregnancy. Similarly, secretion of pancreatic juice is independent of pancreatic innervation; the regulating agents are blood-borne substances (cholecystekinin and secretin) produced by cells in the wall of the small intestine. Generally, in such humorally regulated effector systems, activation is more delayed and more

prolonged than in neurally regulated systems. *See* PANCREAS; UTERUS.

Finally, some effector systems are hybrid in the sense that both nerves and humors regulate their functions. The smooth muscles of arterioles contracts in response to either nerve stimulation or epinephrine, a substance secreted into the bloodstream by the adrenal medulla. Secretion of hydrochloric acid by the gastric mucosa is increased by activation of the vagus nerve or by the presence in the blood of histamine, a substance found in many tissues of the body. Effector systems with both neural and humoral regulation are never completely paralyzed by denervation but may be deficient in reaction patterns when the quick integrated activation provided by neural regulation is essential. For example, following extensive vascular denervation, humoral agents may maintain sufficient arteriolar constriction to sustain the blood pressure in static situations. However, the normal capacity to increase arteriolar constriction to offset the gravitational effects of rising from the prone to standing positions is permanently lost, with the result that such postural changes may lead to inadequate cerebral blood flow and consequent fainting. *See* HISTAMINE; NERVOUS SYSTEM (VERTEBRATE); PSYCHOLOGY. Theodore C. Ruch; Harry D. Patton

Efficiency

The ratio, expressed as a percentage, of the output to the input of power (energy or work per unit time). As is common in engineering, this concept is defined precisely and made measurable. Thus, a gear transmission is 97% efficient when the useful energy output is 97% of the input, the other 3% being lost as heat due to friction. A boiler is 75% efficient when its product (steam) contains 75% of the heat theoretically contained in the fuel consumed. All automobile engines have low efficiency (below 30%) because of the total energy content of fuel converted to heat; only a portion provides motive power, while a substantial amount is lost in radiator and car exhaust. *See* AUTOMOTIVE ENGINE; BOILER; POWER.

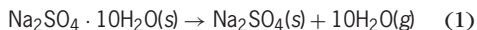
In such simple cases the value is clear. However, in some others it can be difficult to calculate exactly. For example, the efficiency of the process for converting corn to alcohol for use as automobile fuel can be computed as the ratio of the heat value of the alcohol to the heat value of all the energy used to produce it, including the fuel for the tractor to plow and harvest the cornfield and even the fuel used to create the steel and fabricate the tractor. The question is then, how much to include in the overall efficiency determination. *See* SIMPLE MACHINE. F. R. E. Crossley

Efflorescence

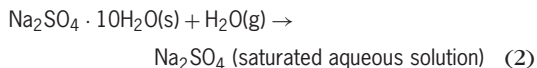
The spontaneous loss of water (as vapor) from hydrated crystalline solids. The thermodynamic requirement for efflorescence is that the partial

pressure of water vapor at the surface of the solid (its dissociation pressure) exceed the partial pressure of water vapor in the air.

A typical efflorescent substance is Glauber's salt, $\text{Na}_2\text{SO}_4 \cdot 10\text{H}_2\text{O}$. At 25°C (77°F) the dissociation pressure for the process in reaction (1) is 2586.4 pas-



cals (0.776 in. Hg), 81% of the saturation vapor pressure of pure water at this temperature. In a sufficiently humid atmosphere Glauber's salt also can deliquesce by the process shown in reaction (2).



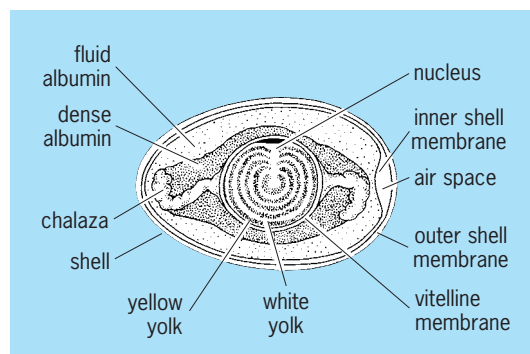
The vapor pressure of the saturated solution is 2919.27 Pa (0.876 in. Hg), 92% of the vapor pressure of pure water. Thus, Glauber's salt at 25°C (77°F) is stable in atmospheres having relative humidities of 81–92%; below 81% it effloresces; above 92% it deliquesces.

The spontaneous loss of water normally requires that the crystal structure be rearranged, and consequently, efflorescent salts usually go to microcrystalline powders when they lose their water of hydration. See DELIQUESCENCE; PHASE EQUILIBRIUM; VAPOR PRESSURE.

Robert L. Scott

Egg (fowl)

A single, large, living, female sex cell enclosed in a porous, calcareous shell through which gases may pass. Although they vary in size, shape, and color, the eggs of chickens, ducks, geese, and turkeys are essentially the same in structure and content (see *illus.*). Inward from the shell are the outer and inner shell membranes which are also permeable to gases. The membranes are constructed to prevent rapid evaporation of moisture from the egg but to allow free entry of oxygen, which is necessary for life. Air begins to penetrate the shell soon after the egg is laid, and it tends to accumulate in a space between the two membranes at the large end of the egg. See CELL (BIOLOGY).



Egg of a bird. (After L. P. Sayles, ed., *Biology of the Vertebrates*, 3d ed., Macmillan, 1949)

The inner shell membrane surrounds a mass of fluid albumin which, in turn, encloses a body of dense albumin; these two types of protoplasm constitute the so-called egg white. The central part of the egg is occupied by the yolk, which contains the vital egg nucleus and its associated parts. The yolk consists of alternating layers of yellow and white yolk. The yolk, enclosed by the vitelline membrane, is held in place by the chalaza which is anchored at each end of the egg and prevents undue mechanical disturbance. See CELL NUCLEUS; YOLK SAC.

When a sperm nucleus fuses with the egg nucleus, the process is called fertilization. Within a few hours or less, the fertilized egg begins a series of cell divisions and differentiations which result in the formation of the embryo. The embryo then undergoes further cell modification and eventually develops into the young of the species. In time the pressure created by the growth of the embryo causes the shell to rupture, and the young hatches. As the young emerges from the shell, it carries with it a part of the food and water originally in the yolk on which it can subsist for a few days. After this initial period, the young must have access to food and water. See CELL DIVISION; FERTILIZATION (ANIMAL); MITOSIS.

For a discussion of the steps involved in the development of various kinds of fertilized eggs into mature embryos see EMBRYOLOGY; OVUM; REPRODUCTION (ANIMAL); POULTRY PRODUCTION. James F. Ferry

Egg water balance. The water vapor content of the environment is usually less than that found inside eggs, so water vapor is lost continually throughout incubation. Eggshell conductance, that is, the inverse resistance of the shell to gas movement, is directly related to egg mass, and inversely related to incubation time. The relationship between eggshell conductance and daily water loss, egg mass, and incubation time predicts that all bird eggs should lose about 15% of their initial mass as water vapor during incubation and that the vapor pressure gradient across the shell is similar for all bird species and relatively constant throughout incubation. The fractional water loss is apparently obligatory since eggs prevented from losing water show increased embryonic mortality and deformity. The obligatory nature of egg water loss may be related to the observation that the embryo inflates its lung just prior to hatching with the gas in the air space of the egg. The volume of the air space is related to the quantity of water loss since the egg shell is rigid.

Respiratory gas exchange. While the daily water loss from avian eggs during incubation appears to be relatively constant, such is not the case for the daily oxygen consumption (O_2 flux) and carbon dioxide production (CO_2 flux). The metabolic activity increases throughout incubation as the embryo develops and increases in mass, so both the O_2 and CO_2 exchange must similarly increase. Since the gas conductance of the shell is constant (all gases diffuse through the same pores), the pressure of O_2 inside the egg must fall and the pressure of CO_2 must rise if the embryo is to gain O_2 and lose metabolically

produced CO₂. The egg O₂ exchange of a variety of bird species measured just prior to hatching appears to increase in proportion to egg mass in much the same way that the eggshell gas conductance (as measured with water vapor) increases with increase in egg mass. These relationships allow one to predict that the gas pressures of O₂ and CO₂ (in the air space) are the same for all bird embryos just prior to hatching and are in fact similar to values one might expect to see in adult bird lungs. It would appear that one function of the eggshell is to prepare the embryo for transition to lung breathing at birth. See RESPIRATION.

Ralph A. Ackerman

Bibliography. R. W. Burley and D. V. Vadehra, *The Avian Egg: Chemistry and Biology*, 1989; B. M. Carlson, *Patten's Foundations of Embryology*, 6th ed., 2000.

Eggplant

A warm-season vegetable (*Solanum melongena*) of Asiatic origin belonging to the plant order Polemoniales (formerly Tubiflorales). Eggplant is grown for its usually egg-shaped flesh fruit (see *illus.*), and is



Eggplant (*Solanum melongena*), cultivar Black Magic. (Joseph Harris Co., Rochester, New York)

eaten as a cooked vegetable. Cultural practices are similar to those used for tomatoes and peppers; however, eggplant is more sensitive to low temperatures. Popular purple-fruited varieties (cultivars) are Black Beauty and a number of hybrid varieties; fruits of other colors, including white, brown, yellow, and green, are used chiefly for ornamental purposes. Harvesting generally begins 70–80 days after planting. Florida and New Jersey are important eggplant-producing states. See PEPPER; SOLANALES; TOMATO.

H. John Carew

Ehrlichiosis

Ehrlichiosis is a tick-borne infection that often is asymptomatic but also can produce an illness ranging from a few mild symptoms to an overwhelming multisystem disease. Ehrlichiosis is included with those infections that are said to be emerging, either because they have been recognized only recently or because they were previously well known but now are occurring more frequently.

History. Prior to 1986, the bacteria of the genus *Ehrlichia* were of interest primarily to veterinarians. These microorganisms have a global distribution and are capable of producing a hemorrhagic disease in dogs and febrile illnesses of horses, cattle, sheep, and bison. A constant feature in all these circumstances is the predilection of the *Ehrlichia* to enter the cellular elements of the blood. The specific *Ehrlichia* species determined whether mononuclear cells, granulocytes, or platelets were invaded. Human illness associated with these organisms was first recognized in southern Japan in the 1950s. Because these illnesses seemed to be a local phenomenon, they received little attention elsewhere.

However, in 1986 a 51-year-old man, after planting some trees while vacationing in Arkansas, discovered some ticks on his neck and removed them uneventfully. Ten days later, he developed fever accompanied by a serious systemic illness that eluded specific diagnosis and required a prolonged hospitalization before his recovery. Careful microscopic examination of a slide preparation of his peripheral blood revealed distinctive inclusion bodies in his lymphocytes and monocytes. When reviewed by microbiologists and pathologists, these inclusions were identified as ehrlichiae. This illness was the first case of ehrlichiosis diagnosed in the United States. To date, over 400 cases of ehrlichiosis have been diagnosed in 30 states, and ehrlichiosis now is considered to be the most common tick-borne infection in this country.

Epidemiology. Human ehrlichiosis is caused by two distinct species: *E. chaffeensis* and an unnamed ehrlichial species. *Ehrlichia chaffeensis* infects primarily mononuclear blood cells; the disease produced by this species is referred to as human monocytic ehrlichiosis. The other ehrlichial species invades granulocytic blood cells, causing human granulocytic ehrlichiosis. The latter organism closely resembles *E. equi*, a species that infects horses.

Both of these ehrlichia species are transmitted to humans by the bite of infected ticks; however, the two species appear to have somewhat different geographic distributions. *Ehrlichia chaffeensis* occurs most commonly in the south-central and southeastern states, where it is associated primarily with the Lone Star tick (*Amblyomma americanum*); it is also transmitted by the common dog tick (*Dermacentor variabilis*). The agent of human granulocytic ehrlichiosis is found in the upper midwestern states of Wisconsin and Minnesota, as well as in several northeastern states. This agent seems to be transmitted principally by the deer tick (*Ixodes scapularis*).

Although ticks (the vector) are the mode of transmission of ehrlichial infections to humans, the ticks must acquire the ehrlichial organisms from animal sources (the reservoir hosts). The animal reservoirs for these organisms have not yet been confirmed, but white-tailed deer are a likely reservoir host of *E. chaffeensis* in the southeastern United States. See IXODIDES.

Clinical signs and pathogenesis. The forms of the disease caused by the two ehrlichial species are indistinguishable. Illness occurs most often during April–September, corresponding to the period when ticks are most active and humans are pursuing outdoor activities. Most individuals relate a history of tick exposure and have been engaged in work or recreation in wooded areas where they encountered ticks.

While obtaining its blood meal by biting an individual, the infected tick inoculates ehrlichiae into the skin. The microorganisms then spread throughout the body via the bloodstream. *Ehrlichia* must achieve an intracellular residence in order to survive and multiply. *Ehrlichia chaffeensis* seeks out tissue macrophages in the spleen, liver, and lymph nodes. It is suspected that the agent of human granulocytic ehrlichiosis has a tropism for granulocytic precursor cells in the bone marrow. Once the organism has entered the cell, it multiplies there for a period of time without producing symptoms (the incubation period). In ehrlichiosis, the incubation period can last from 1 to 3 weeks after exposure to the infected tick. Thereafter, individuals develop fever, chills, headache, and muscle pains. Gastrointestinal symptoms such as nausea, vomiting, and loss of appetite also are common. Laboratory abnormalities regularly include anemia, low white blood cell and platelet counts, and abnormal liver function. More severely ill individuals also may manifest abnormalities of the central nervous system, lungs, and kidneys. Elderly patients are more likely to become seriously ill than children and young adults. Fatality rates up to 4% have been documented. Because the clinical presentation is nonspecific, the diagnosis of ehrlichiosis may not be immediately apparent. Prolonged intervals between the onset of illness and the administration of appropriate therapy can lead to more severe disease symptoms and a greater risk of fatality. See CLINICAL MICROBIOLOGY; HEMORRHAGE; INFECTION.

Diagnosis and treatment. An important clue to the diagnosis of human granulocytic ehrlichiosis is the recognition of cytoplasmic vacuoles filled with ehrlichiae (morulae) in circulating neutrophils. Careful examination of stained smears of peripheral blood often yields such findings in human granulocytic ehrlichiosis. In contrast, similar detection of infection in circulating monocytes is a distinctly rare event in disease caused by *E. chaffeensis*. Laboratory diagnosis usually is made by detecting an increase in species-specific antibodies in serum specimens obtained during the acute and convalescent phases of the illness. However, such serologic testing is of no use in establishing the diagnosis before treatment is initiated. In an attempt to shorten this diagnostic

lag, some research laboratories are offering to test blood specimens by using polymerase chain reaction (PCR) technology. Unfortunately, this method is not yet widely available. Therefore, therapy must be initiated on clinical suspicion.

Ehrlichiosis closely resembles another tick-borne illness, Rocky Mountain spotted fever, except that the rash, characteristic of spotted fever, is usually absent or modest. Hence, ehrlichiosis has been referred to as spotless fever. Fortunately, both diseases can be treated with tetracycline antibiotics. Most individuals respond to tetracycline therapy within 48–72 h. See RICKETTSIOSES.

Prevention. The avoidance of tick bites is fundamental to preventing ehrlichiosis. Sometimes this can best be accomplished by staying out of tick-infested areas. It has been documented that the regular use of insect repellent by individuals who spend time in grassy or wooded areas offers substantial protection against tick bites. It also is advisable to wear long-sleeved shirts and trousers that can be tucked into boots. Individuals should examine themselves for ticks periodically. Ticks attached to the skin can be removed with tweezers by grasping the tick as close to the point of attachment as possible and pulling gently and steadily.

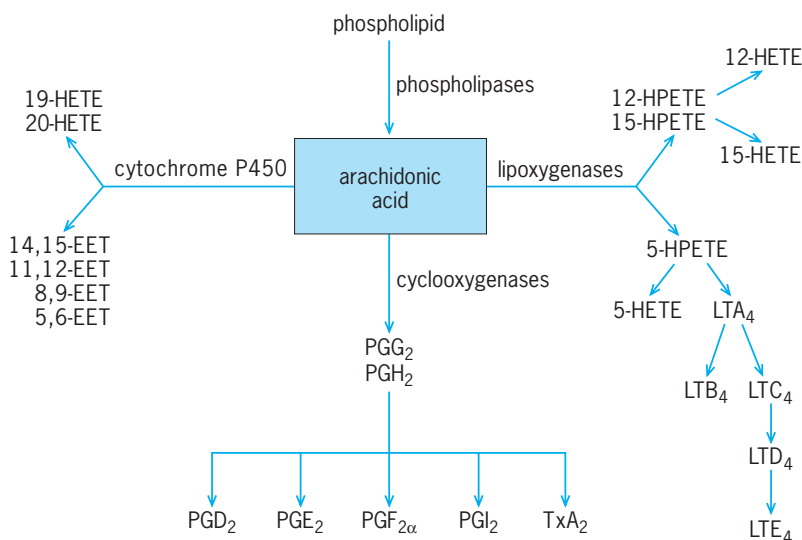
Steven M. Standaert; William Schaffner

Bibliography. J. E. Dawson et al., The interface between research and the diagnosis of an emerging tick-borne disease, human ehrlichiosis due to *Ehrlichia chaffeensis*, *Arch. Int. Med.*, 156:137–142, 1996; J. S. Dumler and J. S. Bakken, Ehrlichial diseases of humans: Emerging tick-borne infections, *Clin. Infect. Dis.*, 20:1102–1110, 1995; S. M. Standaert et al., Ehrlichiosis in a golf-oriented retirement community, *N. Engl. J. Med.*, 333:420–425, 1995; D. H. Walker and J. S. Dumler, Emergence of the ehrlichioses as human health problems, *Emerging Infect. Dis.*, 2:18–29, 1996.

Eicosanoids

A family of naturally occurring, biologically active substances derived from 20 carbon polyunsaturated fatty acids such as arachidonic acid. Members of the family include prostaglandins, thromboxanes, leukotrienes, and epoxyeicosatrienoic acids (EETs). Functioning as local hormones, they are synthesized and exert their biological actions in the same tissue. Many have physiological and pathological effects on the cardiovascular, pulmonary, reproductive, and digestive systems.

Biosynthesis and nomenclature. Eicosanoids are formed from polyunsaturated 20-carbon essential fatty acids. These fatty acids include 8,11,14-eicosatrienoic acid, 5,8,11,14-eicosatetraenoic acid (arachidonic acid), and 5,8,11,14,17-eicosapentaenoic acid. Arachidonic acid is the most abundant precursor of the eicosanoids in humans. It is esterified to phospholipids and other complex lipids in the cell membrane; thus, the cellular concentration of free arachidonic acid is low.



Enzymology of arachidonic acid metabolism. Structures and abbreviations for the compounds are given in Table 1.

Hormones such as histamine, bradykinin, and angiotensin II or physical or chemical stimuli increase the biosynthesis of eicosanoids by releasing arachidonic acid from the membrane lipids and making it available to the eicosanoid-synthesizing enzymes (see **illus.**). These stimuli activate phospholipase A₂ or phospholipase C either directly or by increasing intracellular calcium. Phospholipase A₂ hydrolyzes arachidonic acid from phospholipids, usually phosphatidylcholine or phosphatidylethanolamine. Phospholipase C converts phosphatidylinositol to a 1,2-diglyceride. The sequential action of diglyceride lipase and monoglyceride lipase on the diglyceride releases arachidonic acid. The free arachidonic acid is reesterified into cellular lipids or metabolized by cyclooxygenases, lipoxygenases, or cytochrome P450 monooxygenases. These enzymes add molecular oxygen or oxygens to the fatty acid.

Metabolites of cyclooxygenases. The cyclooxygenases, also called prostaglandin endoperoxide synthases, metabolize unsaturated fatty acids to the prostaglandins and thromboxanes. There are two forms of the enzyme. Cyclooxygenase-1 is a constitutive enzyme found in most cells, whereas cyclooxygenase-2 is an inducible enzyme that becomes elevated in response to mediators of inflammation such as cytokines. Both enzymes perform the same two enzymatic reactions: endoperoxide synthesis and peroxide reduction. The endoperoxide synthase activity oxygenates and cyclizes the arachidonic acid, producing the cyclic endoperoxide prostaglandin G₂ (PGG₂) [**Table 1**]. The peroxidase activity reduces PGG₂ to PGH₂. Both PGG₂ and PGH₂ are unstable intermediates that are further metabolized by isomerases to PGE₂ and PGD₂, by prostacyclin synthase to PGI₂, and by thromboxane synthase to TXA₂. While most cells contain cyclooxygenase, cells vary in their complement of endoperoxide-metabolizing enzymes. As a result, cells make different prostaglandins. For example, vascular endothelial cells contain prostacyclin

synthase and not thromboxane synthase. As a result, they synthesize PGI₂ but not TXA₂. In contrast, platelets contain only thromboxane synthase and thus make predominantly TXA₂. This selective synthesis regulates the prostaglandins produced and limits their sites of action. The cyclooxygenases are inhibited by nonsteroidal antiinflammatory drugs such as aspirin, indomethacin, and ibuprofen. By inhibiting cyclooxygenases, these drugs inhibit the synthesis of all of the prostaglandins and thromboxanes and block all of their biological actions. PGE₂ and PGD₂ are inactivated enzymatically in the lung and liver. TXA₂ and PGI₂ are chemically unstable and degrade nonenzymatically to the biologically inactive TXB₂ and 6-keto PGF_{1α}, respectively. See **ASPIRIN**.

The cyclooxygenase metabolites fall into categories based on the substituents on the cyclopentane ring (**Table 1**). PGG₂ and PGH₂ are cyclic endoperoxides, PGE₂ and PGD₂ are hydroxy ketones, PGI₂ has a double ring structure, and PGF_{2α} is a diol. TXA₂ has a six-membered oxane ring. These categories of prostaglandins are further divided by the number of double bonds in the side chain, which is indicated by a subscript 1, 2, or 3. The prostaglandins produced from arachidonic acid have the subscript 2.

Metabolites of lipoxygenases. Lipoxygenases are a family of enzymes that catalyze the addition of a single molecule of oxygen into the unsaturated fatty acid to form a lipid hydroperoxide (**Table 1**). They are generally constitutive enzymes that depend on the availability of free arachidonic acid. This reaction involves hydrogen abstraction from an allylic carbon adjacent to two cis-double bonds. Since arachidonic acid has four cis-double bonds, several isomeric hydroperoxyeicosatetraenoic acids (HPETEs) may be synthesized, commonly 5-, 12-, and 15-HPETEs. The lipoxygenases place the hydroperoxy group on specific carbons. Cells differ in the lipoxygenases that they contain. For example, platelets contain only 12-lipoxygenase, which catalyzes the formation of only 12-HPETE, whereas leukocytes contain both 5-lipoxygenase and 12-lipoxygenase. 5-Lipoxygenase catalyzes the formation of only 5-HPETE. Like PGG₂ and PGH₂, HPETEs are unstable intermediates that may be further metabolized. They are reduced either enzymatically or nonenzymatically to the hydroxyeicosatetraenoic acid (HETE). Of particular importance is the 5-lipoxygenase pathway which leads to the production of the leukotrienes (**Table 1**). LTA₄ synthase catalyzes the rearrangement of 5-HPETE to the unstable epoxide, LTA₄. There are two possible enzymatic fates for LTA₄: hydrolysis by LTA hydrolase to LTB₄, or conjugation with glutathione by LTC synthase giving LTC₄. The glutathione group of LTC₄ may be sequentially metabolized by specific peptidases. Removal of the glutamic acid from the glutathione group gives LTD₄, and further removal of the glycine group gives LTE₄. LTB₄, LTC₄, LTD₄, and LTE₄ are all biologically active.

Metabolites of cytochrome P450. Cytochrome P450 is a family of greater than 100 isozymes (any two or more enzymes chemically distinct but functionally

TABLE 1. Structures of eicosanoids					
Name	Abbreviation	Structure	Name	Abbreviation	Structure
Arachidonic acid (5,8,11,14-eicosatetraenoic acid)	AA		5-Hydroperoxy-eicosatetraenoic acid	5-HPETE	
Prostaglandin G ₂	PGG ₂		5-Hydroxy-eicosatetraenoic acid	5-HETE	
Prostaglandin H ₂	PGH ₂		Leukotriene A ₄	LTA ₄	
Prostaglandin E ₂	PGE ₂		Leukotriene B ₄	LTB ₄	
Prostaglandin D ₂	PGD ₂		Leukotriene C ₄	LTC ₄	
Prostaglandin F _{2α}	PGF _{2α}		Leukotriene D ₄	LTD ₄	
Thromboxane A ₂	TXA ₂		Leukotriene E ₄	LTE ₄	
Thromboxane B ₂	TXB ₂		11,12-Epoxyeicosatrienoic acid	11,12-EET	
Prostaglandin I ₂	PGI ₂		20-Hydroxyeicosatetraenoic acid	20-HETE	
6-Ketoprostaglandin F _{1α}	6-ketoPGF _{1α}		Anandamide		

the same) that metabolize endogenous compounds such as fatty acids and cholesterol and exogenous compounds such as drugs and xenobiotics. Arachidonic acid is metabolized by ω - and ω -1-hydroxylases or monooxygenases to 19-HETE and 20-HETE and by epoxygenases to EETs (Table 1). The ω -hydroxylations differ from the lipoxygenase reaction: the HETE is formed directly without the

HPETE intermediate, and the ω -hydroxylase does not require the cis-double bonds. These hydroxylations are catalyzed by CYP1A and CYP4A isozymes. The epoxygenases add an oxygen across one of the four double bonds of arachidonic acid, giving 14,15-, 11,12-, 8,9-, or 5,6-EET. Epoxygenations result from CYP2B, CYP2C, and CYP2J isozymes. These enzymes are constitutively found in cells of the liver, kidney,

pituitary gland, and blood vessels. They are regulated by the availability of free arachidonic acid.

Other pathways. The free-radical peroxidation of arachidonic acid and arachidonic acid containing phospholipids gives rise to the isoprostanes. The isoprostanes associated with phospholipids may be released by phospholipases. These eicosanoids are comparable to the prostaglandins in structure, and some also have actions similar to the prostaglandins. Since their synthesis is independent of cyclooxygenase, their formation is not inhibited by nonsteroidal anti-inflammatory drugs.

Arachidonic acid and other fatty acids are coupled to ethanolamine by an amide bond to form *N*-acyl ethanolamides. For arachidonic acid, arachidonyl ethanolamide or anandamide is produced (Table 1). The synthetic pathway for anandamide and related *N*-acyl ethanolamides is unclear. Two pathways are proposed: hydrolysis of anandamide from *N*-arachidonylphosphatidylethanolamine by phospholipase D or direct condensation of arachidonic acid with ethanolamine. Anandamide is degraded by a fatty acid amidohydrolase, giving arachidonic acid and ethanolamine. The amidohydrolase acting in reverse may represent the synthetic pathway. The brain and heart are major sources of these compounds.

Physiological and pathological actions. The actions of the eicosanoids are diverse, affecting most organ systems (Table 2). As a result, when an eicosanoid is administered systemically in pharmacological doses, many of these actions are observed. However, phys-

ologically and pathologically, eicosanoids are local hormones. They are synthesized and exert their action in a tissue, and only certain eicosanoids are made by the cells of a tissue. As a result, the action of an eicosanoid is usually discrete and related to the physiological demands of the tissue.

Mechanism of action. The prostaglandins, leukotrienes, TXA₂, and anandamide act by binding to and activating membrane receptor proteins. These receptors are of the seven-membrane spanning, serpentine receptor subgroup and are coupled to guanine nucleotide binding proteins. Distinct receptors have been identified for many eicosanoids (PGD₂, PGE₂, PGF_{2α}, PGI₂, TXA₂, LTB₄, LTC₄, LTD₄/LTE₄) and anandamide. In many cases, subgroups of receptors have been identified. The existence of receptor subgroups explains how some eicosanoids such as PGE₂ have multiple actions. The prostaglandin receptors are termed prostanoid receptors, and they are named after the primary prostaglandin that activates the receptor. For example, the PGI₂ prostanoid receptor is termed the IP receptor. The other receptors are named in an analogous manner. Receptor subgroups are indicated by numbered subscript. The receptors are coupled to stimulation of adenylyl cyclase in the case of IP, DP, and EP₂ receptors. Adenylyl cyclase is inhibited by anandamide via a cannabinoid receptor and PGE₂ via the EP₃ receptor. The receptors for the leukotrienes, TXA₂, PGF_{2α}, and PGE₂ (EP₁) are coupled to activation of phospholipase C and increases in intracellular calcium. Receptors have not been identified for 20-HETE or the EETs. Their

TABLE 2. Physiological actions of arachidonic acid metabolites

Eicosanoid	Biological effect	Physiologic or pathologic role
PGE ₂	Vasodilation Decreases gastric acid secretion Pain sensitization Bronchodilation	Protects against intense vasoconstriction Protects gastric mucosa Pain, inflammation
PGF _{2α}	Contracts gastrointestinal smooth muscle	Peristalsis, diarrhea
PGI ₂	Contracts uterine smooth muscle Contracts gastrointestinal smooth muscle Vasodilation Inhibits platelet aggregation Decreases gastric acid secretion Renin release	Parturition, dysmenorrhea Peristalsis, diarrhea Protects against intense vasoconstriction Protects against thrombosis Protects gastric mucosa Blood pressure regulation
TXA ₂	Vasoconstriction Causes platelet aggregation Contracts most smooth muscle Bronchoconstriction	Hemostasis, angina pectoris Thrombosis, angina pectoris Bronchial asthma
PGD ₂	Inhibits platelet aggregation Vasodilation	
LTC ₄ /LTD ₄	Pulmonary vasoconstriction Contracts gastrointestinal smooth muscle Constricts peripheral airways Increase mucous secretion Leakage in microcirculation	Bronchial asthma Bronchial asthma Inflammation
LTB ₄	Neutrophil and eosinophil chemotaxis Leakage in microcirculation Causes neutrophil aggregation	Inflammation Inflammation Inflammation
12-HETE/12-HPETE	Neutrophil chemotaxis	Inflammation
20-HETE	Vasoconstriction	Myogenic tone
EETs	Vasodilation	Regulates tissue blood flow
Anandamide	Releases pituitary hormones Hypothermia Antinociception Catalepsy Cerebral vasodilation	

mechanisms of action have not been clearly determined.

Cyclooxygenase metabolites. The prostaglandins and TXA_2 have several actions that are necessary for maintaining normal homeostasis and organ integrity. These effects are mediated by their constitutive synthesis through cyclooxygenase-1. PGE_2 and PGI_2 are potent vasodilators. These prostaglandins are synthesized in the blood vessel in response to vasoconstrictors such as angiotensin II and antagonize the effect of the vasoconstrictor. This action protects the organ from intense vasoconstriction and maintains nutritive blood flow. In the absence of these prostaglandins, injury to the organ occurs. These prostaglandins also protect the gastric mucosa. They inhibit gastric acid secretion and may exert other cellular protective effects. In their absence, gastric ulcers develop. TXA_2 is made principally by platelets. It causes vasoconstriction and promotes platelet aggregation. These actions prevent excessive blood loss when a blood vessel is damaged or severed. PGI_2 is made by the blood vessel wall and has the opposite effects. It inhibits platelet aggregation and prevents intravascular platelet aggregates from forming and obstructing blood flow. The balance between the proaggregatory or vasoconstrictor action of TXA_2 and the antiaggregatory or vasodilator action of PGI_2 is critical to cardiovascular homeostasis. For example, an increase in the effect of TXA_2 over PGI_2 has been implicated in ischemic heart disease. $\text{PGF}_{2\alpha}$ contracts uterine smooth muscle. With menstruation and destruction of the uterine lining, arachidonic acid is released and metabolized to $\text{PGF}_{2\alpha}$. Contraction of the uterine smooth muscle by $\text{PGF}_{2\alpha}$ contributes to the pain associated with menstruation. In pregnancy, an increase in prostaglandin synthesis is a major determinant of the onset of labor.

In inflammatory processes, immune cells migrate into the site of inflammation and release cytokines. These inflammatory mediators induce the expression of cyclooxygenase-2 in the immune cells and cells of the tissue. The increased synthesis of PGE_2 and PGI_2 mediates the vasodilation and pain sensitization associated with inflammation and arthritic conditions. Inhibitors of cyclooxygenase-1 and -2 such as aspirin and ibuprofen relieve the signs and symptoms of inflammation. However, they also block the constitutive, homeostatic functions of the prostaglandins by blocking cyclooxygenase-1, which can result in renal failure and gastric ulcers. Selective cyclooxygenase-2 inhibitors such as celecoxib are anti-inflammatory without the deleterious effects on homeostatic functions of the prostaglandins.

Lipoxygenase metabolites. The leukotrienes are synthesized by leukocytes, and they also serve as inflammatory mediators. The leukocytes are immune cells and become sequestered and activated at a site of inflammation. When activated, they release leukotrienes. LTB_4 is chemotactic and promotes leukocyte aggregation. Thus, it attracts other inflammatory cells to the site of inflammation and sequesters them. LTC_4 , LTD_4 , and LTB_4 increase capillary permeability, causing the leakage of fluid and protein from the vascula-

ture into the tissue. This contributes to the swelling associated with inflammation. LTC_4 and LTD_4 also cause bronchiolar constriction and increase mucus formation in the airways. These actions are thought to contribute to bronchial asthma. Inhibitors of 5-lipoxygenase and leukotriene receptor antagonists are useful in treating inflammation and bronchial asthma.

Cytochrome P450 metabolites. The EETs and 20-HETE are involved in the regulation of vascular tone and organ blood flow. In the blood vessel, EETs are synthesized by the endothelial cells in response to vasodilator hormones such as acetylcholine and bradykinin. The EETs diffuse to the adjacent smooth muscle cells and exert their action. They open potassium channels and hyperpolarize the smooth muscle cell membrane. This results in vasodilation. Thus, EETs serve as mediators of vasodilation. 20-HETE is made by smooth muscle cells in response to increases in intravascular pressure. It inhibits potassium channel activity and depolarizes the cell membrane. Depolarization promotes the influx of calcium and vasoconstriction. 20-HETE is a mediator of myogenic tone in blood vessels. In the blood vessel wall, 20-HETE and the EETs are counterregulatory mechanisms that control blood flow.

Other pathways. Anandamide is a naturally occurring agonist of the cannabinoid receptor. It has biological actions identical to Δ^9 -tetrahydrocannabinol, the active ingredient in marijuana, but is shorter-acting. Anandamide is synthesized and degraded by discrete regions of the brain and causes hypothermia (decreased body temperature), antinociception (decreased pain response), catalepsy (rigidity and decreased response to external stimuli), and cerebral vasodilation. It remains to be determined if anandamide functions as a chemical mediator in the brain. See ARTERIOSCLEROSIS; ARTHRITIS; ASTHMA; GOUT; MARIJUANA; PAIN; ULCER. William B. Campbell

Bibliography. C. Denzlinger, Biology and pathophysiology of leukotrienes, *Crit. Rev. Oncology/Hematology*, 23:167-223, 1996; C. J. Hillard and W. B. Campbell, Biochemistry and pharmacology of arachidonyl ethanolamide, a putative endogenous cannabinoid, *J. Lipid Res.*, 38:2382-2398, 1997; J. C. McGiff, Cytochrome P-450 metabolism of arachidonic acid, *Annu. Rev. Pharmacol.*, 31:339-369, 1991; W. L. Smith, L. J. Marnett, and D. L. DeWitt, Prostaglandin and thromboxane biosynthesis, *Pharmacol. Ther.*, 49:153-179, 1991; C. S. Williams and R. N. Dubois, Prostaglandin endoperoxide synthase: Why two isoforms?, *Amer. J. Physiol.*, 270:G393-G400, 1996.

Eigenfunction

One of the solutions of an eigenvalue equation. A parameter-dependent equation that possesses nonvanishing solutions only for particular values (eigenvalues) of the parameter is an eigenvalue equation, the associated solutions being the eigenfunctions (sometimes eigenvectors). In older usage the

terms characteristic equation and characteristic values (functions) are common. Eigenvalue equations appear in many contexts, including the solution of systems of linear algebraic equations (matrix equations), differential or partial differential equations, and integral equations. The importance of eigenfunctions and eigenvalues in applied mathematics results from the widespread applicability of linear equations as exact or approximate descriptions of physical systems. However, the most fundamental application of these concepts is in quantum mechanics where they enter into the definition and physical interpretation of the theory. Only linear eigenvalue equations will be discussed. See EIGENVALUE (QUANTUM MECHANICS); ENERGY LEVEL (QUANTUM MECHANICS); NONRELATIVISTIC QUANTUM THEORY; QUANTUM MECHANICS.

Matrix equations. An orienting example will first be presented. A common problem, for instance in the theory of angular momentum, is to find the column vectors v which are solutions of the matrix, Eq. (1), where M is an N -by- N matrix with compo-

$$M \cdot v = \lambda v \quad (1)$$

nents M_{ij} which are assumed to be known, v is a column vector of length N with components v_i , and λ is the parameter. All quantities can in general be complex.

For example, the nonvanishing solutions to Eq. (2)

$$\begin{pmatrix} -1 & 3 \\ 2 & 0 \end{pmatrix} \begin{pmatrix} v_1 \\ v_2 \end{pmatrix} = \lambda \begin{pmatrix} v_1 \\ v_2 \end{pmatrix} \quad (2)$$

may be sought. This equation can be rewritten as $M' \cdot v = \vec{0}$ for $M' = (M - \lambda 1_2)$. (Here, 1_2 is the two-dimensional unit matrix, and $\vec{0}$ is the zero vector.) The condition for such an equation to have nontrivial solutions for the vector, v , is that the determinant, $\det M'$, of the matrix M' is zero. Otherwise the inverse of M' would exist, implying that the only allowed solution is $v = (M')^{-1} \vec{0} = \vec{0}$. Thus, Eq. (3) is required.

$$\det \begin{pmatrix} -1 - \lambda & 3 \\ 2 & -\lambda \end{pmatrix} = (1 + \lambda)\lambda - 6 = 0 \quad (3)$$

The solutions are $\lambda = 2$ and $\lambda = -3$, the eigenvalues of the matrix equation. The two associated eigenvectors $v^{(1)}$ and $v^{(2)}$ are discovered by substituting these values into Eq. (2). For the eigenvalue $\lambda = 2$, the relation $v_1^{(1)} = v_2^{(1)}$ is found. This does not set the overall scale of $v^{(1)}$. Similarly, for $\lambda = -3$ the relation $v_1^{(2)} = -3v_2^{(2)}/2$ is found. The scale of the eigenvectors can be fixed by imposing the additional condition that the generalized dot product $(v^{(i)})^* \cdot v^{(i)} = 1$ for $i = 1, 2$.

The general N -by- N matrix equation can similarly be rewritten as Eq. (4). The condition for the existence of nonvanishing solutions for v is Eq. (5). This

$$(M - \lambda 1_N)v = \vec{0} \quad (4)$$

$$\det (M - \lambda 1_N) = 0 \quad (5)$$

determinant, when expanded, gives an N th-order al-

gebraic equation for λ of the form of Eq. (6). By the

$$\lambda^N + a_1 \lambda^{N-1} + a_2 \lambda^{N-2} + \dots + a_N = 0 \quad (6)$$

fundamental theorem of algebra, this equation has N solutions λ_i , $i = 1, \dots, N$, the eigenvalues of M , if λ is allowed to be complex. The solutions can be degenerate; in other words, a given value, λ , can appear more than once as a solution.

A particularly important special case is when the matrix M is self-adjoint (also known as hermitian), namely, when M is equal to its complex-conjugate transpose, as in Eq. (7). The eigenvalues of M are

$$M^\dagger \equiv M^{T*} = M \quad (7)$$

then guaranteed to be real, and the eigenvectors are orthonormal, as in Eq. (8), where $\delta^{ij} = 0$ when $i \neq j$

$$(v^{(i)})^* \cdot v^{(j)} = \delta^{ij} \quad (8)$$

and $\delta^{ij} = 1$ when $i = j$. For k degenerate eigenvalues, only a k -dimensional subspace is uniquely specified, but it is always possible to choose k such orthonormal eigenvectors, one for each eigenvalue. See ANGULAR MOMENTUM; MATRIX CALCULUS; MATRIX THEORY; RIGID-BODY DYNAMICS.

Linear operators. The general setting for linear eigenvalue equations is the notion of a linear operator \hat{L} acting in a vector space V endowed with an inner product. A complex vector space is a set of elements (vectors), $\{|u\rangle, |v\rangle, \dots\}$, that have the operations of addition and scalar multiplication defined, such that $\alpha|u\rangle + \beta|v\rangle$ is an element of V for any elements, $|u\rangle$ and $|v\rangle$, of V , and any complex numbers, α and β . Here, the notation introduced by P. Dirac in the context of quantum mechanics is employed. The operations of addition and scalar multiplication must further satisfy properties such as associativity, commutativity, and the distributive law.

A linear operator \hat{L} is a rule that takes an arbitrary vector $|u\rangle$ in V into a, generally different, vector $|\hat{L}u\rangle$ in V , and that satisfies Eq. (9) for any vectors, $|u\rangle$

$$\hat{L}(\alpha|u\rangle + \beta|v\rangle) = \alpha|\hat{L}u\rangle + \beta|\hat{L}v\rangle \quad (9)$$

and $|v\rangle$, in V , and complex numbers α and β . The usual action of a matrix on a vector by matrix multiplication satisfies this linearity property. In this general setting, eigenvalue equations again take the form $\hat{L}|v\rangle = \lambda|v\rangle$. The total set of eigenvalues of a linear operator \hat{L} acting on V is called the spectrum of \hat{L} . For differential operators this spectrum can contain continuous parts.

An inner product on the vector space V , denoted $\langle u|v\rangle$, is a complex-valued function of two arbitrary vectors, $|u\rangle$ and $|v\rangle$, in V which satisfies Eqs. (10)

$$\langle u|u\rangle \geq 0, \langle u|u\rangle = 0 \text{ only for } |u\rangle = |0\rangle \quad (10a)$$

$$\langle u|v\rangle = \langle v|u\rangle^* \quad (10b)$$

$$\langle u | \alpha v + \beta w \rangle = \alpha \langle u|v\rangle + \beta \langle u|w\rangle \quad (10c)$$

for any vectors, $|u\rangle$, $|v\rangle$, and $|w\rangle$, in V and complex

numbers α and β . Here, $|0\rangle$ is the zero vector in V . The inner product generalizes the usual dot product for three-dimensional vectors in physical space.

A basis for an N -dimensional complex vector space V is any set of N vectors $|e^{(i)}\rangle$ in V such that any vector $|u\rangle$ in V can be uniquely written as the sum in Eq. (11), where the N coefficients can be complex.

$$|u\rangle = \sum_{i=1}^N \alpha_i |e^{(i)}\rangle \quad (11)$$

For a general vector space, the number of basis elements can be (countably) infinite, with $N = \infty$, or the basis elements can even be labeled by a continuous parameter, in which case the dimension of the vector space is uncountably infinite and sums, such as in Eq. (11), need to be replaced by integrals. In both these latter cases, subtle questions of the convergence of various expressions arise.

Given any such basis and an inner product, a matrix representation L_{ij} of any abstract linear operator \hat{L} can be constructed, by defining the components as in Eq. (12). The notion of the adjoint \hat{L}^\dagger of an operator \hat{L} can be generalized by the definition in Eq. (13)

$$L_{ij} = \langle e^{(i)} | L e^{(j)} \rangle \quad (12)$$

$$\langle u | L^\dagger v \rangle = \langle u | L v \rangle^* \quad (13)$$

for any $|u\rangle$ and $|v\rangle$ in V , and the matrix adjoint is then given by $L^\dagger_{ij} = L^*_{ji}$, in accordance with Eq. (7). A self-adjoint operator satisfies the property $\hat{L}^\dagger = \hat{L}$ for all elements in the vector space V . Just as in the case of a self-adjoint matrix, the eigenvalues of a such an operator are purely real, and the eigenvectors can always be chosen to be orthonormal.

Since the eigenvectors $\{|v^{(i)}\rangle\}$ of any self-adjoint operator \hat{L} acting on V are a set of N orthonormal vectors in the N -dimensional space V , they form a basis for V and any vector $|w\rangle$ can be expanded as in Eq. (14), where the generally complex numbers

$$|w\rangle = \sum_{i=1}^N |v^{(i)}\rangle \langle v^{(i)} | w \rangle \quad (14)$$

$\langle v^{(i)} | w \rangle$ play the role of the α_i in Eq. (11). Since $|w\rangle$ is arbitrary, this equation implies the important consequence, Eq. (15), where $\hat{1}$ is the identity operator

$$\sum_{i=1}^N |v^{(i)}\rangle \langle v^{(i)} | = \hat{1} \quad (15)$$

on the space V (namely, the operator that takes any vector to itself). Equation (15) is known as the completeness relation for the eigenvectors of \hat{L} .

The relation (14), or equivalently (15), implies a very useful representation for any self-adjoint operator \hat{M} . If $|w\rangle$ is expanded in the set of eigenvectors $|v^{(i)}\rangle$ associated with \hat{M} , then, using the definition $\hat{M}|v^{(i)}\rangle = \lambda_i|v^{(i)}\rangle$, along with the fact that $|w\rangle$ is arbitrary, then gives Eq. (16). This is a simple,

$$\hat{M} = \sum_{i=1}^N \lambda_i |v^{(i)}\rangle \langle v^{(i)} | \quad (16)$$

finite-dimensional example of the spectral theorem, stating that the action of an operator may be represented purely in terms of its eigenvalue spectrum and eigenvectors. For matrices, this result says that any self-adjoint matrix can be put into a purely diagonal form by choosing the eigenvalues as a basis, the diagonal entries being the eigenvalues of the matrix. See HILBERT SPACE; LINEAR ALGEBRA; OPERATOR THEORY.

Differential eigenvalue equations. Eigenvalue problems also appear very frequently in the solution of linear differential equations. An example is the motion of a uniform vibrating string with fixed ends at $x = 0, L$. The partial differential equation that describes the displacement $y(x, t)$ of the string from its equilibrium position (taken to be $y = 0$) is Eq. (17).

$$\frac{\partial^2 y(x, t)}{\partial t^2} - c^2 \frac{\partial^2 y(x, t)}{\partial x^2} = 0 \quad (17)$$

where c is the wave velocity on the string. Along with this equation, appropriate initial and boundary conditions must be specified. The fixed end-point condition implies $y(0, t) = y(L, t) = 0$ for all t . A typical initial condition is that the string is held stationary with some particular displacement and then released.

The solution of the problem is greatly facilitated by the principle of superposition, which states that, for two solutions, $y^{(1)}(x, t)$ and $y^{(2)}(x, t)$, of a linear equation, the sum in expression (18) is also a solution

$$\alpha y^{(1)}(x, t) + \beta y^{(2)}(x, t) \quad (18)$$

for arbitrary constants α and β . The general strategy to solve such a problem is therefore to decompose $y(x, t)$ into sum of simple motions, the eigenfunctions of the problem, whose time development can straightforwardly be solved for, and then form a suitable linear combination such that the initial conditions are satisfied.

Explicitly, elementary solutions of the form $y(x, t) = b(x) \cos \omega(t - t_0)$ are studied. Substituting this into Eq. (17) results in Eq. (19). This is of the

$$\frac{d^2 b(x)}{dx^2} = -\frac{\omega^2}{c^2} b(x) \quad (19)$$

form of an eigenvalue equation with ω^2 as the parameter. The set of (differentiable) functions $\{f(x)\}$ on the interval $[0, L]$ which vanish at the end points forms a vector space V (of infinite dimension), and the differential operator d^2/dx^2 is a linear operator acting on V . An inner product on this space can be defined by Eq. (20), the adjoint operator $(d^2/dx^2)^\dagger$

$$\langle f | g \rangle \equiv \int_0^L dx f(x)^* g(x) \quad (20)$$

can be defined with respect to this inner product by applying Eq. (13), and d^2/dx^2 can be shown to be a self-adjoint operator. The eigenvalues of this operator must therefore be real, and the associated eigenfunctions satisfy orthonormality and completeness relations.

The general solution to Eq. (19) is Eq. (21), with

$$b(x) = A \cos(\omega x/c) + B \sin(\omega x/c) \quad (21)$$

A, B arbitrary constants. Imposing the boundary conditions at $x = 0, L$ gives $A = 0$ and $B \sin(\omega L/c) = 0$. A nontrivial solution ($B \neq 0$) therefore fixes the frequency ω to be one of the discrete set given by Eq. (22), the eigenvalues of the problem. The associ-

$$\omega_n = \frac{\pi n c}{L} \quad \text{for } n = 1, 2, \dots \quad (22)$$

ated eigenfunctions are $b_n(x) = \sin(\omega_n x/c)$. In the absence of fixed end-point boundary conditions (that is, for a string of infinite length) the possible values of ω would not be restricted, and the spectrum of eigenvalues would be continuous.

A linear combination, Eq. (23), which, by the prin-

$$y(x, t) = \sum_{n=1}^{\infty} a_n \sin(\omega_n x/c) \cos \omega(t - t_0) \quad (23)$$

ciple of superposition, is still a solution to the original equation (17), is now constructed in such a way that the initial conditions are satisfied. The analog of the completeness relations, Eqs. (11) and (15), appropriate for this case imply that any continuous function $y_0(x)$ representing the initial displacement [satisfying $y_0(0) = y_0(L) = 0$] on the interval $[0, L]$ can be expanded in such a series (Fourier's theorem). The orthogonality property then leads to the formula (24) for the coefficients. Equation (24), together

$$a_n = \frac{2}{L} \int_0^L dx \sin(\pi n x/L) y_0(x) \quad (24)$$

with Eqs. (22) and (23), constitutes a full solution of the problem. See DIFFERENTIAL EQUATION; FOURIER SERIES AND TRANSFORMS; INTEGRAL TRANSFORM; ORTHOGONAL POLYNOMIALS.

John March-Russell

Bibliography. G. B. Arfken and H.-J. Weber (eds.), *Mathematical Methods for Physicists*, 4th ed., 1995; R. Geroch, *Mathematical Physics*, 1985; E. Merzbacher, *Quantum Mechanics*, 3d ed., 1997; R. Shankar, *Principles of Quantum Mechanics*, 2d ed., 1994.

Eigenvalue (quantum mechanics)

One of the values of the parameter in an eigenvalue equation for which the equation has a solution. A linear eigenvalue equation is a parameter-dependent equation of the form (1) that possesses nonvanish-

$$\hat{L}|v\rangle = \lambda|v\rangle \quad (1)$$

ing solutions, $|v^{(n)}\rangle$, only for particular values, λ_n , of the parameter. These are the eigenvalues, the associated solutions being called the eigenfunctions (sometimes eigenvectors). Here \hat{L} is a linear operator acting in a generalized vector space, V , with an inner product. The concepts of eigenvalue and eigenfunction are of fundamental importance in quantum mechanics. The relation between the mathematical formulation of the theory and the possible results

of a physical measurement is couched in terms of the eigenvalues and eigenvectors of a privileged set, $\{\hat{L}\}$, of linear operators which represent all possible physical observables. See EIGENFUNCTION.

Wave functions and operators. The mathematical formulation and physical interpretation of quantum mechanics is encapsulated in a number of fundamental postulates. The first is that all possible information regarding the physical state of any system at a given time is encoded in a generalized vector of a vector space, V . In most quantum-mechanics applications these vectors are complex-valued functions, $\psi(x)$, of the coordinates, x , of the system, termed wave functions. Only this case is considered in the following. Physical observables, L , are represented by self-adjoint linear operators acting on this function space. For instance, the momentum of a particle moving in one dimension is represented by the differential operator in Eq. (2), where \hbar is Planck's constant over 2π .

$$\hat{p} = -i\hbar \frac{d}{dx} \quad (2)$$

Furthermore, any self-adjoint linear operator, \hat{L} , has a set of eigenfunctions, $\phi^{(n)}(x)$ (for simplicity, assumed indexed by a discrete label n), in terms of which any wave function, $\psi(x)$, in V may be expanded, as in Eq. (3). Here the coefficients a_n are

$$\psi(x) = \sum_n a_n \phi^{(n)}(x) \quad (3)$$

complex numbers. The sense of Eq. (3) is that both sides are equivalent, at least as far as all physical measurements are concerned. This is known as the expansion postulate.

Measurement of observables. A particularly important aspect of the theory is that the only possible result of a physical measurement of an observable, L , is one of the eigenvalues, λ_n , of the operator \hat{L} corresponding to L , as in Eq. (4). Since the result of a

$$\hat{L}\phi^{(n)}(x) = \lambda_n \phi^{(n)}(x) \quad (4)$$

physical measurement is necessarily a real number, the set of operators allowed as representatives of observables must be restricted by the requirement that their eigenvalues are purely real. Self-adjoint operators have this property, and this is the reason why representatives of observables are restricted to this class. In addition, the probability of finding the result λ_n when the system is initially in an arbitrary state $\psi(x)$ is given by Eq. (5), where the brackets rep-

$$\text{Prob}(\lambda_n) = |\langle \phi^{(n)} | \psi \rangle|^2 = |a_n|^2 \quad (5)$$

resent the inner product on the space V , given by Eq. (6). After measuring a particular value, λ_n , the

$$\langle \phi^{(n)} | \psi \rangle = \int dx (\phi^{(n)}(x))^* \psi(x) \quad (6)$$

new state of the system is given purely by $\psi(x) = \phi^{(n)}(x)$. Only when the wave function of the system is a single eigenfunction, $\phi^{(n)}(x)$, of some observable,

\hat{L} , can the system be said to definitely have the property λ_n .

Eigenvalue spectra. Unlike the case of classical mechanics where dynamical quantities such as energy or angular momentum can take on continuously varying values, the spectrum of eigenvalues is very often discrete, or contains a discrete portion. The values of the associated dynamical quantities are then said to be quantized. For instance, the energy of an electron bound in an atom can be only one of a discrete set of energy eigenvalues. *See* ENERGY LEVEL (QUANTUM MECHANICS).

Interpretation of superposition. Although Eq. (3) used the fact that the superposition principle applies to quantum mechanics, its interpretation is quite different from that in classical mechanics. For example, if two functions describing the displacement of an elastic string are superposed, $a_1y_1(x) + a_2y_2(x)$, then this new string configuration has, in general, quite different physical properties from a string with either displacement $y_1(x)$ or $y_2(x)$. However, in quantum mechanics the superposition $a_1\phi^{(1)}(x) + a_2\phi^{(2)}(x)$ of two eigenfunctions has the interpretation that sometimes the system has exactly the physical properties of $\phi^{(1)}(x)$, other times that of $\phi^{(2)}(x)$, with probabilities $|a_1|^2$ and $|a_2|^2$ respectively.

Commutation of operators. If the energy operator (the hamiltonian), \hat{H} , of the system commutes with an operator \hat{L} , that is, $\hat{H}\hat{L} = \hat{L}\hat{H}$, then a wave function that is an eigenfunction $\phi^{(n)}(x)$ of \hat{L} stays so as the system dynamically evolves. Thus, the associated eigenvalue, λ_n , is a constant of the motion, and can be thought of as labeling the state of the system. The term quantum number is then used for λ_n . Further, it is possible to find a wave function which is simultaneously the eigenvector of two different operators if and only if they commute. A fundamental property of quantum mechanics is that the maximal number of attributes that can be simultaneously, and precisely, known about any system is given by the eigenvalues of the largest possible commuting set of operators that also includes the hamiltonian \hat{H} , in other words, the maximal set of quantum numbers. For example, the operators representing momentum and position do not commute, and it is not possible to know both simultaneously with arbitrary accuracy.

Although these postulates imply a picture of the physical world quite different from that in classical physics, the consequences of these postulates, together with the rest of quantum mechanics, have been verified with extraordinary precision. *See* NONRELATIVISTIC QUANTUM THEORY; QUANTUM MECHANICS.

John March-Russell

Bibliography. S. Gasiorowicz, *Quantum Physics*, 2d ed., 1995; E. Merzbacher, *Quantum Mechanics*, 3d ed., 1997; R. Shankar, *Principles of Quantum Mechanics*, 2d ed., 1994.

Einsteinium

A chemical element, Es, atomic number 99, a member of the actinide series in the periodic table. It is not found in nature but is produced by artificial nu-

1																	18
H	2															He	
3	4															10	
Li	Be															Ne	
11	12															18	
Na	Mg	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	Ar
19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36
K	Ca	Sc	Ti	V	Cr	Mn	Fe	Co	Ni	Cu	Zn	Ga	Ge	As	Se	Br	Kr
37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52	53	54
Rb	Sr	Y	Zr	Nb	Mo	Tc	Ru	Rh	Pd	Ag	Cd	In	Sn	Sb	Te	I	Xe
55	56	71	72	73	74	75	76	77	78	79	80	81	82	83	84	85	86
Cs	Ba	Lu	Hf	Ta	W	Re	Os	Ir	Pt	Au	Hg	Tl	Pb	Bi	Po	At	Rn
87	88	103	104	105	106	107	108	109	110	111	112	113					
Fr	Ra	Lr	Rf	Db	Sg	Bh	Hs	Mt	Ds	Rg							

lanthanide series	57	58	59	60	61	62	63	64	65	66	67	68	69	70
	La	Ce	Pr	Nd	Pm	Sm	Eu	Gd	Tb	Dy	Ho	Er	Tm	Yb

actinide series	89	90	91	92	93	94	95	96	97	98	99	100	101	102
	Ac	Th	Pa	U	Np	Pu	Am	Cm	Bk	Cf	Es	Fm	Md	No

clear transmutation of lighter elements. All isotopes of einsteinium are radioactive, decaying with half-lives ranging from a few seconds to about 1 year. *See* ACTINIDE ELEMENTS; PERIODIC TABLE; RADIOACTIVITY.

Einsteinium is the heaviest actinide element to be isolated in weighable form. The metal is chemically reactive, is quite volatile, and melts at 860°C (1580°F); one crystal structure is known. *See* TRANSURANIUM ELEMENTS.

Stanley G. Thompson; Glenn T. Seaborg

Bibliography. S. Hofmann, *On Beyond Uranium: Journey to the End of the Periodic Table*, 2002; J. J. Katz and G. T. Seaborg, *The Chemistry of the Actinide Elements*, 2d ed., 1986; L. Morss and J. Fuger (eds.), *Transuranium Elements: A Half Century*, 1992; G. T. Seaborg, *Modern Alchemy: Selected Papers of Glenn T. Seaborg*, 1994; G. T. Seaborg and W. D. Loveland, *The Elements Beyond Uranium*, 1990.

El Niño

In general, an influx of warm water to the central and eastern equatorial Pacific Ocean off the coast of Peru and Ecuador, with a return period of 4–7 years. El Niño events come in various strengths: weak, moderate, strong, very strong, and extraordinary. The size of an event can be determined using various criteria: the amount of warming of sea-surface temperatures in the central and eastern Pacific from their average condition; the areal extent of that warm water anomaly; and the length of time that the warm water lingers before being replaced by colder-than-average sea-surface temperatures in this tropical Pacific region.

Under normal conditions the winds blow up the west coast of South America and then near the Equator turn westward toward Asia. The surface water is piled up in the western Pacific, and the sea level there is several tens of centimeters above average, while the sea level in the eastern Pacific is below average. As the water is pushed toward the west, cold water from the deeper part of the ocean along the Peruvian coast wells up to the surface to replace it. This cold water is rich with nutrients, making the coastal upwelling region along western South America among the most productive fisheries in the world.

Every 4–7 years those winds tend to die down and sometimes reverse, allowing the warm surface waters that piled up in the west to move back toward the eastern part of the Pacific Basin. With reduced westward winds, the surface water also heats up. Sea level drops in the western Pacific and increases in the eastern part of the basin. An El Niño condition can last for 12–18 months or longer, before the westward-flowing winds start to pick up again.

Limiting the impact. El Niño causes shifts in regional and local climate conditions around the globe. For many locations, forecasting El Niño provides useful information about how local weather conditions might be affected during the lifetime of the phenomenon. This lead time enables societies to prepare for the impacts of these extreme events.

The 1997–1998 and 1982–1983 events are considered the two biggest El Niños in the past 100 years. Both surprised the scientific community, but in different ways. The 1982–1983 event was the most damaging up to that time; no one had forecast it. The 1997–1998 event was larger and developed earlier

and faster than expected. While El Niño's onset was not forecast, many of its worldwide impacts were, enabling societies to prepare for the onslaught. The 1982–1983 event emphasized to scientists and policy makers the importance of understanding and forecasting El Niño. The 1997–1998 event exemplified how many sectors of society (for example, agriculture, mining, construction, fisheries, hydroelectric power, and health) use El Niño forecasts to improve their decision-making processes.

El Niño is considered the second biggest climate-related influence on human activities, after the natural flow of the seasons. Although the phenomenon is at least thousands of years old, its impacts on global climate have only recently been recognized. Due to improved scientific understanding and forecasting of El Niño's interannual process, societies can prepare for and reduce its impacts considerably.

Michael H. Glantz

Tropical Pacific Ocean. The tropical Pacific Ocean is 15,000 km (9300 mi) wide. Over the very warm water at its western end, there is heavy rainfall, low atmospheric surface pressure, and large-scale atmospheric ascent. At the eastern end the water is usually cool, there is little rainfall, atmospheric surface pressure is high, and the air is descending (Fig. 1a). This east-west asymmetry of sea-surface temperature in the Pacific is quite remarkable, for the Earth rotating about its axis sees the Sun with no east-west asymmetry.

The proximate cause of this east-west difference in the sea-surface temperature is the direction of the surface wind, which in the tropics normally blows from east to west. If there were no wind, the sea-surface temperature would be uniformly warm in the east-west direction, and it would rapidly drop at deeper levels in the ocean. The region of rapid temperature change with depth is called the thermocline (Fig. 1a), and in the absence of winds would remain flat from east to west.

The thermocline can be thought of as marking the depth of the cold water. When the thermocline is close to the surface, cold water is near the surface; when the thermocline is deep, cold water is far from the surface. The normal westward wind tilts the thermocline, and deepens it in the west and shallows it in the east (Fig. 1a). The westward wind tends to bring cold water to the surface in the east (where the cold water is near the surface), and to bring warm water to the surface in the west (where the cold water is far from the surface). The west-to-east decrease of sea-surface temperature is therefore due to the cold water upwelled in the east by the westward winds. In the absence of westward winds, the temperature of the tropical Pacific would be uniform from east to west and would have about the same temperature that the western Pacific has now.

Superimposed on this east-to-west temperature difference is an annual cycle that has east Pacific sea-surface temperatures coolest in August and warmest in April, in contrast to midlatitudes, which have January and July as the extreme months. In general, the warm pool in the west moves northwest

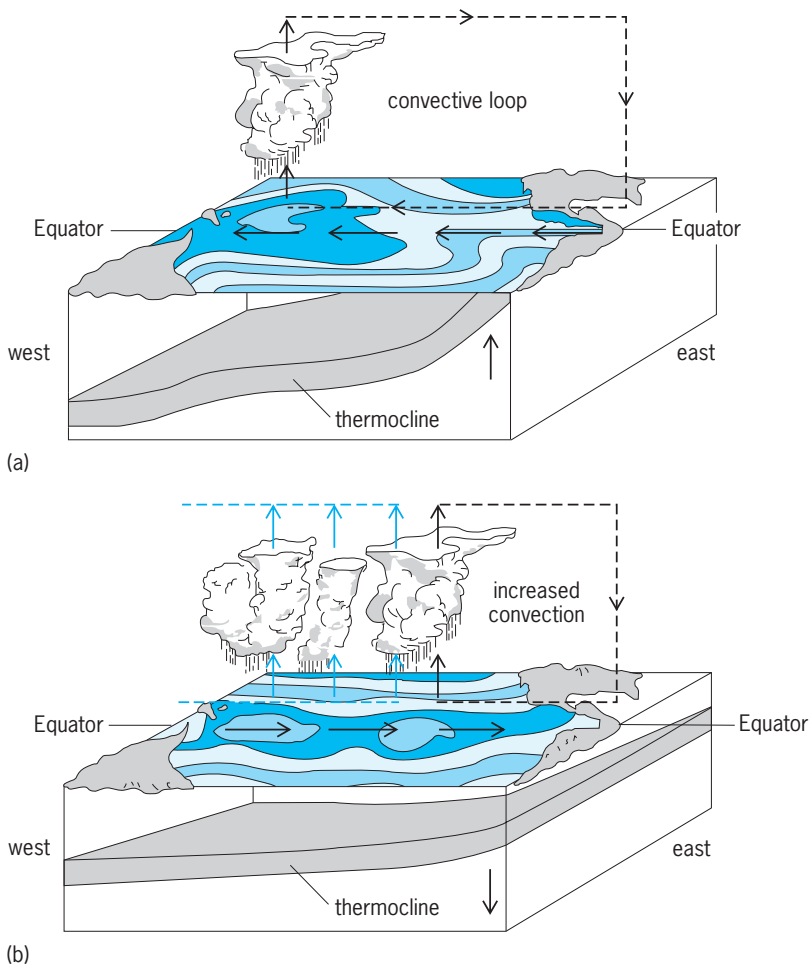


Fig. 1. Conditions of the tropical atmosphere–ocean system over the equatorial Pacific. (a) Normal, in which the thermocline is highly tilted and approaches the surface in the east Pacific, the eastern sea-surface temperature is cool, the western sea-surface temperature is warm, and rainfall remains in the far west. (b) El Niño (warm phase of ENSO) in which the thermocline flattens, the east and central Pacific warms, and the rainfall moves into the central Pacific.

during northern summer and southeast during northern winter. Throughout this complicated annual cycle, the basic pattern is maintained, with warm water in the west and its concomitant heavy rainfall, low pressure, and upward atmospheric motion. See PACIFIC OCEAN; SEAWATER; TROPICAL METEOROLOGY; UPWELLING.

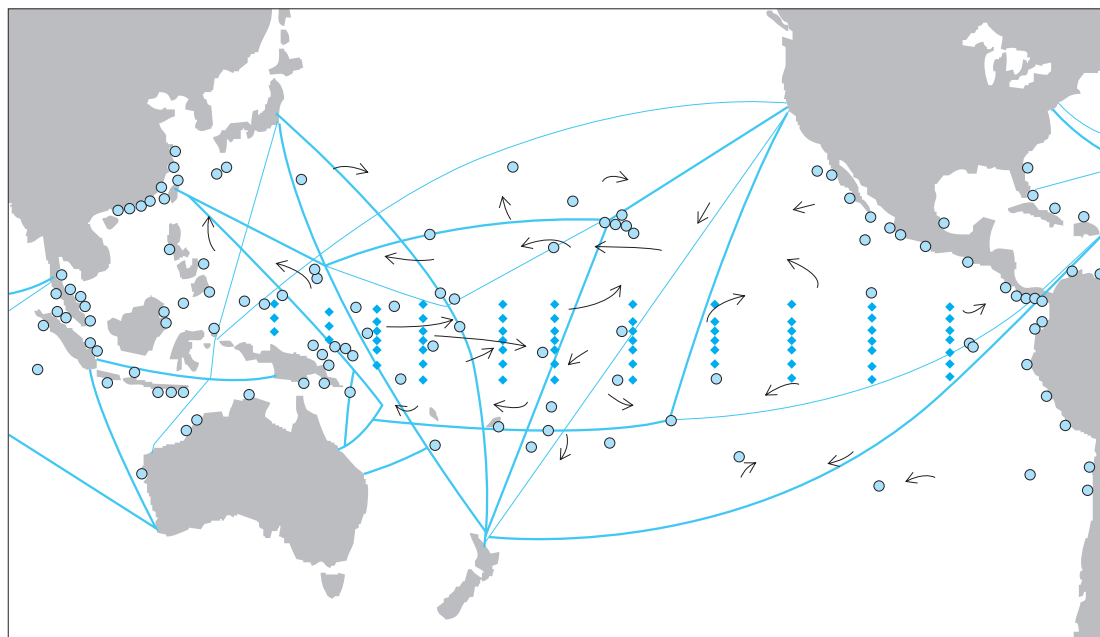
Temperature, pressure, and precipitation cycles.

Every 4–7 years, the temperature pattern of the tropical Pacific changes dramatically. It is as if all the processes that normally maintain the climatic conditions no longer occur and the tropical Pacific becomes uniformly warm, which would be expected in the absence of winds. This sporadic failure of the normal climate and its transition to a warmer Pacific has acquired the name El Niño (Fig. 1*b*). Occasionally the opposite also occurs, where the eastern Pacific becomes cooler than normal, rainfall decreases still more, atmospheric surface pressure increases, and the westward winds become stronger. This irregular cyclic swing of warm and cold phases in the tropical Pacific is called the El Niño Southern Oscillation (ENSO). The Southern Oscillation is the oscillation of the difference of surface pressure between the east and the west as the sea-surface temperature changes. When the eastern tropical Pacific is warm (El Niño or warm-phase conditions), the pressure in the east decreases and the pressure in the west increases, so the anomalous pressure difference (the Southern Oscillation Index) is negative. Precipitation actively participates in the cycle since regions of persistent precipitation lie over the warmest waters.

See CLIMATOLOGY; MARITIME METEOROLOGY; PRECIPITATION (METEOROLOGY).

Observing ENSO. Prior to 1979, observations of equatorial Pacific climate were limited to surface measurements from islands and merchant ships; beginning in 1979, satellite estimates of surface and atmospheric-column quantities became available. Aside from a small number of research cruises, the interior of the ocean remained a mystery. As a result of the international Tropical Ocean Global Atmosphere (TOGA) program (1985–1995), an observing system specifically designed to reveal and monitor those features then believed most relevant to ENSO was deployed, and is maintained to this day (Fig. 2).

The centerpiece of the system is a set of 70 buoys moored to the bottom of the ocean by anchors (usually discarded railroad wheels) and connected to a surface torus by 4–5 km (2.5–3 mi) of line. The torus contains a suite of instruments that measure near-surface quantities (winds, humidity, and sea-surface temperature). Hung on the line beneath the surface and extending downward to about 500 m (1600 ft) is a set of instruments that measure temperature and pressure (with some measuring ocean currents) and relay the information to computers on the torus. Once a day the information is relayed to satellites and then becomes freely available by the Global Telecommunication System, whose users are the weather services of the world, and by the Internet to other users and to the public. Additional observations by ships, surface buoys drifting with the currents, and tide



Key:

- ◆ moored buoys
- ◀ drifting buoys
- tide gauge stations
- volunteer observing ships

Fig. 2. ENSO Observing System. Diamonds are the moored buoys, arrows are drifting buoys, points are tide gauges, and lines are ships.

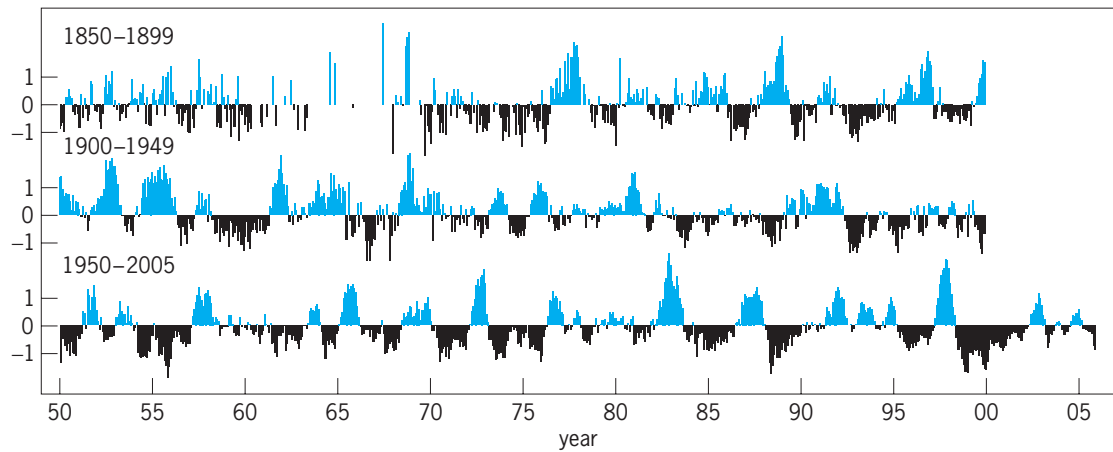


Fig. 3. Time series of an index of ENSO consisting of the monthly average temperature anomaly in the equatorial Pacific from 6°S to 6°N, 90°W to 180°W. The global mean sea-surface temperature deviation has been subtracted for each month to remove the effect of the gradual warming of the Earth's temperature, which is approximately 0.8°C (1.44°F) over the record. There are no observations for many of the months in the earlier periods, especially the 1860s, which accounts for their more ragged appearance.

gauges at islands complete the system. The ENSO observing system has been in full operation since 1995 and provides an unprecedented view of the evolution of ENSO. In particular, it was available to observe what is arguably the largest warm phase of ENSO in the last 150 years, the event of 1997–1998. *See INSTRUMENTED BUOYS; SATELLITE METEOROLOGY.*

History of ENSO. Instrumental records of surface atmospheric pressure and sea-surface temperatures document ENSO variability back to the 1870s and, with less confidence, to the 1850s. The monthly average central and eastern equatorial Pacific sea-surface temperature exhibits persistent, large-magnitude (>1°C or 1.8°F) deviations from the annual cycle (**Fig. 3**) that are the temporal signature of ENSO. The 1982–1983 and 1997–1998 warm phases of ENSO were characterized by temperature deviations in excess of 2°C (3.6°F). It also can be seen that the frequency of the warm and cold episodes varies throughout the record and, in particular, from decade to decade. The years since the 1997–1998 warm episodes have been, on average, colder than the long-term mean.

ENSO-related precipitation and temperature variations produce changes and hence proxy records of ENSO variability in many geophysical and biological systems. Occasional measurements of South American glacier masses corroborate some of the major ENSO episodes of the past 100 years. For the period back to 1500, records of sailing transit times, newspaper and literary accounts of rainfall and flooding, descriptions of mortality of marine organisms and birds, and other sources have been used to chronicle large-magnitude ENSO warm episodes. Centennial to millennial ENSO variability can be inferred from lake deposit and tree growth bands, and coral, shellfish, and glacier isotope concentrations.

The latter proxy records have been used to document ENSO back to 130,000 years ago (hereafter 130 Ka), and include epochs in which the solar forcing was different from today. ENSO variability ap-

pears to have been diminished during 6 Ka, a period when precession of the Earth's orbit enhanced the boreal summer equatorial solar forcing (output) and diminished the boreal winter solar forcing. ENSO episodes tend to develop during the summer, and the enhanced forcing is thought to have inhibited ENSO occurrences by differentially warming the western equatorial Pacific, which leads to increased precipitation in the west, and westward surface winds, stronger upwelling, and colder sea surface temperatures in the east. The global glacier extent in 6 Ka was not too different from today, at least in comparison with the epochs centered on 20 Ka (last glacial maxima) and 140 Ka, during which the ice volumes were five times the current values. The coral record suggests that ENSO variability was diminished relative to today during these icy epochs. Changes in the eccentricity of the Earth's orbit significantly changed the solar forcing during these periods, and the mechanisms by which these slow insolation changes influence ENSO are a subject of current research. *See CLIMATE HISTORY; INSOLATION.*

Teleconnections. The core of the ENSO is in the tropical Pacific, but the phenomenon influences the climate throughout the tropics and the extratropical latitudes of the winter hemisphere. Such planetary-scale atmospheric patterns are referred to as teleconnections and, in the case of ENSO, they are due to changes in the distribution of tropical Pacific precipitation.

Typical teleconnections during the boreal winter of a warm ENSO cycle are shown in **Fig. 4**. Within the tropics, warm ENSO is associated with enhanced precipitation in the central and eastern equatorial Pacific, torrential precipitation in the northern coastal desert of Peru, and diminished precipitation (and sometimes droughts) over Indonesia, northern Australia, northeast Brazil, and southern Africa. The expansion of the region of persistent precipitation into the tropical Pacific also leads to warming of the entire tropical troposphere by about

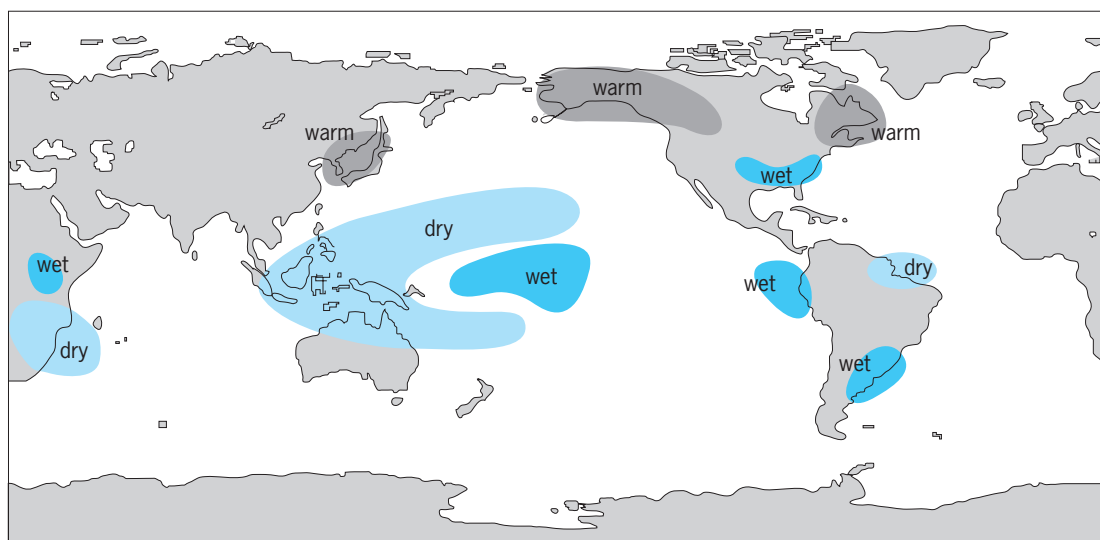


Fig. 4. Teleconnection pattern of warm phases of ENSO to the globe.

1°C (1.8°F). Precipitation deviations of opposite sign are observed during cold ENSO cycles, with the exception of the Peru coast, where it only rains only during warm ENSO. Enhanced atmospheric subsidence during warm ENSO warms the tropical Indian Ocean (Fig. 5b, upper panel) and inhibits the normal February through May rains in northeast Brazil as well as the precipitation that maintains the equatorial Andean glaciers. Warm ENSO is associated with more hurricanes in the eastern Pacific and less in the Caribbean/Atlantic, and vice versa during cold ENSO cycles.

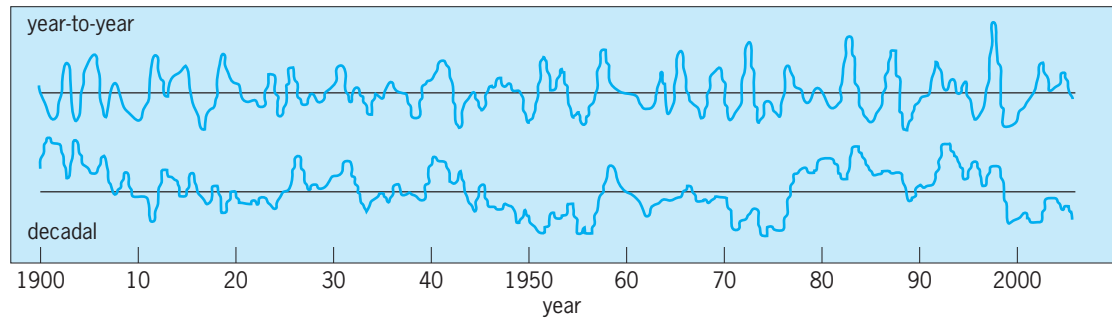
ENSO influences the higher latitudes primarily through changes in the preferred paths of storms, called storm tracks, in the winter hemisphere. Winter storms (weather) and their preferred locations are associated with the mean position of the jet streams, and these storms move heat north and south and produce precipitation in the middle latitudes. The rearrangement of storm tracks by ENSO produces deviations in the seasonal mean climate and in the distribution of extreme daily weather events over North America (Fig. 4), with diminished storminess over Alaska and Canada, and enhanced precipitation over the southern United States. Similar relationships exist in the Southern Hemisphere during austral winter with, for example, a southward shift of precipitation from the Andes' Altiplano at 15°S to the Chilean coast during ENSO warm episodes.

The year-to-year, decadal, and secular variability of the tropical Pacific climate all influence the teleconnections and climate impacts on the Northern Hemisphere. Tropical Pacific climate variability on both the year-to-year and longer time scales influence surface temperature and precipitation over the contiguous United States in the regions shown in Fig. 4. The extratropical temperature and precipitation deviations are different for warm and cold ENSO phases and reflect a nonlinear relationship between the storm tracks and the planetary-scale atmospheric circulation in the extratropics. Although the telecon-

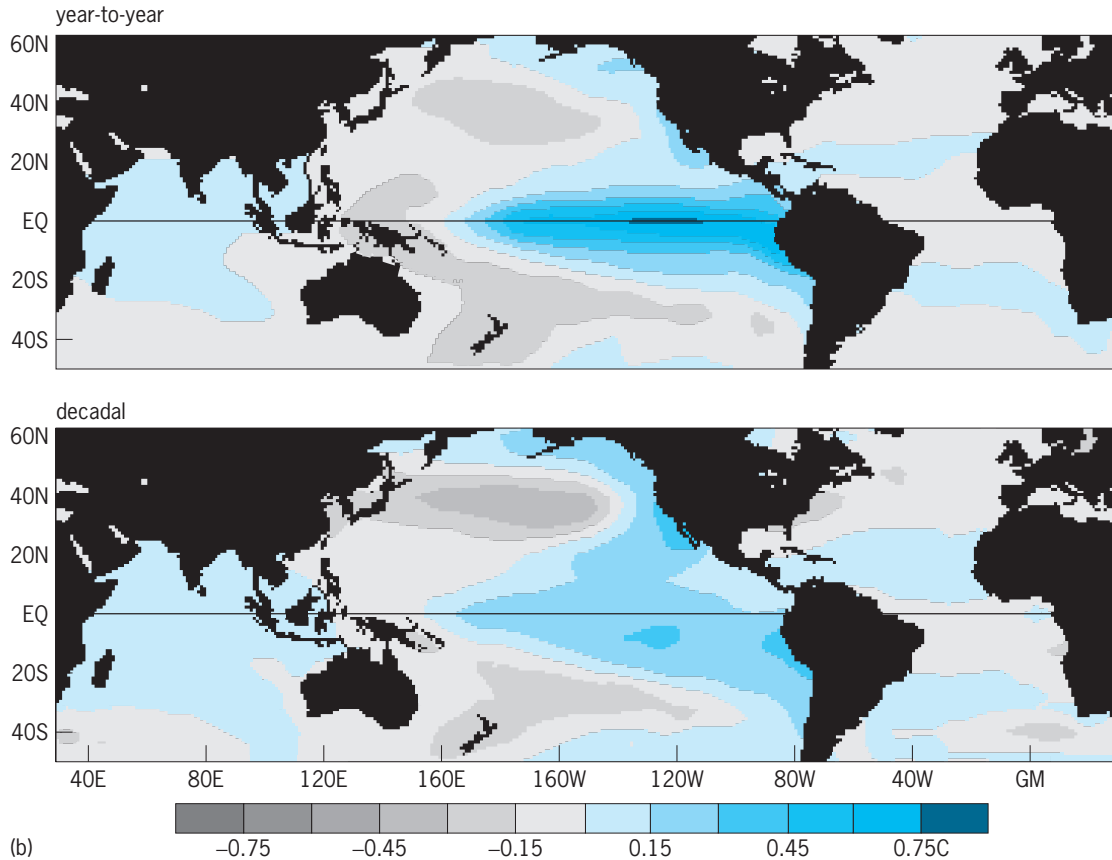
nections from the Pacific are primarily into the winter hemisphere, there is observational and modeling evidence that cool surface-temperature deviations in the tropical Pacific can contribute to persistent summertime droughts in the central United States.

Decadal ENSO. Tropical Pacific climate variability is dominated by the year-to-year swings of the ENSO cycle that are shown in Fig. 3. There are also small-amplitude changes in sea-surface temperatures over broad regions of the Pacific, including the tropics, on time scales longer than the year-to-year variability. Figure 5a presents a pair of time series that documents tropical Pacific temperature variability.

Year-to-year ENSO temperature deviations (Fig. 5a upper) are characterized by deviations in excess of 0.5°C (0.9°F), centered on the Equator in the central and eastern Pacific, with equatorial upwelling playing a dominant role in producing these temperature changes. Deviations of comparable magnitude and like sign are found along the coasts of Ecuador and Peru, and are due to changes in upwelling. Smaller-magnitude (<0.15°C or 0.27°F) warm deviations occur in the tropical Indian Ocean, and cold deviations (<-0.15°C) are found in the western equatorial Pacific, North Pacific (30°S–50°N), and South Pacific (20°S–50°S). Longer, decadal time-scale surface-temperature deviations (Fig. 5b lower) are superficially similar to those for the year-to-year variability, but the extratropical and tropical Pacific deviations are comparable in magnitude (>0.2°C or 0.36°F) on longer time scales, and the central and eastern equatorial temperature deviations extend over a broader latitudinal scale than for the year-to-year variability. For comparison, the secular warming of the oceans, which has been removed from these analyses, is approximately 0.8°C (1.44°F) over the same period. The decadal variability is associated with cooler eastern Pacific conditions in the mid-1940s through late 1970s and warmer conditions in the 1980s and 1990s (Fig. 5a lower). This decadal oscillation is related to the Pacific Decadal Oscillation,



(a)



(b)

Fig. 5. Time series of (a) tropical Pacific ENSO and ENSO-like variations and (b) the characteristic global patterns that go with these series.

which has effects similar to those of ENSO on the United States.

The existence of decadal (and longer time scale, say global warming) modulations of ENSO call into question the definition of ENSO anomalies with respect to a single unchanging background climate.

Theories of ENSO. It is impossible to understand ENSO by considering the atmosphere or the ocean alone. The normal tropical Pacific sea-surface temperature (warm in the west, cold in the east) can be understood in terms of the wind blowing westward over the surface of the ocean. Similarly, the normal atmosphere over the tropical Pacific (heavy rainfall in the west, westward winds, low pressure, and ascending air) can be understood in terms of the sea-surface temperature. But to predict the future evolution of ENSO, the winds would have to be

known in order to predict the sea-surface temperature, and the sea-surface temperature would have to be known in order to predict the winds.

The only way to simulate the consistent evolution of the atmosphere and ocean is to make coupled numerical models of them. A coupled model is one in which the atmospheric component sees and drives the oceanic component at the same time as the oceanic component sees and drives the atmospheric component. The way to predict the future evolution of ENSO is to start with the atmosphere-ocean system in its current state and allow the coupled system to evolve in a mutually consistent manner. See CLIMATE MODELING.

The first such coupled model was constructed in the mid-1980s by S. E. Zebiak and M. A. Cane. It was a simplified model in that only the upper layers

of the ocean were modeled and, more importantly, the annual cycle of sea-surface temperature and surface winds was specified. In this sense, the model was an anomaly model, calculating anomalies with respect to the specified annual cycle. Its importance lay in the fact that it was able to simulate ENSO in ways thought to be relevant. This advance led to a relatively simple explanation of ENSO and made possible the first successful long-term prediction of sea-surface temperature in the eastern Pacific by means of a numerical model of the coupled atmosphere-ocean system. See CLIMATIC PREDICTION.

The mechanism of ENSO that proved to be operative in the simplified Zebiak-Cane model relied on instability of the coupled system: when water in the eastern Pacific was warm, eastward winds were forced at the western edge of the warm water. These eastward winds further warm the water by a number of different mechanisms: by advecting warm water from the warmer West Pacific; by deepening the thermocline and therefore quelling the upwelled cold water; and by weakening the westward winds and therefore weakening the upwelling itself. As the instability sets on and the warm water grows, planetary (Rossby) waves are excited that contain the seeds of the reversal of the instability. They propagate westward, are reflected at the western edge of the Pacific, and return to counter the growing warming in the eastern Pacific. The combined sequence of events is called the delayed oscillator mechanism. The period and amplitude of the ENSO cycle depends on the strength of the coupling, the size of the Pacific, and the delay of the Rossby waves. This mechanism is not complete, however. It explains the regular ENSO cycle but not its irregularity as evidenced in Fig. 3. This irregularity has variously been explained by the addition of random noise to the cycle, the nonlinear interaction of the ENSO cycle with the annual cycle, the nonlinear interaction of the ENSO cycle with other growing instabilities, and by many other possibilities. See OCEAN WAVES.

A completely different class of theories is that the interaction between the atmosphere and the ocean is not unstable. How can a stable system exhibit ENSO? Although the theory is mathematically recondite, there is a class of dynamics (nonnormal) that can exhibit a growth of disturbances followed by decay. Since all ENSO phases first grow and then decay, this has to be allowed as an explanation of ENSO. A stable coupled atmosphere-ocean system perturbed by random noise has been shown to have properties and predictive power currently indistinguishable from those of unstable theories.

Modeling ENSO. Since ENSO is a coupled atmosphere-ocean phenomenon, the system must be composed of a numerical model of the atmosphere suitably coupled to a numerical model of the ocean. If we want also to model the effects of ENSO on the rest of the world, the models must be global and include land, ice and snow, and so on. The model must be run for a very long time, and the statistics of the annual cycle and ENSO are recorded and compared to historical data such as

in Fig. 3. This has been done to very good effect for simplified models where the annual cycle is specified and only anomalies from the annual cycle are calculated. But for fully coupled numerical models where the only input is from the Sun and all else is calculated, the simulation of ENSO has not been so good. The problem is that, due to systematic biases, the simulation of the annual cycle in the tropical Pacific is generally poorly done, and some of the higher-frequency phenomena (such as the intraseasonal oscillations) are also poorly done. These issues generally are agreed to be the dominant problem with the simulation of ENSO. This makes it difficult to examine the changes of ENSO with global warming, but it does not preclude effective predictions since the predictions are initialized by observations. While the systematic biases in the model may tend to pull the prediction back to unrealistic states over time, there will generally be a time in which the predictions are pretty good. It would still be of great value for both short-term predictions (a season or so in advance) and global-warming simulations to improve the models in the manner indicated.

Todd Mitchell; Edward S. Sarachik

Bibliography. K. Achuta Rao and K. R. Sperber, Simulation of the El Niño Southern Oscillation, *Clim. Dyn.*, 19:191–209, 2002; A. G. Barnston et al., Long-lead seasonal forecasts—where do we stand?, *Bull. Amer. Meteorol. Soc.*, 75:2097–2114, 1994; A. G. Barnston et al., Multimodel ensembling in seasonal climate forecasting at IRI, *Bull. Amer. Meteorol. Soc.*, 84:1783–1796, 2003; M. A. Cane, S. E. Zebiak, and S. C. Dolan, Experimental forecasts of El Niño, *Nature*, 321:827–832, 1986; H. R. Diaz and U. Markgraf (eds.), *El Niño: Historical and Paleoclimatic Aspects of the Southern Oscillation*, 1993; L. Goddard et al., Current approaches to seasonal-to-interannual climate predictions. *Int. J. Climatol.*, 21:1111–1152, 2001; T. N. Palmer et al., Development of a European multimodel ensemble system for seasonal-to-interannual prediction (DEMETER), *Bull. Amer. Meteorol. Soc.*, 85:853–872, 2004; S. G. H. Philander, *El Niño, La Niña, and the Southern Oscillation*, 1990; C. F. Ropelewski and M. S. Halpert, North American precipitation and temperature patterns associated with the El Niño/Southern Oscillation (ENSO), *Mon. Weath. Rev.*, 114:2352–2362, 1986; J. M. Wallace et al., On the structure and evolution of ENSO-related climate variability in the tropical Pacific: Lessons from TOGA, *J. Geophys. Res. Oceans*, 103:14,241–14,260, 1998; S. E. Zebiak and M. A. Cane, A model El Niño/Southern Oscillation, *Mon. Weath. Rev.*, 115:2262–2278, 1987; Y. Zhang, J. M. Wallace, and D. S. Battisti, ENSO-like interdecadal variability 1900–93, *J. Clim.*, 10:1004–1020, 1997.

Elasmobranchii

The subclass within the Chondrichthyes (cartilaginous fishes) that includes the sharks (Euselachii) and the skates and rays (Batoidei). The other subclass within Chondrichthyes, according to traditional

classifications, is the Holocephali (chimaera, or ratfishes). It is probable that both groups arose independently during the Silurian or Early Devonian from a group of extinct armored fishes, the Placodermi. The elasmobranchs are distinguished by separate gill openings, amphistylic or hyostylic jaw suspension, and sensory ampullae (of Lorenzini) in the head region. Characters shared with the holocephalans include a variably calcified cartilaginous endoskeleton, placoid scales, urea-retention mechanism, clasper organs in the male for internal fertilization, and the absence of an air (swim) bladder. *See* CHIMAERIFORMES; CHONDRICHTHYES; PLACODERMI; SCALE (ZOOLOGY); SWIM BLADDER.

Evolutionary development. The history of the elasmobranchii can best be understood in terms of three successive evolutionary levels: cladodont, hybodont, and modern. Because of their incomplete fossil record, it is impossible to work out the detailed phylogenetic relationships of the extinct and living groups.

Cladodonts. The cladodont sharks first appeared in the Devonian; except for the specialized pleuracanth (Fig. 1a), which persisted into the Triassic, they had disappeared by the end of the Pennsylvanian. The cladodonts all possess amphistylic jaw suspension. Their teeth have a pointed central cusp, two or more lateral cusps, and a flattened disklike base. The notochord was persistent and was not replaced by centra.

The basal part of the pectoral fin skeleton is variable in design but never tribasal. The radial cartilages

of the pectoral and pelvic fins are not divided and extend more or less to the fin margin. The pelvic plates are separated. Claspers are usually present. The dorsal fins are preceded by spines. The caudal fin is heterocercal, nearly equilobate, and the radials of the hypochordal lobe are unsegmented. The body scales are usually multilobed, each lobe having a separate pulp cavity. Some cladodont sharks possess relatively large semicircular canals; in others these are the same size as in modern sharks. *See* COPULATORY ORGAN; EAR (VERTEBRATE).

Cladodonts were primarily pelagic predators, capable of seizing and tearing prey. The relatively stiff fins suggest that their maneuverability was more limited than in modern sharks. Representative cladodonts include the Devonian genera *Cladodus* and *Cladoselache* (Fig. 1b). *See* DEVONIAN.

Hybodonts. The hybodonts first appeared in the Mississippian Period. Although their skeleton resembles that of the cladodonts, they were more advanced in having only a tribasal pectoral fin skeleton, divided and shortened paired fin radials, clearly differentiated anal fin, and ribs and hemal arches along the entire length of the unreduced notochord. The acquisition of the tribasal pectoral fin, present in all later sharks, and a more flexible hypochordal lobe in the caudal fin apparently enabled these sharks to maintain their equilibrium more effectively than the cladodonts.

Some hybodonts had cladodont-like teeth and presumably fed on swimming prey, but much of the adaptive radiation at this level involved the development of crushing or grinding dentitions related mainly to bottom feeding. The late Palaeozoic genus *Ctenacanthus* is a good example of a conservative hybodont; the Mesozoic *Hybodus* (Fig. 1c) has a more advanced fin structure and a crushing dentition. *Ctenacanthus* lived from the Mississippian into the Permian, and *Hybodus* from the Triassic through the Cretaceous.

Modern elasmobranchs. Modern sharks arose during the Jurassic Period long before the hybodonts had become extinct in the Late Cretaceous. The transition from the hybodont to the modern level is poorly documented by fossils. However, there are three groups of aberrant (part hybodont, part modern) living sharks. Heterodontidae (Port Jackson sharks), Chlamydoselachidae (frilled sharks), and Hexanchidae (six-gill and seven-gill sharks). All except the Port Jackson sharks are rare deep-water forms, and each group has a somewhat different combination of hybodont and modern shark characters. These archaic sharks suggest the kind of evolutionary experimentation that took place during the rise of the modern sharks. *See* JURASSIC; SELACHII.

The skates and rays (batoids) are mostly bottom-dwelling elasmobranchs that possess many features distinct from the sharks. Present evidence suggests that the batoids arose during the Jurassic from some benthic hybodont group.

Modern features. Morphologically the modern shark level may be defined by the following features:

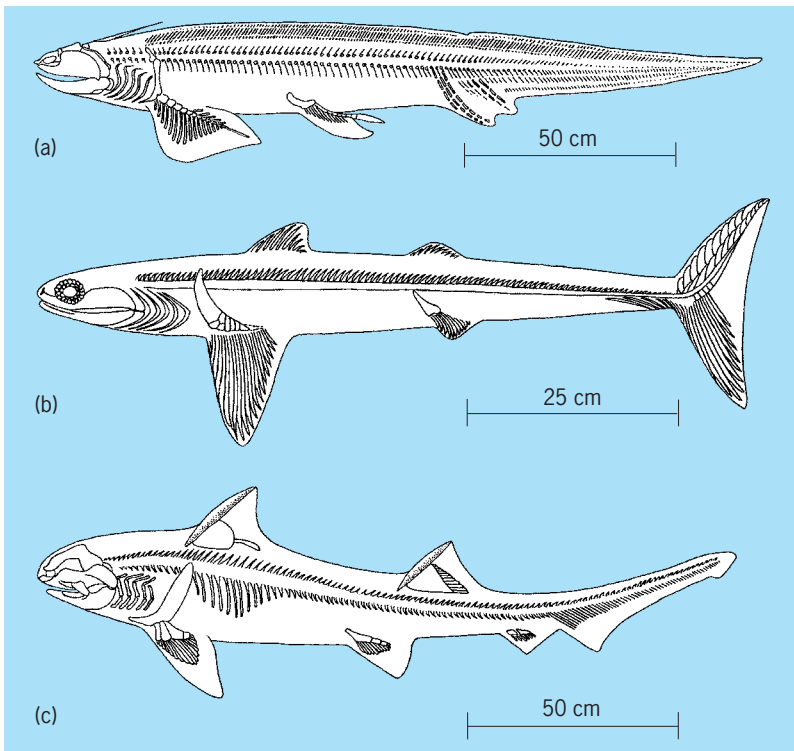


Fig. 1. Lateral views of fossil sharks. (a) *Pleurocanthus*, a specialized cladodont. (b) *Cladoselache*, a cladodont. (c) *Hybodus*, a hybodont.

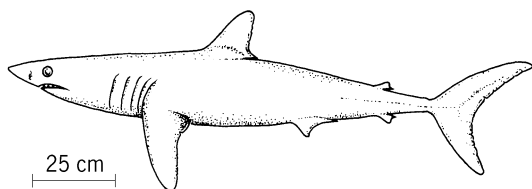


Fig. 2. Mackerel shark (*Isurus*). (After H. B. Bigelow and W. C. Schroeder, *Fishes of the Western North Atlantic*, pt. 1, 1948)

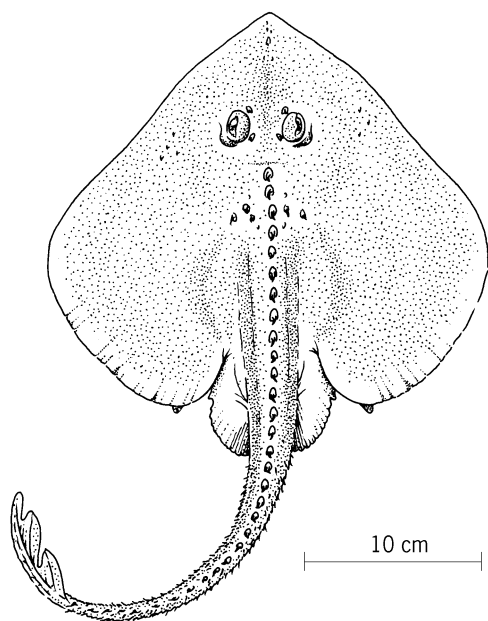


Fig. 3. Skate (*Raja*). (After H. B. Bigelow and W. C. Schroeder, *Fishes of the Western North Atlantic*, pt. 2, 1953)

hyostylic jaw suspension, shortened jaws, jaw protrusion mechanism, complete replacement of the notochord by calcified centra, fusion of the pelvic plates, and presence of only single-lobed placoid scales. Although the rostrum was short and blunt in the cladodonts, it became elongated in some hydodonts and, with a few exceptions, it projects well beyond the mouth in all modern sharks (Fig. 2).

Skates and rays differ from the sharks in having a depressed body, modification in the mode of gill ventilation, freer hyostylic jaw suspension, elaboration of the pavement dentition for crushing or grinding, enlargement of the pectoral fins, and disappearance of the anal fin. With the pectoral fins assuming the main role in locomotion, the posterior part of the body, including the caudal fin, became reduced to whiplike proportions (Fig. 3). See BATOIDEI; SELACHII. Bobb Schaeffer

Bibliography. H. B. Bigelow and W. C. Schroeder, *Fishes of the Western North Atlantic*, pt. 1, 1948, pt. 2, 1953; C. E. Bond, *Biology of Fishes*, 2d ed., 1997; E. S. Herald, *Living Fishes of the World*, 1961; F. H. Pough et al., *Vertebrate Life*, 5th ed., 1998; B. Schaeffer, Comments on elasmobranch evolution, in P. W. Gilbert et al. (eds.), *Sharks, Skates and Rays*, 1967.

Elasticity

The property whereby a solid material changes its shape and size under the action of opposing forces, but recovers its original configuration when the forces are removed. The theory of elasticity deals with the relations between the forces acting on a body and the resulting changes in configuration, and is important in many branches of science and technology, for instance, in the design of structures, in the theory of vibration and sound, and in the study of the forces between atoms in crystal lattices.

Elastic constants. The forces acting on a body are expressed as stresses and measured as force per unit area. Thus if a bar $ABCD$ of square cross section (Fig. 1a) is fixed at one end and subjected to a force F uniformly distributed over the other end DC , the stress is $F/(DC)^2$. This stress causes the bar to become longer and thinner and to assume the shape $A'B'C'D'$. The strain is measured by the ratio (change in length)/(original length), that is, by $(B'C' - BC)/(BC)$. According to Hooke's law, stress is proportional to strain, and the ratio of stress to strain is therefore a constant, in this case Young's modulus, denoted by E , so that $E = F(BC)/(DC)^2 (B'C' - BC)$. See HOOKE'S LAW; STRESS AND STRAIN; YOUNG'S MODULUS.

Poisson's ratio σ is defined as the ratio of lateral strain to longitudinal strain, so that $\sigma = BC(DC - D'C')/DC (B'C' - BC)$. The bar of Fig. 1a is in a state of tension, and the stress is tensile; if the force F were reversed in direction, the stress would be compressive. Stresses of this type are called direct or normal stresses; a second type of stress, known as tangential or shear stress, is illustrated in Fig. 1b. In this case, the configuration $ABCD$ becomes $ABC'D'$, with the shear forces F acting in the directions AB and CD . The shear strain is measured by the angle θ , and if the body is originally a cube, the shear stress is $F/(DC)^2$. The ratio of stress to strain, $F/(DC)^2\theta$, is the shear or rigidity modulus G , which measures the

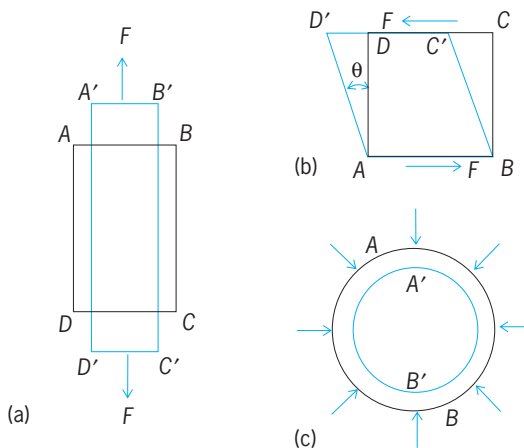


Fig. 1. Stresses on a bar. (a) Direct or normal stress. (b) Tangential or shear stress. (c) Change in volume with no change in shape. (All deformations are exaggerated.)

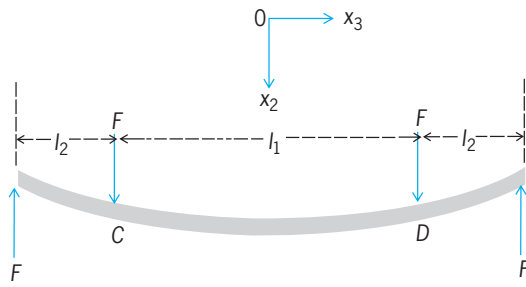


Fig. 2. Flexure and torsion in a bar.

resistance of the material to change in shape without change in volume.

A further elastic constant, the bulk modulus k , measures the resistance to change in volume without changes in shape, and is illustrated in Fig. 1c. The original configuration is represented by the circle AB , and under a hydrostatic (uniform) pressure P , the circle AB becomes the circle $A'B'$. The bulk modulus is then $k = Pv/\Delta v$, where $\Delta v/v$ is the volumetric strain. The reciprocal of the bulk modulus is the compressibility.

Determination of values. The elastic constants may be determined directly in the way suggested by their definitions; for instance, Young's modulus can be determined by measuring the relative extension of a rod or wire subjected to a known tensile stress. Less direct methods are, however, usually more convenient and accurate. Prominent among these are the dynamic methods involving frequency of vibration and velocity of sound propagation. The elastic constants can be expressed in terms of frequency of (or velocity in) regularly shaped specimens, together with the dimensions and density, and by measuring these quantities, the elastic constants can be found. See ULTRASONICS.

The elastic constants can also be determined from the flexure and torsion of bars. As an illustration, consider a bar AB (Fig. 2) of breadth b (in the x_1 direction) and depth d (in the x_2 direction) supported by forces F at the ends, and loaded symmetrically by forces F at points C and D . Over the portion CD there is a uniform bending moment $M = Fl_2$, and the theory of bending shows that the portion CD is bent into the arc of a circle such that Eq. (1) applies, where R

$$R = \frac{EI}{M} \quad (1)$$

is the radius of curvature, E is Young's modulus, and I is the moment of inertia of cross section, equal to $bd^3/12$ for a rectangular cross section. The longitudinal stress at the lower face of the bar is tensile, and at the upper face, compressive. The middle plane of the bar is free of stress, and is the neutral axis. The stress at a distance x_2 from the neutral axis is shown in Eq. (2).

$$T = \frac{Ex_2}{R} = \frac{Mx_2}{I} \quad (2)$$

It is thus possible to determine E from Eq. (1) by measuring I , R , and M ; conversely, if E is known, the stress may be determined from Eq. (2). See LOADS, TRANSVERSE.

Practical limitations. In practice, stress is only proportional to strain, and the strain is only completely recoverable within certain limits, called the elastic limits of the material. The stress below which the strain is completely recoverable is sometimes called the limit of perfect elasticity, and the stress up to which Hooke's law is obeyed is sometimes called the proportional limit or limit of linear elasticity. Above the elastic limits, the material is subject to time-dependent effects, and as the stress is further increased, the ultimate strength of the material is approached. See PLASTICITY; STRENGTH OF MATERIALS.

Theory of elasticity. In classical elasticity theory, it is assumed that the strains are always small; Hooke's law is therefore obeyed; the strains are completely recoverable and, moreover, are superposable, so that the strain produced by the joint action of two or more stresses is the sum of the strains produced by them individually.

In order to develop the theory, it is necessary to specify the stresses and strains more closely. Figure 3 shows the stress components T_{ij} (where i, j may take the values 1, 2, or 3) acting on the faces of a cube, parallel to the coordinate axes x_1, x_2, x_3 . The first suffix indicates the direction of the stress component and the second the direction of the normal to the plane under consideration. Stresses of the type T_{11} are normal stresses, and of the type T_{12} , shear stresses. The conditions for zero rotation of the cube are $T_{12} = T_{21}, T_{13} = T_{31}, T_{23} = T_{32}$, and there are therefore six independent stress components.

In addition to the stresses T_{ij} , body forces proportional to volume (for instance, forces due to the weight of the body) may also be acting. If the stresses T_{ij} vary with position, application of Newton's second law leads to Eq. (3) for the x_1 direction where ρ

$$\frac{\partial T_{11}}{\partial x_1} + \frac{\partial T_{12}}{\partial x_2} + \frac{\partial T_{13}}{\partial x_3} + X_1 = \rho f_1 \quad (3)$$

is the density, f_1 is the acceleration, and X_1 the body force component per unit volume along x_1 , together with two similar equations for the x_2 and x_3 directions. If $f_1 = f_2 = f_3 = 0$, these equations become the equations of equilibrium, and if, further, $X_1 = X_2 = X_3 = 0$, they become the equations of equilibrium in the absence of body forces. The preceding

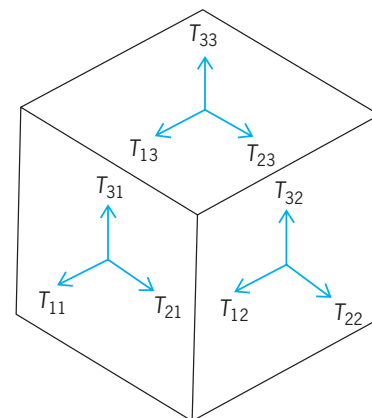


Fig. 3. Stress components acting on the faces of a cube.

equations are important in many branches of elastic theory and, for example, provide a starting point in the study of vibrating bodies and of the twisting of cylinders and prisms with cross sections of various shapes. See TORSION.

The components of strain are specified in a similar way to the stresses. There are six independent strain components: S_{11} , S_{22} , S_{33} , S_{23} , S_{13} , and S_{12} . If, as a result of strain, the coordinates of a point x_1 , x_2 , x_3 become $x_1 + u_1$, $x_2 + u_2$, $x_3 + u_3$, the quantities u_1 , u_2 , and u_3 are the components of the displacement vector, and the strain components are as displayed in Eq. (4), so that, for example the relations shown by Eq. (5) would hold true.

$$S_{ij} = \frac{1}{2} \left(\frac{\partial u_i}{\partial x_j} + \frac{\partial u_j}{\partial x_i} \right) \quad (4)$$

$$S_{11} = \frac{\partial u_1}{\partial x_1} \quad S_{12} = \frac{1}{2} \left(\frac{\partial u_1}{\partial x_2} + \frac{\partial u_2}{\partial x_1} \right) \quad (5)$$

By eliminating the displacements from these equations, the so-called compatibility equations are obtained with three of the type shown in Eq. (6) and three of the type shown in Eq. (7).

$$\frac{\partial^2 S_{22}}{\partial x_3^2} + \frac{\partial^2 S_{33}}{\partial x_2^2} = 2 \frac{\partial^2 S_{23}}{\partial x_2 \partial x_3} \quad (6)$$

$$\frac{\partial^2 S_{11}}{\partial x_2 \partial x_3} = \frac{\partial}{\partial x_1} \left(-\frac{\partial S_{23}}{\partial x_1} + \frac{\partial S_{13}}{\partial x_2} + \frac{\partial S_{12}}{\partial x_3} \right) \quad (7)$$

The stresses and strains have so far been denoted by two suffixes. This is essential if the methods of tensor analysis are to be applied to elasticity problems, but for many purposes a single suffix notation is adequate. The change from a two- to a one-suffix notation for the stresses is simply $T_{11} = T_1$, $T_{22} = T_2$, $T_{33} = T_3$, $T_{23} = T_4$, $T_{13} = T_5$, $T_{12} = T_6$. The change of notation for the strains is $S_{11} = S_1$, $S_{22} = S_2$, $S_{33} = S_3$, $2S_{23} = S_4$, $2S_{13} = S_5$, $2S_{12} = S_6$; the factor 2 is required to make the strains S_4 , S_5 , and S_6 conform with the usual definition of shear strain (Fig. 1b).

Hooke's law generalized. Hooke's law may be generalized to the statement that each stress component is proportional to each strain component, equivalent to the six equations (8a) which may be written more concisely as Eq. (8b), where the summation extends

$$\begin{aligned} T_1 &= c_{11}S_1 + c_{12}S_2 + c_{13}S_3 + c_{14}S_4 + c_{15}S_5 + c_{16}S_6 \\ &\dots \dots \dots \\ T_6 &= c_{61}S_1 + c_{62}S_2 + c_{63}S_3 + c_{64}S_4 + c_{65}S_5 + c_{66}S_6 \end{aligned} \quad (8a)$$

$$T_q = \sum_r c_{qr} S_r \quad (8b)$$

over $r = 1, 2, 3, 4, 5$, and 6. The elastic constants c_{qr} are termed the stiffnesses; there are altogether 36 of them but they are subject to the reciprocal relations $c_{qr} = c_{rq}$ imposed by thermodynamic requirements, and the number is thus reduced to 21.

Additional relations can be derived from the three assumptions that the interatomic forces act along the lines joining the centers of atoms in the lattice, that

the atoms are situated at centers of symmetry, and that the lattice is initially at zero stress. These relations, called Cauchy relations, are $c_{23} = c_{44}$, $c_{13} = c_{55}$, $c_{12} = c_{66}$, $c_{14} = c_{56}$, $c_{25} = c_{46}$, $c_{45} = c_{36}$ and, if true, would reduce the number of stiffnesses to 15. Experiment shows, however, that they are not true in general; nevertheless, their investigation provides an indication of the extent to which these three assumptions hold in any particular case.

The generalized Hooke's law can also be written to express the strains in terms of the stresses given in Eq. (9), in which the quantities s_{qr} are the elas-

$$S_q = \sum_r s_{qr} T_r \quad (r = 1, 2, 3, 4, 5, 6) \quad (9)$$

tic compliances. If the six simultaneous equations of Eq. (8) are solved for the strains, the compliances are obtained in terms of the stiffnesses as in Eq. (10), where Δc is the determinant shown in Eq. (11) and

$$s_{qr} = \frac{\Delta c_{qr}}{\Delta c} \quad (10)$$

$$\Delta c = \begin{vmatrix} c_{11} & c_{12} & c_{13} & c_{14} & c_{15} & c_{16} \\ c_{12} & c_{22} & c_{23} & c_{24} & c_{25} & c_{26} \\ c_{13} & c_{23} & c_{33} & c_{34} & c_{35} & c_{36} \\ c_{14} & c_{24} & c_{34} & c_{44} & c_{45} & c_{46} \\ c_{15} & c_{25} & c_{35} & c_{45} & c_{55} & c_{56} \\ c_{16} & c_{26} & c_{36} & c_{46} & c_{56} & c_{66} \end{vmatrix} \quad (11)$$

Δc_{qr} is the cofactor obtained by deleting the row and column containing c_{qr} from the determinant Δc .

The 21 stiffnesses (or compliances) of the generalized Hooke's law describe the elastic behavior of a material belonging to the triclinic crystal system. The existence of symmetry elements reduces the number of independent elastic constants in the other crystal systems to the following numbers: monoclinic, 13; orthorhombic, 9; tetragonal, 7 or 6; trigonal, 7 or 6; hexagonal, 5; and cubic, 3. Materials belonging to all of these systems are anisotropic, and the elastic properties depend upon direction within the material. If the properties are independent of direction, the material is isotropic and its elastic behavior is completely described by two independent stiffnesses (or compliances). See CRYSTALLOGRAPHY.

The stress-strain relations, referred to the principal axes in the orthorhombic, hexagonal, cubic, and isotropic systems, are given in the table. The equations involving the compliances are completely analogous, with S and T interchanged, and s_{qr} written for c_{qr} , except where $T_q = 1/2 (c_{11} - c_{12}) S_q$, in which case $S_q = 2(s_{11} - s_{12}) T_q$.

Rochelle salt is an example of an orthorhombic crystal; materials which, although not crystalline, possess the same symmetry and matrix of elastic constants as orthorhombic crystals are said to be orthotropic. Wood and plywood are materials of this description, and orthotropic elastic theory has also been applied to laminated plastics and reinforced concrete.

Single-crystal zinc, cobalt, magnesium, and ice are hexagonal materials; they are transversely isotropic because the properties are independent of direction

Stress	Strain			
	Orthorhombic	Hexagonal	Cubic	Isotropic
$T_1 =$	$c_{11}S_1 + c_{12}S_2 + c_{13}S_3$	$c_{11}S_1 + c_{12}S_2 + c_{13}S_3$	$c_{11}S_1 + c_{12}S_2 + c_{12}S_3$	$c_{11}S_1 + c_{12}S_2 + c_{12}S_3$
$T_2 =$	$c_{12}S_1 + c_{22}S_2 + c_{23}S_3$	$c_{12}S_1 + c_{11}S_2 + c_{13}S_3$	$c_{12}S_1 + c_{11}S_2 + c_{12}S_3$	$c_{12}S_1 + c_{11}S_2 + c_{12}S_3$
$T_3 =$	$c_{13}S_1 + c_{23}S_2 + c_{33}S_3$	$c_{13}S_1 + c_{13}S_2 + c_{33}S_3$	$c_{12}S_1 + c_{12}S_2 + c_{11}S_3$	$c_{12}S_1 + c_{12}S_2 + c_{11}S_3$
$T_4 =$	$c_{44}S_4$	$c_{44}S_4$	$c_{44}S_4$	$(c_{11} - c_{12})S_4/2$
$T_5 =$	$c_{55}S_5$	$c_{44}S_5$	$c_{44}S_5$	$(c_{11} - c_{12})S_5/2$
$T_6 =$	$c_{66}S_6$	$(c_{11} - c_{12})S_6/2$	$c_{44}S_6$	$(c_{11} - c_{12})S_6/2$

in all of the planes normal to the hexagonal axis.

Single-crystal copper, gold, silver, nickel, and the alkali halides (for example, sodium chloride) are important cubic materials. The stress-strain equations are derived from those of the orthorhombic system by superimposing the condition that the three principal directions are all equivalent. This does not mean that the properties are independent of direction; for example, the compliance s_{11}' in an arbitrary direction is given by Eq. (12), where a_1, a_2, a_3 are the

$$s'_{11} = s_{11} - 2(s_{11} - s_{12} - s_{44}/2) \cdot (a_1^2 a_2^2 + a_2^2 a_3^2 + a_3^2 a_1^2) \quad (12)$$

cosines of the angles between the arbitrary direction and the cubic axes. This equation shows that s'_{11} depends on orientation unless $s_{44}/2 = s_{11} - s_{12}$. See PHOTOELASTICITY.

R. F. S. Hearmon
 Bibliography. A. P. Boresi, R. J. Schmidt, and O. M. Sidebottom, *Advanced Mechanics of Material*, 5th ed., 1993; S. F. Borg, *Fundamentals of Engineering: Elasticity*, 1990; J. E. Marsden and T. J. Hughes, *Mathematical Foundations of Elasticity*, 1983, reprint 1994; S. P. Timoshenko and J. N. Goodier, *Theory of Elasticity*, 3d ed., 1970.

Electret

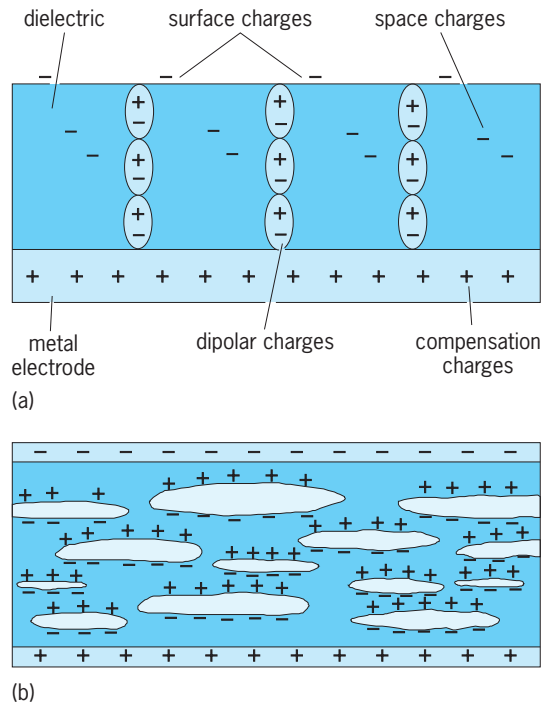
A solid dielectric with a quasipermanent electric moment. Electrets may be classified as real-charge electrets and dipolar-charge electrets. Real-charge electrets of present importance are either solid dielectrics with charges of one polarity at or near one side of the dielectric and charges of opposite polarity at or near the other side, or cellular dielectrics with charges of opposite polarity on the two sides of the voids. The dipolar-charge electrets are dielectrics with aligned dipolar charges. Some dielectrics are capable of storing both real and dipolar charges. Examples of charge arrangements in electrets are shown in the illustration. See POLARIZATION OF DIELECTRICS.

Modern electrets used in research and in applications are often films ranging in thickness from the submicrometer scale for silicon-based materials to sheets of 5–50- μm (0.2–2-mil) thickness (foil electrets) for suitable polymer materials. They are frequently metallized on one or both sides, depending on the intended use. In many applications of the solid real-charge electrets, the external fields of the

charged dielectric are used and the films are metallized on only one side to allow the electric field to emanate from the nonmetallized side. In applications of the cellular real-charge electrets and the dipolar electrets, the piezoelectric effect generated by such samples is used and the films are metallized on both surfaces. See PIEZOELECTRICITY.

Solid real-charge electrets. In these electrets, the intrinsic conductivity and carrier mobility must be small enough that excess charges can persist for long periods of time. Suitable materials are the fluorocarbons polytetrafluoroethylene (PTFE) and its copolymer fluoroethylenepropylene (FEP), as well as polypropylene (PP), polycarbonate (PC), polyimide (PI), polyethylene terephthalate (PETP), and inorganic materials such as silicon dioxide and aluminum oxide. PETP can also assume a considerable dipole polarization. See POLYMER.

Charging of real-charge electrets, metallized on one side, can be achieved by exposure of the nonmetallized side to a corona discharge; by application



Schematic cross sections of electrets. (a) Solid electret disk metallized on one side with real charges (surface and space charges) and dipolar charges. (b) Cellular electret metallized on both sides with charges of opposite polarities on the upper and lower surfaces of the voids.

of an electric field at elevated temperatures through an electrode opposing the nonmetallized side of the dielectric; or by contacting the nonmetallized side with a biased, wet electrode. Another charging technique, allowing greater control than other methods, is electron injection with a partially penetrating electron beam. In all cases, carriers of one polarity are deposited onto or injected into the insulator and trapped, and a compensation charge of equal magnitude but opposite sign flows into the sample electrode. Thus, the charge arrangement exhibits an electric moment. Another method, applied to photoconductors metallized on both sides with transparent electrodes, consists in the application of an electric field during irradiation with ultraviolet or visible light. Carrier generation by the light and charge separation by the field yield the desired electric moment, resulting in a photoelectret. *See* CORONA DISCHARGE; PHOTOCONDUCTIVITY.

PTFE and FEP are electronegative and therefore hold negative charges better than positive ones. Real-charge densities of up to 5 millicoulombs/m² (3.2 microcoulombs/in.²) can be stored in these polymers. Electrets are stable over long periods of time. At the customarily used charge densities of about 0.2 millicoulombs/m² (130 nanocoulombs/in.²), corresponding to surface potentials of approximately 250 V in a 25- μ m-thick (1-mil) sample, the time constant of the decay of the surface potential of PTFE and FEP is about 200 years if storage is at room temperature and low relative humidity under shielded conditions. This time constant drops to about 10 years at 50°C (122°F) and 90% relative humidity.

Silicon dioxide (SiO₂) electrets are preferably made of a thermally grown silicon dioxide layer, usually less than 1 μ m (0.04 mil) thick, on a silicon substrate. Negative charges deposited onto the free oxide surface show better stability than negative volume charges or positive charges. Surface-charging techniques, such as the corona method, are therefore utilized. After charging, a compensation charge of opposite polarity is located at the silicon-silicon dioxide interface. Charge densities are generally chosen in the range 10–20 mC/m² (6–12 μ C/in.²), corresponding to surface potentials of 300–600 V on a 1- μ m-thick layer. Time constants of the charge decay exceed those for PTFE and FEP.

Cellular real-charge electrets. These newer electrets with piezoelectric properties consist of a cellular dielectric, typically a polymer such as PP, with voids of the order of a few micrometers height and lateral dimensions of several tens of micrometers. Such voids can be generated by biaxially stretching a polymer film filled with small inorganic particles. During the stretching, microcracks form around the particles and develop into flat voids. These can be enlarged in the thickness direction by a pressure treatment, consisting of an increase of the pressure of the surrounding gas such that gas penetrates into the voided structure. A subsequent pressure release inflates the material. Thereafter, the sample is metallized on one of its surfaces.

Charging of such a sample is achieved by expo-

sure of its unmetallized surface to a corona discharge. The deposited charges together with their countercharges on the metallization generate a high internal field that will cause breakdowns in the voids and thus deposition of charges of opposite sign on the internal void surfaces. Finally, the sample is metallized on the other surface. The material now exhibits a large piezoelectric d_{33} coefficient that is very stable at room temperature. Problematic is the thermal depolarization at elevated temperatures found in the PP samples. Other materials with higher depolarization temperatures are now under study. The samples also show hysteresis behavior as a function of the applied field and are therefore often referred to as ferroelectrets.

Dipolar-charge electrets. In these electrets, the relaxation time of the molecular reorientation has to be large enough that the dipole alignment persists for long periods of time. Most prominent among such materials are poly(vinylidene fluoride) [PVDF] and vinylidene fluoride-trifluoroethylene copolymers [VDF-TrFE], which exhibit piezoelectric and pyroelectric properties linked to the dipole polarization. Inorganic piezoelectric crystals and ceramics also belong to this category but are not discussed here. *See* PIEZOELECTRICITY; PYROELECTRICITY.

The preferred technique for the polarization of dipolar electrets is also the above corona method applied to one-sided metallized samples. Again, the injected charges set up a large field within the dielectric that aligns the dipoles. Another method consists in the direct application of a high field to a two-sided metallized sample at an elevated temperature. After the polarization process, the samples exhibit sizable d_{31} and d_{33} coefficients.

In the case of PVDF polarization densities of up to 0.1 C/m² (65 μ C/in.²) can be achieved by corona charging at room temperature with fields of up to 400 MV/m (10 MV/in.) on samples prestretched three to five times their original length. The stretching makes the material, which is semicrystalline, recrystallize in a polar form and also reorients the lamellar crystals so that the dipoles are more precisely aligned by the applied field.

PVDF has a relatively stable polarization at room temperature. At 100°C (212°F), the polarization decays initially (within minutes) to half its original value and stabilizes thereafter. Similar stability is achieved for the piezoelectric activity of this material.

Comparison of piezoelectric properties. Typical d_{33} coefficients for cellular PP and solid PVDF are compared in the **table** with values for the conventional piezoelectric materials quartz and lead zirconate

Comparison of piezoelectric properties		
Piezoelectric material	d_{31} , pC/N	d_{33} , pC/N
Quartz	0.73 (d_{14})	2.3 (d_{11})
PZT-5	170	370
PVDF	25	20
Cellular PP (expanded)	1	600

titanate (PZT-5). Since the electret materials are polymers, they have certain other properties such as mechanical softness and flexibility that distinguish them from the conventional piezoelectrics. Particularly noteworthy is the huge d_{33} coefficient of cellular PP. The electrets have already found numerous commercial applications.

Applications. Important commercial applications of real-charge electrets are in electroacoustic and electromechanical transducers, air filters, electret dosimeters, and adhesive posters. Of particular significance are electret microphones, which are made in quantities of more than 2×10^9 annually and account for 90% of the microphones produced. Also of interest are biological applications based on the blood compatibility of charged polymers or on their favorable influence on wound or fracture healing. Commercial applications of dipolar electrets are in piezoelectric transducers for airborne and waterborne sound and in pyroelectric detectors. See AIR FILTER; DIELECTRIC MATERIALS; DOSIMETER; ELECTRET TRANSDUCER; ELECTRICAL INSULATION; MICROPHONE.

Gerhard M. Sessler

Bibliography. S. Bauer, R. Gerhard-Multhaupt, and G. M. Sessler, Ferroelectrets: Soft electroactive foams for transducers, *Phys. Today*, 57(2): 37–43, 2004; Institute of Electrical and Electronics Engineers, *IEEE 12th International Symposium on Electrets*, 2005; V. Kestelman, L. Pinchuk, and V. Goldade, *Electrets in Engineering*, Kluwer Academic, 2000; G. M. Sessler and R. Gerhard-Multhaupt (eds.), *Electrets*, vols. 1 and 2, 3d ed., Laplacian Press, 1999.

Electret transducer

A device for the conversion of acoustical or mechanical energy into electrical energy, and vice versa, which utilizes a quasipermanently charged or polarized dielectric material (electret). Examples are certain microphones, headphones, hydrophones, and ultrasonic devices. Depending on the electret material and the usage, either the external electric field of real-charge electrets or the piezoelectric activity of dipolar electrets is used. In the simplest implementation, a transducer based on real-charge electrets consists of a metal backplate (first electrode) covered by a mechanically tensioned diaphragm. The diaphragm is a film electret carrying a metal coating (second electrode) on the side facing away from the backplate. Provisions are made to maintain a shallow air gap between electret and backplate. The air gap is occupied by an electrostatic field originating from the electret charges. Upon acoustical or mechanical deflection of the diaphragm, such a device generates an electrical output signal between its two electrodes; similarly, application of an electrical signal results in diaphragm deflections. Such electret devices are therefore self-biased electrostatic or condenser transducers. They thus exhibit all the advantages of this transducer class, such as wide dynamic range and flat response over a frequency range of several decades, without requiring the external bias necessary in conventional transducers of this kind. On the

other hand, a transducer based on the piezoelectric activity of a dipolar electret may consist of just the electret film covered by two electrodes. Compression or bending of the film causes again an electrical output signal between the electrodes, while the application of a voltage to the electrodes results in mechanical action. See ELECTRET; ELECTROSTATICS; TRANSDUCER.

Film electret. Suitable polymer films are electrically charged or polarized to produce an external electric field or piezoelectric activity, respectively. Materials capable of quasipermanently trapping charge carriers (that is, storing charges over periods of years or decades) include polytetrafluoroethylene (PTFE), its copolymer polyfluoroethylene (FEP or Teflon[®]), and polypropylene (PP), while substances capable of carrying a quasipermanent dipole polarization include poly(vinylidene fluoride) [PVDF] and some of its copolymers. A polymer of considerable interest is cellular PP, which can carry charges on the surfaces of its internal voids and therefore exhibits a very large piezoelectric activity. The cellular PP loses its charges at elevated temperatures and is therefore suitable only for room-temperature use. Before charging or poling, all the materials are either metallized on one side or backed up with a metal electrode.

Electret microphones. The most widely used electret transducer is the electret microphone. This is a condenser microphone consisting of two electrodes separated by an electret and an air gap. The electret generates an electric field in the air gap. If one of the electrodes is designed as a diaphragm that vibrates in response to an incident sound wave, an output signal proportional to the sound pressure is generated.

Such microphones are designed as transducers with a backplate electret or a diaphragm electret. The backplate-electret microphones, which are the more common transducers, consist of a (6–12- μm) (0.25–0.5-mil) Teflon layer adhering to a metallic or metal-coated backplate and charged to 200–400 microcoulombs/ m^2 (130–260 nanocoulombs/ in^2), corresponding to an external bias of about 100 V. The surface of this electret is placed next to a metallic or metallized diaphragm, leaving a shallow air layer of about 20 μm (0.8 mil) thickness. The stiffness of the air layer can be decreased (and thus the sensitivity of the microphone can be improved) by connecting the air layer to a larger cavity by means of small holes through the backplate. The electrical output of the microphone is taken between the backplate, which is insulated from the outer case, and the metal part of the diaphragm. The output is fed into a high-impedance preamplifier. The diaphragm-electret microphones typically consist of a 12–25- μm -thick (0.5–1-mil) Teflon or PTFE film, metallized on one side and charged on the other side to charge densities similar to the above-mentioned values. A backplate and an intermediate air gap complete the system. These transducers are simpler and less expensive than the backplate-electret microphones but mechanically not as stable. This is attributed to stress relaxation of the tensioned Teflon diaphragm, which

causes small changes in sensitivity.

At frequencies below the resonance frequency, the acoustical properties of electret microphones are largely governed by the restoring force on the diaphragm. Since the mechanical tension of the diaphragm is generally kept at a relatively low value (about 10–50 newtons/m or 0.7–3.5 lb/ft), the restoring force is determined by the compressibility of the air layer. Controlling the restoring force by the air layer is advantageous because changes in tension due to stress relaxation thus have only a minor effect on the sensitivity, which is largely independent of transducer area.

Under open-circuit conditions, a displacement d of the diaphragm of an electret microphone causes a frequency-independent output voltage given by Eq. (1), where E is the (constant) electric field in

$$v = Ed = \frac{\sigma D_1 d}{\varepsilon_0(D_1 + \varepsilon D_2)} \quad (1)$$

the air layer of the transducer; σ is the charge density of the electret projected onto its surface; D_1 and D_2 are the thicknesses of electret film and air layer, respectively; ε_0 is the permittivity of free space; and ε is the relative dielectric constant of the electret material. As in conventional electrostatic transducers, the displacement is proportional to the applied pressure in a wide frequency band extending from a lower cutoff given by a pressure-equalization leak in the back cavity to an upper cutoff determined by the resonance frequency. The voltage response for constant sound pressure is frequency-independent in this range.

Typical electret microphones designed for the audio-frequency range have constant sensitivities of 10–50 mV/pascal (1–5 mV/microbar or 70–350 V/lb·in.⁻²) in the frequency range 20 to 15,000 Hz. Nonlinear distortion is less than 1% for sound-pressure levels below 140 dB, and the impulse response is excellent, owing to the flat amplitude and phase characteristics. Other properties of electret microphones are their low sensitivity to vibration, owing to the small diaphragm mass, and their insensitivity to magnetic fields. Compared with conventional electrostatic transducers, electret microphones have the following advantages: they do not require a dc bias; they have three times higher capacitance per unit area, resulting in a better signal-to-noise ratio; and they are not subject to destructive arcing between foil and backplate in humid atmospheres and under conditions of water condensation.

Various electret microphones for operation at infrasonic and ultrasonic frequencies, covering the range from 0.001 Hz to 200 MHz, have been designed by properly positioning the upper and lower cutoff frequencies discussed above. Furthermore, transducers with directional characteristics such as cardioid, bidirectional, toroidal, and second-order unidirectional have been developed. See MICROPHONE.

Electret headphones. Another group of electret transducers made of real-charge electrets are headphones. The above microphone designs can also be utilized in this case. These transducers produce a cer-

tain amount of nonlinear distortion due to the partially quadratic dependence of the force F per unit area, exerted on the diaphragm, on the applied ac voltage v , as can be seen from Eq. (2). High-voltage

$$F = \frac{-\frac{1}{2}(D_1\sigma + \varepsilon\varepsilon_0v)^2}{\varepsilon_0(D_1 + \varepsilon D_2)^2} \quad (2)$$

sensitivities have been achieved by making the air-gap thickness small. The distortion can be eliminated by the use of push-pull transducers and monocharge electrets.

For a push-pull transducer, a diaphragm consisting of two electrets, each metallized on one side and with these metal layers in contact, is sandwiched between two perforated metal electrodes, forming a symmetrical system. Application of a signal voltage \bar{v} in antiphase to the two electrodes causes the electrodes to exert forces F_1 and F_2 on the diaphragm. These forces are given by Eq. (2) if $\pm\bar{v}$ is substituted for v . The net force $F = F_1 - F_2$ on the diaphragm is determined by Eq. (3), which is linear in \bar{v} , indicat-

$$F = \frac{-2\bar{v}\sigma\varepsilon D_1}{(D_1 + \varepsilon D_2)^2} \quad (3)$$

ing the absence of nonlinear distortion. The linearity, however, is maintained only if both air-layer thicknesses remain at the same value D_2 . This is possible only if D_2 is relatively large, which makes such systems less sensitive than the above-described single-ended transducers.

This transducer can be simplified and improved by using an asymmetric system in which a single electret coated on one side with a low-conductivity layer is used. The diaphragm holds an excess charge because of the presence of an induction charge on the front electrode. Thus, application of an ac signal causes a net force between diaphragm and plates that is linear if the resistivity of the coating (or the resistor R) prevents charge equalization during a period of lowest-frequency signal applied.

An interesting modification of the above push-pull transducers can be achieved by using nonmetallized monocharge electrets, which differ from customary electrets in that they carry only a single-polarity charge and no compensation charge. It is possible to form these electrets by electron injection or by the use of wet-contact charging methods. Since a compensating charge is absent, monocharge electrets have a strong external field that is independent of electrode distance. Thus, if such electrets are used in transducer configurations, large air gaps of the order of 1 mm (0.04 in.) can be used without loss in field strength. This in turn results in highly efficient transducers that can operate at large displacement amplitudes. See EARPHONES.

Electromechanical transducers. Apart from their use in electroacoustic transducers, such as microphones and earphones, real-charge electrets have been applied to electromechanical transducers. Examples are touch or key transducers, impact-sensitive line transducers, relay switches, and optical display panels.

Relatively close in their design principles to the

electroacoustic transducers described above are the first two of these groups. The touch or key transducers depend on the manual deflection of an electret relative to an electrode in a microphonelike structure. Large-voltage signals can be generated with such devices. Line transducers consist essentially of a coaxial cable with polarized dielectric in which the center conductor and the shield serve as electrodes. Mechanical excitation resulting in a deformation of the shield at any point along the length of such a cable produces an electrical output signal.

Very different design principles are used in the relay switches and the optical display panels. The relay switches utilize the external fields of electrets mounted in bistable arrangements to open or close contacts. Such devices have power requirements 1000 times smaller than equivalent electromagnetic systems; thus considerable power conservation over present magnetic components is possible. A similar design is used in the optical display panels where illuminated channels are opened and closed by hinged and opaque electrets moved by electrostatic forces. Although relatively high voltages are required, the necessary currents are small and flow only during the switching operation. *See* ELECTRONIC DISPLAY; RELAY.

Hydrophones based on PVDF. Many types of hydrophones based on the piezoelectric polymer PVDF and its copolymers have been described. Among these are hydrophones for the sonic and near-ultrasonic frequency ranges based on the use of thick sheets or cylinders of the piezoelectric material, as well as hydrophones for the higher ultrasonic range, generally designed as membrane or needle hydrophones.

The lower-frequency hydrophones can be relatively simple devices consisting of an electroded film of the piezoelectric polymer that is encapsulated for waterproofing. If the dimensions of the hydrophone are small compared to the wavelength of the sound, equal sound pressure in all directions prevails and the output signal of the transducer is given by the hydrostatic piezoelectric constant. More sophisticated designs include stiffening plates to avoid resonances due to extensional modes of the polymer film, two-layer constructions to improve shielding, and cylindrical designs. Arrays or transducers with shaped elements allow directional characteristics to be obtained. Hydrophones based on these principles generally operate in the frequency range from less than 1 kHz to several MHz.

The membrane hydrophones for the higher ultrasonic range, up to 100 MHz, also consist of one or two peripherally clamped thin sheets of piezoelectric polymer with a diameter of several centimeters, poled in a tiny central section of less than 1 mm (0.04 in.) in diameter, which is the actual sensing area. These hydrophones have the advantage of being more or less transparent to sound, thus allowing sound pressure to be measured without disturbing the sound field by generating unwanted reflections. In the needle hydrophones, a metallic needle serves as the support and inner electrode for a sur-

rounding layer of the piezoelectric polymer, which is covered on its outer side with another electrode. The outside diameter of the needle is generally less than 1 mm (0.04 in.) to 0.3 mm (0.012 in.), thus permitting again measurements without perturbation of the sound field. The frequency range of needle hydrophones can extend up to 80 MHz. *See* HYDROPHONE.

Piezoelectric cellular transducers. The very large piezoelectric thickness coefficient of charged cellular PP can be used to sense sound, force, impact, vibration, or motion. For example, cellular PP has been used in microphones of very simple design. Such transducers consist essentially of only the PP film and suitable housing and shielding. If stacked layers of film are used, the microphone sensitivity increases proportionally to the number of layers. Using this technique, microphones consisting of six layers have been built that show sensitivities of 15 mV/Pa, comparable to the sensitivity of conventional electret microphones. Other features of these transducers are relatively flat frequency response, low noise level, and low distortion. The major advantage is that a minute and critical air gap, as used in electret microphones, is not needed. Cellular PP can also be used for sound generation. Since the generated sound pressure increases with frequency, sound sources are particularly suitable for the ultrasonic range, below the resonance frequency of the PP film that is in the region above 100 kHz.

Important electromechanical transducers based on cellular PP are piezoelectric floor mats, keyboards, control panels, pickups for musical instruments, and biomedical sensors. The floor mats are step- or impact-sensitive and can detect the motion of patients in hospitals, of intruders in industrial plants or military installations, or of vehicles on highways. Because of the inexpensive material, large floor areas may be covered. Push-buttons of cellular PP for keyboards and control panels have the advantage of high sensitivity and mechanical flexibility. The biomedical sensors are used to monitor, for example, heartbeat, breathing, and forces on limbs in humans and animals.

Applications. Because of their favorable properties, simplicity, and low cost, electret transducers have been used in many applications as research tools and in the commercial market.

Among the research applications are microphones for use in acousto-optic spectroscopy, applied to the detection of air pollution and to the study of reaction kinetics of gases and optical absorption of solids. Because of the favorable noise performance of electret microphones, the detection threshold for air pollutants has been lowered by more than an order of magnitude. Other applications of electret microphones have been in aeronautics and shock-tube studies, in which the low vibration sensitivity of these transducers is crucial. The wide frequency range of electret transducers, discussed above, made possible their application in infrasonic atmospheric studies and in ultrasonic investigations of liquids and solids. A broad field of application is in noise

monitoring and noise abatement. Examples are all kinds of mobile and stationary noise measurements as well as active noise reduction. In addition, audio and ultrasonic arrays of electret microphones have been used in near-field acoustic holography. Yet other uses of electromechanical electret transducers are in such diverse areas as microforce detection, vibration analysis, and leak detection in space stations. For example, microforce sensors based on PVDF are used in the micromanipulation and microassembly of MEMS (micro-electro-mechanical systems) structures. In all these applications, the simplicity and reliability of electret transducers are of importance. *See* ACOUSTICAL HOLOGRAPHY; ACTIVE SOUND CONTROL; INFRASOUND; MICRO-ELECTRO-MECHANICAL SYSTEMS (MEMS); NOISE MEASUREMENT; PHOTOACOUSTIC SPECTROSCOPY; SHOCK TUBE; ULTRASONICS.

Of all commercial applications of electret devices, the electret microphone for amateur, professional, studio, and telephone use is most prominent. About 2×10^9 such microphones are produced annually, accounting for about 90% of the entire output of microphones. Almost half of these are used in stationary and mobile telephones. Other uses of electret microphones are in high-fidelity work, such as audio studios, and in MP-3 players, sound-level meters, noise dosimeters, camcorder, and toys. Miniature microphones are used in hearing aids, in telephone operator headsets, and in communication headsets, often as noise-canceling (first-order gradient) transducers. The success of the electret microphone is primarily due to its acoustic quality and low cost. It is noteworthy that, owing to their low vibration sensitivity, electret microphones were the first transducers to be built directly into widely used tape recorders. In addition to electret microphones, there are single-backplate and push-pull electret headphones, also with monocharge electrets. The latter are among the highest-quality headsets available. *See* HEARING AID; MAGNETIC RECORDING; SOUND RECORDING; TELEPHONE.

Apart from their use in electroacoustics, electret devices have found many commercial applications in the underwater, ultrasonic, and electromechanical fields. The PVDF hydrophones in the varieties discussed above for the sonic and ultrasonic frequency ranges have been commercially available for a long time and, because of their many advantages, enjoy an undisputed position as the preferred devices in the field. The same holds for the ultrasonic devices for airborne sound based on space-charge and piezoelectric polymers. In the electromechanical field, transducers based on PVDF and on cellular piezoelectric PP have been introduced commercially. In particular, step-sensitive floor mats, touch-sensitive keyboards and control panels, and biomedical force sensors are available.

Gerhard M. Sessler

Bibliography. Institute of Electrical and Electronics Engineers, *IEEE 12th International Symposium on Electrets*, 2005; G. M. Sessler and R. Gerhard-Multhaupt (eds.), *Electrets*, vols. 1 and 2, 3d ed., Laplacian Press, 1999.

Electric charge

A basic property of elementary particles of matter. One does not define charge but takes it as a basic experimental quantity and defines other quantities in terms of it. The early Greek philosophers were aware that rubbing amber with fur produced properties in each that were not possessed before the rubbing. For example, the amber attracted the fur after rubbing, but not before. These new properties were later said to be due to "charge." The amber was assigned a negative charge and the fur was assigned a positive charge.

According to modern atomic theory, the nucleus of an atom has a positive charge because of its protons, and in the normal atom there are enough extranuclear electrons to balance the nuclear charge so that the normal atom as a whole is neutral. Generally, when the word charge is used in electricity, it means the unbalanced charge (excess or deficiency of electrons), so that physically there are enough "nonnormal" atoms to account for the positive charge on a "positively charged body" or enough unneutralized electrons to account for the negative charge on a "negatively charged body."

The rubbing process mentioned "rubs" electrons off the fur onto the amber, thus giving the amber a surplus of electrons, and it leaves the fur with a deficiency of electrons. *See* ELECTROSTATICS.

In line with the previously mentioned usage, the total charge q on a body is the total unbalanced charge possessed by the body. For example, if a sphere has a negative charge of 1×10^{-10} coulomb, it has 6.24×10^8 electrons more than are needed to neutralize its atoms. The coulomb is the unit of charge in the SI system of units. *See* COULOMB'S LAW; ELECTRICAL UNITS AND STANDARDS.

The surface charge density σ on a body is the charge per unit surface area of the charged body. Generally, the charge on the surface is not uniformly distributed, so a small area ΔA which has a magnitude of charge Δq on it must be considered. Then σ at a point on the surface is defined by the equation below.

$$\sigma = \lim_{\Delta A \rightarrow 0} \frac{\Delta q}{\Delta A}$$

The subject of electrostatics concerns itself with properties of charges at rest, while circuit analysis, electromagnetism, and most of electronics concern themselves with the properties of charges in motion. *See* CAPACITANCE; ELECTRIC CURRENT; ELECTRICITY.

Ralph P. Winch

Electric contact

A part, in an electrical switching device, made of conducting material, for the purpose of closing, opening, or changing the conductive path of an electrical circuit. To open or close a circuit, an electric contact is made to come in contact with or separate from its mating part. Devices embodying contacts for these

purposes are electric switches, relays, contactors, and circuit breakers. Contacts may be actuated directly or through a linkage that is driven either manually, mechanically, electromagnetically, hydraulically, or pneumatically.

An electric contact is also an essential part of a variable resistor. In such an application the contact or wiper provides a moving connection to the resistive element. For such service the contact material is chosen for low abrasion combined with absence of a high-resistance film such as would be formed by an adherent oxide. *See* POTENTIOMETER; RHEOSTAT.

Requirements for contacts differ markedly depending on applications. For contacts in relays and low-power applications, reliability in completing a circuit may be of the utmost importance. Surface films and contaminants are not tolerable. Noble metals such as gold and platinum are sometimes sealed hermetically for this purpose. For the purpose of carrying high continuous currents, electrical resistance must be kept to a minimum. Contacts with high silver content are usually used. High contact forces, up to hundreds of pounds, may be applied when the contacts are in the closed position. For contacts that must interrupt high currents, refractory materials such as tungsten and molybdenum are often used. These materials also possess antiwelding characteristics, which are essential for contacts that must close in on short-circuit currents and later open under normal operating force.

On many occasions all the aforementioned properties are desired in one material. Special alloys such as silver cadmium oxide or refractory materials impregnated with silver have been developed for this purpose.

A related electrical part is the brush on rotating machines. Pairs of brushes provide continuity between the external circuit and slip rings or commutator segments on the rotating portion of the machine. Similar arrangements of brushes and slip rings connect to continuously rotating radar antenna mounts. Graphite, because of its inherent lubricating nature, is a usual constituent of such sliding contacts for high-current operation; corrosion-resistant metals or metal alloys are used for low-current applications. *See* CIRCUIT (ELECTRICITY); ELECTRIC ROTATING MACHINERY.

Thomas H. Lee

Bibliography. ARS Electronics, *Receiving Tube Characteristics and Socket Connection Guide*, 1981; CES Industries Inc. Staff, *Contactors Sensor Operation*, 1982; R. Holm and E. Holm, *Electric Contacts*, 4th ed., 1967; Institute of Electrical and Electronics Engineers, *1998 IEEE 44th Holm Conference on Electrical Contacts*, 1998.

Electric current

The net transfer of electric charge per unit time. It is usually measured in amperes. The passage of electric current involves a transfer of energy. Except in the case of superconductivity, a current always heats the medium through which it passes. For a

discussion of the heating effect of a current *see* JOULE'S LAW.

On the other hand, a stream of electrons or ions in a vacuum, which also may be regarded as an electric current, produces no local heating. Measurable currents range in magnitude from the nearly instantaneous 10^5 or so amperes in lightning strokes to values of the order of 10^{-16} A, which occur in research applications.

All matter may be classified as conducting, semi-conducting, or insulating, depending upon the ease with which electric current is transmitted through it. Most metals, electrolytic solutions, and highly ionized gases are conductors. Transition elements, such as silicon and germanium, are semiconductors, while most other substances are insulators.

Electric current may be direct or alternating. Direct current (dc) is necessarily unidirectional but may be either steady or varying in magnitude. By convention it is assumed to flow in the direction of motion of positive charges, opposite to the actual flow of electrons. Alternating current (ac) periodically reverses in direction.

Conduction current. This is defined as the transfer of charge by the actual motion of charged particles in a medium. In metals the current is carried by free electrons which migrate through the spaces between the atoms under the influence of an applied electric field. Although the propagation of energy is a very rapid process, the drift rate of the individual electrons in metals is only of the order of a few centimeters per second. In a superconducting metal or alloy the free electrons continue to flow in the absence of an electric field after once having been started. In electrolytic solutions and ionized gases the current is carried by both positive and negative ions. In semiconductors the carriers are the limited number of electrons which are free to move, and the "holes" which act as positive charges.

Displacement current. When alternating current traverses a capacitor, there is no physical flow of charge through the dielectric (insulating material, gas, or vacuum), but the effect on the rest of the circuit is as if there were a continuous flow. Energy can pass through the capacitor by means of the so-called displacement current, and this displacement current generates a surrounding magnetic flux, just as an orthodox current would. James Clerk Maxwell introduced the concept of displacement current in order to make complete his theory of electromagnetic waves. *See* ALTERNATING CURRENT; CAPACITANCE; CONDUCTION (ELECTRICITY); DIELECTRIC MATERIALS; DIRECT CURRENT; DISPLACEMENT CURRENT; ELECTRIC INSULATOR; ELECTRICAL RESISTANCE; FREE-ELECTRON THEORY OF METALS; SEMICONDUCTOR; SUPERCONDUCTIVITY.

John W. Stewart

Electric distribution systems

Systems that comprise those parts of an electric power system between the subtransmission system and the consumers' service switches. It includes distribution substations; primary distribution feeders;

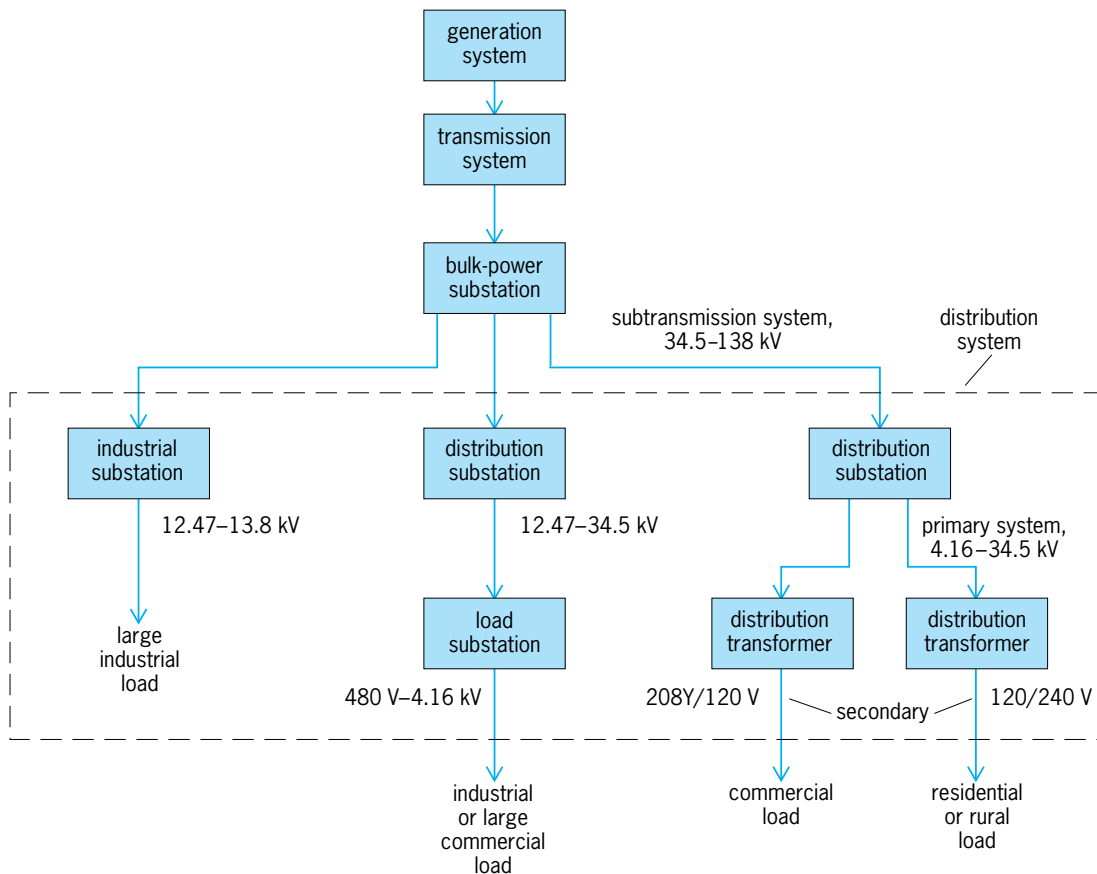


Fig. 1. Overview of the power system from generation to consumer's switch.

distribution transformers; secondary circuits, including the services to the consumer; and appropriate protective and control devices. Sometimes, the subtransmission system is also included in the definition.

The subtransmission circuits of a typical distribution system (Fig. 1) deliver electric power from bulk power sources to the distribution substations. The subtransmission voltage is usually between 34.5 and 138 kV. The distribution substation, which is made up of power transformers together with the necessary voltage-regulating apparatus, bus-bars, and switchgear, reduces the subtransmission voltage to a lower primary system voltage for local distribution. The three-phase primary feeder, which usually operates at voltages from 4.16 to 34.5 kV, distributes electric power from the low-voltage bus of the substation to its load center, where it branches into three-phase subfeeders and three-phase and occasionally single-phase laterals. Most of the three-phase distribution system lines consist of three-phase conductors and a common or neutral conductor, making a total of four wires. Single-phase branches (made up of two wires) supplied from the three-phase mains provide power to residences, small stores, and farms. Loads are connected in parallel to common power-supply circuits. See ALTERNATING CURRENT.

Subtransmission. The subtransmission system is that part of the power system which delivers power from bulk power sources, such as large transmission substations. The subtransmission may be made up of

overhead open-wire construction on wood poles or underground cables. Although the voltages of these circuits can occasionally range from 12.47 to 345 kV, the majority are at the 69-, 115-, and 138-kV voltage levels. There is a trend toward the use of the higher voltages as a result of the increasing use of higher transmission voltages. The subtransmission-system designs vary from simple radial systems to a subtransmission network. The major considerations affecting the design are cost and reliability. See CONDUCTOR (ELECTRICITY); ELECTRIC POWER TRANSMISSION.

Substation. The distribution substation is an assemblage of equipment for the purpose of switching, regulating, and changing the supply voltage from subtransmission to primary distribution. More important substations are designed so that the failure of a piece of equipment in the substation or one of the subtransmission lines to the substation will not cause an interruption of power to the customer. The location of a substation is dictated by the voltage levels, voltage regulation considerations, subtransmission costs, substation costs, and the costs of primary feeders, mains, and distribution transformers. It is also restricted by esthetic considerations. A typical substation may have power transformers, circuit breakers, disconnecting switches, station buses and insulators, current-limiting reactors, shunt reactors, current transformers, potential transformers, capacitor voltage transformers, coupling capacitors,

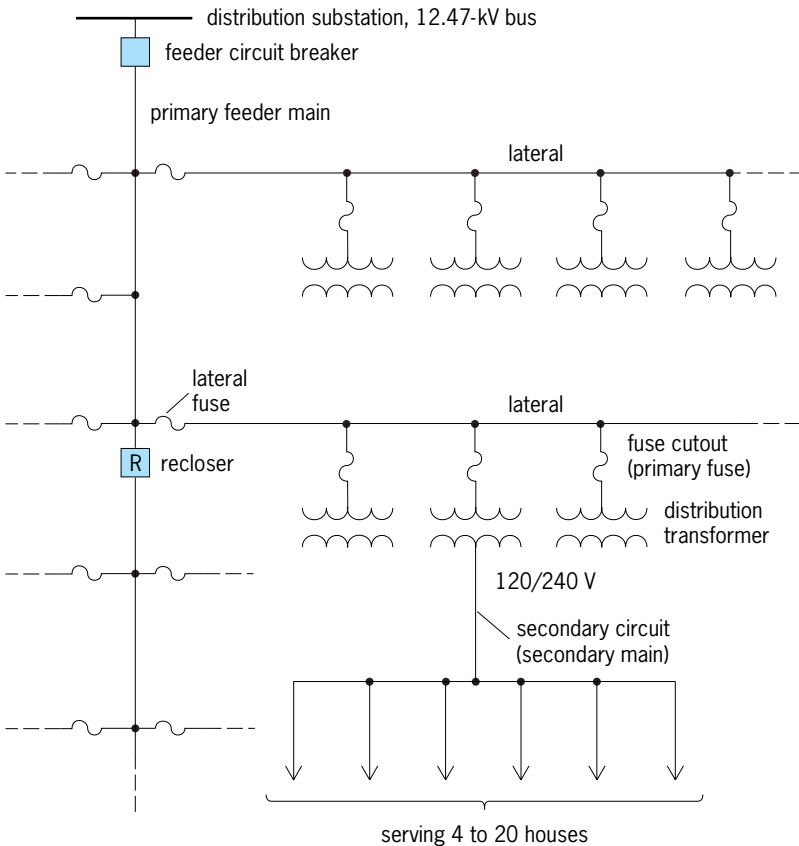


Fig. 2. Typical power distribution feeder.

protective relays, lightning arresters, station batteries, and other equipment. See ELECTRIC POWER SUBSTATION.

Primary system. The part of the electric utility system which is between the distribution substation and distribution transformers is called the primary system. It is made up of circuits known as primary feeders or distribution feeders. A feeder includes the main or main feeder, which usually is a three-phase four-wire circuit, and branches or laterals, which usually are single-phase or three-phase circuits, tapped off the main (Fig. 2). A feeder is usually sectionalized by means of reclosing devices at various locations so as to remove as little as possible of a faulted circuit and to hinder service to as few consumers as possible. This goal is achieved through the coordination of all the fuses and reclosers.

The primary system leaving the substation is most often in the 11,000–15,000-V range. A specific voltage used is 12,470-V line-to-line and 7200-V line-to-neutral (conventionally written 12,470Y/7200 V). Some utilities still use a lower voltage, such as 4160Y/2400 V. However, the use of primary distribution circuits in the 25- and 35-kV classes is increasing; all are four-wire systems. Single-phase loads are connected line-to-neutral on the four-wire systems. Congested and heavy-load locations in metropolitan areas as well as new residential areas are customarily served by underground primary feeders made up of radial three-conductor cables.

The primary system design is affected by the cost, importance of the load, and the required level of service continuity and reliability. The simplest and least expensive, and thus the most common, is the radial design. The low reliability of this design is improved by a modification called the loop-type primary design, in which the feeder loops through the feeder load area and returns to the substation bus. The most expensive and reliable design is the primary-network type, in which a system of interconnected feeders is supplied by a number of substations. Such a system supplies a load from several directions and consequently provides the maximum reliability and quality of service.

Secondary system. This is the part of the electric power system between the primary system and the consumer's property. The secondary distribution system includes distribution transformers, secondary circuits (secondary mains), consumer services (or service drops), and meters to measure consumer energy consumption. Generally, the secondary systems are designed in single phase for areas of residential customers and in three phase for areas of industrial or commercial customers with high load densities. Secondary voltages are provided by distribution transformers that are connected to the primary system, and they usually are associated with utilization voltages. Residential and most rural customers are supplied by 120/240-V single-phase three-wire systems. Commercial and small industrial customers are supplied by 208Y/120-V or 480Y/277-V three-phase four-wire systems. The secondary voltage also supplies multiple street lights.

Spot networks. These are used in downtown areas for high-rise buildings with extremely high load densities, and sometimes also for areas of industrial or commercial customers. Often, large shopping centers are supplied by spot networks. In such a network, all transformers and protecting apparatus are placed at the same location.

Good voltage. Essentially, good voltage means that the average voltage level is correct, that variations are within acceptable voltage limits, and that sudden and momentary changes in voltage level do not cause objectionable light flicker. Of course, the utilization voltage of the customer changes with changing load on the system as well as the location of the customer on the system. However, voltage variation is usually less than 5% at the customer's meter.

Good continuity. Service continuity is the provision of uninterrupted electric power to the customer. Good continuity means that such service is almost always provided.

Back-up systems. To increase the service reliability for critical loads, such as hospitals, computer centers, and crucial industrial loads, back-up systems, such as emergency generators or batteries, with automatic switching devices are provided. Computer installations are backed up by uninterruptible power supplies (UPS). See UNINTERRUPTIBLE POWER SUPPLY.

Overhead construction. Most existing distribution systems in residential, industrial, and rural areas are

of overhead construction. The distribution transformers are mounted near the tops of poles, and bare primary and secondary conductors are strung from pole to pole.

Underground construction. Almost all residential developments now being built are served by underground residential distribution (URD) systems using underground cables for esthetic reasons. The URD system is free from service outages and accompanying repair expenses caused by lightning, rain, ice, sleet, snow, and wind storms. Also, the probability of hazard caused by building fires and by motor vehicles and other foreign objects coming into contact with the system elements is far less than what might be expected on overhead systems. Other costs, such as the cost of tree trimming, are totally eliminated by using the URD systems. In general, underground distribution systems cost more than comparable overhead systems. However, greater public interest in the esthetic appearance of residential communities and reductions in the cost of underground equipment and installations have made the use of URD systems popular. Residential underground systems normally do not have the refinements and operating advantages found in conventional downtown commercial underground systems. Extensive duct banks, transformer vaults, access holes, and submersible equipment are not required in residential systems.

Cables. In the past the most frequently used underground cable was stranded copper with an oil-impregnated paper insulation wrapped in a lead sheath. It is expensive to buy, install, and splice, but it has a long service life. It is available in single- and three-phase cable and in voltages up to 46 kV. Practical residential underground distribution was made possible by the development of direct-buried underground cable that used a relatively inexpensive plastic insulation, typically cross-linked polyethylene (XLPE), and practical elbow connectors. The insulated cable has stranded conductors surrounding a bare concentric neutral. Other insulating plastics include high-molecular-weight polyethylene (HMWP), sometimes called high-density track-resistant polyethylene (HDPE), and ethylene propylene rubber (EPR). The use of such plastics as the protective sheath for cables has made practical their direct burial in the ground without the need for ducts and access holes. Splicing procedures for such cables are also simple and less costly. Long lengths of cable capable of being plowed directly into the ground or placed in narrow and shallow trenches naturally reduce installation and maintenance costs. Further savings may be realized if other facilities, such as telephone, cable television, gas, and water pipes, are installed simultaneously.

System design. In general, the heavy three-phase feeders are located overhead along the periphery of a residential development, and the laterals to the pad-mount transformers are buried about 40 in. (100 cm) deep. The lateral cable conductors come up from underground into the high-voltage compartment of the distribution transformer. The secondary service lines then run to the individual dwellings at a depth of

about 24 in. (60 cm), and come up into the dwelling meter through a conduit. The service conductors run along easements and do not cross adjacent property lines.

Faults. The frequency of faults is lower on underground systems than on overhead systems. Also, faults are not so likely to “bunch up” due to storm conditions. However, faults are much more difficult and time-consuming to find, isolate, and repair on underground systems. Service-restoration requirements have dictated primary-system designs which operate as a normally open loop. In the event of a cable fault, such a design facilitates the location and isolation of the failure and faster service restoration to all customers on the unfaulted segment of the primary loop. A well-designed URD system should permit ease of sectionalizing in order to isolate a faulted portion, and should be flexible enough to accept load growth with a minimum number of changes. The looped primary and single-phase banked secondary system provide these requirements at the lowest cost. The banked secondary system takes better advantage of the diversity among loads and allows the use of transformers with a smaller kilovolt-ampere capacity and possibly fewer transformers. It also offers better average voltage conditions at the loads and less voltage dip due to motor starting for a given transformer rating, conductor size, and spacing between transformers than the simple-radial system offers. *See* VOLT-AMPERE.

Transformers. The transformers used for the URD system are hermetically sealed against moisture, including their bushings and terminals. The terminals may contain one or more insulated disconnecting elbows, enabling the disconnection of the transformer and the sectionalizing of the primary circuits. The transformers may be completely buried or installed on ground-level pads or semiburied pads. Pad-mounted transformers are the predominant type of transformer in URD systems. They are usually installed on concrete slabs, or pads. They may utilize “dead-front” configuration with separable insulated connectors, or elbows. Many other combinations of pad-mount construction and accessory equipment are available. Residential subsurface transformers are used much less frequently. They are installed in relatively tight-fitting vaults with the cover grating of the vault at ground level. Direct-buried transformers are not popular because of the problem of location of accessory equipment and its operation under adverse weather conditions such as snow and ice. *See* TRANSFORMER.

Dispersed storage and generation. Electric distribution systems are undergoing major changes to accommodate the development of alternative sources for generating electric energy and enhanced opportunities for small power producers and cogenerators. Usually, these generators are small (typically ranging from 10 kW to 10 MW and connectable to either side of the meter) and can be economically connected only to the distribution system. If properly planned and operated, dispersed storage and generation may provide benefits to distribution

systems by reducing capacity requirements, improving reliability, and reducing losses. Dispersed-storage-and-generation technologies include hydro-electric, diesel generators, wind-electric systems, solarelectric systems, batteries, storage space and water heaters, storage air conditioners, hydroelectric pumped storage, photovoltaics, and fuel cells. See ELECTRIC POWER SYSTEMS; ENERGY SOURCES.

Turan Gonen

Bibliography. American National Standards Institute, *Preferred Voltage Ratings for A-C Systems and Equipment: Guide for (EEI R-6-1949) (IEC38)*, ANSI C84.1, 1989; M. T. Bishop, A. G. Jones, and W. F. Israel, Overcurrent protection alternatives for underground distribution systems, *IEEE Trans. Power Delivery*, 10(1):69-77, January 1995; T. Gonen, *Electric Power Distribution System Engineering*, 1986; A. J. Pansini, *Guide to Electrical Power Distribution Systems*, 1992; W. L. Weeks, *Transmission and Distribution of Electrical Energy*, Harper and Row, New York, 1981.

Electric field

A condition in space in the vicinity of an electrically charged body such that the forces due to the charge are detectable. An electric field (or electrostatic field) exists in a region if an electric charge at rest in the region experiences a force of electrical origin. Since an electric charge experiences a force if it is in the vicinity of a charged body, there is an electric field surrounding any charged body.

Field strength. The electric field intensity (or field strength) \mathbf{E} at a point in an electric field has a magnitude given by the quotient obtained when the force acting on a test charge q' placed at that point is divided by the magnitude of the test charge q' . Thus, it is force per unit charge. A test charge q' is one whose magnitude is small enough so it negligibly alters the field in which it is placed. The direction of \mathbf{E} at the point is the direction of the force \mathbf{F} on a positive test charge placed at the point. Thus, \mathbf{E} is a vector point function, since it has a definite magnitude and direction at every point in the field, and its defining equation is (1).

$$\mathbf{E} = \frac{\mathbf{F}}{q'} \quad (1)$$

Ralph P. Winch

Principle of superposition. The field at a distance of r meters from a point charge of q coulombs in vacuum is given by Eq. (2), where \mathbf{E} is the field

$$\mathbf{E} = \frac{q}{4\pi\epsilon_0 r^2} \quad (2)$$

strength in volts per meter and $\epsilon_0 = 8.85 \times 10^{-12}$ farad/meter is the permittivity of free space. \mathbf{E} is a vector, directed along the radius vector from the point charge; positive \mathbf{E} is directed away from a positive charge or toward a negative charge. For an assembly of charges, the resultant field is, by the principle of superposition, the vector sum of the field

components due to the individual charges. This summation may be performed directly, but in practical cases Gauss' theorem often affords a more powerful and convenient method. See GAUSS' THEOREM.

A. E. Bailey

Electric displacement. Electric flux density or electric displacement \mathbf{D} in a dielectric (insulating) material is related to \mathbf{E} by either of the equivalent equations (3), where \mathbf{P} is the polarization of the

$$\mathbf{D} = \epsilon_0 \mathbf{E} + \mathbf{P} \quad \mathbf{D} = \epsilon \mathbf{E} \quad (3)$$

medium, and ϵ is the permittivity of the dielectric which is related to ϵ_0 by the equation $\epsilon = \epsilon_r \epsilon_0$, ϵ_r being the relative permittivity of the dielectric. In empty space, $\mathbf{D} = \epsilon_0 \mathbf{E}$. The units of \mathbf{D} are coulombs per square meter.

In addition to electrostatic fields produced by separations of electric charges, an electric field is also produced by a changing magnetic field. The relationship between the \mathbf{E} produced and the rate of change of magnetic flux density $d\mathbf{B}/dt$ which produces it is given by Faraday's law of induced electromotive forces (emfs) in Eq. (4), where $d\mathbf{s}$ is a vector element

$$\oint \mathbf{E} \cdot d\mathbf{s} = - \int_A \frac{d\mathbf{B}}{dt} \cdot d\mathbf{A} \quad (4)$$

of path length directed along the path of integration in the general sense of \mathbf{E} . Thus $\oint \mathbf{E} \cdot d\mathbf{s}$ is the emf induced in this closed path of integration. The area of the surface bounded by the path of integration is A , and the direction of $d\mathbf{A}$, an infinitesimal vector element of this area, is the direction of the thumb of the right hand when the fingers encircle the path of integration in the general sense of \mathbf{E} . The right side of Eq. (4) is seen to be the negative of the time rate of change of the magnetic flux linking the path of integration chosen for the left side.

In an electrostatic field, $\oint \mathbf{E} \cdot d\mathbf{s}$ is always zero. See ELECTRIC CHARGE; ELECTROMAGNETIC INDUCTION; POTENTIALS.

Ralph P. Winch

Bibliography. D. Halliday, R. Resnick, and J. Walker, *Fundamentals of Physics*, 7th ed., 2005; E. M. Purcell, *Electricity and Magnetism*, 2d ed., 1985; H. D. Young and R. A. Freedman, *Sears and Zeman-sky's University Physics*, 11th ed., 2003.

Electric filter

A transmission network used to selectively modify the components of a signal according to their frequencies. In most cases a filter is used to enhance signals of desired frequencies while suppressing signals of undesired frequencies. An ideal filter would pass only desired frequencies while completely suppressing all unwanted frequencies, without any dispersion in time of the frequencies. Unfortunately, ideal filters are impossible to achieve.

Applications. Electric filters are used in most electronic communication systems. Whether communication is over wire, free space, or optical fiber, multiple channels of information can be multiplexed on different frequency bands. Unwanted signals

and noise are introduced along the communications path. The main function of electric filters is to separate the desired signal or channel from all others and from any noise or interference. For example, an AM radio receiver may have a low-pass filter after the antenna to separate the AM frequency band from higher frequency bands and, elsewhere in the radio, a band-pass filter to select the desired station out of the AM band. See ELECTRICAL COMMUNICATIONS; MULTIPLEXING AND MULTIPLE ACCESS; RADIO RECEIVER.

Although electronic filters are commonly thought of as devices for conferring selectivity to communication paths, they are used in almost every part of electronic equipment, such as the damping element in phase-locked loops, cleanup devices for frequency sources, and pulse expansion and compression devices for radar. One of the simplest and most common filters is the bypass capacitor used to restrict high-frequency electronic noise. See ELECTRICAL NOISE; PHASE-LOCKED LOOPS; RADAR.

Transfer function. The performance of a filter is, as is any transmission network, characterized by a transfer function. A transfer function is the ratio of the output signal to the input signal. In many filter types, including those composed of lumped elements, the transfer function can be expressed as a ratio of two polynomials, shown in Eq. (1). Here H is the trans-

$$H(s) = \frac{b_0 + b_1s + b_2s^2 + \dots + b_ns^n}{1 + a_1s + a_2s^2 + \dots + a_ns^n} \quad (1)$$

fer function, a_i and b_i are the coefficients of the two polynomials, and $s = j\omega$, ω being the frequency in radians and $j = \sqrt{-1}$. The zeros of the denominator of the transfer function are called the poles or natural frequencies of the filter. The zeros of the nu-

merator are referred to as the transmission zeros or occasionally finite attenuation poles. Attenuation is the reciprocal of the transfer function or $1/H(s)$. Both the poles and the zeros of the filter can be complex and, for stability considerations, the poles must lie in the left half of the complex s -plane. See COMPLEX NUMBERS AND COMPLEX VARIABLES.

The characteristics of a filter can be analyzed from its transfer function. The transfer function can be expressed in a polar form, shown in Eq. (2), where

$$H(j\omega) = |H(j\omega)|e^{j\theta(\omega)} \quad (2)$$

$|H(j\omega)|$ is the magnitude and $\theta(\omega)$ is the phase function. The square of the magnitude is equal to the ratio of the output power to maximum available power. The power gain of the filter, expressed in decibels, is given in Eq. (3). Loss would be defined as the

$$\alpha(\omega) = 20 \log_{10} |H(j\omega)| \quad \text{dB} \quad (3)$$

negative of gain. For example, the gain function of a band-pass filter approximates the ideal amplitude response, wherein the gain is constant within the passband and frequencies outside the passband are completely suppressed (Fig. 1a). The phase function of a typical band-pass filter decreases as the frequency is increased through the passband (Fig. 1b). Another important parameter is group delay (Fig. 1c), defined in Eq. (4). Group delay is used to character-

$$\tau(\omega) = \frac{d\theta(\omega)}{d(\omega)} \quad (4)$$

ize how different frequencies are dispersed in time. Digital data are often transmitted by using shifts in frequency or phase. If the delay variation is too great,

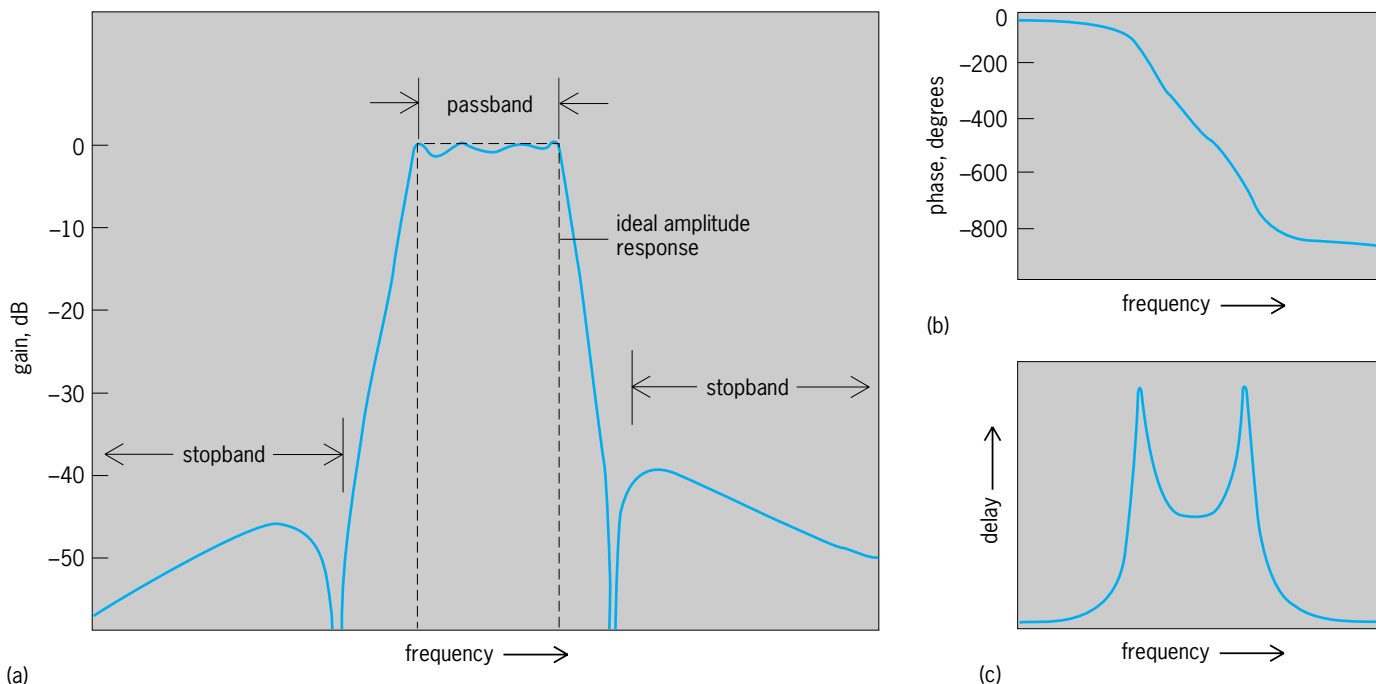


Fig. 1. Characteristics of a typical band-pass filter. (a) Gain characteristics, compared with those of an ideal filter. (b) Phase plot. (c) Delay plot.

these shifts will be detected incorrectly. See DECIBEL; GAIN; RESPONSE.

Characteristics. Filters are characterized in multiple ways. The expression low-pass, Butterworth, *LC* describes a filter. The descriptor low-pass indicates the relation of the passed to the rejected frequencies. Butterworth describes the type of polynomials in the transfer function. *LC* indicates the construction method. This filter is made of inductors (*L*'s) and capacitors (*C*'s).

Filter types. Filters are classified by the relationship of the frequencies that are selectively passed, referred to as the passband, to those which are attenuated, referred to as the stopband. An ideal low-pass filter passes all frequencies below a specified cutoff frequency and rejects those above. A high-pass filter does the opposite. An ideal band-pass filter will pass a band of frequencies while rejecting all others (Fig. 1*a*); a band-reject filter will reject a band of frequencies and pass all others.

An all-pass filter passes all frequencies but does, however, modify the time delay characteristics. It normally corrects delay distortions caused by other sections of a communication path.

All the above classifications are based on frequency-domain considerations. In addition, there are two terms that apply to the time-response characteristics of a filter. A finite impulse response (FIR) filter, when exposed to a change in input, will settle to a steady state within a finite amount of time. An infinite impulse response (IIR) filter will continue oscillating in a decaying manner forever.

A further consideration in classifying a filter is whether the frequency response is constant in time or varies. If it varies with time, as a function of the input signal, it is known as an adaptive filter. This type of filter finds use in speech and image enhancement and echo cancellation.

Approximations. Realizable filters are only approximations to the ideal. An ideal filter has instantaneous cutoff from the passband to the stopband. It passes all the desired frequencies with equal amplitude and delay at each frequency. Actual filters have a band of frequencies where the filter transitions from the passband to the stopband. The passband amplitude may not be flat but can have ripple (Fig. 1*a*). In the stopband of a filter that includes attenuation poles, the amplitude response will flare back up between those attenuation poles. The filter can also be dispersive, where different frequencies have different time delays.

The transfer function is the ratio of two polynomials. At least one of these polynomials would require an infinite number of terms to meet the instantaneous cutoff requirements of the ideal filter. However, there are several classes of polynomial that do approximate the ideal filter. These functions are used to synthesize actual filters.

The Chebyshev approximation to the ideal filter is a polynomial function optimized so that for a given stopband rejection the maximum variation in the passband is a minimum. The numerator of the transfer function of a Chebyshev filter would be 1, and

it would have equal ripple in the passband; that is, all the peaks in the amplitude would be of the same height and all the dips of the same depth. An elliptic-integral filter, also known as a Cauer filter, has both equal ripple in the passband and equal minima of attenuation in the stopband; that is, the heights of the flarebacks between the attenuation poles are all the same. Both the numerator and denominator polynomials of the transfer function are of order greater than zero. The passband has the same criteria as the Chebyshev and, in addition, the minimum attenuation in the stopband is maximized. For polynomials of a given order, Chebyshev and Cauer filters provide the least ripple for the maximum stopband selectivity. Higher ripple can be traded for an increase in selectivity. However, an increase in delay variation and a degradation in stability will also occur. Increased stability and lower delay variation is obtained by allowing the ripple to decrease. By decreasing it in a Chebyshev design to zero, where the derivatives of $H(s)$ at center frequency become zero, a Butterworth filter is obtained. The Butterworth polynomial is the most stable. Changes in its coefficients or the element values used to realize Butterworth response cause the least changes in ripple of all the polynomial types.

In many applications the time-domain characteristics are as important as those of the frequency domain. In radar the key information is provided by pulse peaks. In many digital systems the information is encoded in the phase of a signal. If all the frequencies passing through a filter have the same time delay, the signal shape will remain unchanged and the encoded information will be recoverable. A filter with a gaussian-function amplitude response will have a linear phase and a constant delay characteristic. The normalized transfer function for a low-pass filter is given by Eq. (5). This function is not a polynomial

$$|H(j\omega)| = e^{-\omega^2/2} \quad (5)$$

and not physically realizable. Approximations that try to match the amplitude or delay response of this function are known collectively as gaussian filters. The Bessel filters, also known as Thomson filters, are members of this group. They have maximally flat delay in the sense that the derivatives of delay at center frequency vanish. Other members of this group try to match the gaussian function only in the passband portion of their response. This characteristic allows increased selectivity in the stopband. Others try to provide constant delay in the Chebyshev sense by having equal delay ripple. See DISTRIBUTION (PROBABILITY).

Filter realization. Because filters are used over wide frequency and bandwidth ranges and with such varying performance criteria, many methods have been devised for creating a filter function (see table).

A term used to measure the quality of components in passive filters is Q . The quantity Q is defined as the ratio of energy lost in one cycle to total energy stored. It is a measure of how much loss an element has. An inductor with an inductance of L henrys and

Application ranges of various filter types		
Filter type	Frequency range	Bandwidth (% of center frequency)
LC	0–5 GHz	>0.5
Crystal	5 kHz–300 GHz	0.01–5
Ceramic	300 kHz–50 MHz	0.5–20
Mechanical	1 kHz–700 kHz	0.02–10
Surface-acoustic-wave (SAW) resonator	30 kHz–2 GHz	0.01–0.1
SAW transversal	30 kHz–2 GHz	0.2–50
Active	<100 kHz	>0.5
Switched capacitor	<40 kHz	>0.5
Digital	<10 MHz	>0.1
Microwave	0.5 GHz–200 GHz	>0.5
Ceramic dielectric resonator (CDR)	0.7 GHz–5 GHz	>0.1

series parasitic resistance of R ohms will have a Q equal to $\omega L/R$, where ω is the angular frequency. Filter elements with low Q values increase the loss of a filter and limit the minimum bandwidth that may be achieved. See Q (ELECTRICITY).

LC and RC filters. These filters are based on lumped element inductors (L 's), capacitors (C 's), and resistors (R 's). The transfer function for a typical RC filter (Fig. 2a), which contains one pole at $-1/RC$ on the real axis and no zeros, is given in Eq. (6). Because the

$$H(s) = \frac{1}{1 + sCR} \quad (6)$$

poles of RC filters can be located only on the negative real axis, only low-pass responses can be achieved. To achieve complex poles and zeros, inductors must be used (Fig. 2b). Inductors have relatively low Q 's, and at low frequencies are bulky. The Q 's of typical inductors range from ten to a few thousand, and the larger the Q , the larger the inductor. Many of the other filter types were created to avoid the problems with inductors.

Most LC filters are designed as ladder structures with resistive termination. A ladder structure consists of alternating series and shunt arms (Fig. 3a). The poles and zeros of the transfer function are provided by resonant circuits within the arms. The ability to control these resonances gives the ladder structure stability in both the passband and stopband. An alternate to a ladder would be a multipath structure. To achieve high rejection in the stopband requires that the amplitude and phase be closely controlled at the points where the paths join in order that the signals from the separate circuits will cancel. Therefore components must have tighter tolerances than those needed to realize a ladder design. Therefore,

ladder structures tend to predominate in LC filter design. In crystal-filter design, however, because of the closeness of the impedance pole and zero of crystals, multipath (mainly lattice) structures are common. See RESONANCE (ALTERNATING-CURRENT CIRCUITS).

In a ladder structure a transfer-function pole occurs if a series arm has an impedance zero or a shunt arm has an impedance pole. Transfer-function zeros occur where a series impedance arm has a pole, or a shunt arm a zero. A capacitor connected in series or parallel with an inductor would be the simplest circuit for creating either an impedance zero or a pole. An impedance structure that contains both a finite-frequency pole and a zero consists of a capacitor (C_1), in series with an inductor capacitor, (L_1) and a second capacitor (C_2), connected in parallel (Fig. 3b). This circuit also contains a pole at zero frequency and a zero at infinite frequency.

An LC structure is created by choosing a transfer function, of the form shown in Eq. (1), to meet the filter's specification. This function may be a Chebyshev or other standard type or one created by the designer that is optimum for the specification. By finding the zeros of the numerator and denominator polynomials of the transfer function, the numbers of transfer-function poles and zeros at zero, finite, and infinite frequency are determined. At this point an LC circuit structure can be synthesized. By factoring out poles or zeros sequentially from the transfer function, the type of arm and component values required can be determined. As the process continues, the ladder gains arms as the transfer function describing the residual portion of the circuit is reduced in complexity.

Acoustic filters. This class of filters includes crystal, ceramic, mechanical, and surface-acoustic-wave (SAW) filters. These devices convert electrical energy

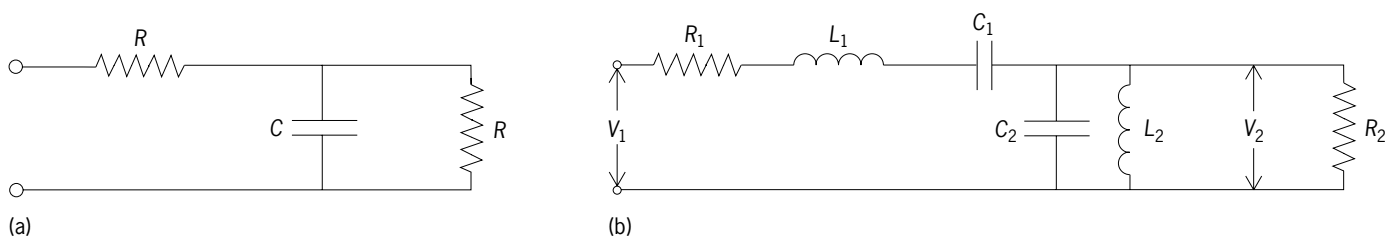


Fig. 2. Passive, lumped-element filters. (a) RC low-pass filter. (b) LC band-pass filter.

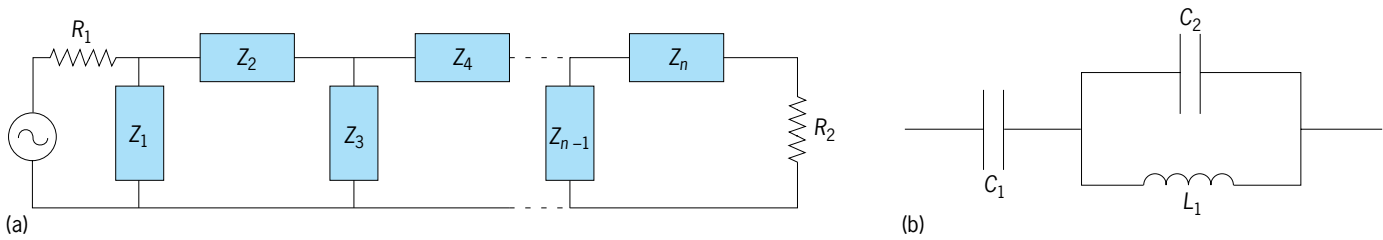


Fig. 3. Ladder structure. (a) Overall structure with alternating series and shunt arms. (b) LC circuit whose impedance contains both a finite-frequency pole and a zero. Such circuits are used in the series and shunt arms of the ladder structure.

to mechanical vibrations, process the signal acoustically, and then convert the energy back to an electrical form. The equations describing a mechanically vibrating resonator, where energy is cycled between kinetic motion and stress, match those of an inductor and capacitor (LC) attached in parallel, where energy is cycled between the electric field of the capacitor and the magnetic field of the inductor. However, the mechanical resonators have much higher Q 's and better stability than the LC circuit. With the addition of a transducer to convert electrical energy to acoustic, the LC circuits can be replaced with mechanical resonators. See RESONANCE (ACOUSTICS AND MECHANICS); TRANSDUCER.

SAW filters can be built in both transversal structures and resonator structures. Crystal, ceramic, and mechanical filters are realized only as resonator filters. Acoustic resonator filters use mechanical resonators coupled together either electrically or acoustically. Each resonator provides either a pole or a zero in the transfer function. In addition, the piezoelectric coupling of mechanical to electrical energy provides a second resonance that may be used to create addi-

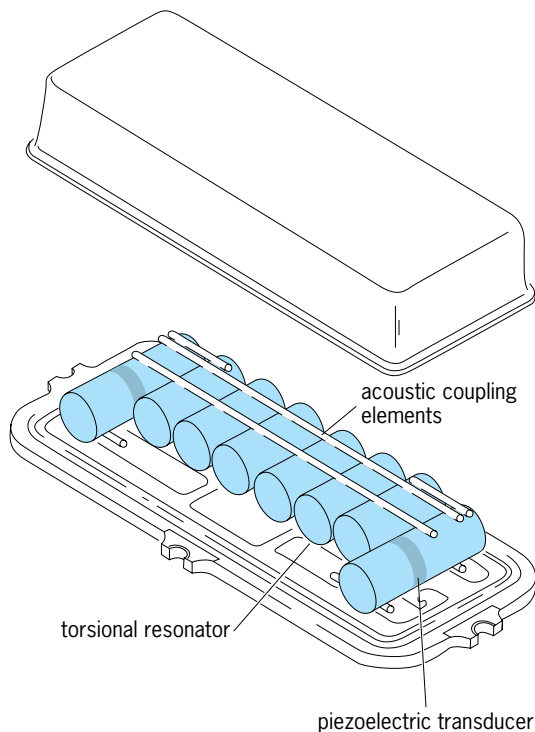


Fig. 4. Mechanical filter.

tional poles and zeros. See HELMHOLTZ RESONATOR; PIEZOELECTRICITY.

The difference between crystal and ceramic filters is mostly in the material that they use for their resonators. In crystal filters the resonators are made from quartz. It is a piezoelectric material with excellent temperature stability and high Q . Ceramic filters use other piezoelectric materials, mainly a lead-zirconate-titanate (PZT) ceramic, for their resonators. These materials have much higher coupling coefficients than quartz and therefore can achieve much wider bandwidths, but have lower Q 's and poorer temperature stability.

Mechanical filters (Fig. 4) use metal for their resonant elements. Coupling between elements is done mechanically. Piezoelectric or other transducer types are used at the input and output to convert electrical energy to acoustic energy. Mechanical filters have temperature stabilities and Q 's between those of quartz and ceramic filters.

SAW resonators trap surface acoustic waves between two reflecting structures. Interdigital fingers are used to couple electrical energy piezoelectrically into and out of the acoustic wave. Transversal SAW filters use a transmitting and receiving set of piezoelectric interdigital fingers. One or both of the finger sets are weighted to create a desired impulse response. The spacings between interdigital fingers act as time delays for the acoustic wave. These time delays and weighting of the fingers are analogous to processes used in digital filters. Because of this analogy, many of the algorithms that are used for transversal digital filters have been applied to SAW filters. See SURFACE-ACOUSTIC-WAVE DEVICES.

Other filter realizations. Many techniques are used to create filters. Inductors are replaced with transistor networks in active filters, discussed below, to reduce size and cost. By using an analog-to-digital converter the transfer function can be created mathematically by a digital processor. At high frequencies, transmission lines and waveguide structures replace lumped elements. See DIGITAL FILTER; INTEGRATED-CIRCUIT FILTER; MICROWAVE FILTER; SWITCHED CAPACITOR CIRCUIT.

Donald P. Havens

Active filters. An active filter comprises resistors, capacitors, and active elements such as operational amplifiers. It is also referred to as an active- RC filter.

Advantages and disadvantages. Active filters can realize the same filter characteristics as passive ones comprising resistor, capacitor, and inductor elements.

They have, however, several advantages over their passive counterparts:

1. Active filters can provide gain, and are frequently used to simultaneously match filtering (frequency-determining) and gain specifications.

2. They are readily implemented in integrated-circuit technology, whereas the inductor element of passive filters is not readily realized. As a result, the active filter is inexpensive, and is attractive for its small size and weight. In addition, it is readily included with other signal-processing functions on a single integrated circuit. See INTEGRATED CIRCUITS.

3. The design of active filters is considerably simpler than that of passive ones. In addition, it is easy to provide for variability, which can be used to change filter characteristics by electrical input signals.

The active filter also has some disadvantages:

1. Since the active filter contains electronic components, it requires a power supply, which adds to the complexity of the realization. The electronic components also place restrictions on the level of the signals that can be applied to the filter and on the noise component that the filter may add to the filtered signal. See ELECTRONIC POWER SUPPLY.

2. The mathematical process by which the active filter produces filtering characteristics in general requires the use of internal feedback. When this feedback is positive, the resulting filter may be very sensitive to lack of precision in component values, and the effects of aging and environmental conditions. See FEEDBACK CIRCUIT.

Single-amplifier filters. One important classification of active filters is according to the number of gain elements (amplifiers) that are used in their design. Historically, the first practical filter used a single amplifier of positive low gain, in the range of 1–5. The amplifier was realized by an operational amplifier and two feedback resistors, R_A and R_B . The resulting configuration could be used for low-pass, high-pass, and band-pass applications. Other filtering functions could be obtained by additional circuit complexity. The circuit is usually referred to as a Sallen-Key filter (Fig. 5a). Since the Sallen-Key filter uses positive feedback, the sensitivity problem mentioned above usually restricts the use of this filter to low- Q applications.

A more frequently utilized single-amplifier filter uses an inverting operational amplifier directly as the active element. The filter is called an infinite-gain, single-amplifier filter (Fig. 5b). Similar filter configurations may be used to realize high-pass and band-pass filter characteristics. Although the design equations for the infinite-gain, single-amplifier filter are more complex than those for the Sallen-Key one, the use of negative feedback which is implicit in this realization provides for a far more practical realization with considerably lower sensitivities.

The single-amplifier, second-order filter realizations described above can also be used to realize higher-order filter functions. [The order of a filter is the number n in Eq. (1), giving its transfer function.] To do this, separate filter sections are used to realize second-order factors, $T_i(s)$ of the overall transfer

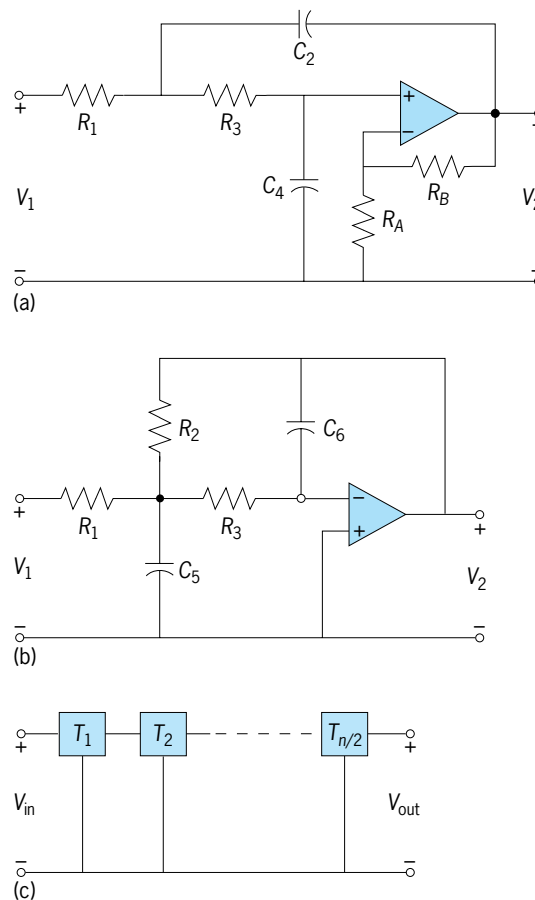


Fig. 5. Single-amplifier filters. (a) Low-pass Sallen-Key filter. (b) Low-pass, infinite-gain, single-amplifier filter. (c) Cascade of second-order filter realizations.

function, and the resulting sections are connected in a cascade configuration in which the output of one filter section is used as the input of the next one (Fig. 5c). There are $n/2$ sections labeled T_i , and since each is of second order, an n th-degree function results.

Multiple-amplifier filters. Another type of second-order filter, which can be used either by itself or as a building block in a cascade, is the multiple-amplifier filter. Such a filter uses more than one active element (operational amplifier), typically three. Although these filters have a more complex structure than the single-amplifier ones, and usually have a greater number of passive RC elements as well, they have excellent sensitivity properties and are easy to design. They are the most important general-purpose filter structures.

One of the best-known multiple-amplifier active filters is the state-variable filter, also known as the KHN (from the names of its inventors, W. J. Kerwin, L. P. Huelsman, and R. W. Newcomb) [Fig. 6a]. Three outputs, V_{HP} , V_{BP} , and V_{LP} , for high-pass, band-pass, and low-pass transfer functions respectively, are available. The versatility of the state-variable filter is improved if a fourth operational amplifier is provided as a summing amplifier. By summing the high-pass, low-pass, and band-pass outputs, a biquadratic filter function can be realized. This includes notch filters,

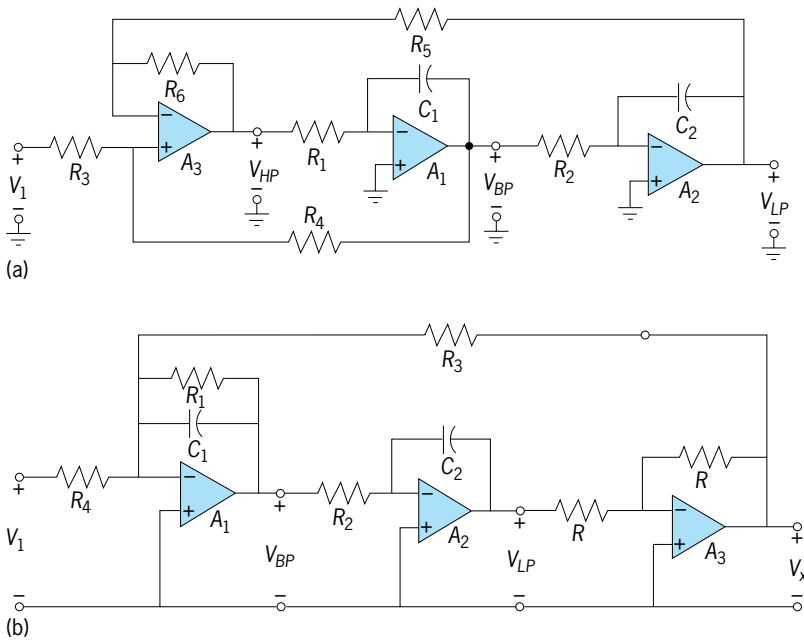


Fig. 6. Multiple-amplifier filters. (a) State-variable (KHN) filter. (b) Tow-Thomas filter.

elliptic (Cauer) filters, and all-pass filters. Realizations of higher than second order can be produced by cascading appropriately designed second-order realizations.

Another versatile multiple-amplifier filter, the Tow-Thomas filter (Fig. 6b), has the advantage of being even easier to design than the state-variable filter, but has the disadvantage of not providing a high-pass output in its basic configuration.

Direct realizations. A completely different class of active filters than the single- and multiple-amplifier ones described above is referred to as direct realizations. This class uses a passive *RLC* filter realization as a prototype, then directly derives an active-*RC* configuration from the prototype. Since passive filters can be of any order, the resulting class of active filters can also have any desired order. These high-order active filters are said to be directly realized in contrast to the cascade higher-order realizations de-

scribed above. There are two main classes of direct realization filters.

The first type of direct realization is one in which the inductor element is replaced by a simulated inductor. The other elements in the circuit remain the same. This technique is especially applicable to the realization of high-pass filters in which the inductor elements all appear with one terminal connected to ground. This arrangement simplifies the connection of the power supplies for the active elements. An example of a simulated inductor realization is a fifth-order high-pass filter (Fig. 7a). In a direct active-*RC* realization of the filter (Fig. 7b), each of the inductors is simulated by an active-*RC* circuit consisting of four resistors, a capacitor, and two operational amplifiers.

Another example of an element-replacement direct realization is one in which the elements are first subjected to an impedance transformation of $1/s$, in which $R \rightarrow C$, $L \rightarrow R$, and C becomes an element with an impedance $1/Cs^2$. This element is called a frequency-dependent negative resistor (FDNR), and is indicated in circuit diagrams by a four-bar symbol. Such a transformation leaves the voltage transfer function invariant but permits the realization of a low-pass filter with grounded active-*RC* elements. An example is a fifth-order *RLC* low-pass filter (Fig. 8), in which two active-*RC* circuits, each consisting of two capacitors, three resistors, and two operational amplifiers, realize the FDNRs.

The second type of direct realization is one in which the voltages and currents of a passive *RLC* filter are simulated by the voltages in a circuit consisting of damped and undamped integrators. The resulting configuration has a series of feedback paths and is called a leapfrog filter.

Computer-aided design. The design of various active filters (as well as their passive counterparts) can be simplified by the use of various computer-aided design programs which operate on a personal computer. See COMPUTER-AIDED DESIGN AND MANUFACTURING.

Parasitics. In general, active filters are designed with the assumption that the active elements are ideal (no parasitics). For example, the operational amplifier is assumed to have infinite gain, infinite input

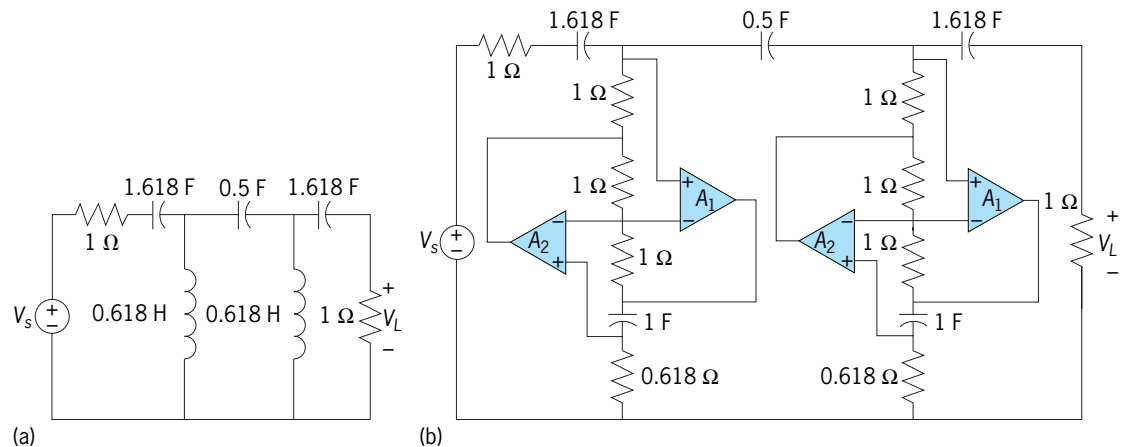


Fig. 7. Simulated-inductor realization of fifth-order high-pass filter. (a) *RLC* prototype, normalized for cutoff frequency of 1 rad/s (0.159 Hz). (b) Direct active-*RC* realization.

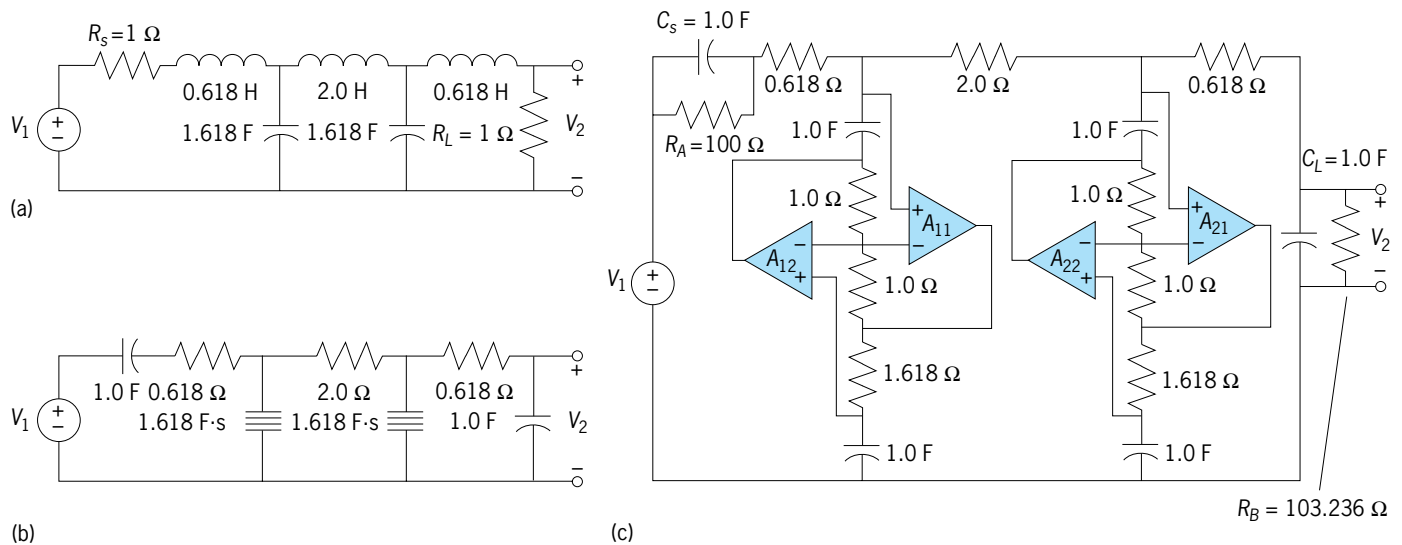


Fig. 8. FDNR realization of fifth-order low-pass filter. (a) RLC prototype, with normalized low-pass function with cutoff frequency of 1 rad/s (0.159 Hz). (b) Filter after impedance transformation. (c) Final realization, in which active-RC circuits realize FDNRs.

impedance, zero output impedance, and an infinite frequency range (gain-bandwidth). All practical filter realizations must be evaluated for the effect that the nonideal parasitics have on actual filter performance. One of the most troublesome of the parasitics is the operational-amplifier gain bandwidth. Typically, for a given application, tuning charts may be developed which can be used to provide compensation for such operational-amplifier limitations. See AMPLIFIER; CIRCUIT (ELECTRONICS); OPERATIONAL AMPLIFIER.

Lawrence P. Huelsman

Bibliography. H. Blinichoff and A. Zverev, *Filtering in the Time and Frequency Domains*, 1976, reprint 1986; M. Hasler and J. Neiryneck, *Electric Filters*, 1986; L. P. Huelsman, *Active and Passive Analog Filter Design*, 1993; S. Mitra and C. Kurth, *Miniature and Integrated Filters*, 1989; R. Schaumann and M. Van Valkenburg, *Design of Analog Filters: Passive, Active-RC, and Switched Capacitor*, 2d ed., 2000.

Electric furnace

A chamber heated to high temperature by electricity. The furnace consists of a refractory shell to resist the high temperatures attained within the chamber and minimize the heat losses to the surrounding area, plus a power source and electrical circuit to provide the heat. The operating temperatures range from a few hundred degrees Celsius to the melting points of refractory metals. The heating mechanism may be resistance, induction, microwave, or an electric arc in the largest installations. Electric furnaces are used primarily in the metallurgical industries for heat-treating, melting, and smelting operations. Resistance furnaces are used for heat treating. Induction furnaces are used to melt relatively low-melting alloys for casting, while arc furnaces are used for melting and smelting applications. The

range of furnace powers, sizes, and geometries varies immensely. Powers range from a few kilowatts to 100 megawatt-amperes in a single unit. The furnace volume may be less than a cubic meter to more than 300 m³. Both cylindrical and rectangular geometries are common. See ELECTRIC HEATING; ELECTROMETALLURGY; FURNACE CONSTRUCTION; HEAT TREATMENT (METALLURGY); INDUCTION HEATING; METALLURGY; REFRACTORY; RESISTANCE HEATING.

Electric-arc furnaces, in which the prime mechanism of heat transfer is from an electric arc, may be of the direct-arc or indirect-arc type. In the direct-arc furnace, an arc is struck between the charge being melted (usually scrap steel) and an electrode. In the indirect-arc type, an arc is struck between adjacent electrodes, and this arc then heats the charge. The arc may also be submerged below the surface of the charge. In this case, some of the current may heat the charge by resistance heating through the charge, while the rest of the current passes through the arc. Submerged-arc furnaces are among the largest units in operation and are used primarily for the smelting of ferroalloys and silicon metal, as well as phosphorus, calcium carbide, calcium silicon, copper, and copper-nickel mattes. See ARC HEATING.

Both alternating current and direct current are used to operate electric-arc furnaces. Historically, alternating-current devices with three or six carbon electrodes were used. A cylindrical geometry with the electrodes arranged in a triangle is common for three electrode systems, while the largest smelting furnaces are usually rectangular with the six electrodes installed in a line. Direct-current furnaces are cylindrical with a single axial electrode acting as cathode, while the furnace charge and bottom becomes the anode. Submerged-arc furnaces are continuously charged with feed, and there is a continuous discharge of gas from the furnace roof. Molten metal and/or slag are removed by tapping at discrete time intervals. Melting furnaces have a removable

roof for charging and often a means of tilting the furnace to remove the molten product. Although most arc furnaces use graphite electrodes, some smaller-capacity plasma furnaces use nonconsumable water-cooled metal electrodes, which allow for better control of the furnace atmosphere.

Presently one-third of the world's steel production and 50% of the United States' steel production comes from electric-arc furnaces. The largest units have graphite electrodes 60–80 cm (24–31 in.) in diameter and a tap weight of 80–100 metric tons (88–110 tons) of liquid steel. The active power may be up to 100 megawatts. A typical FeNi submerged-arc smelting furnace would have six electrodes and an operating load of 75 MVA. See STEEL MANUFACTURE.

Richard J. Munz

Bibliography. N. Barcza, The application of dc plasma-arc technology to extractive metallurgical processing, *High Temp. Mater. Proc.*, 1:1–15, 1997; F. Fereday, *Commentary: Submerged Arc Furnaces*, TechCommentary published by CMP, EPRI Center for Materials Production, Carnegie Mellon Research Institute, 15230-2950, 1996; J. Kunze and R. Degel, History, current status and new trends of submerged arc furnace technology for ferroalloy metals, *Steel GRIPS*, 1(3):164–169, 2003; D. Neuschütz, Metallurgical applications of high power arc heating systems, *High Temp. Mater. Proc.*, 4(3):309–321, 2000.

Electric heating

Methods of converting electric energy to heat energy by resisting the free flow of electric current. Electric heating has several advantages: it can be precisely controlled to allow a uniformity of temperature within very narrow limits; it is cleaner than other methods of heating because it does not involve any combustion; it is considered safe because it is protected from overloading by automatic breakers; it is quick to use and to adjust; and it is relatively quiet. For these reasons, electric heat is widely chosen for industrial, commercial, and residential use.

Types of electric heaters. There are four major types of electric heaters: resistance, dielectric, induction, and electric-arc.

Resistance heaters. Resistance heaters produce heat by passing an electric current through a resistance—a coil, wire, or other obstacle which impedes current and causes it to give off heat. Heaters of this kind have an inherent efficiency of 100% in converting electric energy into heat. Devices such as electric ranges, ovens, hot-water heaters, sterilizers, stills, baths, furnaces, and space heaters are part of the long list of resistance heating equipment. See RESISTANCE HEATING.

Dielectric heaters. Dielectric heaters use currents of high frequency which generate heat by dielectric hysteresis (loss) within the body of a nominally non-conducting material. These heaters are used to warm to a moderate temperature certain materials that have low thermal conducting properties; for example, to soften plastics, to dry textiles, and to work

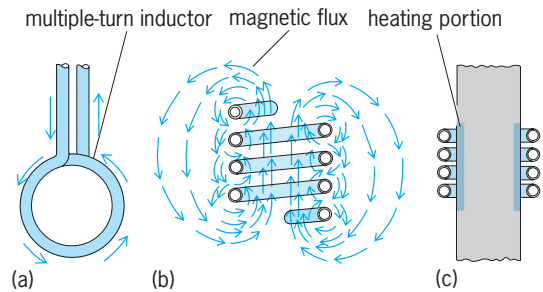


Fig. 1. Induction heater. Arrows indicate current flow. (a) End-on view. (b) Magnetic fields. (c) Heating pattern.

with other materials like rubber and wood. See DIELECTRIC HEATING.

Induction heaters. Induction heaters produce heat by means of a periodically varying electromagnetic field within the body of a nominally conducting material. This method of heating is sometimes called eddy-current heating and is used to achieve temperatures below the melting point of metal. For instance, induction heating is used to temper steel, to heat metals for forging, to heat the metal elements inside glass bulbs, and to make glass-to-metal joints. **Figure 1** shows the magnetic fields and heating pattern produced by an induction heater. See INDUCTION HEATING.

Electric-arc heaters. Electric-arc heating is really a form of resistance heating in which a bridge of vapor and gas carries an electric current between electrodes. The arc has the property of resistance. Electric-arc heating is used mainly to melt hard metals, alloys, and some ceramic metals (**Fig. 2**). See ARC HEATING.

General design features. All electrical parts must be well protected from contact by operators, work materials, and moisture. Terminals must be enclosed within suitable boxes, away from the high heat zone, to protect the power supply cables. Repairs and replacements should be possible without tearing off heat insulations.

Resistance heaters are often enclosed in pipes or tubes suitable for immersion or for exposure to difficult external conditions. Indirect heating is done by circulating a heat transfer medium, such as special oil or Dowtherm (liquid or vapor), through jacketed vessels. This permits closer control of heating-surface temperature than is possible with direct heating.

Some conducting materials can be heated by

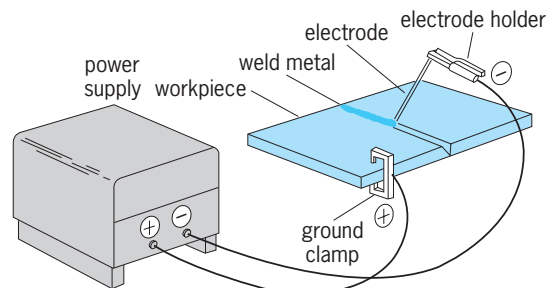


Fig. 2. Typical electric-arc heater.

passing electric current through them, as is done in the reduction of aluminum. Some conducting liquids can be heated by passing an electric current between immersed electrodes. Heat is produced by the electrical resistance of the liquid.

The supply of necessary electric power for large heating installations necessitates consultation with the utility company. The demand, the power factor of the load, and the load factor all affect the power rates. Large direct-current or single-phase alternating-current loads should be avoided. Polyphase power at 440–550 V permits lower current and reduced costs. *See* FURNACE CONSTRUCTION.

Electric heating for houses. Electricity is one choice for heating houses, but with only a 35% efficiency rate, electricity has been a less attractive option than the direct use of gas and oil for heating homes.

One type of electric heater for houses is the heat pump. The heat pump can be used alone or to supplement the output of direct resistance heating (heat strips). The heat pump works on a reversed refrigeration cycle; the pump pulls the heat from the cold outside air and forces it into the house. The efficiency of this process is directly affected by the ambient temperature and the desired interior temperature; the larger the difference in these temperatures, the lower the efficiency of the pump. Generally, however, the efficiency of the heat pump is higher than the efficiency of resistance heating, so that the efficiency of the heat pump is actually higher than 100%. *See* HEAT PUMP.

Since the heat pump performs both heating and cooling and is more efficient than other forms of electric heat, it can be an economical way to heat houses.

Common electric heating systems in houses are central heating employing an electric furnace with forced air circulation; central heating employing an electric furnace with forced water circulation; central heating using radiant cables; electrical duct heaters; space (strip) heaters which use radiation and natural convection for heat transfer; and portable space heaters.

Mo-Shing Chen

Bibliography. C. J. Erikson, *Handbook of Electric Heating for Industry*, 1994; D. G. Fink and H. W. Beaty (eds.), *Standard Handbook for Electrical Engineers*, 14th ed., 1999; L. Miles, *Heat Pumps: Theory and Service*, 1993; M. Orfeuill, *Electric Process Heating*, 1987; H. J. Sauer, *Heat Pump Systems*, 1990.

Electric insulator

A material that blocks the flow of electric current across it. Insulators are to be distinguished from electrolytes, which are electronic insulators but ionic conductors. *See* ELECTROLYTE; IONIC CRYSTALS.

Electric insulators are used to insulate electric conductors from one another and to confine electric currents to specified pathways, as in the insulation of wires, electric switchgear, and electronic compo-

nents. They provide an electrical, mechanical, and heat-dissipation function. The electrical function of an insulator is characterized by its resistivity, its dielectric strength (breakdown voltage), and its dielectric constant (relative permittivity). Insulators can be solid, liquid, or gaseous (see **table**).

Resistivity. The resistivity of a material is a measure of the electric current density that flows across it in response to an applied electric field. The electric current density is the sum of contributions from all species of mobile charge carriers, including electrons and ions. Each contribution is proportional to the concentration of the mobile species, its charge, and its mean (drift) velocity in response to the applied field. Solid and liquid insulators have direct-current resistivities of 10^{10} ohm-meters at room temperature as compared to 10^{-8} $\Omega\cdot\text{m}$ for a good metal or 10^{-3} $\Omega\cdot\text{m}$ for a fast ionic conductor. Insulators minimize the concentrations of mobile charge carriers by having a filled valence band separated from an empty conduction band by a large energy gap and by having all equivalent lattice sites either occupied or vacant of ionic species. An energy gap greater than 5 eV can be obtained in either ionically or covalently bonded crystals and molecules. Oxides may have the advantage of being inert to further oxidation, but they are more subject to short-range, if not long-range, ionic motion. Protons introduced by the absorption or adsorption of water are the most prevalent source of ionic conduction. Impurities and crystalline defects are sources of mobile electrons accessible at lower temperatures and electric fields. *See* BAND THEORY OF SOLIDS; ELECTRICAL RESISTANCE; ELECTRICAL RESISTIVITY; SEMICONDUCTOR.

Dielectric strength. A material used in the electrical industry for insulation, capacitors, or encapsulation may be exposed to large voltage gradients that it must withstand for the operating life of the system. Power lines operate at transmission voltages in excess of 1 MV, for example. Failure occurs if an electric short-circuit develops across the material. Such a failure is called dielectric breakdown. The breakdown voltage gradient, expressed in kilovolts per millimeter, is a measure of the dielectric strength. Dielectric breakdown of a solid is destructive; liquids and gases are self-healing.

Breakdown is initiated by the injection of mobile charge carriers into the conduction or valence band of the insulator. There are three basic types: intrinsic, thermal, and avalanche. Intrinsic breakdown occurs where electrons are injected across a metal-insulator interface. It depends primarily on the applied field gradient, it occurs in less than 1 microsecond, and it is independent of sample shape. Thermal breakdown is associated with thermal excitation of electrons into a conduction band; resistive losses of the excited electrons lead to thermal runaway. This mechanism depends primarily on the temperature rather than the applied field gradient, and is dependent on both sample shape and the rate of application of the field. Avalanche breakdown occurs where collisions between mobile and stationary electrons multiply

Characteristics of electric insulators*

Insulator	Breakdown voltage, kV/mm	Resistivity, $\Omega\text{-m}$,	Relative permittivity (k)
Solid			
Paper	20	10^{10}	3–4
Nylon	16	10^{13}	3–4
Glass	20–40	10^{10} – 10^{14}	5
Rubber	16–50	10^{13} – 10^{14}	3
Mica	25–200	10^{11} – 10^{15}	5
Poly(tetrafluoroethylene)	60	10^{14} – 10^{15}	2.1
Porcelain	10–20	10^{16}	5–7
Polystyrene	25	10^{16}	2.5
Liquid			
Saturated hydrocarbons [†]	140–200	10^{17} – 10^{18}	
Aromatic hydrocarbons [†]	160–280	10^{17} – 10^{18}	
Transformer and silicon oils [‡]	15		
Gaseous [§]			
Air	3		
Carbon dioxide	2.7		
Sulfur hexafluoride	7.5		

* At 20°C (68°F) and atmospheric pressure (101.3 kPa)

[†] Physically and chemically purified.

[‡] Commercial high purity.

[§] Between two electrodes 25 cm (10 in.) in diameter.

the concentration of mobile electrons. It generally occurs in conjunction with intrinsic or thermal breakdown.

In addition, there are three pseudo-types of breakdown of great practical importance: discharge, electrochemical, and mechanical. Dielectric discharge is associated with a gaseous breakdown at the surface or within included pores of the material. Electrochemical breakdown is a result of chemical reactions that, over time, reduces the dielectric strength. Mechanical breakdown is precipitated by cracks or other defects that distort the applied field. *See* BREAKDOWN POTENTIAL; ELECTRICAL BREAKDOWN; ELECTRICAL CONDUCTION IN GASES.

Dielectric constant. An insulator is also known as a dielectric. The dielectric constant k is defined as the ratio of the capacitance of a flat-plate condenser, or capacitor, with a dielectric between the plates to that with a vacuum between the plates; this ratio is also the relative permittivity of the dielectric. The dielectric constant is a measure of the ability of the insulator to increase the stored charge on the plates of the condenser as a result of the displacement of charged species within the insulator. Displacement of a charge q through a distance δ creates an electric dipole of moment $p = q\delta$. The displacement of charge takes time, so in an alternating-current circuit $k = k' - ik''$ is a complex number; a measure of the losses introduced by the charge displacements is the loss tangent, $\tan \delta \equiv k''/k'$. Low dielectric constants are desirable for the insulation of alternating-current circuits.

In ceramics and glasses, four primary short-range motions contribute to the total polarization of the dielectric: electronic, atomic, dipolar, and interfacial. Electronic polarization is due to a shift of the valence-electron cloud relative to the atomic nucleus. It gives a resonance absorption peak at optical frequencies. Atomic polarization is due to a displacement of posi-

itive and negative ions with respect to each other. It gives a resonance in the infrared range 10^{12} – 10^{13} Hz that may be quite broad where several types of ions are present. Dipolar, or orientational, polarization is a perturbation of the thermal motion of ionic or molecular dipoles to give a net dipolar orientation in the direction of the applied field. The preferred orientation of permanent dipoles (Stevens deformation polarization), which may be particularly important in liquid and gaseous insulators, occurs in the 10^{11} – 10^{12} -Hz frequency range and the switching of dipoles between two equivalent equilibrium orientations in a solid occurs at lower frequencies, 10^3 – 10^6 Hz. Interfacial, or space-charge, polarization occurs where the motions of mobile charges are blocked by internal barriers such as two-phase interfaces. This mechanism produces a large capacitance and dielectric constant at low frequencies. Orientational and interfacial polarizations introduce losses in insulators from 10^{-3} to 10^9 Hz; atomic and electronic polarizations in glasses and ceramics are important for optical communication. *See* CAPACITANCE; DIELECTRIC MATERIALS; FERROELECTRICS; OPTICAL COMMUNICATIONS; PERMITTIVITY; POLARIZATION OF DIELECTRICS.

Mechanical properties. Successful application of solid insulating materials also depends on their mechanical properties. Insulation assemblies commonly must withstand thermal-expansion mismatch, tension, compression, flexing, or abrasion as well as a hostile chemical-thermal environment. The introduction of cracks promotes the penetration of moisture and other contaminants that promote failure, and the presence of pores may cause damaging corona discharge on the surface of a high-voltage conductor. As a result, composite materials and engineered shapes must be tailored to meet the challenges of a particular operational environment. *See* CORONA DISCHARGE.

For example, glasses and varnishes are used as sealants, and oil is used to impregnate high-voltage, paper-insulated cables to eliminate air pockets. Porcelain is a commonly used material for the suspension of high-voltage overhead lines, but it is brittle. Therefore, a hybrid insulator was developed that consists of a cylindrical porcelain interior covered by a mastic sealant and a silicone elastomer sheath heat-shrunk onto the porcelain core. The circular fins of the outer sheath serve to shed water. However, twisted-pair cables insulated with poly(tetrafluoroethylene) are used for high-speed data transmission where a small dielectric constant of the insulator material is needed to reduce signal attenuation. See CONDUCTOR (ELECTRICITY); ELECTRIC POWER TRANSMISSION; FLUORINE; PORCELAIN; SILICONE RESINS.

Liquid and gaseous insulators. Liquid or gas insulation provides no mechanical strength, but it may provide a cheap, flexible insulation not subject to mechanical failure. Biphenyls are used as insulating liquids in capacitors; alkyl benzenes in oil-filled cables; and polybutenes for high-pressure cables operating at alternating-current voltages as high as 525 kV. Sulfur hexafluoride (SF_6) is a nonflammable, nontoxic electron-attracting gas with a breakdown voltage at atmospheric pressure more than twice that of air. Fluorocarbons such as C_2F_6 and C_4F_8 as well as the Freons are also used, and breakdown voltages have been increased significantly in gas mixtures through a synergistic effect. Used as an insulating medium in high-voltage equipment at pressures up to 10 atm (1 megapascal), sulfur hexafluoride can reduce the size of electrical substations by a factor of 10 over that of air-insulated equipment. Enclosure of metal cable in a metallic conduit filled with sulfur hexafluoride gas has been used to shield nearby workers from exposure to high electric fields. See CIRCUIT BREAKER; FLUOROCARBON; POLYCHLORINATED BIPHENYLS.

Thermal conductivity. Finally, the ability to transfer heat may be an overriding consideration for the choice of an electric insulator. Electrical machines generate heat that must be dissipated. In electronic devices, the thermal conductivity of the solid substrate is a primary consideration. Where mechanical considerations permit, circulating liquid or gaseous insulation is commonly used to carry away heat. Liquids are particularly good transporters of heat, but they are subject to oxidation. In transformers, for example, the insulators are generally mineral or synthetic oils that are circulated, in some places with gaseous nitrogen to inhibit oxidation, to carry away the heat generated by the windings and magnetic core. Silicone fluids, while flammable, have a higher fire point than mineral oil and a lower dissipation factor. Circulating gaseous sulfur hexafluoride is also used for transformer insulation. Hydrogen is a gaseous insulator used to cool rotating electrical machines because it has a heat-conducting capacity over seven times that of air. See ELECTRIC ROTATING MACHINERY; ELECTRICAL INSULATION; HYDROGEN; TRANSFORMER.

John B. Goodenough

Bibliography. R. Bartnikas, *Electrical Insulating Liquids in Dielectrics*, vol. 3, 1994; R. Bartnikas et al. (eds.), *Engineering Dielectrics*, vol. 1, 1979, vol. 2A, 1983, vol. 3, 1987; C. Johnson, *Handbook of Electrical and Electronic Technology*, 1994; F. H. Kreuger, *Industrial High Voltage*, vol. 1: *Electric Fields, Dielectrics, Constructions*, 1991; W. T. Shugg, *Handbook of Electrical and Electronic Insulating Materials*, 2d ed., 1995.

Electric organ (biology)

An effector organ found in six different groups of fishes; output is an electric pulse (**Table 1**; **Fig. 1**). Voltages large enough to aid in prey capture or predator deterrence are produced by various strongly electric fishes. These include the electric eel (*Electrophorus electricus*) from South America; the electric catfish (*Malapterurus electricus*) from Africa; the family of electric rays, Torpedinidae, which are widely distributed in the world's oceans; and possibly the stargazer genus, *Astroscopus*, of the western Atlantic. Weakly electric fishes emit a lower voltage, the energy source for an active electrosensory system that monitors electrical impedance in the environment. These weak signals also serve in intra- and interspecific communication. There are three groups of weakly electric fishes. First, the South American knifefishes, the Gymnotiformes, are a large and diverse group of several families that also include *Electrophorus*. Second, the electrically active African Mormyriiformes are composed of the numerous species of the family Mormyridae and the single species *Gymnarchus niloticus* in the family Gymnarchidae. Finally, many species of skates and rays of the family Rajidae occur in marine waters worldwide.

The strongly electric fishes are remarkable for the high voltage and power of their discharges (**Table 2**). For example the electric eel generates pulses in excess of 500 V. In contrast, a large torpedo generates a smaller voltage, about 50 V in air, but the maximum current is larger, and the peak pulse power can exceed 1 kW (or about 1 hp). The electric organs may constitute a substantial fraction of the body mass. The weakly electric fishes have much smaller organs emitting pulses of a few tenths of a volt to several volts; such amplitudes are still large compared to those recorded outside ordinary, non-electric fishes.

Electric organ discharge is explicable in terms of the properties of ordinary excitable cells, that is, those of nerve and muscle. The single generating cells of electric organs are called electrocytes. They are modified striated muscle fibers that have lost the ability to contract, although they retain muscle filaments to varying degrees, and they produce electric signals as muscle fibers do.

Electrocyte operation. The basic operation of an electrocyte is as follows. Consider a flattened cell which contains a high concentration of potassium ions (K^+), a low concentration of sodium ions (Na^+), and an equivalent number of unspecified anions. Let

TABLE 1. Groups of electric fishes*

Common name	Family	Genera and species	Strength of organ discharge	Distribution
Skates, ordinary rays	Rajidae	<i>Raja</i> , many species, a number of other genera not known to be electric	Weak	Marine, cosmopolitan
Electric rays, torpedos	Torpedinidae	About 20 genera	Strong, up to 60 V or 1 kW, some perhaps weak	Marine, cosmopolitan
Mormyrids, elephant-nosed fish (many lack enlarged chin or snout)	Mormyridae	About 11 genera, several with many species	Weak	Fresh-water, Africa
<i>Gymnarchus</i>	Gymnarchidae	1 species, <i>Gymnarchus niloticus</i>	Weak	Fresh-water, Africa
Gymnotid eels, electric eel, knifefishes	Electrophoridae	1 species, <i>Electrophorus electricus</i>	Strong, more than 500 V	Fresh-water, South America
	Gymnotidae	1 species, <i>Gymnotus carapo</i>	Weak	Fresh-water, South America
	Sternopygidae	About 8 genera, a number of species	Weak	Fresh-water, South America
Electric catfish	Sternarchidae (Apteronotidae)	About 5 genera, a number of species	Weak	Fresh-water, South America
	Malapteruridae	1 species, <i>Malapterurus electricus</i>	Strong, more than 300 V	
Stargazers	Uranoscopidae	1 electric genus, <i>Astroscoopus</i> , several species	Strong (?), about 6 V in air from small animals	Marine, western Atlantic

*From M. V. L. Bennett, Electric organs, in W. S. Hoar and D. J. Randa II (eds.), *Fish Physiology*, vol. 5, pp. 347–491, Academic Press, 1971.

it be bathed in a solution which is high in Na^+ , and low in K^+ . Initially its membrane is specifically permeable to potassium ions (K^+) over its entire surface, and in this condition the membrane potential is uniform and no currents flow (Fig. 2a). The cell interior will be at some negative voltage, at which the tendency for K^+ to diffuse out and down its concentration gradient is just balanced by the negativity of the interior, tending to prevent the positively charged ions from diffusing out. This “equilibrium” potential is about 90 mV inside negative for most excitable cells. When the cell is activated, one of its membrane faces becomes much more permeable to Na^+ . The potential at which this membrane would be at equilibrium in terms of Na^+ flow is about 50 mV inside positive, but the cell interior is negative. Thus, Na^+ flows inward through the activated face, making the interior more positive. The other membrane

face that still is specifically permeable to K^+ is displaced from its equilibrium potential, and as a result K^+ flows outward through it. In effect, the two faces act as batteries oriented in series (Fig. 2b). All operate just this simply; apposed membranes have different permeabilities during activity, allowing ions to flow down their concentration gradients and thereby generating electric currents. The proximate sources of energy for organ discharge are the differences in ionic concentrations between interior and exterior of the electrocytes. During activity these concentrations become more alike; between periods of activity, metabolic energy is expended to restore the concentration differences. See BIOPOTENTIALS AND IONIC CURRENTS.

During activity, currents crossing all of the electrolyte membranes in series flow around the outside of the electric organ through the surrounding body

TABLE 2. Discharge characteristics of electric fishes

Species	Approximate discharge voltage, V	Duration, ms	Frequency	Number of electrocytes in series
<i>Electrophorus electricus</i> , electric eel	700	2	0–250	6000
<i>Malapterurus electricus</i> , electric catfish	300	1.5	0–200	?
<i>Torpedo nobiliana</i> , electric ray	60	5	0–150	1000
<i>Astroscoopus γ-graecum</i> , stargazer	6	5	0–150	200
<i>Gnathonemus petersii</i> , mormyridform	5	0.3	2–100	100
<i>Gymnarchus niloticus</i> , mormyridform	1	2	250	140
<i>Gymnotus carapo</i> , gymnotiform knifefish	1	1.5	50–150	90
<i>Sternarchus albifrons</i> , sternarchid (apteronotid) knifefish	1	1	800	90
<i>Sternopygus macrurus</i> , gymnotiform knifefish	1	10	50 male, 150 female	90
<i>Raja erinacea</i> , rajid skate	2	100	0–5	200

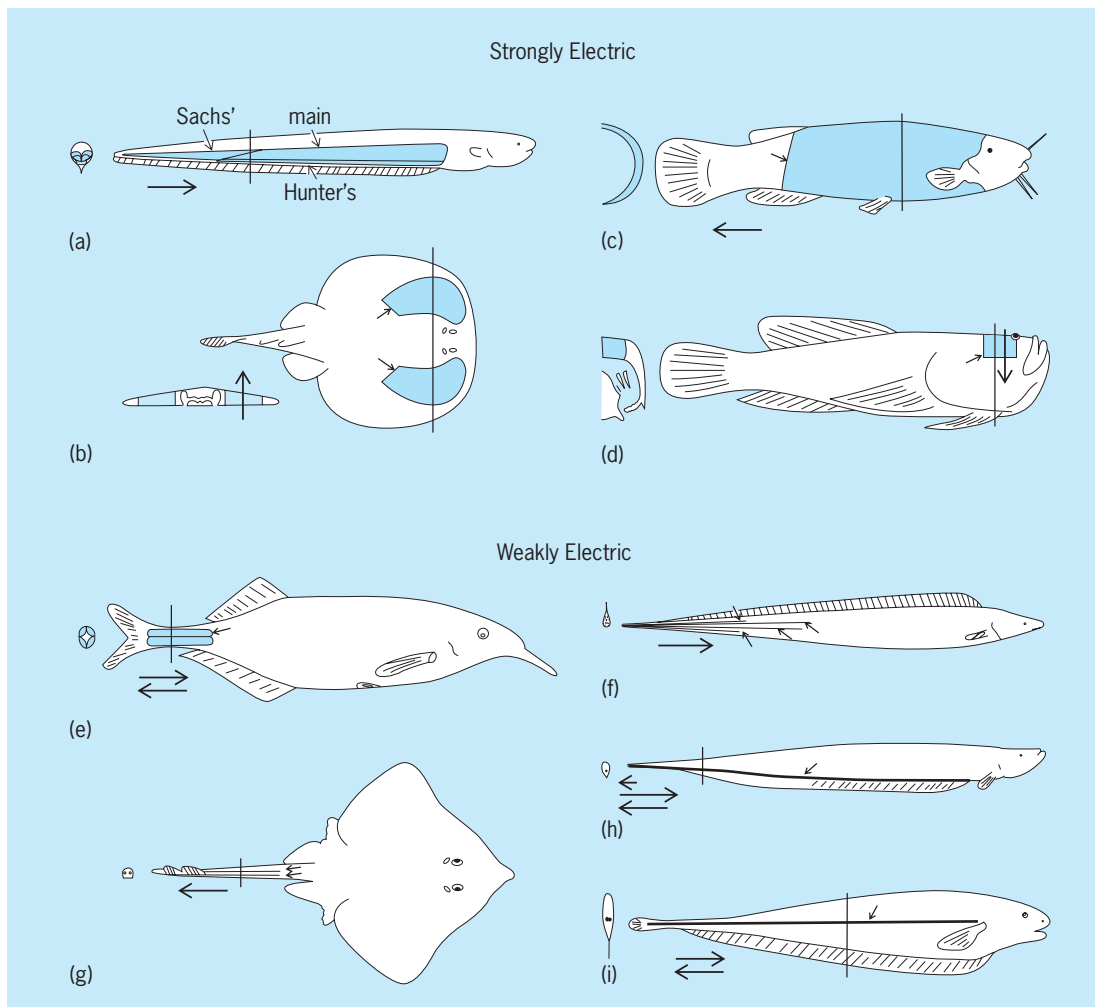


Fig. 1. Representative electric fishes: (a) *Electrophorus*; (b) *Torpedo*; (c) *Malapterurus*; (d) *Astroscopus*; (e) *Gnathonemus*; (f) *Gymnarchus*; (g) *Raja*; (h) *Gymnotus*; (i) *Sternarchus*. Electric organs are shaded or solid and indicated by small arrows. Cross-sectional views through the organs are also shown at the levels indicated by the solid lines. The large arrows indicate the direction of active current flow through the organs; more than one arrow indicates successive phases of activity, and the relative lengths are proportional to relative amplitudes of the phases. (After W. S. Hoar and D. J. Randall, eds., *Fish Physiology*, vol. 5, Academic Press, 1971)

and water, producing a field, similar in form to that around a bar magnet (Fig. 2c). A voltage drop occurs across any object in the field. If the objects differ in conductivity from water, they distort the pattern of current flow, and these distortions can be detected by the electrosensory system (see below).

Ions enter or leave a cell by passing through specific protein channels in the cell membrane. The extracellular fields generated by electrocytes are significantly larger than those generated by ordinary excitable cells, but the difference is not in the nature of the ionic permeabilities or the magnitude of the membrane voltages. Rather it is in the shape of the cells and the distribution and density of ionic channels. The maximum voltage generated by a single electrocyte is not much over 0.1 V. The extraordinarily large voltages that are generated by strongly electric organs result from the arrangement of many elements in series (Fig. 2c) and in parallel. For example, in the electric eel many thousands of cells occur in series (Table 2). In the torpedos there are hundreds of cells in series.

Electrocyte membranes. The arrangement of specific kinds of membrane differs in different electrocytes. In the strongly electric fishes, electrocytes are always flattened (hence the terms electroplaque and electroplax, Fig. 3). The simplest are those of the electric ray and stargazer. One face is densely innervated and generates only postsynaptic potentials (PSPs), as at the nerve-muscle junction. The PSPs result from a permeability increase to Na^+ and K^+ , so that the membrane is at electrical equilibrium near 0 V instead of inside positive. The PSPs are produced by action of a neurotransmitter substance, acetylcholine. The transmitter is secreted by the innervating nerves. It combines briefly with receptor molecules in the electrocyte membrane, causing the ion channels to open; then it is inactivated by the enzyme acetylcholinesterase, which is found at high concentration at the synapses. The uninnervated face of these electrocytes is inexcitable; its permeability does not change when a potential is applied. At rest the uninnervated face has a much lower electric resistance than the innervated face. During

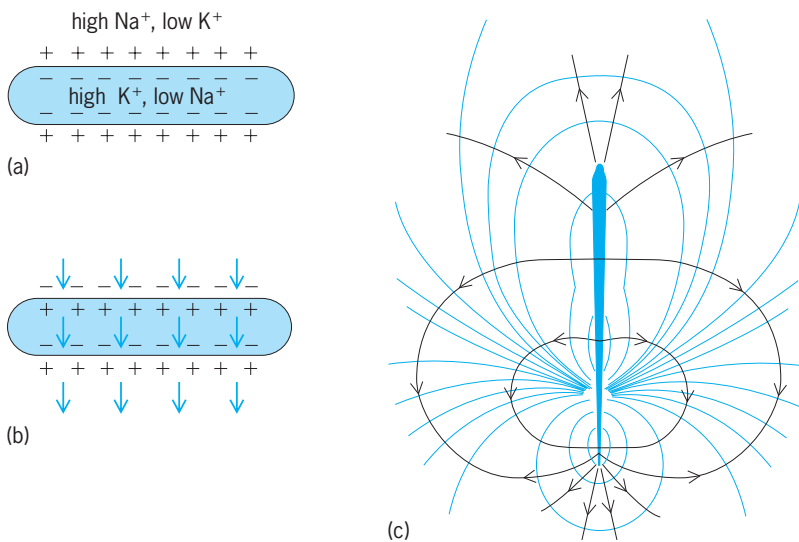


Fig. 2. Current flow through a single electrocyte and fields around an electric fish. (a) Cross section of a flattened cell at rest. Each face is specifically permeable to K^+ , and the cell interior is at a negative potential. Because equal potentials are opposed, no current flows. (b) During activity that greatly increases the permeability to Na^+ in the upper face, this membrane moves toward inside positivity. As a result, two batteries act in series and current flows. (c) Flow of current (black lines) and isopotentials (colored lines) around a weakly electric gymnotid during a head positive phase of organ discharge. (After W. Heiligenberg, *Principles of Electrolocation and Jamming Avoidance in Electric Fish*, Springer, 1977)

PSPs, however, the resistance of the innervated face greatly decreases and the two faces have comparable resistances. These characteristics maximize the effectiveness of the cells as generating elements.

The electrocytes of the electric eel (and of several weakly electric relatives) are slightly more complex. In these cells, as in those of the electric ray and stargazer, one face is innervated and one is not, the uninnervated face being of low resistance and inexcitable. In addition to PSP channels, the innervated membrane has channels specifically permeable to Na^+ ; they open whenever the cell interior becomes about 10 mV or so more positive than its resting potential (depolarization). The PSPs generated by the innervated membrane cause some of these channels to open, greatly increasing the permeability to Na^+ . The initial opening of Na^+ channels allows an influx of Na^+ , which makes the interior yet more positive, causing a further opening of Na^+ channels, and so forth; the response is said to be regenerative. The positive feedback action continues until virtually all the Na^+ channels are open, and an impulse resembling the action potential of axons and muscle fibers is generated. In contrast to the electrically opened Na^+ channels, the PSP generation channels are little affected by membrane potential. In the electric eel the density of PSP channels is much lower than in the electric ray, and most of the change in membrane potential results from the electrically opened Na^+ channels. Two mechanisms can contribute to termination of the impulse. In the eel the Na^+ channels simply close after a brief period, and the membrane returns to the K^+ equilibrium potential. In others there are additional K^+ channels that are opened by depolarization somewhat more slowly than the Na^+ chan-

nels, and these speed the return of the membrane potential to its resting value.

The three classes of response, mediated by the PSP channels and by the Na^+ and K^+ channels opened by depolarization, can be mixed together in a number of different arrangements. In some weakly electric fishes, both innervated and uninnervated membranes generate impulses. First, current flows in one direction through the cell as the innervated membrane generates its impulse; this current then excites the uninnervated face, though with some delay, producing an impulse with current flow in the opposite direction. Thus, a biphasic external potential results. In the electric catfish, for example, one face generates a briefer impulse than the other. The resting resistances of the two faces are about equal, and the role of the briefer impulse appears to be to increase the K^+ permeability of that face, thus making the resistances of the two faces comparable during the later part of the longer lasting impulse, thereby increasing current flow. A similar phenomenon occurs in some skates, except that there are no electrically opened Na^+ channels. PSPs generated in the innervated face depolarize the uninnervated face, opening K^+ channels in it, and thus increasing current flow.

In several weakly electric fishes (for example, *Gymnarchus*) the uninnervated membrane has a very low permeability to K^+ or to any other ions. This membrane does have a large electrical capacity that permits significant current flow. Because the current flow is capacitative, the discharge has no direct current (dc) component. The absence of a dc component can be related to the functioning of the discharge in the electrosensory system.

Discharge. In all strongly electric fishes the activity of individual electrocytes is highly synchronized, thus maximizing the voltage and current produced. The discharges are monophasic pulses which are presumably more effective in stimulating prey or predator (A and A' in Fig. 4). In weakly electric fishes the discharges are more diverse. There are a few with discharges like the strongly electric fishes, but others generate diphasic or triphasic pulses that need not be synchronous throughout the electric organ (B and B' in Fig. 4). Triphasic discharges can result from out-of-phase addition of diphasic responses from two sets of electrocytes (*Gymnotus*). Alternatively a single electrocyte can generate a triphasic response (some mormyrids). The shape of electrocytes is also more variable in weakly electric fishes. The cells can be drum-shaped (*Gymnotus*, Fig. 3) or spindle-shaped (*Sternopygus*, Fig. 3). The innervation can occur on a thin process or stalk from the main body of the cell (*Malapterurus*, Fig. 3; some gymnotids). The stalks may also be multiple (mormyrids, Fig. 3; some gymnotids) and may even turn around and pass through holes in the flattened part to the opposite side before receiving their synapses (some mormyrids). Propagated activity in these stalks introduces a small component in overall organ discharge.

A number of gymnotids have a separate small "accessory" electric organ in the head region.

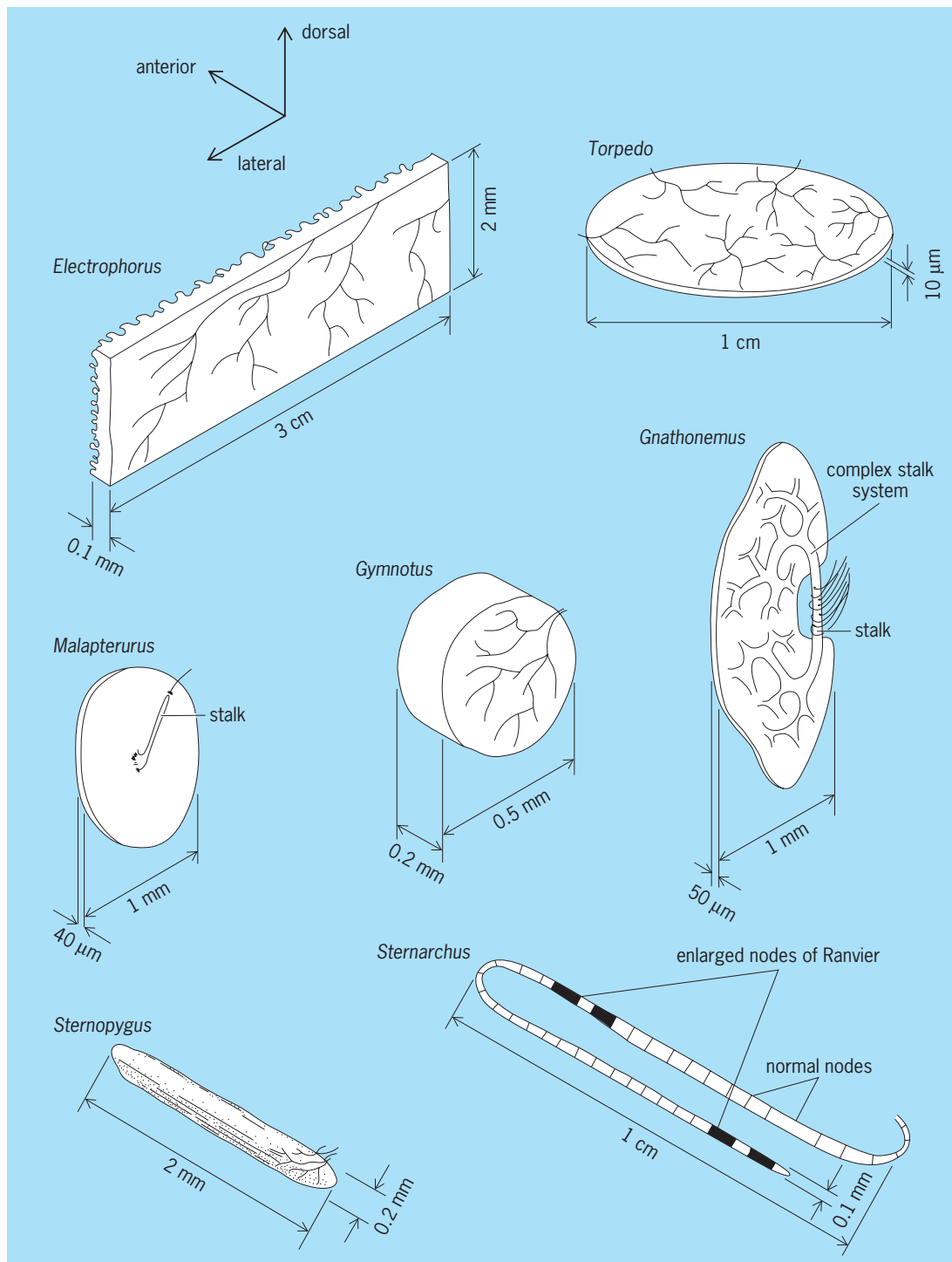


Fig. 3. Representative electrocytes. Body axes are indicated by the arrows except that the *Torpedo* electrocyte is shown ventral surface up. Innervation is indicated by branching, except for *Malapterurus* in which a single nerve fiber makes a single contact on the stalk tip. Also, the electrocyte of *Sternarchus* is a nerve fiber.

Electroreceptors are also concentrated on the head, and presumably the accessory organ functions in detecting localized objects in this region.

In strongly electric fishes the shape of the electric organ can be related to the conductivity of the environment (Fig. 1). In the marine forms, the torpedos and the stargazers, the organ is short and wide, particularly in the torpedos. In the fresh-water forms, the electric catfish and eel, the organs are long and

relatively narrower. In the marine, as compared to fresh-water forms, fewer cells occur in series and more in parallel, and the voltages are lower and the currents higher in the higher-conductance medium.

In weakly electric fishes, including the marine rajids, the organs are generally long and slender. They have fewer than a dozen cells in parallel and of the order of 100 in series (Table 2).

In view of the high power that can be generated

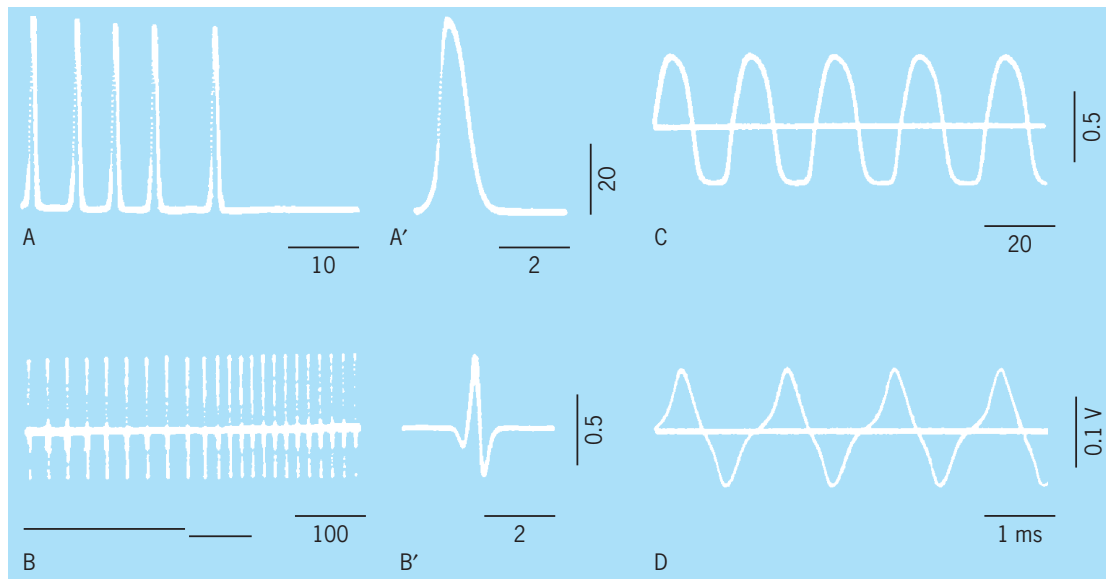


Fig. 4. Patterns of electric organ discharge. (A) An electric catfish; potentials recorded between head and tail in a small volume of water with head negativity upward. Pulses attain a maximum frequency of 190/s. (A') Single pulse. (B–D) Weakly electric gymnotids immersed in water; discharges recorded between head and tail, head positivity upward. (B) *Gymnotus*; pulses are emitted at a basal frequency of approximately 35/s. Tapping the side of the fish at the time indicated by the downward step in the lower trace causes an acceleration up to about 65/s. (B') Faster sweep showing the triphasic pulse shape. (C) *Sternopygus*; pulse frequency is about 55/s. The horizontal line indicates the zero potential level. (D) *Sternarchus*; pulse frequency is about 800/s. The horizontal line indicates the zero potential level. (After W. S. Hoar and D. J. Randall, eds., *Fish Physiology*, vol. 5, Academic Press, 1971)

by strongly electric fishes, it is not surprising that they discharge the organs only infrequently, and generally in response to an external stimulus (A in Fig. 4). The fresh-water weakly electric fishes, however, emit pulses continually (B–D in Fig. 4), and the range of frequency is wide, from a few pulses per second to more than a thousand. The frequencies and discharge patterns are species-specific and differ considerably. Two broad types can be defined, pulse species and continuous-wave species. The pulse species emit brief pulses separated by a much longer interval; nonspecific sensory stimulation usually causes an increase in the frequency of discharge (B and B' and Fig. 4). Continuous-wave species emit monophasic or biphasic pulses separated by intervals of about the same duration as the pulses (C and D in Fig. 4). These species generally modify their discharge frequency very little, although all show small changes when presented with a stimulus at a frequency close to their own. These changes, termed jamming avoidance responses, are presumed to reduce interference in electrolocation.

The discharges of rajids under normal conditions have not been described. They are evoked with difficulty by handling and consist of long-lasting, repetitive, asynchronous discharges of the generating cells.

Neural control. The neural control of electric organ discharge has been studied in representative species of most groups. Generally the frequency is set in a command nucleus, a compact group of neurons in the brain stem. The command neurons integrate various neural inputs, both excitatory and inhibitory, and the group of neurons generates a highly synchronous impulse. Each synchronous discharge is the command signal that is responsible for each elec-

tric organ discharge. Synchronization is mediated by electrical synapses that interconnect the cells. In the continuously discharging fishes the command neurons are probably spontaneously active, but they still receive sensory and other inputs that modulate discharge frequency. The command signal may be relaxed by one or two intervening groups of cells, the ultimate group innervating the electrocytes. (In the sternarchids the electric organ actually consists of the axons of the final neurons; Fig. 3.) Since different parts of the electric organ are at significantly different distances from the brain, compensatory mechanisms are present to adjust for the variations in nerve conduction and thus to ensure synchrony of organ discharge. For example, nerves to nearer parts may take more devious pathways, and they may conduct impulses at lower velocities.

Electrosensory systems. A knowledge of the electrosensory systems is integral to understanding the biological significance of electric organs. A number of groups of fishes, both electric and nonelectric, have receptors specialized for the detection of electric fields. These groups include all the elasmobranchs; the primitive bony fishes, lungfishes, bichirs, sturgeons, and paddlefish; the catfishes; and the two teleostean groups of fresh-water electric fishes. Electrosensory systems can be divided into two categories: those activated by fields from extrinsic sources, or passive systems, and those activated by the fishes' electric organs, or active systems. The passive systems generally are sensitive to voltages of low frequency such as are generated by respiratory movements in fishes, muscle or skin lesions, and geomagnetic sources. They can be extraordinarily sensitive. For example, elasmobranchs can respond to

voltage gradients as low as 0.01 microvolt/cm or 1 V uniformly distributed over 1000 km. The active systems utilize the electric organ discharge to test the conductivity of the surrounding medium, thereby revealing changes or inhomogeneities. The active systems are also involved in intra- and interspecific communication.

The electroreceptors are cutaneous organs distributed over the body surface. They are innervated by cranial nerves of the acousticolateralis system. The gymnotids and mormyrids have separate sets of receptors for active and passive systems; the electric organ discharges are often of high frequency with little dc component, and the two receptor systems detect electric fields over different frequency bands. The Mormyriiformes have a third set of receptors, very sensitive to high frequencies, that appears to be specialized for detecting the electric discharges of other fishes. The diversity of organ discharges in terms of both frequency and pulse shape may be partly accounted for in terms of interspecific communication and interaction.

In the skate the electric organ discharge is sufficiently long-lasting to be capable of activating the passive electrosensory system. It is unknown whether their discharges are utilized for electrolocation or for communication.

Evolution. The electrosensory systems allow an explanation of the evolution of strongly electric organs. Charles Darwin supposed that a weakly electric organ would have been an intermediate stage in the gradual development of a strongly electric organ, but could see no use for a weak organ. Later, Hans Lissman, who initially described the active electrosensory system, proposed that they provide a function for Darwin's hypothesized intermediate stage. This notion is supported by the electric eel and at least one torpedinid; both have strongly and weakly electric organs, and the eel is known to utilize the weak pulses for electrolocation in the same way as the related weakly electric gymnotids. The electric catfish has a passive electrosensory system, as do most nonelectric catfish. No weakly electric catfish are known, and the hypothetical intermediate leading to *Malapterurus* is either extinct or has not yet been discovered. The electrosensory capabilities of the stargazer have not been adequately investigated.

Electrocution. A question often asked about the strongly electric fishes is how they keep from electrocuting themselves. When immersed in water, they do not appear to be at all affected by their own discharges, although their electroreceptors are undoubtedly activated. In air, when the organ is not loaded down by the conductivity of the water, they often do twitch when they discharge their organs, but as would be expected, they are much more resistant than other fishes. One factor is that their nerves and central nervous system are well insulated by many layers of connective and fatty tissue. Whether their brains are more resistant to stunning by a given current density has not been determined; certainly the excitability of individual cells is no different from that in animals in general.

Source of data. Electric fishes have proved useful in defining extremes of excitability properties. Historically they served to validate Luigi Galvani's proposal of animal electricity, and in modern times have provided, and are continuing to provide, a wealth of biophysical data of general and comparative interest. They also supply important starting materials for the biochemical characterization of excitable membranes, for example, the acetylcholine receptor at the synapse. Because evolution has led to large volumes of tissue with high concentrations of the relevant macromolecules, the electric organs provide a rich source for isolation and analysis. See ACETYLCHOLINE; SYNAPTIC TRANSMISSION.

Michael V. L. Bennett

Electric power generation

The production of bulk electric power for industrial, residential, and rural use. Although limited amounts of electricity can be generated by many means, including chemical reaction (as in batteries) and engine-driven generators (as in automobiles and airplanes), electric power generation generally implies large-scale production of electric power in stationary plants designed for that purpose. The generating units in these plants convert energy from falling water, coal, natural gas, oil, and nuclear fuels to electric energy. Most electric generators are driven either by hydraulic turbines, for conversion of falling water energy; or by steam or gas turbines, for conversion of fuel energy. Limited use is being made of geothermal energy, and developmental work is progressing in the use of solar energy in its various forms. See GENERATOR; PRIME MOVER.

Fossil-fuel and hydroelectric plants, and various types of advanced power sources, are discussed in detail below. For a discussion of nuclear power plants see NUCLEAR POWER; NUCLEAR REACTOR.

General Considerations

An electric load (or demand) is the power requirement of any device or equipment that converts electric energy into light, heat, or mechanical energy, or otherwise consumes electric energy as in aluminum reduction, or the power requirement of electronic and control devices. The total load on any power system is seldom constant; rather, it varies widely with hourly, weekly, monthly, or annual changes in the requirements of the area served. The minimum system load for a given period is termed the base load or the unity load-factor component. Maximum loads, resulting usually from temporary conditions, are called peak loads, and the operation of the generating plants must be closely coordinated with fluctuations in the load. The peaks, usually being of only a few hours' duration (**Figs. 1 and 2**), are frequently served by gas or oil combustion-turbine or pumped-storage hydro generating units. The pumped-storage type utilizes the most economical off-peak (typically 10 P.M. to 7 A.M.) surplus generating capacity to

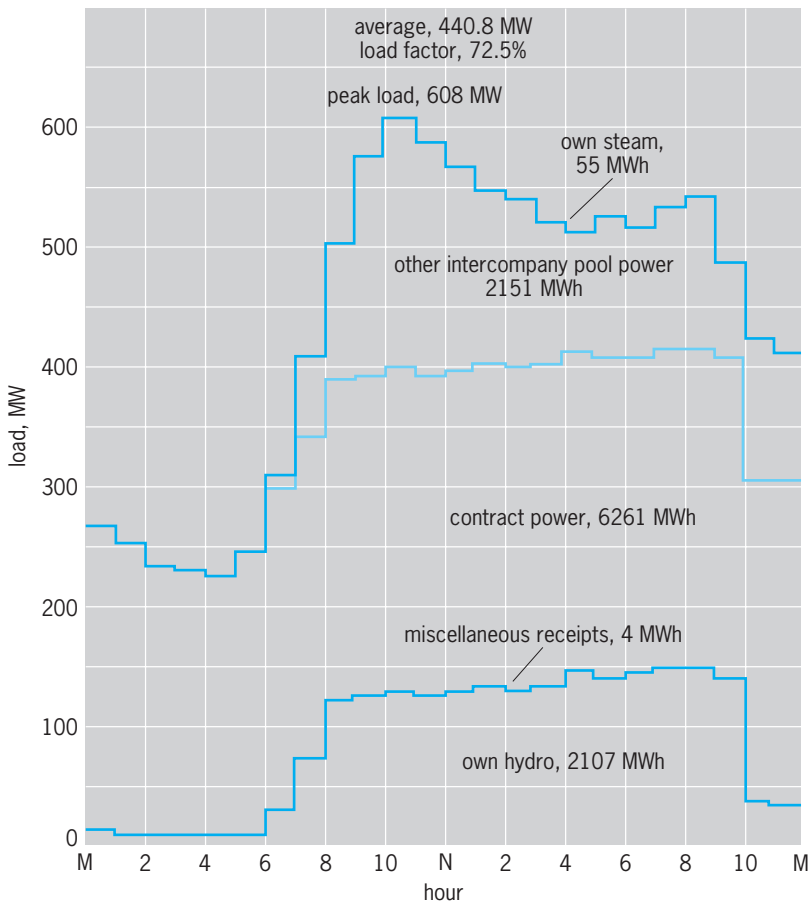


Fig. 1. Load graph indicates net system load of a metropolitan utility for typical 24-h period (midnight to midnight), totaling 10,578 MWh. Such graphs are made to forecast probable variations in power required.

released through hydraulic turbine generators during peak periods. This type of operation improves the capacity factors or relative energy outputs of base-load generating units and hence their economy of operation.

Actual variations in the load with time are recorded, and from these data load graphs are made to forecast the probable variations of load in the future. A study of hourly load graphs (Figs. 1 and 2) indicates the generation that may be required at a given hour of the day, week, or month, or under unusual weather conditions. A study of annual load graphs and forecasts indicates the rate at which new generating stations must be built; they are an inseparable part of utility operation and are the basis for decisions that profoundly affect the financial requirements and overall development of a utility.

Generating unit sizes. The size or capacity of electric utility generating units varies widely, depending upon type of unit; duty required, that is, base-, intermediate-, or peak-load service; and system size and degree of interconnection with neighboring systems. Base-load nuclear or coal-fired units may be as large as 1200 MW each, or more. Intermediate-duty generators, usually coal-, oil-, or gas-fueled steam units, are of 200 to 600 MW capacity each. Peaking units, combustion turbines or hydro, range from several tens of megawatts for the former to hundreds of megawatts for the latter. Hydro units, in both base-load and intermediate service, range in size up to 825 MW.

The total installed generating capacity of a system is typically 20 to 30% greater than the annual

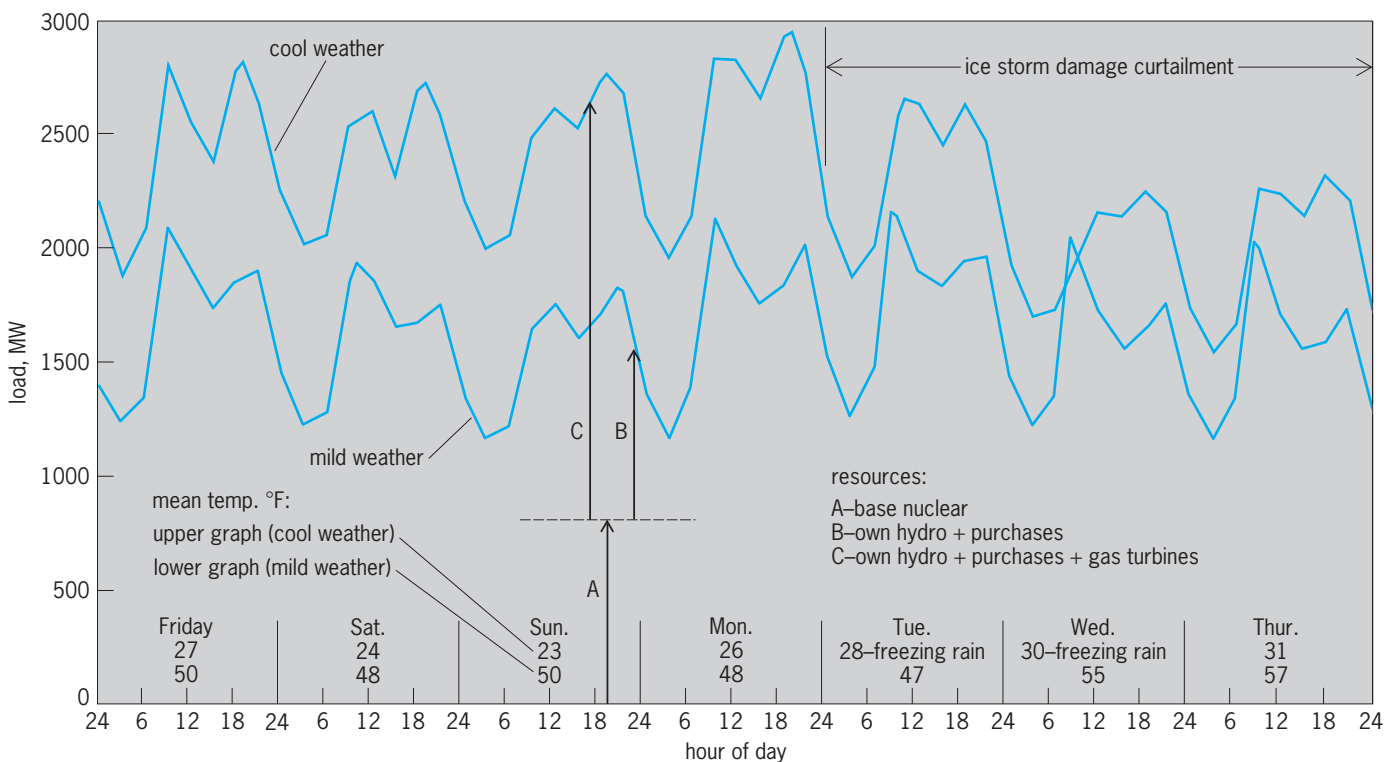


Fig. 2. Examples of northwestern United States electric utility weekly load curves showing same-year weather influence.

predicted peak load in order to provide reserves for maintenance and contingencies.

Power-plant circuits. Both main and accessory circuits in power plants can be classified as follows:

1. Main power circuits to carry the power from the generators to the step-up transformers and on to the station high-voltage terminals.
2. Auxiliary power circuits to provide power to the motors used to drive the necessary auxiliaries.
3. Control circuits for the circuit breakers and other equipment operated from the plant's control room.
4. Lighting circuits for the illumination of the plant and to provide power for portable equipment required in the upkeep and maintenance of the plant. Sometimes special circuits are installed to supply the portable power equipment.
5. Excitation circuits, which are so installed that they will receive good physical and electrical protection because reliable excitation is necessary for the operation of the plant.
6. Instrument and relay circuits to provide values of voltage, current, kilowatts, reactive kilovoltamperes, temperatures, pressures, flow rates, and so forth, and to serve the protective relays.
7. Communication circuits for both plant and system communications. Telephone, radio, transmission-line carrier, and microwave radio may be involved.

It is important that reliable power service be provided for the plant itself, and for this reason station service is usually supplied from two or more sources. To ensure adequate reliability, auxiliary power supplies are frequently provided for start-up, shutdown, and communication services.

Generator protection. Necessary devices are installed to prevent or minimize other damage in cases of equipment failure. Differential-current and ground relays detect failure of insulation, which may be due to deterioration or accidental overvoltage. Overcurrent relays detect overload currents that may lead to excessive heating; overvoltage relays prevent insulation damage. Loss-of-excitation relays may be used to warn operators of low excitation or to prevent pulling out of synchronism. Bearing and winding overheating may be detected by relays actuated by resistance devices or thermocouples. Overspeed and lubrication failure may also be detected.

Not all of these devices are used on small units or in every plant. The generator is immediately deenergized for electrical failure and shut down for any over-limit condition, all usually automatically. *See* ELECTRIC PROTECTIVE DEVICES.

Voltage regulation. This term is defined as the change in voltage for specific change in load (usually from full load to no load) expressed as percentage of normal rated voltage. The voltage of an electric generator varies with the load and power factor; consequently, some form of regulating equipment is required to maintain a reasonably constant and predetermined potential at the distribution stations or load centers. Since the inherent regulation of most alternating-current (ac) generators is rather

poor (that is, high percentagewise), it is necessary to provide automatic voltage control. The rotating or magnetic amplifiers and voltage-sensitive circuits of the automatic regulators, together with the exciters, are all specially designed to respond quickly to changes in the alternator voltage and to make the necessary changes in the main exciter or excitation system output, thus providing the required adjustments in voltage. A properly designed automatic regulator acts rapidly, so that it is possible to maintain desired voltage with a rapidly fluctuating load without causing more than a momentary change in voltage even when heavy loads are thrown on or off.

In general, most modern synchronous generators have excitation systems that involve rectification of an ac output of the main or auxiliary stator windings, or other appropriate supply, using silicon controlled rectifiers or thyristors. These systems enable very precise control and high rates of response. *See* SEMICONDUCTOR RECTIFIER; VOLTAGE REGULATOR.

Generation control. Computer-assisted (or on-line controlled) load and frequency control and economic dispatch systems of generation supervision are being widely adopted, particularly for the larger new plants. Strong system interconnections greatly improve bulk power supply reliability but require special automatic controls to ensure adequate generation and transmission stability. Among the refinements found necessary in large, long-distance interconnections are special feedback controls applied to generator high-speed excitation and voltage regulator systems.

Synchronization of generators. Synchronization of a generator to a power system is the act of matching, over an appreciable period of time, the instantaneous voltage of an alternating-current generator (incoming source) to the instantaneous voltage of a power system of one or more other generators (running source), then connecting them together. In order to accomplish this ideally the following conditions must be met:

1. The effective voltage of the incoming generator must be substantially the same as that of the system.
2. In relation to each other the generator voltage and the system voltage should be essentially 180° out of phase; however, in relation to the bus to which they are connected, their voltages should be in phase.
3. The frequency of the incoming machine must be near that of the running system.
4. The voltage wave shapes should be similar.
5. The phase sequence of the incoming polyphase machine must be the same as that of the system.

Synchronizing of ac generators can be done manually or automatically. In manual synchronizing an operator controls the incoming generator while observing synchronizing lamps or meters and a synchroscope, or both. The operator closes the connecting switch or circuit breaker as the synchroscope needle slowly approaches the in-phase position.

Automatic synchronizing provides for automatically closing the breaker to connect the incoming machine to the system, after the operator has properly

adjusted voltage (field current), frequency (speed), and phasing (by lamps or synchroscope). A fully automatic synchronizer will initiate speed changes as required and may also balance voltages as required, then close the breaker at the proper time, all without attention of the operator. Automatic synchronizers can be used in unattended stations or in automatic control systems where units may be started, synchronized, and loaded on a single operator command. See ALTERNATING-CURRENT GENERATOR; PHASE-ANGLE MEASUREMENT; SYNCHROSCOPE.

Eugene C. Starr

Fossil-Fuel Plants

Fossil fuels are of plant or animal origin and consist of hydrogen and carbon (hydrocarbon) compounds. The most common fossil fuels are coal, oil, and natural gas. The less common ones include peat, oil shale, and biomass (wood and so forth), as well as various waste or by-products such as steel mill blast furnace gas, coke-oven gas, and refuse-derived fuels. Fossil-fuel electric power generation uses the combustion heat energy from these fuels to produce electricity. See COAL; COKE; FOSSIL FUEL; NATURAL GAS; OIL SHALE; PEAT; PETROLEUM.

Steam power plants. A fossil-fuel steam power plant operation essentially consists of four steps (Fig. 3): (1) Water is pumped at high pressure to a boiler, where (2) it is heated by fossil-fuel combustion to produce steam at high temperature and pressure. (3) This steam flows through a turbine, rotating an electric generator (connected to the turbine shaft) which converts the mechanical energy to electricity. (4) The turbine exhaust steam is condensed by using cooling water from an external source to

remove the heat rejected in the condensing process. The condensed water is pumped back to the boiler to repeat the cycle. Figure 3 also shows features to increase cycle efficiency, including preheating of the boiler feedwater by using steam extracted from the turbine, and reheating the high-pressure turbine exhaust steam before it enters the intermediate-pressure turbine.

Steam cycle. Modern fossil-fuel power plants are based on a steam cycle first proposed by W. J. M. Rankine in 1908. The basic Rankine cycle underwent conceptual evolution through the 1950s that improved cycle efficiency to current levels but which also increased the complexity of building and operating a fossil-fuel power plant. Commercial development lagged behind conceptual cycle development because of the unavailability of suitable materials and fabrication techniques to accommodate the higher steam pressures and temperatures required to increase cycle efficiency. **Table 1** summarizes this steam cycle evolution with respect to turbine inlet steam pressure, turbine heat rate, and turbine cycle and overall thermal efficiencies, along with related available main steam piping materials from the 1930s through the 1950s. See RANKINE CYCLE; THERMODYNAMIC CYCLE; VAPOR CYCLE.

During the 1960s and 1970s, there was no significant cycle efficiency improvement, primarily due to the advent of nuclear power and the abundant supply of low-cost fossil fuel. In the 1980s, the interest in cycle efficiency improvement was renewed due to the uncertain future of nuclear energy and escalating fossil-fuel costs. However, due to the return of relative cost and price stability of fossil fuels since the late 1980s, interest in cycle efficiency

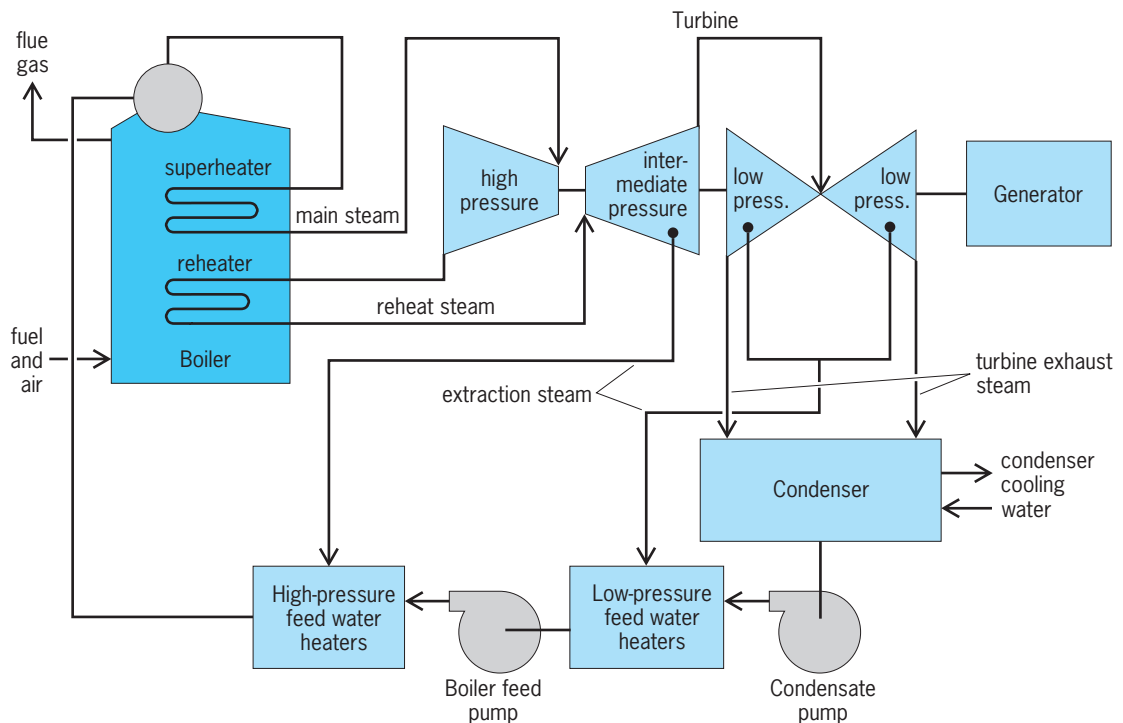


Fig. 3. Steam power plant flow diagram.

TABLE 1. Steam cycle development in fossil-fuel power plant

Years	Steam cycle, psig/°F/°F* (MPa/°C/°C)	Main steam piping material	Turbine cycle [†] heat rate, Btu/kWh (MJ/kWh)	Thermal efficiency (%)	
				Turbine cycle [†]	Overall [‡]
Late 1930s	1250/950 (8.7/510)	Chromium-molybdenum alloy (0.2% C, 0.5% Cr, 0.5% Mo)	9350 (9.86)	36.50	29.2
Mid 1940s	1450/1000/1000 (10.1/538/538)	Chromium-molybdenum alloy (0.2% C, 1.0% Cr, 0.5% Mo)	8150 (8.60)	41.9	33.5
Late 1940s	2000/1050/1050 (13.9/566/566)	Chromium-molybdenum alloy (0.2% C, 1.0% Cr, 0.5% Mo)	7700 (8.12)	44.3	35.5
Early 1950s	2400/1100/1050	Chromium-molybdenum	7500 (7.91)	45.5	36.4

*Main steam pressure/main steam temperature/reheat steam temperature.
[†] Cycle heat rate and thermal efficiency are at rated load for a condenser backpressure of 1.43 in. Hg (36.3 mmHg or 4.84 kPa) absolute and a condensate temperature of 90° F (32.2° C).
[‡] Based on a boiler efficiency of 85% and a power plant auxiliary power requirement of 5% of plant output.

improvement has diminished significantly. In fact, few if any large steam power plants with main steam pressure, main steam temperature, and reheat steam temperature higher than 2400 psig, 1100 °F, and 1050 °F (16.6 MPa, 593 °C, and 566 °C) respectively (the level achieved in the early 1950s) have been built since 1990. With the rapid development of Asian and Pacific rim countries including China, there is renewed interest in smaller-size 200–300-MW units, with emphasis on reducing capital cost.

Plant size and facilities. A given power plant’s size is determined by the utility company’s need for power as dictated by the electrical demand growth forecast. The additional generating capacity needed is matched to commercially available turbine-generator capacity sizes. The boiler is selected to suit the turbine steam flow requirements, and the remaining accessory equipment is sized to serve the needs of the resulting boiler-turbine steam cycle.

A typical large fossil-fuel power plant consists of several major facilities and equipment (Fig. 4), including fuel handling and processing, boiler (including furnace), turbine and electric generator, con-

denser and condenser heat removal system, feedwater heating and pumping system, flue gas-cleaning system, and plant controls and control system.

1. *Fuel handling.* Fuel type determines the fuel-handling system requirements. Natural gas is normally delivered by pipeline directly and fired in the furnace, needing no special handling or storage. For coal and oil, which are usually delivered by rail, ship, or pipelines, the fuel-handling system includes fuel unloading, storage, reclaiming from storage, and fuel preparation for combustion (such as coal pulverizing and pneumatic conveying to the burners, and oil heating and atomizing). Coal is usually stored in open stockpiles and oil in closed tanks. Normally, 30 to 180 days of fuel requirement is stored at the plant site to maintain operation during fluctuations in fuel delivery.

2. *Boiler.* The fuel type also determines the size, design, and operation of the furnace and boiler, which become larger and more complex (in ascending order) for gas, oil, and coal. For coal, the size and complexity also increase with the decline in coal rank and quality (in ascending order)

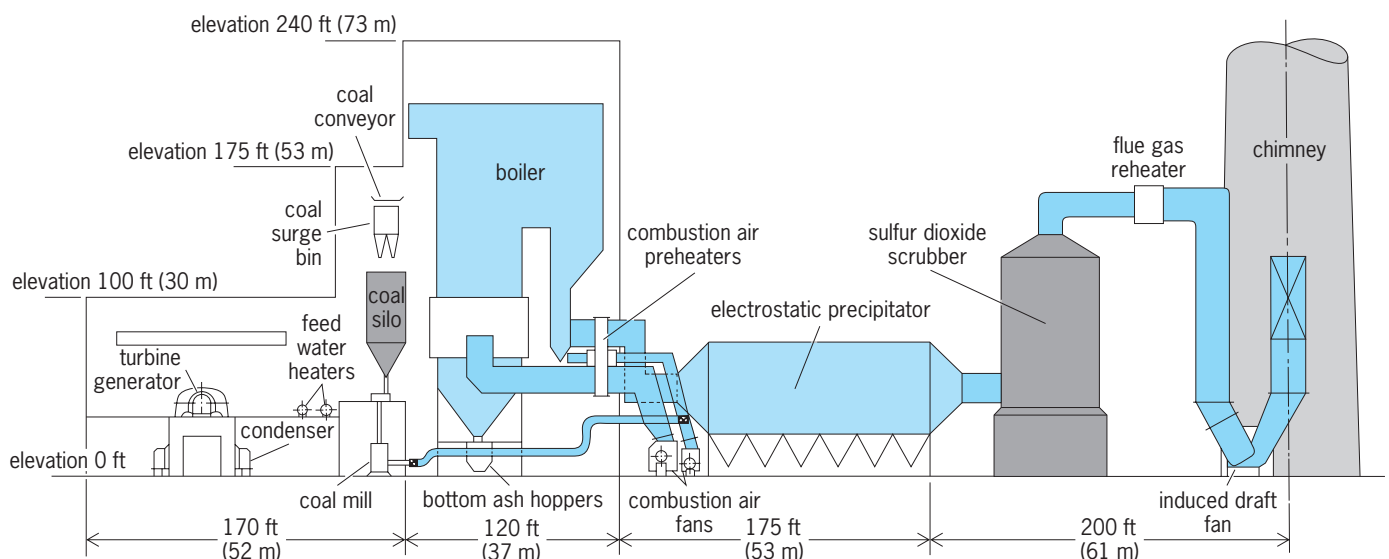


Fig. 4. Schematic cross section of a 600-MW coal-fired steam power plant.

for bituminous, subbituminous, lignite, and brown coal. Specialized boiler designs are required to handle peat, shale, and the various waste fuels. One evolving technology, fluidized-bed boilers, can effectively burn a wide range of low-quality solid fuels in the same boiler. Units of up to 200 MW are commercially proven and available. Most operating plants are in the size range of 120–180 MW. *See* BOILER; FIRE-TUBE BOILER; FLUIDIZED-BED COMBUSTION; STEAM-GENERATING FURNACE; STEAM-GENERATING UNIT; WATER-TUBE BOILER.

3. *Turbines.* Modern power plants, depending on the size and steam conditions, have single or multiple-casing turbines with high-, intermediate-, and low-pressure sections. High-pressure steam from the boiler after expanding through one or more turbines is exhausted to a condenser under vacuum, where the steam is condensed by cooling water. *See* STEAM TURBINE.

4. *Condenser.* About 50% of the fuel heat input to the boiler is ultimately rejected as waste heat in the condenser. Condenser heat-removal-system selection depends on the geography, climatology, and water availability at each plant site. There are basically two types of heat rejection systems: open (or once-through) cooling and closed cooling. The open system is the most simple and normally used where water is abundantly available. It takes water from an ocean, lake, or river and discharges the heated water to a downstream location. The closed system is used when it is determined not to return heated water to the water source due to environmental or other constraints. The closed system rejects heat to the atmosphere by using a secondary heat rejection system such as a cooling tower or pond. However, closed-system operation involves evaporating a portion of the water circulated, requiring a water supply capable of restoring the quantity of water evaporated. Dry cooling (with air condensers) is less efficient, but is used in arid areas where the water supply is limited. The temperature of the cooling medium (water or air) at the condenser inlet determines the cycle thermal efficiency by establishing the steam condensation temperature. *See* COOLING TOWER; STEAM CONDENSER; VAPOR CONDENSER.

5. *Feedwater heating.* The condensed steam is pumped back to the boiler inlet to complete the cycle, with heat from turbine extraction steam transferred along the way to increase cycle efficiency. Condensate pumps deliver feedwater from the condenser through low-pressure feedwater heaters to the boiler feed pumps, which increase the feedwater pressure and pump it through high-pressure feedwater heaters to the boiler inlet. A modern power plant typically has six feedwater heaters, which can be increased in number to improve cycle efficiency, but with offsetting increases in plant complexity and cost.

6. *Flue gas and liquid waste.* Fossil-fuel power plants produce gaseous and liquid wastes which are treated prior to discharge to minimize effects on the environment. In many countries, environmental reg-

ulations define the extent of flue gas and liquid waste treatment required.

Major pollutants in the flue gas include particulates (except for natural gas fuel), sulfur dioxide, and nitrogen oxides. Incombustible matter in coal and oil yields particulates upon combustion which are collected in devices such as electrostatic precipitators and fabric filters (baghouses). These devices are capable of removing more than 99% of the particulate matter. *See* AIR FILTER; ELECTROSTATIC PRECIPITATOR.

A variety of chemical processes is available for the removal of sulfur dioxide from the flue gas. Flue gas desulfurization systems use alkaline chemicals such as lime, limestone, sodium, or magnesium compounds in solution, slurry, or dry form to react with and remove sulfur dioxide. The choice depends on the fuel sulfur content, performance requirement, capital and operating costs, and market for the desulfurization by-products. The throw-away wet limestone slurry process is the most widely used in the United States due to the limited market for by-products. A similar process producing gypsum as the by-product is common in Europe and Japan. *See* GAS ABSORPTION OPERATIONS.

Nitrogen oxides emission can be reduced by combustion modifications which limit their formation during combustion. These techniques involve reduction of flame temperature by prolonging the combustion time (staged combustion) and reducing the combustion air quantity (low-nitrogen oxides burners).

Processes to remove nitrogen oxides after their formation involve injection of ammonia into the hot flue gas in the presence or absence of catalysts. They are referred to as the selective catalytic reduction and noncatalytic reduction processes, respectively. The selective catalytic reduction process is widely used in Japan and Europe, particularly in Germany, in both large-scale electric utilities and in small-scale industrial boilers. The noncatalytic reduction process is used in the United States in some small-scale industrial boilers. Typical nitrogen oxide reductions with the selective catalytic reduction process range from 70 to 90% and with the noncatalytic reduction process from 30 to 70%.

Waste liquid pollutants include suspended solids, toxic metal compounds, dissolved chlorine, and high-acidity or high-alkalinity, high-temperature, and oily compounds. Treatment prior to discharge or reuse usually involves one or more of the following: neutralization, sedimentation, clarification, filtration, dechlorination, and incineration. *See* CLARIFICATION; FILTRATION; SEDIMENTATION (INDUSTRY).

7. *Plant controls.* The power plant size and steam cycle pressure and temperature determine plant controls and control systems requirements. In ascending order of complexity and sophistication, these include pneumatic analog controls, electronic analog controls, hard-wired digital controls, and distributed digital control with microprocessors and minicomputers. Modern power plants typically use pneumatic local controls and electronic analog remote

TABLE 2. Typical cost distribution of major systems and equipment for a coal-fired steam power plant*

Systems and equipment	% of capital cost
Boiler (including furnace)	30
Turbine and electric generator	18
Flue gas desulfurization system	16
Particulate collection system	7
Coal-handling equipment	5
Condenser and cooling power	5
Ash-handling equipment	3
Pumps	3
Chimney	2
Water treatment	1
Other mechanical equipment	10

*Excluding land, site preparation, and building costs.

controls with electromechanical or solid-state interlocking devices and variable-speed motor drives. *See* CONTROL SYSTEMS; DIGITAL CONTROL; DISTRIBUTED SYSTEMS (CONTROL SYSTEMS); MICROCOMPUTER; MICROPROCESSOR; PROCESS CONTROL.

A typical cost distribution by major systems and equipment for a coal-fired steam power plant is given in **Table 2**.

Gas turbine plants. Power plants with gas turbine-driven electric generators are often used to meet short-term peaks in electrical demand. They are generally small (up to 150–200 MW), have a low thermal efficiency, require clean fuels such as natural gas or light oils, and consequently are more expensive to operate. However, these plants have a relatively low capital cost and can be used in areas with low power demand and limited water availability. Diesel engine-driven electric generators are also used under similar conditions. *See* DIESEL ENGINE.

Gas turbine power plants use atmospheric air as the working medium, operating on an open cycle where air is taken from and discharged to the atmosphere and is not recycled. In a simple gas turbine plant (**Fig. 5**), air is compressed and fuel is injected into the compressed air and burned in a combustion chamber. The combustion products expand through a gas turbine and exhaust to the atmosphere. Vari-

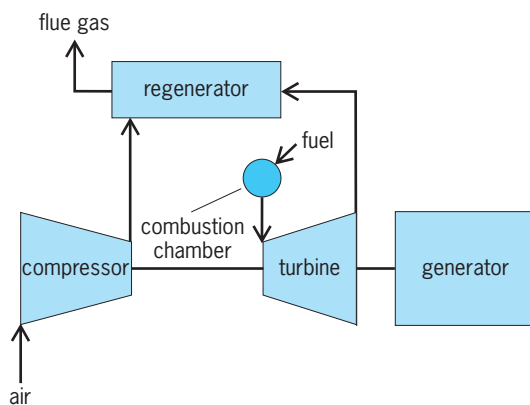


Fig. 5. Gas turbine plant flow diagram.

ations of this basic operation to increase cycle efficiency include regeneration (where exhaust from the turbine is used to preheat the compressed air before it enters the combustion chamber) and reheating (where the combustion gases are expanded in more than one stage and are reheated between stages).

Combined-cycle plants use both gas and steam turbines and offer overall thermal efficiency improvement and operating economy. The gas turbine exhaust heat is used to generate steam for operating steam turbine generators. Combined-cycle plants with overall thermal efficiencies greater than 50%, where operating economies justify it, are commercially proven and available. Some very large plants in the 500–1500-MW size range are in commercial operation. *See* GAS TURBINE; STEAM TURBINE.

Cogeneration is another method of improving overall thermal efficiency, where steam exhausted from the steam turbine generator is used in industrial processes. Substantial energy is saved as heat rejected in the condenser is put to use. *See* COGENERATION.

Advanced concepts. Advanced fossil-fuel power generation concepts such as integrated coal gasification combined cycle and magnetohydrodynamics offer the promise of significant overall thermal efficiency improvement. In the integrated coal gasification combined cycle, gas produced from coal (after proper purification) is burned and expanded through a gas turbine, with the exhaust gas heat used to generate steam for a steam turbine generator. In the magnetohydrodynamic process, an electrically conducting medium such as hot flue gas from a coal combustor (seeded with alkali metal salts to increase the flue gas electrical conductivity) is passed through a magnetic field (magnetohydrodynamic channel) to produce electricity. The magnetohydrodynamic channel exhaust gas is then used to generate steam in a conventional steam power plant. The alkali salts are recovered and reused. Power generation based on these concepts can reach a thermal efficiency of 48–52%, whereas a steam turbine plant has an efficiency of 30–40% and a conventional gas turbine plant 22–28%. *See* COAL GASIFICATION; MAGNETOHYDRODYNAMIC POWER GENERATOR.

Maris T. Fravel; Natarajan Sekhar

Hydroelectric Plants

Hydroelectric generation is an attractive source of electric power because it is a renewable resource and a nonconsumptive use. In the broadest sense, hydroelectric power is a form of solar power; the resource is renewed by the solar cycle in which water is evaporated from the oceans, transported by clouds, and falls as precipitation on the landmasses, and returns through rivers to the ocean, generating power on the way. Hydroelectric power can be defined as the generation of electricity by flowing water; potential energy from the weight of water falling through a vertical distance is converted to electrical energy. The amount of electric power P that can be

generated is given by the equation below, where Q

$$P = \frac{Qb}{k} \quad \text{kW}$$

is the volume flow of water, b is the height through which the water falls, and k is a constant equal to 11.8 when Q is given in cubic feet and b in feet (English units), and equal to 0.102 when Q is in cubic meters and b is in meters (SI units).

Development. In the United States, the first hydroelectric generating station was built in 1882 in Appleton, Wisconsin. The total power output was sufficient to operate some 250 electric lights. Rapid expansion followed, and by 1900 about 300 hydroelectric plants were in operation worldwide. This rapid expansion was a consequence of the development in the nineteenth century of the hydraulic reaction turbine. Also, by the same date, the concept of pumped storage had been introduced in Switzerland, but extensive development in the United States did not occur until the advent of the reversible pump-turbine design in the early 1950s.

Major components. A typical hydroelectric development consists of a dam to divert or store water; waterways such as a forebay, canals, tunnels, and penstocks to deliver the water to the hydraulic turbine and a draft tube, tunnel, or tailrace to return the water to the stream; hydraulic turbines and gover-

nors; generators and exciters; electrical controls to provide protection and to regulate frequency, voltage, and power flow; a powerhouse to enclose the machinery and equipment; transformers and switching equipment; and a transmission line to deliver the power to the load center for ultimate distribution (Fig. 6).

Dam. Dams are functionally important as they regulate the water supply to the plant and also are frequently the largest single element of cost in a hydroelectric development. Dams may be classified according to materials or type of construction, height, or use. See DAM.

Waterways. Waterways may be very short or as much as 10 mi (16 km) or more in length. The type of waterway is selected on the basis of length, pressure, plant arrangement, and economy. Very short waterways are required in installations where the powerhouse and the dam are integral structures; then the waterway is short and the length dictated by the plant arrangement. In other low-head plants, canals and forebays are frequently used. Conversely, very long waterways are also common in the developments where the head is great. Tunnels and penstocks are commonly used in high-head developments to deliver the water to a powerhouse some distance from the impoundment.

Some of the important aspects of the hydraulic design of waterways are provision for dewatering to facilitate inspection and maintenance, provision for water hammer resulting from surges caused by sudden start-up or shutdown, and optimization of hydraulic losses for an economic design. See WATER HAMMER.

An intake structure may be provided at the entrance to the waterway to reduce the hydraulic losses and to provide a closure gate for dewatering the waterway. Surge chambers are required at abrupt vertical changes in the direction of the waterways to provide relief of surges on shutdown and to prevent separation of the water column on start-up. When tunnels are used for the tailrace, surge chambers may be required for surge suppression.

Penstocks are defined as large-diameter pipelines for delivering water to the turbine or the steel-lined part of a tunnel adjacent to the powerhouse. Materials other than steel have been used in penstock construction; in the past, wood staves were common; and fiber glass is now used selectively.

Turbines. Hydraulic turbines are used to drive an electric generator, thereby converting the potential energy of the water to electrical energy. Reaction- or impulse-type turbines are used depending upon the head available, with a reaction turbine being used for heads under 1200 ft (365 m) and an impulse turbine for greater heads. Reaction turbines operate in a closed system from headwater to tailwater. For impulse turbines, only the portion to the turbine is closed and under pressure; from that point, flow is by gravity. As water passes through the reaction turbine, conversion from potential to mechanical energy takes place and, through a shaft,

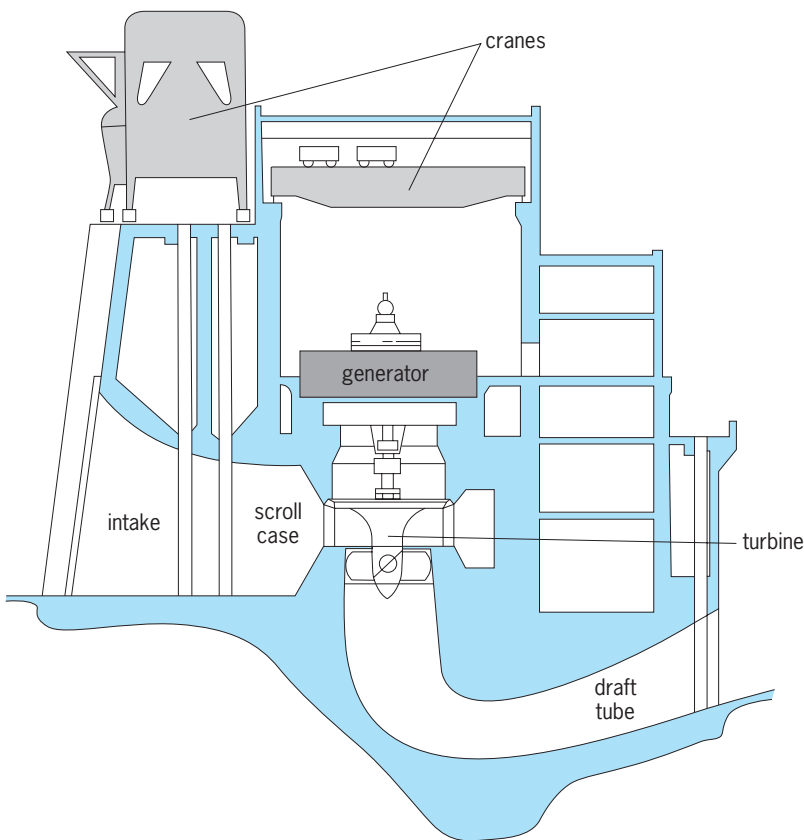


Fig. 6. Cross section of a hydroelectric powerhouse. Arrangement is typical for an adjustable propeller-type reaction turbine. The intake and powerhouse are integral; waterways consist of intake, scroll case, and draft tube. Turbine and generator are set with the axis vertical.

drives the generator. For an impulse turbine, the potential energy of the water is first converted to kinetic energy in high-velocity jets which impinge on buckets at the circumference of the impulse wheel. There the kinetic energy is converted to mechanical energy for driving the generator. Reaction turbines are classified as Francis or propeller types, and propeller types are further categorized as fixed-blade and adjustable-blade, or Kaplan-type. Both classes of reaction turbines are normally installed with the axis vertical. *See* HYDRAULIC TURBINE; IMPULSE TURBINE; REACTION TURBINE.

Generator. In vertical installations, the generator is located above the turbine. When the turbine rotates, so does the rotor of the synchronous generator which produces electrical energy in a conventional way. An operating speed is chosen, as well as a generator with the correct number of poles to develop the desired frequency. From the generator, electricity is carried by buses to switching apparatus and then to transformers, where the voltage is stepped up for transmission. *See* ALTERNATING-CURRENT GENERATOR; HYDROELECTRIC GENERATOR.

Bulb units. Bulb units, which take their name from their shape, are an adaptation of the propeller unit and differ from the conventional unit in two significant respects: the alignment is horizontal, and the generator is totally enclosed within the bulb and mounted in the water passage. Water flows to the unit in a horizontal conduit and passes in the annular space between the bulb and the walls of the conduit. As the water passes the bulb, it drives the propeller and exits from the unit through a divergent horizontal water passage. Flow is controlled in a manner similar to the vertical units by wicket gates.

Classification of developments. Hydroelectric developments may be classified according to purpose, use, location of powerhouse, and head. The descriptive classifications commonly used are the run-of-river or storage, high- or low-head, base- or peak-load, multipurpose or single-purpose, and conventional or pumped-storage. Powerhouses may be indoor or outdoor, aboveground (surface) or underground. A number of these classifications may apply to any specific project such as conventional, run-of-river with surface, outdoor powerhouse, or pumped-storage with an underground powerhouse.

Multipurpose projects. A multipurpose project has a number of uses, some of which may be governed by competing criteria. Hydroelectric generation is generally compatible with other uses because it is non-consumptive. Such projects normally require large reservoirs to accommodate the multiple uses. The most frequent purposes are irrigation, power, flood control, navigation, water supply, and recreation. Some adjustment must be made in either the amount of or the timing of the generation for compatibility with other uses. On major river systems in the United States, multipurpose projects have been developed in conjunction with navigation, flood control, and power. In the semiarid states, irrigation has played a dominant role in multipurpose developments.

Pumped storage. Pumped storage is a process for converting large quantities of electrical energy to potential energy by pumping water to a higher elevation where it can be stored indefinitely, then released to pass through hydraulic turbines and generate electrical energy on demand. Storage is desirable, as the consumption of electricity is highly variable according to the time of day or week, as well as seasonally. Consequently, there is excess generating capacity at night and on weekends. This excess capacity can be used to generate energy for pumping, hence storage. Normally, pumping energy can be obtained cheaply at night or on weekends, and its value will be upgraded when used for daytime peak loads.

In a typical operation, water is pumped at night or on weekends from a lower to an upper reservoir, where it is stored. The water can be retained indefinitely without deterioration or significant loss. During the daylight hours when the loads are greatest, stored water is released to flow from the upper to the lower reservoir through hydraulic turbines which generate electricity. No water is consumed in either the pumping or the generating cycles. A hydroelectric pumped-storage development is similar to a typical hydro installation, except that it has two reservoirs of essentially equal size situated to maximize their difference in elevation. Also, since it is a closed system, it need not be located directly on a large stream but must have a source for initial filling and to make up for the losses of evaporation and seepage. The second principal difference is the pump-turbine which is capable of rotation in either direction and acts as a pump in one direction and a turbine in the other. *See* ENERGY STORAGE; WATER-POWER.
Dwight L. Glasscock

Advanced Power Sources

Interest has developed in alternative or advanced power sources. Most of this attention has focused on solar or solar-related technologies, although several additional options including geothermal and fusion have been examined. Genuine interest in alternatives developed following the oil embargo in 1973.

The supplies of fossil and nuclear fuels are classified as nonrenewable and therefore are gradually being depleted. The advanced power sources of interest are those which are renewable and linked to the Sun, naturally occurring steam or hot water in the Earth's crust, or fusion. Since all solar sources suffer from intermittent availability, they require special consideration. As these various options are considered, three significant problems must be addressed. First, availability and the primary form of occurrence of the resource must be assessed. Second, a conversion system must be developed to transform the source into a usable form of energy, typically electricity. Third, the electricity production source must be effectively and efficiently appended to the existing electrical supply system.

Solar power. The solar options consisting of photovoltaics, wind, solar thermal process, and biomass will be discussed below. Other solar alternatives,

including solar ponds, passive solar, and ocean-thermal energy conversion, and other ocean power systems such as tide and wave power will not be treated. *See* TIDAL POWER.

Photovoltaics. The collection of solar energy takes many forms, but one of its most desirable configurations is direct conversion to electricity. The device configurations are modular in form, making them convenient from a mass-production viewpoint.

The heart of the photovoltaic system is a thin flat layer of semiconductor material. When the material is struck by sunlight, electrons are freed, producing an electric current. The direct-current (dc) power is passed through a dc load, into a storage battery, or converted to alternating current (ac) for general use in electric utility grids. Typically, individual solar cells are ganged together to form photovoltaic modules about 5 ft² (0.5 m²) in size with a generating capacity of about 50 W. Typically, about half the cost of a solar system lies with the solar cell modules, and the remainder is directed toward power conditioning, electrical wiring, installation, and site preparation.

Typically, silicon or gallium arsenide are used to fabricate solar cells, although other semiconductor materials are being developed. Silicon technology is the most advanced because it is the least expensive, and takes many different forms including single-crystal, polycrystal, and amorphous configurations. The efficiency of these configurations ranges from about 6 to about 14% in module form. Amorphous silicon offers significant potential because of low manufacturing cost.

One way to advance the power yield from photovoltaics is to use concentrator systems. Special lenses and mirrors are used to focus sunlight on the cells, thereby raising the output of a module. Concentrator systems result in larger electric current which increases cell losses but improves overall economic efficiency of the system. Concentrator technology can raise module efficiencies to about 20%. Photovoltaic systems have potential application at the household level; however, the best near-term application is in array configurations for utility application. *See* SOLAR CELL.

Wind. The use of wind energy dates back to sailing ships and windmills. Today the interest in wind is for electricity generation by wind turbines. Wind energy conversion systems consist of several major subsystems. Among the most important is the mechanical system and electromechanical rotary converter. A wind energy conversion system is designed to rotate at either constant or variable speed as the wind varies. The variable-speed system usually offers high wind-collection efficiency; however, constant-speed units permit simpler electrical systems.

There are a multitude of aeroturbine designs available, with each striving to enhance the wind to mechanical power conversion effectiveness. A measure of this effectiveness is the power coefficient (C_p). The power coefficient is expressed as a function of the tip speed ratio, that is, the ratio of the blade tip

speed to the wind speed. In variable-speed systems, C_p is tuned to its optimum value, whereas constant-speed systems cannot operate at optimum C_p over the entire wind speed range. This does not suggest that variable speed is a panacea, since it exposes the mechanical drive train and tower to potential resonance problems and requires more power conditioning of the electrical output.

The characteristics of the aeroturbine and the electric generator are not totally compatible; therefore, an element is required to effectively tie these two components. The mechanical interface has two major assignments. First, aeroturbines usually rotate at low speeds, which are not desirable for electrical production. Typically, a transmission is used to convert low-speed, high-torque mechanical power into high-speed, low-torque mechanical power for electrical conversion. Second, the mechanical interface can be used to regulate the drive train shaft stiffness, and thus shape its dynamic performance.

Ultimately the mechanical energy that has been produced by the aeroturbine and conditioned by the mechanical interface must be converted to an electrical form. In general, the type of wind system, constant- or variable-speed, determines the type of electric generator—synchronous, induction, or dc. Variable-speed systems tend to produce initially poor-quality ac or dc electric power, and then use power conditioners in the electrical interface to enhance it. Usually, constant-speed systems produce high-quality electric power at the generator terminals.

Further development of wind turbine technology is required to provide economical systems; however, wind electric systems show promise. Systems under development will produce as much 4–5 MW per machine. Typically, units larger than 200 kW are used by utilities. *See* WIND POWER.

Solar thermal. The process of collecting and concentrating solar energy at a focal point is known as solar thermal. There are two broad classes of solar thermal collection systems: power tower and point-and-line-focus. The power tower is a large central power-generating system which shows significant promise, while the point- and line-focus systems are smaller and would be deployed as so-called dispersed sources. In both cases, solar energy is focused at a point or a line and heat is transferred to a working fluid.

The power tower concept has received significant attention because of its potentially large power production capability. A 10-MW demonstration plant, called Solar One, developed by the U.S. Department of Energy near Barstow, California, consists of 1818 heliostats, each containing 12 separate mirror facets designed to aim their energy at a central receiver on top of a structural steel tower. The heliostats are computer-controlled and adjust during the day to track the Sun so that optimal collection efficiency is maintained. The working fluid in this system is water, which is converted to steam and then passed through

a conventional steam turbine for eventual conversion to electricity. One significant application for the solar thermal system is as a repowering plant. Repowering is a term applied to the concept of adding a solar thermal plant to a conventional power plant to supplement steam production.

The point- and line-focus systems are smaller facilities which produce power at lower temperatures (200–600°F or 100–300°C). A typical point system is a parabolic dish which concentrates sunlight on a receiver a few feet away. Line-focus or trough systems consist of a series of cylindrical or parabolic trough lines with mirrors to collect and concentrate the Sun's radiation. Both the point- and line-focus systems could be used in smaller cogeneration facilities for commercial and industry electric power systems. See SOLAR ENERGY.

Biomass. Biofuel applications, involving consumption of wood products and garbage, have existed for many years. Perhaps the single largest hurdle for this application of wood and other products is the creation of effective distribution and supply systems. In general, these materials are low-density, thus causing significant transportation problems in delivering them to the chosen site for consumption.

An alternative to direct burning of biofuels is conversion to combustible liquids, gases, and solids. One process is classified as thermochemical—pyrolysis and gasification. Pyrolysis is a low-temperature process for forming gases, liquids, and solids by heating them in the absence of air, while gasification is a high-temperature process for producing gaseous fuels in the presence of an oxidant (air). A second class of processes for converting biomass to an alternative form is biochemical, which is better known as fermentation, that is, microbial transformation of organic feed materials that takes place without oxygen and produces alcohols and organic chemicals, such as methane, ethanol, acetic acid, and acetone. See FERMENTATION; PYROLYSIS.

The typical process of converting solar radiation to electricity by using vegetation as the collector is very inefficient (Table 3). It is clear that direct conversion by means other than vegetation is desirable. See BIOMASS.

Geothermal power. The natural emissions of steam (geysers), hot springs, and volcanoes represent po-

tential sources of electricity production and are indicative of the steam and hot water potential embedded in the Earth's crust. Geothermal energy manifests itself in three basic types: hydrothermal, geopressured, and petrothermal. Hydrothermal is the natural production of steam when water is vaporized by coming into contact with hot rock. In geopressured systems, hot water is generated in deep reservoirs embedded in sand and shale formations within the Earth. The hot water is trapped and pressurized and is saturated with natural gas. Petrothermal systems consist of hot rock at the Earth's surface which could have a fluid injected into it and pumped out again, to extract thermal energy. The amount of geothermal potential for electricity production is estimated at 220×10^{15} Btu (2.3×10^{20} joules), which is about three times the United States' energy consumption in 1980.

Presently, most geothermal power-producing facilities are limited to application of geysers or so-called dry steam sources. Since dry steam represents less than 1% of the total geothermal potential, other systems are under development to harness the Earth's natural heat. The largest near-term potential appears to be the hydrothermal resource. The present approach of converting high-temperature (above 410°F or 210°C) water is direct-flash technology. Hot water is extracted from the earth and its pressure is dropped, causing the water to vaporize (flash-boil). The steam is applied to conventional steam turbine technology. The overall efficiency of this process requires improvement since significant energy is lost when flashing the hydrothermal fluid. Since the hot water arrives from the geothermal well as a two-phase mixture of steam and water, a centrifugal process is under development to separate steam and water, capturing natural steam and optimizing conversion of the hot water to steam.

Since much of the geothermal potential occurs at lower temperatures, below 410°F (210°C), a process is being developed to capture this energy. The process is called binary since hot geothermal energy is used to vaporize a secondary fluid such as a hydrocarbon like isobutane or isopentane which boils at a lower temperature. The heat from the geothermal source is transferred to the secondary fluid through a heat exchanger, which in turn produces steam for a conventional steam cycle.

Although the various geothermal sources follow different paths to the production of steam, several common problems exist. Most hydrothermal fluids contain dissolved minerals. As the temperature of the fluid drops, the minerals precipitate out causing scaling in most flow channels and thus reducing the efficiency and effectiveness of the system. The corrosive nature of the brine hydrothermal fluids is a general problem. Perhaps the other general problem is the noncondensable gases found dissolved in the hydrothermal fluids. These gases reduce thermal efficiency, accelerate equipment wear, and cause environmental concerns. See GEOTHERMAL POWER.

TABLE 3. Typical energy losses in using vegetation for generation of electricity

Energy losses	% of energy in incident solar radiation
Due to leaf reflection and ineffective absorption	9
Due to invisible light, unusable for photosynthesis	55
Due to photosynthetic conversion	31
Due to plant respiration and metabolism	2
Due to conversion to electricity	2
Available as electricity	1

Fusion power. The use of nuclear energy as a power source has been confined to nuclear fission or the splitting of atomic nuclei for energy production. In the process of fusion, atomic nuclei merge by overcoming their normal electrostatic repulsion. When the nuclei fuse, energy is released as a result of the freeing of fast-moving nuclear particles. The Sun is a natural fusion reactor and is indicative of the potential from this source of energy.

Atomic mixtures used in fusion, called plasmas, consist of positively and negatively charged particles resulting from heating a gas to a high temperature. Fusion requires that the working plasma have not only extreme temperature to free the nuclei but also a sufficiently dense or compact environment to merge them. As a gas heats up, it will expand unless confined, resulting in a very low-density mixture that will not allow nuclear fusion. A successful fusion effort must therefore combine both high temperature and confinement.

The Sun accomplishes the process of fusion with its own high temperature (15×10^6 °C) and its enormous mass which holds the plasma by gravity at proper densities. On Earth, such gravitational systems are not possible and thus higher temperatures will be necessary. The development of fusion reactors has concentrated on the two problems of temperature and confinement. Heating appears a less formidable problem than confinement. Electric currents, radio-frequency waves, and neutral-beam heaters have been effective approaches to the heating problem. Containment systems have focused on two approaches: magnetic and inertial.

The leading contender in magnetic confinement is a doughnut-shaped device called a tokamak. A magnetic field is used to control the plasma by deflecting the charged particles into a cylindrical or magnetic bottle and preventing them from striking the walls of the containment vessel. Such an occurrence would drop the temperature of the plasma, preventing the fusion process. In inertial confinement, a pulsed-energy source, called a drive, is used to compact and heat the fusion fuel in a single step. This results in rapid burning of the fuel to yield energy release. See ELECTRIC POWER SYSTEMS; NUCLEAR FUSION.

Thomas W. Reddick

Bibliography. P. N. Cheremisinoff, *Air Pollution Control and Design for Industry*, 1993; D. G. Fink and H. W. Beaty (eds.), *Standard Handbook for Electrical Engineers*, 14th ed., 2000; L. L. Freris, *Principles of Wind Energy Conversion Systems*, 1990; General Electric Co., *Updating Excitation Systems*, Publ. GET-6675B, 1986; K. W. Li and A. P. Priddy, *Power Plant System Design*, 1985; P. Longrigg and E. H. Buell, *Electric Power from Renewable Energy Technologies*, Solar Energy Research Institute, May 1983; A. J. Pansini and K. D. Smalling, *Guide to Electric Power Generation*, 1994; J. G. Singer, *Combustion-Fossil Power Systems*, 4th ed., 1993; S. C. Stultz and J. B. Kitto, *Steam: Its Generation and Use*, 40th ed., 1990; R. Van Overstraeten and R. P. Mertens, *Physics, Technology and Use of Photovoltaics*, 1986; C. L. Wadhwa, *Electrical Power*

Systems, 2d ed., 1992; J. Weisman and L. E. Eckart, *Modern Power Plant Engineering*, 1985; J. H. Willenbrock and H. R. Thomas (eds.), *Planning, Engineering and Construction of Electric Power Generation Facilities*, 1980; A. J. Wood and B. J. Wollenberg, *Power Generation, Operation, and Control*, 2d ed., 1996; K. E. Zweibel, *Harnessing Solar Power: The Photovoltaics Challenge*, 1990.

Electric power measurement

The measurement of the time rate at which electrical energy is being transmitted or dissipated in an electrical system. The potential difference in volts between two points is equal to the energy per unit charge (in joules/coulomb) which is required to move electric charge between the points. Since the electric current measures the charge per unit time (in coulombs/second), the electric power p is given by the product of the current i and the voltage v (in joules/second = watts), as in Eq. (1).

$$p = vi \quad (\text{watts}) \quad (1)$$

See ELECTRIC CURRENT; ELECTRICAL UNITS AND STANDARDS; POWER.

Alternate forms of the basic definition can be obtained by using Ohm's law, which states that the voltage across a pure resistance is proportional to the current through the element. This results in Eq. (2),

$$p = i^2 R = \frac{v^2}{R} \quad (\text{watts}) \quad (2)$$

where R is the resistance of the element and i and v are the current through and voltage across the resistive element. Other commonly used units for electric power are milliwatts (1 mW = 10^{-3} W), kilowatts (1 kW = 10^3 W), megawatts (1 MW = 10^6 W), and, in electromechanical systems, horsepower (1 hp = 746 W). See ELECTRICAL RESISTANCE; OHM'S LAW.

These fundamental expressions yield the instantaneous power as a function of time. In the dc case where v and i are each constant, the instantaneous power is also constant. In all other cases where v or i or both are time-varying, the instantaneous power is also time-varying. When the voltage and current are periodic with the same fundamental frequency, the instantaneous power is also periodic with twice the fundamental frequency. In this case a much more significant quantity is the average power defined by Eq. (3), where T is the period of the periodic wave.

$$P = \frac{1}{T} \int_t^{t+T} vi dt \quad (\text{watts}) \quad (3)$$

The average power is usually the quantity of interest since in most cases the electric power is converted to some other form such as heat or mechanical power and the rapid fluctuations of the power are smoothed by the thermal or mechanical inertia of the output

system. The remainder of this article is concerned primarily with average power in systems with periodic voltages and currents.

Sinusoidal ac waves. A number of basic concepts in electric power measurement are directly related to the behavior of systems where the voltages and currents are simple sinusoids.

Average power. In a circuit containing only resistance, the voltage and current waves have the same points of zero crossing and are said to be in phase. The power curve, obtained by directly multiplying the instantaneous values of voltage and current, is a double-frequency sinusoidal wave which is entirely positive. It is clear from the symmetry of the power wave that the average power over one cycle is one-half of the peak power.

When a circuit contains both resistance and inductance, the current lags behind the voltage wave. The power curve has the same double-frequency sinusoidal shape but is shifted downward so that over a portion of the cycle the instantaneous power is negative. During this part of the cycle, power is actually being delivered back to the source from the magnetic field of the inductance. The average of the double-frequency power wave is less than one-half of the peak because of the portion which has negative values.

Quantitative relationships describing such results can be obtained by writing the voltage and current waves as Eqs. (4) and (5), where V_m and I_m are maximum values, V and I are root-mean-square (rms) or effective values, f is the frequency, and ϕ is the angle by which i lags v . The instantaneous power is then given by Eq. (6), which can be rewritten as Eq. (7), from which the average power is given by Eq. (8).

$$v = V_m \sin 2\pi ft = \sqrt{2} V \sin 2\pi ft \quad (\text{volts}) \quad (4)$$

$$i = I_m \sin (2\pi ft - \phi) \\ = \sqrt{2} I \sin (2\pi ft - \phi) \quad (\text{amperes}) \quad (5)$$

imum values, V and I are root-mean-square (rms) or effective values, f is the frequency, and ϕ is the angle by which i lags v . The instantaneous power is then given by Eq. (6), which can be rewritten as Eq. (7), from which the average power is given by Eq. (8).

$$p = vi = 2 VI \sin 2\pi ft \sin (2\pi ft - \phi) \quad (\text{watts}) \quad (6)$$

$$P = VI[\cos \phi - \cos(4\pi ft - \phi)] \quad (\text{watts}) \quad (7)$$

$$P = VI \cos \phi \quad (\text{watts}) \quad (8)$$

See ALTERNATING CURRENT; ALTERNATING-CURRENT CIRCUIT THEORY.

Power factor, apparent power, and reactive power. The average power P in Eq. (8) is also called the real or active power and is distinguished from the simple product VI by the factor $\cos \phi$. The product VI is called the apparent power (which is measured in volt-amperes), and the factor $\cos \phi$ is called the power factor so that, by definition Eq. (9a) holds or, in words, Eq. (9b).

$$\cos \phi = \frac{P}{VI} \quad (9a)$$

$$\text{Power factor} = \frac{\text{average power}}{\text{apparent power}} \quad (9b)$$

One additional quantity is introduced by rewriting Eq. (7) in the alternative form of Eq. (10), which indicates that the single double-frequency sinusoidal

$$p = VI[\cos \phi (1 - \cos 4\pi ft) \\ + \sin \phi \sin 4\pi ft] \quad (\text{watts}) \quad (10)$$

instantaneous power wave of a circuit containing both resistance and inductance can be broken into two double-frequency sinusoidal waves. One of the waves is entirely above the zero line, and the average of this wave is the average power given by Eq. (8). The other wave is symmetric above and below the zero line and represents the instantaneous power flow into and out of the magnetic field of the inductor. This portion of the instantaneous power is termed the reactive power. The peak amplitude is given by Eq. (11), and is measured in volt-amperes-

$$Q = VI \sin \phi \quad (\text{vars}) \quad (11)$$

reactive or vars. The quantity Q is nominally taken as positive for an inductive load. It is evident from Eqs. (8) and (11) that the apparent power, the average power, and the reactive power are related by Eq. (12). See VOLT-AMPERE.

$$V^2 I^2 = P^2 + Q^2 \quad (\text{volt-amperes}) \quad (12)$$

Nonsinusoidal periodic waves. For more complex periodic waves, the voltage and current can be expressed as Fourier series summations consisting of a dc component V_0 , a fundamental component V_1 , and a set of harmonics V_n as in Eq. (13), with the same

$$v = V_0 + \sum_{n=1}^{\infty} \sqrt{2} V_n \sin (2\pi fnt + \phi) \quad (\text{volts}) \quad (13)$$

general form for the current. It is easily shown by the methods used previously that only terms with the same frequency contribute to the power. All other products contribute only reactive power. The average power can be found by repeated application of Eq. (8) for each frequency present in both v and i . See FOURIER SERIES AND TRANSFORMS; HARMONIC (PERIODIC PHENOMENA); NONSINUSOIDAL WAVEFORM.

Low-frequency measurements. The measurement of power in a dc circuit can be carried out by simultaneous measurements of voltage and current by using standard types of dc voltmeters and ammeters. The product of the readings typically gives a sufficiently accurate measure of dc power. If great accuracy is required, corrections for the power used by the instruments should be made. In ac circuits the phase difference between the voltage and current precludes use of the voltmeter-ammeter method unless the load is known to be purely resistive. When this method is applicable, the instrument readings lead directly to average power since ac voltmeters and ammeters are always calibrated in rms values. See AMMETER.

Wattmeter measurements. In power-frequency circuits the most common instrument for power measurement is the moving-coil, dynamometer wattmeter. This instrument can measure dc or ac power by carrying out the required multiplication and averaging on a continuous analog basis. The instrument has four terminals, two for current and two for voltage, and reads the average power directly. It can be built with an accuracy better than 0.25% of full-scale deflection and with frequency response up to about 1 kHz. For harmonics within its frequency range, the reading is independent of the voltage and current waveforms. See WATTMETER.

Digital wattmeters. Electronic wattmeters are available which give a digital indication of average power. Their primary advantage, in addition to minimizing errors in reading the instrument, is that the frequency range can be greatly extended, up to 100 kHz or more, with good accuracy. In addition, such instruments also often provide other measurements such as rms voltage and current, apparent power, reactive power, and power factor. These instruments are of great value in power measurements involving very distorted waveforms as are common at the input of many electronic power supplies. Accurate measurement of power in these situations can be accomplished only with the wide bandwidth of these electronic instruments.

Polyphase power measurements. Measurement of the total power in a polyphase system is accomplished by combinations of single-phase wattmeters or by special polyphase wattmeters which are integrated combinations of single-phase wattmeter elements. A general theorem called Blondel's theorem asserts that the total power supplied to a load over N wires can be measured by using $N - 1$ wattmeters. The theorem states that the total power in an N -wire system can be measured by taking the sum of the readings of N wattmeters so arranged that each wire contains the current coil of one wattmeter. One voltage terminal of each wattmeter is connected to the same wire as its current coil, and the second voltage terminal is connected to a common point in the circuit. If this common point is one of the N wires, one wattmeter will read zero and can be omitted.

The most common polyphase system is the three-phase system, and the two most common connections are the three-phase, three-wire system, in which the sources and loads may be Y-connected or delta (Δ)-connected, and the three-phase, four-wire system, which utilizes Y-connected sources and loads. The fourth wire is the neutral wire, which connects the neutral point of the source to the neutral point of the load.

In one arrangement for measuring the power in a four-wire system, each three-phase wattmeter directly measures the power in one of the load elements. The sum of the three readings is the total power. However, this connection will read the correct total power even if the load is more complex, including, for example, one or more Δ -connected

loads in parallel with the load illustrated in the figure. With multiple loads, only the sum of the meter readings is significant. For three-wire systems, the most common connection is the two-wattmeter arrangement. With this arrangement, again only the sum of the two readings is of significance. Care must be exercised in making the connection (the positive and negative terminals must be connected as shown) since one of the wattmeters will have a negative reading when the power factor is low or the load unbalanced. When this occurs, the negative-reading wattmeter must have its current or voltage coil reversed to obtain the value of the negative meter reading. Thus, careful attention to the positive and negative terminal connections is a necessity in this method.

High-frequency measurements. At frequencies significantly above power frequencies, dynamometer wattmeters become inaccurate and cannot be used. The newer digital wattmeters have usable ranges well above audio frequencies and make accurate audio-frequency power measurements quite feasible. Generally, however, power measurements at higher frequencies are based on indirect methods which do not attempt to carry out the multiplication and integration required for direct electric power measurement.

Voltmeter-based methods. For frequencies up to a few hundred megahertz, the voltage across a standard resistance load can be measured and the power calculated from V^2/R . Instruments which combine the resistive load and voltmeter are called absorption power meters. For the lower frequency end of the range, such instruments often contain an input transformer with taps to provide a range of input resistance values. Adjusting the transformer to achieve a maximum meter reading allows matching the effective load resistance to the source output impedance to maximize the power output from the source.

Diode detector-based methods. A diode may be used as a detector for radio-frequency (rf) power measurement. This type of instrument is simple and easy to use but is less accurate than the thermally based systems described below. The diode is operated on the square-law portion of its characteristic curve, and as a result the voltage on the capacitor C is proportional to the average rf power input. There are a number of limitations to the performance of these instruments. The dc signal at the output is very small, requiring amplification to enable accurate measurement, and care must be taken to avoid errors caused by electronic and thermoelectric effects. Temperature changes affect the diode resistance, and hence the sensitivity and reflection coefficient can vary; the matching resistor, usually 50 ohms, is used to terminate the input source to reduce these temperature effects. See DIODE.

Thermocouple-based methods. A thermocouple consists of two dissimilar metals joined at one end. When the joined end is heated and the other end is at a lower temperature, an electric current is produced (the

thermoelectric or Seebeck effect). The current is proportional to the temperature difference between the two ends. For electric power measurements, the hot end is heated by a resistor supplied from the rf power source to be measured. Modern devices use thin-film techniques to make the thermocouple and resistor to ensure good thermal coupling. The output is a low-level dc signal as in diode detector systems. Since the thermocouple method is affected by drift in the thermocouple and associated electronic circuits, most instruments contain a source arranged to be connected in place of the rf input to calibrate the instrument. Thermocouple instruments generally have a wider range of measurement and are more rugged than most other high-frequency power measurement instruments. *See* THERMOCOUPLE.

Bolometric methods. A bolometer is basically an electric bridge circuit in which one of the bridge arms contains a temperature-sensitive resistor. In principle, the temperature is detected by the bridge circuit. Bolometric bridges use either barreters or thermistors as the temperature-sensitive resistor. A barreter is simply a short, thin wire, usually of platinum, whose resistance changes significantly for small amounts of power dissipation. Since the element must be operated at relatively high temperatures, near its burn-out level, it is easily destroyed by accidental overloads. The thermistor, which has largely replaced the barreter since it is more rugged, is a semiconductor device with a negative temperature coefficient. For rf power measurements, the thermistor is fabricated as a small bead with very short lead wires so that essentially all the resistance is in the bead. *See* BARRETER; BOLOMETER; THERMISTOR.

In power measurements, absolute measurement of the resistance variation is not generally acceptable since the change in resistance can affect the reflection coefficient of the probe. Thus, to expand the useful power range of the instrument, comparison methods are employed. For example, the bridge can be arranged so that, without any input power, a dc or low-frequency bias is applied to the sensing element and the bridge is balanced. The input power is then applied causing an imbalance. The bias is then changed to return the bridge to balance, and the change in the bias power is used to indicate the input power. A major advantage of thermistor-based bolometric techniques is that they operate at high signal levels and shielding requirements are minimized.

Calorimeter methods. The most accurate high-frequency power measurement methods involve calorimetric techniques based on direct determination of the heat produced by the input power. The power to be measured is applied to the calorimeter, and the final equilibrium temperature rise is recorded. Then the input signal is removed, and dc power applied until the same equilibrium temperature is attained. The dc power is then the same as the signal power. Most power measurement calorimeters use a flow-through coolant rather than a static

fluid coolant, and can be built to measure large amounts of power. Calorimeter methods are usually used in standards laboratories rather than in industrial applications. *See* CALORIMETRY.

Sensors and connections. Ideally the power sensor should accept all of the high-frequency power to be measured. There are, however a number of factors that reduce the calibration factor of the sensor below the ideal of 100%. For example, mismatch between the impedance of the rf or microwave source and that of the sensor will result in some reflection of the incident power. Some power may also be lost in the conducting walls of the sensor, and some may be radiated into space. Typically, the calibration factor of a commercial power sensor at various frequencies within its range is supplied by the manufacturer.

The accuracy of high-frequency power measurements is very dependent on correct setup of the instrumentation to minimize losses. Most often, directional couplers and attenuators are used to supply the power meter with the desired signal at the appropriate power level for the instrument. Care must be exercised to assure that all portions of the test circuit are properly matched to avoid reflections and that losses are minimized. *See* DIRECTIONAL COUPLER; ELECTRICAL MEASUREMENTS; IMPEDANCE MATCHING; TRANSMISSION LINES.

Donald W. Novotny

Bibliography. M. Braccio, *Basic Electrical and Electronic Tests and Measurements*, 1978; *Code for Electrical Metering*, ANSI C12, 1982; D. Fink and H. W. Beaty (eds.), *Standard Handbook for Electrical Engineers*, 14th ed., 1999; Instrument Society of America, *ISA Standards and Practices for Instrumentation*, 9th ed., 1988; P. Kantrowitz, G. Kousouri, and L. Zucker, *Electronic Measurements*, 1979; F. F. Mazda, *Electronic Instruments and Measurement Techniques*, 1987.

Electric power substation

An assembly of equipment in an electric power system through which electrical energy is passed for transmission, distribution, interconnection, transformation, conversion, or switching. *See* ELECTRIC POWER SYSTEMS.

Specifically, substations are used for some or all of the following purposes: connection of generators, transmission or distribution lines, and loads to each other; transformation of power from one voltage level to another; interconnection of alternate sources of power; switching for alternate connections and isolation of failed or overloaded lines and equipment; controlling system voltage and power flow; reactive power compensation; suppression of overvoltage; and detection of faults, monitoring, recording of information, power measurements, and remote communications. Minor distribution or transmission equipment installation is not referred to as a substation.

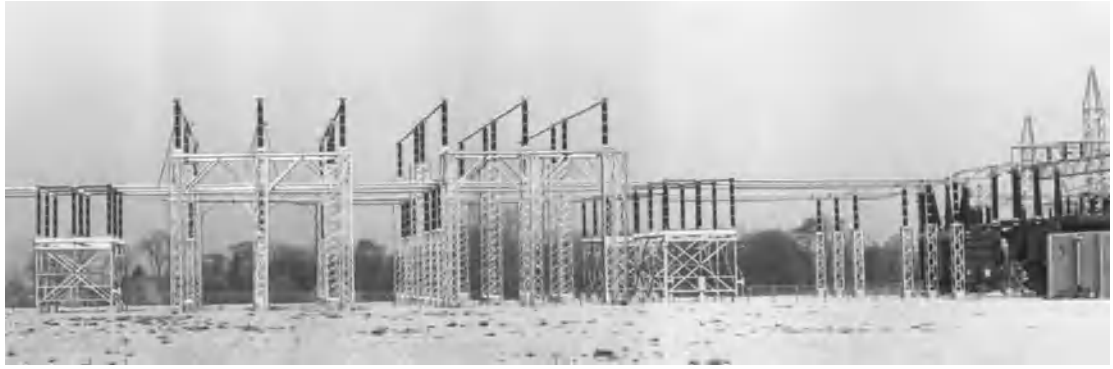


Fig. 1. A typical transmission substation (345 kV).

Classification. Substations are referred to by the main duty they perform. Broadly speaking, they are classified as transmission substations (Fig. 1), which are associated with high voltage levels; and distribution substations (Fig. 2), associated with low voltage levels. See ELECTRIC DISTRIBUTION SYSTEMS.

Substations are also referred to in a variety of other ways, which may be described as follows.

Transformer substations. These are substations whose equipment includes transformers.

Switching substations. These are substations whose equipment is mainly for various connections and interconnections, and does not include transformers.

Customer substations. These are usually distribution substations located on the premises of a larger customer, such as a shopping center, large office or commercial building, or industrial plant.

Converter stations. The main function of converter stations is the conversion of power from ac to dc and vice versa. Converter stations are complex substations required for high-voltage direct-current (HVDC) transmission or interconnection of two ac systems which, for a variety of reasons, cannot be connected by an ac connection. The main equip-

ment includes converter valves usually located inside a large hall, transformers, filters, reactors, and capacitors. Figure 3 shows a ± 250 -kV 500-MW HVDC converter station together with a 245-kV ac substation. The building in the center of the converter station contains all the thyristor valves. On each side of the building against the wall are converter transformers connected to the valves inside. Further away from the building are filters and capacitors. Going away from the building is the dc line, with a typical dc line tower in the uppermost part of the figure. See DIRECT-CURRENT TRANSMISSION.

Air-insulated substations. Most substations are installed as air-insulated substations, implying that the busbars and equipment terminations are generally open to the air, and utilize insulation properties of ambient air for insulation to ground. Modern substations in urban areas are esthetically designed with low profiles and often within walls, or even indoors.

Metal-clad substations. These are also air-insulated, but for low voltage levels; they are housed in metal cabinets and may be indoors or outdoors (Fig. 2).



Fig. 2. View of a residential outdoor distribution substation with double-ended underground supply to transformers which are positioned at either end of the metal-clad switchgear.



Fig. 3. View of ± 250 -kV 500-MW HVDC converter station (upper part) and 245-kV ac substation (lower part). The large building of the converter station contains all the thyristor valves. (General Electric Co.)

Gas-insulated substations. Acquiring a substation site in an urban area is very difficult. Land in many cases is either unavailable or very expensive. Therefore, there has been a definite trend toward increasing use of gas-insulated substations, which occupy only 5–20% of the space occupied by the air-insulated

substations. In gas-insulated substations, all live equipment and bus-bars are housed in grounded metal enclosures, which are sealed and filled with sulfur hexafluoride (SF_6) gas, which has excellent insulation properties. These substations have the appearance of a large-scale plumbing (Fig. 4).

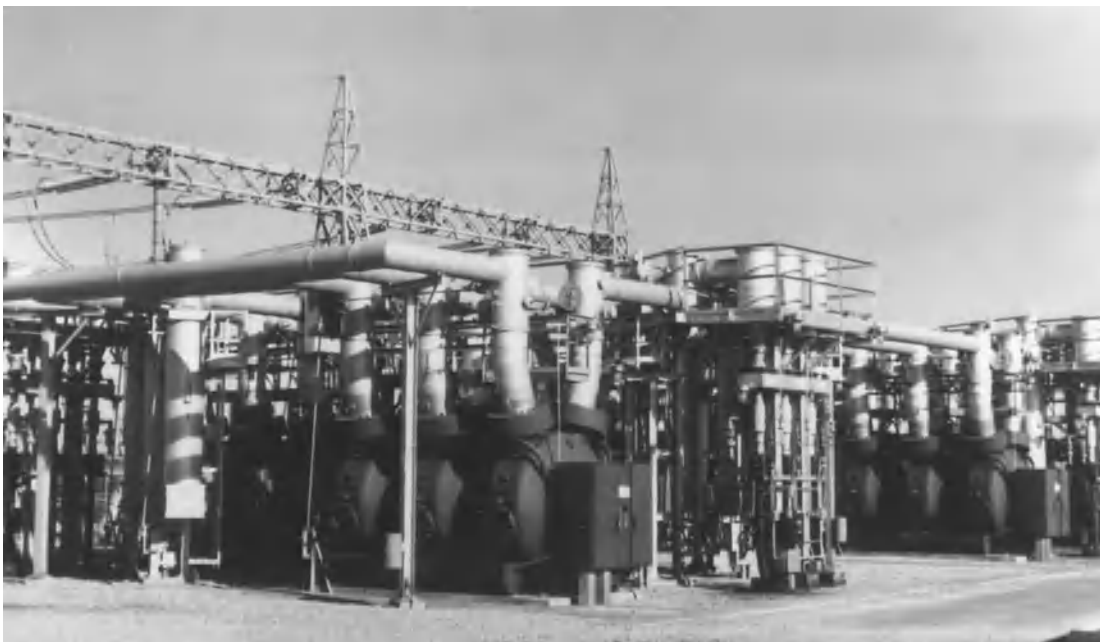


Fig. 4. A 345-kV ac gas-insulated substation. (Gould-BBC).

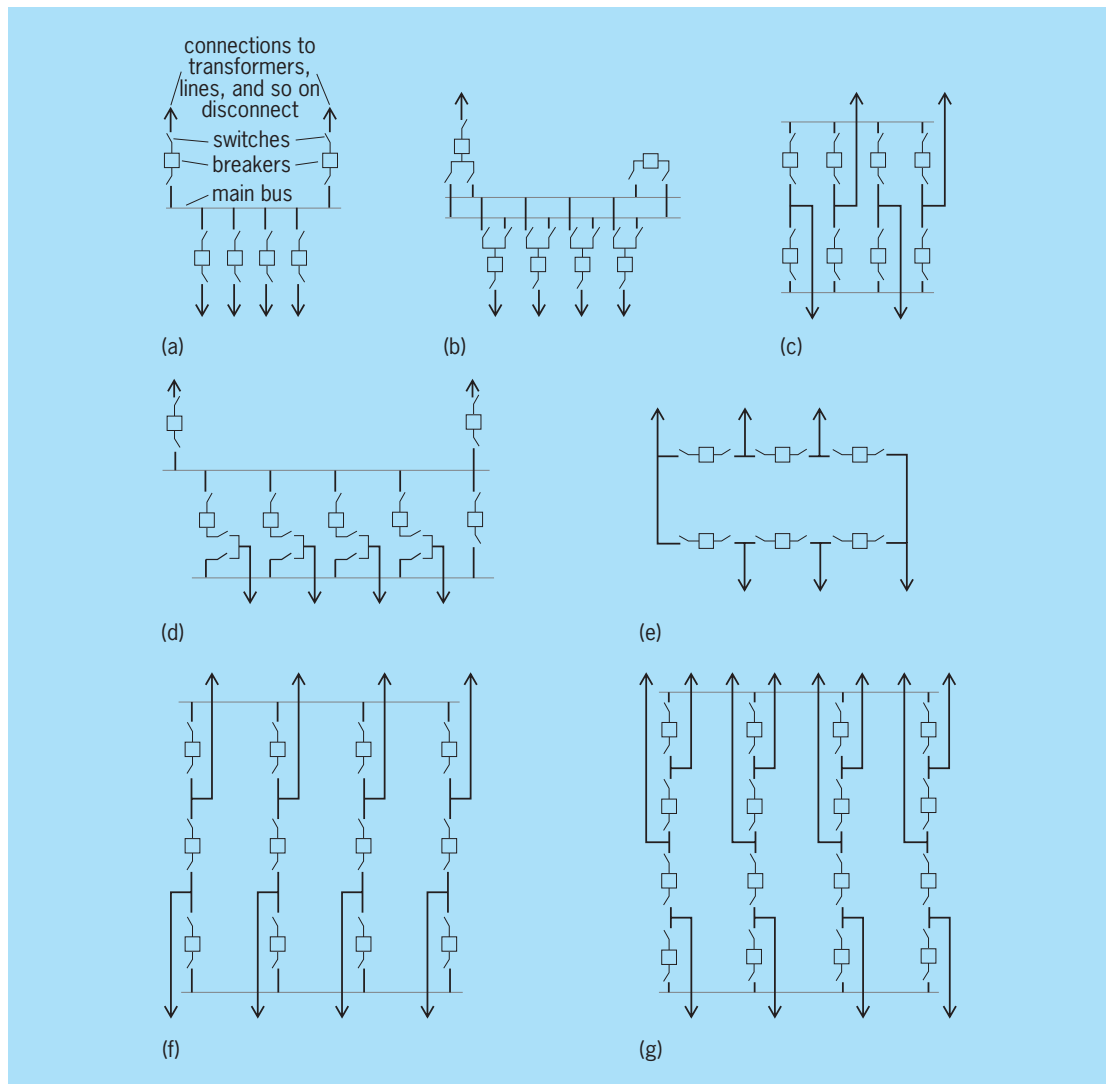


Fig. 5. One-line diagrams of substation switching arrangements. (a) Single bus. (b) Double bus, single breaker. (c) Double bus, double breaker. (d) Main and transfer bus. (e) Ring bus. (f) Breaker-and-a-half. (g) Breaker-and-a-third.

The disk-shaped enclosures in Fig. 4 are circuit breakers.

Mobile substations. For emergency replacement or maintenance of substation transformers, mobile substations are used by some utilities. They may range in size from 100 to 25,000 kVA, and from 2400 to 220,000 V. Most units are designed for travel on public roads, but larger substations are built for travel on railroad tracks. These substations are generally complete with a transformer, circuit breaker, disconnect switches, lightning arresters, and protective relaying.

Substation switching arrangement. An appropriate switching arrangement for “connections” of generators, transformers, lines, and other major equipment is basic to any substation design. There are seven switching arrangements commonly used (Fig. 5). Each breaker is usually accompanied by two disconnect switches, one on each side, for maintenance purposes. Selecting the switching arrangement involves considerations of cost, reliability, maintenance, and flexibility for expansion.

Single bus. This involves one common bus for all connections and one breaker per connection. This is the least costly arrangement, but also the least desirable for considerations of reliability and maintenance.

Double bus, single breaker. This involves two common buses, and each connection of the line equipment and so forth can be made to either bus through one or the other breaker. However, there is only one breaker per connection, and when a breaker is out for maintenance, the connected equipment or line is also removed from operation.

Double bus, double breaker. This arrangement has two common buses and two breakers per connection. It offers the most reliability, ease of maintenance, and flexibility, but is also the most expensive.

Main and transfer bus. This is the single-bus arrangement; however, an additional transfer bus is provided so that a breaker can be taken out of service by transferring its connection to the transfer bus which is connected to the main bus through a breaker between the two buses.

Ring bus. This may consist of four, six, or more breakers connected in a closed loop, with the same number of connection points. A breaker can be taken out of service without also taking out a connection.

Breaker-and-a-half. This arrangement involves two buses, between which three breaker bays are installed. Each three-breaker bay provides two circuit connection points—thus the name breaker-and-a-half.

Breaker-and-a-third. In this arrangement, there are four breakers and three connections per bay.

Substation equipment. A substation includes a variety of equipment. The principal items are listed and briefly described below.

Transformers. These involve magnetic core and windings to transfer power from one side to the other at different voltages. Substation transformers range from small sizes of 1 MVA to large sizes of 2000 MVA. Most of the transformers and all those above a few MVA size are insulated and cooled by oil, and adequate precautions have to be taken for fire hazard. These precautions include adequate distances from other equipment, fire walls, fire extinguishing means, and pits and drains for containing leaked oil. See TRANSFORMER.

Circuit breakers. These are required for circuit interruption with the capability of interrupting the highest fault currents, usually 20–50 times the normal current, and withstanding high voltage surges that appear after interruption. Switches with only normal load-interrupting capability are referred to as load break switches. See CIRCUIT BREAKER.

Disconnect switches. These have isolation and connection capability without current interruption capability. See ELECTRIC SWITCH.

Bus-bars. These are connecting bars or conductors between equipment. Flexible conductor buses are stretched from insulator to insulator, whereas more common solid buses (hollow aluminum alloy tubes) are installed on insulators in air or in gas-enclosed cylindrical pipes. See BUS-BAR.

Shunt reactors. These are often required for compensation of the line capacitance where long lines are involved.

Shunt capacitors. These are required for compensation of inductive components of the load current.

Current and potential transformers. These are for measuring currents and voltages and providing proportionately low-level currents and voltages at ground potential for control and protection.

Control and protection. This includes (a) a variety of protective relays which can rapidly detect faults anywhere in the substation equipment and lines, determine what part of the system is faulty, and give appropriate commands for opening of circuit breakers; (b) control equipment for voltage and current control and proper selection of the system configuration; (c) fault-recording equipment; (d) metering equipment; (e) communication equipment; and (f) auxiliary power supplies. See ELECTRIC PROTECTIVE DEVICES; RELAY; VOLTAGE REGULATOR.

Many of the control and protection devices are solid-state electronic types, and utilize digital tech-

niques with microprocessors. Most of the substations are fully automated locally with a provision for manual override. The minimum manual interface required, along with essential information on status, is transferred via communications channels to the dispatcher in the central office.

Other items which may be installed in a substation include: phase shifters, current-limiting reactors, dynamic brakes, wave traps, series capacitors, controlled reactive compensation, fuses, ac to dc or dc to ac converters, filters, and cooling facilities.

Substation grounding and shielding. Good substation grounding is very important for effective relaying and insulation of equipment; however, the safety of the personnel is the governing criterion in the design of substation grounding. It usually consists of a bare wire grid, laid in the ground; and all equipment grounding points, tanks, support structures, fences, shielding wires and poles, and so forth, are securely connected to it. The grounding resistance is reduced to be low enough that a fault from high voltage to ground does not create such high potential gradients on the ground, and from the structures to ground, to present a safety hazard. Good overhead shielding is also essential for outdoor substations, so as to virtually eliminate the possibility of lightning directly striking the equipment. Shielding is provided by overhead ground wires stretched across the substation or tall grounded poles. See GROUNDING; LIGHTNING AND SURGE PROTECTION. Narain G. Hingorani

Bibliography. J. Arrilaga, *High Voltage Direct Current Transmission*, 2d ed., 1998; W. M. Flanagan, *Handbook of Transformer Design and Applications*, 2d ed., 1993; C. H. Flurscheim (ed.), *Power Circuit Breaker Theory and Design*, rev. ed., 1982; A. C. Franklin and D. P. Franklin, *The J&P Transformer Book*, 11th ed., 1983.

Electric power systems

Complex assemblages of equipment and circuits for generating, transmitting, and distributing electrical energy. The elements of a power system (Fig. 1) form complex networks which require energy control centers to monitor and regulate their operation. Various sources of primary energy, such as coal, nuclear fission, hydro power, geothermal, and wind, can be used to drive the electric generator in Fig. 1.

Since electrical energy plays a central role in industrialized societies, the reliability of electric power systems is a critical factor in their planning, design, and operation, and considerable effort is directed toward quantifying performance. In addition to the provision of electrical energy, attention is given to the myriad of items that consume this energy, that is, the power system load.

Principal Elements

Electrical systems require generating stations to produce electric power, transmission systems to carry

it to areas of consumption, substations to transform it for use in industrial and residential areas, and distribution systems to carry it to customers (Fig. 1). Coordinated interconnections between power systems play a critical role in ensuring an acceptable degree of reliability.

Generation. Electricity in the large quantities required to supply electric power systems is produced in generating stations, commonly called power plants. Such generating stations, however, should be considered as conversion facilities in which the chemical or nuclear energy of fuel (coal, oil, gas, or uranium); the kinetic energy of falling water, wind, or sea waves; or light from the Sun is converted to electricity. See ELECTRIC POWER GENERATION; POWER PLANT.

Steam stations. Most of the electric power used in the United States is obtained from generators driven by steam turbines. Units of 650, 800, and 950 MW are commonplace for new fossil-fuel-fired stations, with 1300 MW the largest in service. The most commonly installed units in nuclear stations are 800–1100 MW. See STEAM TURBINE.

Coal is an important fuel for more than half of steam turbine generation. Other fuels include natural gas and heavy fuel oil; the remainder is generated from the nuclear energy of slightly enriched uranium. As nuclear units have come into commercial operation, the contribution of uranium to the electrical energy supply has risen to about 15% of the total fuel generated. However, uranium's share should decrease, because no additional nuclear plants are slated for service and instead some existing nuclear power stations are slated for closure.

Combustion of coal produces sulfur dioxide and nitrogen oxides in the stack gases. To reduce these emissions below those permitted by the Environmental Protection Agency, modern coal-fired stations use either low-sulfur coal (that is, less than 0.3% sulfur by weight) or scrubbers that react the sulfur dioxide with a reagent to permit its removal. Nitrogen oxides are controlled by controlling combustion temperatures. Research is being conducted on clean coal technologies to minimize the production of carbon dioxide. See AIR POLLUTION.

Nuclear steam stations in the United States are mostly of the water-cooled and moderated types in which the heat of a controlled nuclear reaction is used to convert ordinary or "light" water into steam to drive a conventional turbine generator. See NUCLEAR POWER; NUCLEAR REACTOR.

Hydroelectric plants. Waterpower supplies about 10% of the electric power consumed in the United States, but this share will decline because very few sites remain undeveloped where sufficient water drops far enough in a reasonable distance to drive large hydraulic turbines. Much of the very small share of the planned additional hydroelectric capability will be used at existing plants to increase their effectiveness in supplying peak power demands, and as a quickly available source of emergency power. Special incen-

tives from the federal government have made development of small hydroelectric plants of 15 MW or less financially attractive to private developers. See HYDRAULIC TURBINE; HYDROELECTRIC GENERATOR.

Some hydroelectric plants actually draw power from other generating facilities during light system-load periods, to drive their turbines in a reverse mode to pump water from a river or lake into an artificial reservoir at a higher elevation. From there, the water can be released through the hydraulic turbines when the power system needs additional generation. These pumped-storage installations consume about 35% more energy than they return to the power system and, accordingly, cannot be considered as primary energy sources. Their use is justified, however, by their ability to convert surplus power that is available during low-demand periods into prime power to serve system needs during peak-demand intervals, a need that otherwise would require building more generating stations for operation during the relatively few hours of high system demand. See ENERGY STORAGE; WATERPOWER.

Combustion turbine plants. Gas-turbine-driven generators, now commonly called combustion turbines because of the use of light oil as fuel, have gained wide acceptance as an economical source of additional power for heavy-load periods. In addition, they offer the fastest erection time and the lowest investment cost per kilowatt of installed capability. Offsetting these advantages, however, is their relatively less efficient consumption of more costly fuel. Combustion turbine units, even in the largest rating (100 MW), offer extremely flexible operation and can be started and run up to full load in as little as 10 minutes. Thus they are useful as emergency power sources, as well as for operating during the few hours of daily load peaks.

Combustion turbines have an additional role besides applying peaking or emergency power. Some installations use their exhaust gases to heat boilers that generate steam to drive steam turbine generators. Such combined-cycle units offer fuel economy comparable to that of modern steam plants and at considerably less cost per kilowatt. In addition, because only part of the plant uses steam, the requirement for cooling water is considerably reduced. However, wide acceptance is inhibited by the government's restrictions on light fuel oil for them. This barrier will be resolved by the development of systems for fueling combustion-turbine installations with gas derived from coal, natural gas, or by direct firing with pulverized coal. See COAL GASIFICATION; COGENERATION; GAS TURBINE.

Internal combustion plants. Internal combustion engines of the diesel type drive generators in many small power plants. In addition, they offer the ability to start quickly for operation during peak loads or emergencies. However, their small size, commonly about 2 MW per unit (although a few approach 10 MW), has limited their use. Such installations account for less than 1% of the total power-system generating capability in the United States, and make an

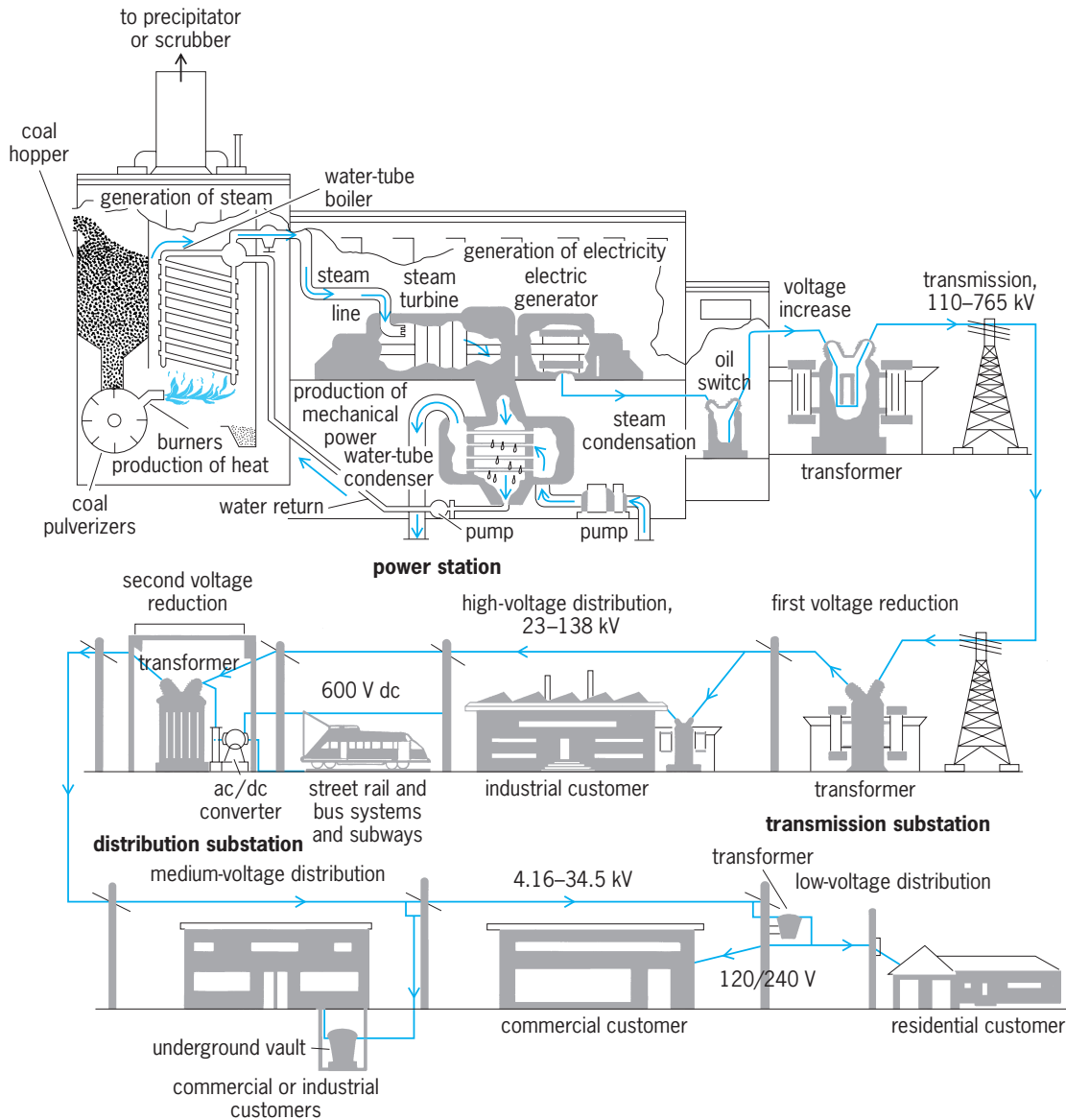


Fig. 1. Major steps in the thermal generation, transmission, and distribution of electricity.

even smaller contribution to total electric energy consumed. At present they are used primarily as standby units to provide electricity in case of emergencies and loss of electric grid. See DIESEL ENGINE; INTERNAL COMBUSTION ENGINE.

Renewable resources. Utilities actively seek to develop generating resources that do not consume fuel, including geothermal steam, wind, Sun, ocean waves, tides, and biomass-powered capacity. Of these, wind power is now the fastest growing segment of the power generation industry. Wind energy technology worldwide has grown 20% annually since 2000, and it is projected that by 2020 wind energy will meet 6% of the electric needs of the United States, matching the share of hydro power. See BIOMASS; ENERGY SOURCES; GEOTHERMAL POWER; SOLAR ENERGY; WIND POWER.

Three-phase output. Because of their simplicity and efficient use of conductors, three-phase 60-Hz

alternating-current systems are used almost exclusively in the United States. Consequently, power-system generators are wound for three-phase output at a voltage usually limited by design features to a range from about 11 kV for small units to 30 kV for large ones. The output of modern generating stations is usually stepped up by transformers to the voltage level of transmission circuits used to deliver power to distant load areas. See TRANSFORMER.

Transmission. The transmission system carries electric power efficiently and in large amounts from generating stations to consumption areas. Such transmission is also used to interconnect adjacent power systems for mutual assistance in case of emergency and to gain for the interconnected power systems the economies possible in regional operation. Interconnections have so expanded that most of the generation east of the Rocky Mountains regularly operates in parallel and in synchronism. More than 90% of all

Line-to-line voltage, kV	Capability, MVA
115 ac	60
138 ac	90
230 ac	250
345 ac	600
500 ac	1200
765 ac	2500
800 dc*	1500

* Bipolar line with grounded neutral.

generation in the United States and Canada, exclusive of Alaska and Hawaii, is or can be linked.

Transmission circuits are designed to operate up to 765 kV, depending on the amount of power to be carried and the distance to be traveled. The permissible power loading of a circuit depends on many factors, such as the thermal limit of the conductors and their clearances to ground, the voltage drop between the sending and receiving end and the degree to which system service reliability depends on it, and how much the circuit is needed to hold various generating stations in synchronism. A widely accepted approximation to the voltage appropriate for a transmission circuit is that the permissible load-carrying ability varies as the square of the voltage (see table).

Alternating-current transmission voltages have increased many times since transmission began as a distinct function around 1890 (Fig. 2). Lines at 500 and 765 kV were introduced in the 1960s. The very high capacity of lines at 765 kV limits their use, but 2454 mi (3950 km) are in service in the United States.

Transmission engineers have anticipated even higher voltages, of 1100 to 1500 kV, but they are fully aware that this objective may prove too costly in space requirements and funds to gain wide acceptance. Experience gained at 500 and 765 kV verifies that the prime requirement no longer is insulating the lines to withstand lightning discharges, but insulating them to tolerate voltage surges caused by the operation of circuit breakers. Audible noise levels, especially in rain or humid conditions, are high, requiring wide buffer zones. Environmental challenges have been brought on the basis of possible negative biological effects of the electrostatic field produced under EHV lines, although research has not shown any such effects. See ELECTRIC PROTECTIVE DEVICES; LIGHTNING AND SURGE PROTECTION.

Experience has indicated that, within about 10 years after the introduction of a new voltage level for overhead lines, it becomes necessary to begin connecting underground cable. The first 500-kV cable in the United States was placed in service in 1976.

Another approach to high-voltage long-distance transmission is high-voltage direct current (HVDC), which offers the advantages of less costly lines, lower transmission losses, and insensitivity to many system problems that restrict alternating-current systems. Its greatest disadvantage is the need for costly equipment for converting the sending-end power to direct current, and for converting the receiving-end direct-current power to alternating current for distribution to consumers. However, the development of solid-state converter valves made overhead HVDC economical for distances more than about 400 mi (650 km), but only 25–30 mi (40–50 km) in underground construction. Where systems that are out of synchronism must be interconnected, much shorter distances may be economic. An extreme example of this is at Eel River, New Brunswick, Canada, where back-to-back converters connect Canadian and United States systems. Starting in the late 1950s with a 65-mi (105-km) 100-kV system in Sweden, HVDC has been applied successfully in a series of special cases around the world, each one for a higher voltage and greater power capability. The first such installation in the United States was put into service in 1970. It operates at 800 kV line to line, and is designed to carry a power interchange of 1440 MW over an 841-mi (1354-km) overhead tie line between the winter-peaking Northwest Pacific coastal region and the summer-peaking southern California area. These HVDC lines perform functions other than power transfer, however. The Pacific Intertie is used to stabilize the parallel alternating-current transmission lines, permitting an increase in their capability; and back-to-back converters with no tie line between them are used to tie together two systems in Nebraska that otherwise could not be synchronized. The first urban installation of this technology was energized in 1979 in New York. See DIRECT-CURRENT TRANSMISSION.

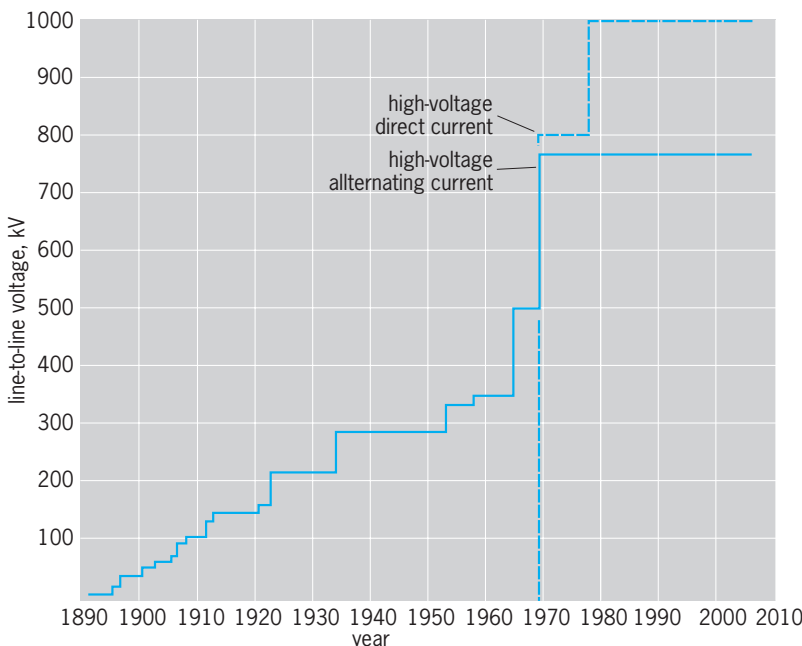


Fig. 2. Growth of alternating-current transmission voltages from 1890 onwards.

In addition to these high-capability circuits, every large utility has many miles of lower-voltage transmission, usually operating at 110 to 345 kV, to carry bulk power to numerous cities, towns, and large industrial plants. These circuits often include extensive lengths of underground cable where they pass through densely populated areas. See ELECTRIC POWER TRANSMISSION; TRANSMISSION LINES.

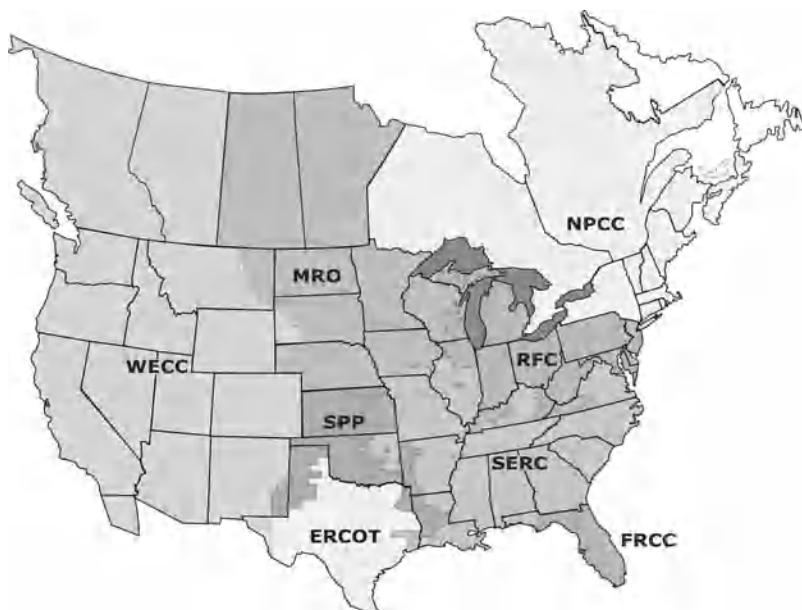
Interconnections. As systems grow and the number and size of generating units increase, and as transmission networks expand, higher levels of bulk-power-system reliability are attained through properly coordinated interconnections among separate systems. This practice began around 1930.

Most of the electrical utilities in the contiguous United States and a large part of Canada operate as members of power pools, and these pools in turn are interconnected into one gigantic power grid, known as the North American Power Systems Interconnection. The operation of this interconnection, in turn, is coordinated by the operating committee of the North American Electric Reliability Council (NERC). Each individual utility in such pools operates independently, but has contractual arrangements with other members in respect to generation additions and scheduling of operation. Their participation in a power pool affords a higher level of service reliability and important economic advantages.

The Northeast blackout of November 9, 1965, stemmed from the unexpected trip-out of a key transmission circuit carrying emergency power into Canada and cascaded throughout the northeastern states to cut off electrical service to some 30,000,000 people. It spurred the utilities into a chain reaction affecting the planning, construction, operation, and control procedures for their interconnected systems. They soon organized regional coordination councils to cover the entire contiguous United States and four Canadian provinces. Their objective was to further improve reliability of the planning and operation of their generation and transmission facilities. Currently, there are eight regional reliability councils (Fig. 3).

Then, in 1968, the North American Electric Reliability Council (NERC) was established to serve as an effective body for the collection, unification, and dissemination of various reliability criteria for use by individual utilities in meeting their planning and operating responsibilities (Fig. 3).

Increased interconnection capability among power systems reduces the required generation reserve of each of the individual systems. In most utilities the loss-of-load probability (LOLP) is used to measure the reliability of electrical service, and it is based on the application of probability theory to unit-outage statistics and load forecasts. A common LOLP criterion is 1 day in 10 years when load may exceed generating capability. The LOLP decreases (that is, reliability increases) with increased interconnection between two areas until a level is reached which depends upon the amount



Regional Reliability Councils

ERCOT, Electric Reliability Council of Texas
 FRCC, Florida Reliability Coordinating Council
 MRO, Midwest Reliability Organization
 NPCC, Northeast Power Coordinating Council
 RFC, Reliability First Corporation
 SERC, SERC Reliability Corporation
 SPP, Southwest Power Pool Inc.
 WECC, Western Electricity Coordinating Council

Fig. 3. Areas served by the eight regional reliability councils, coordinated by the North American Electric Reliability Council, that guide the planning, coordination, and operation of generation and transmission facilities. (North American Electric Reliability Council, <http://www.nerc.com>)

of reserve, unit sizes, and annual load shape in each area.

Traditionally, systems were planned to withstand all reasonably probable contingencies, and operators seldom had to worry about the possible effect of unscheduled outages. Reserve margin, which is the excess of capacity over load at maximum annual peak, while difficult for every system, is generally considered acceptable at 25%. Operators' normal security functions were to maintain adequate generation online and to ensure that such system variables as line flows and station voltages remained within the limits specified by planners. However, stronger interconnections, larger generating units, and rapid system growth spread the transient effects of sudden disturbances and increased the responsibilities of operators for system security.

System security is concerned with service continuity at standard frequency and voltage levels. The system is said to be insecure if a contingency would result in overloading some system components, in abnormal voltage levels at some stations, in change of system frequency, or in system instability, even if there is adequate capability as indicated by some reliability index. The concepts of power system

reliability and security are discussed in more detail below.

Substations. Power delivered by transmission circuits must be stepped down in facilities called substations to voltages more suitable for use in industrial and residential areas. On transmission systems, these facilities are often called bulk-power substations; at or near factories or mines, they are termed industrial substations; and where they supply residential and commercial areas, distribution substations.

Basic equipment in a substation includes circuit breakers, switches, transformers, lightning arresters and other protective devices, instrumentation, control devices, and other apparatus related to specific functions in the power system. *See* ELECTRIC POWER SUBSTATION.

Distribution. That part of the electric power system that takes power from a bulk-power substation to customers' switches, commonly about 35% of the total plant investment, is called distribution. This category includes distribution substations, subtransmission circuits that feed them, primary circuits that extend from distribution substations to every street and alley, distribution transformers, secondary lines, house service drops or loops, metering equipment, street and highway lighting, and a wide variety of associated devices.

Primary distribution circuits usually operate at 4160 to 34,500 V line to line, and supply large commercial institutional and some industrial customers directly. Since about 1950, by far the majority of lines constructed have been in the 15-kV class. Primary lines may be overhead open wire on poles, spacer or aerial cable, or underground cable. Legislation in more than a dozen states requires that all new services to developments of five or more residences be put underground. The bulk of existing lines are overhead, however, and will remain so for the foreseeable future.

At conveniently located distribution transformers in residential and commercial areas, the line voltage is stepped down to supply low-voltage secondary lines, from which service drops extend to supply the customers' premises. Most such service is at 120/240 V, but other common voltages are 120/208 and 125/216 V, and for larger commercial and industrial buildings, 240/480, 265/460, or 277/480 V. These are classified as utilization voltages. *See* ELECTRIC DISTRIBUTION SYSTEMS.

Electrical utility industry. In the United States, which has the third highest per capita use of electricity in the world and more electric power capability than any other nation, electrical systems as measured by some criteria are the country's largest industry.

The utility industry in the United States is pluralistic in the nature of its ownership. Investor-owned utilities comprise about 78% of total installed capacity, and about the same percentage of customers served. Publicly owned utilities (that is, municipal, state, and power-district utilities) own and operate about 10% of the installed capacity and serve almost 14% of the customers; cooperatives

own 3% of the capacity and serve about 10% of the customers; and federal agencies, while serving no measurable percentage of customers, own almost 10% of installed capacity. Electrical utilities are among the most capital-intensive of industries, primarily because of the huge investment in generating units.

Introduction of deregulation over the past few years has resulted in a significant reorganization of the electric utility industry. Instead of vertically integrated utilities being engaged in all three aspects—generation, transmission, and distribution—the industry is being organized horizontally with independent power producers (IPPs), independent system operators (ISOs), and independent distribution companies (DISCOs). The ISOs manage transmission of power from the IPPs to the DISCOs to meet the contractual obligations of the IPPs and DISCOs. This structure is supplemented by electricity markets with bidding mechanisms for the sale and purchase of electricity, and long-term and short-term contracts between the IPPs and the DISCOs. In this reorganization, the ISO link is still regulated and has the responsibility of maintaining the reliability of electricity supply.

William C. Hayes; Om P. Malik

Power System Operation

The operation and control of the generation-transmission-distribution grid is quite complex because this large system has to operate in synchronism and because many different organizations are responsible for different portions of the grid. In North America and Europe, many public and private electric power companies are interconnected, often across national boundaries. Thus, many organizations have to coordinate to operate the grid, and this coordination can take many forms, from a loose agreement of operational principles to a strong pooling arrangement of operating together.

Power-system operations can be divided into three stages: operations planning, real-time control, and after-the-fact accounting. The main goal is to minimize operations cost while maintaining the reliability (security) of power delivery to customers. When all the utilities were vertically integrated—that is, the operator was also the owner of the generation and transmission—centrally optimizing the generation production was a main duty of the operator. Many regions of the world have deregulated generation ownership, and in those areas the market transactions determine generation production while the network operator is mainly concerned with the reliability of the grid. Operations planning requires the forecasting of load demand in the next few hours, days, weeks, or months that allows the market to determine the schedules and the operator to approve these schedules if they meet the reliability requirements. The market determines the buying and selling of power for the long and short terms, and after the markets clear (which involves matching the buy and sell bid prices) the

generating companies can schedule hydro resources, fossil fuels, and equipment maintenance over many weeks, and the commitment (startup and shutdown) of generating units over many hours. The operator has to run contingency (what-if) studies using these schedules to ensure that all reliability criteria are met or else the buy-sell transactions have to be adjusted to ensure that the reliability criteria are met.

Real-time control of the system is required to respond to the actual demand of electricity, which always deviates somewhat from the forecast, and any unforeseen contingencies (equipment outages). Maintaining reliability of the system so that a possible contingency cannot disrupt power supply is an integral part of the real-time control function.

After-the-fact accounting tracks actual delivery of energy between organizations so that billing can be generated.

For loosely coordinated operation of the grid, each balancing authority or utility takes responsibility for the operation of its own portion while exchanging all relevant information with its partners. For pool-type operations, a group of utilities sets up a pool where the operational decisions may be made centrally and then implemented by each utility. All of this requires significant data communication as well as engineering computation within a utility as well as between utilities. The use of modern computers and communications technology makes this possible, and the heart of system operations in a utility is the energy control center.

The monitoring and control of a power system from a centralized control center became desirable quite early in the development of electric power systems, when generating stations were connected together to supply the same loads. As electrical utilities interconnected and evolved into complex networks of generators, transmission lines, distribution feeders, and loads, the control center became the operations headquarters for each utility. Since the generation and delivery of electrical energy are controlled from this center, it is referred to as the energy control center or energy management system.

Although these centers have greatly advanced in technology, their basic control objective of maintaining reliability and security remains the same. In a vertically integrated utility the economic goal of minimizing the cost of supply is achieved by optimizing the dispatch and scheduling of its generators, whereas in deregulated regions the electricity markets optimize the supply costs. Increasingly, the generation is being privatized, and the generation companies have to compete in the electricity markets to sell their supply. The computerized tools needed to run the market and those used by the generation companies to determine their bid strategies are becoming more sophisticated. The rest of this section, however, describes the tools that are used to operate the grid itself in a reliable manner that avoids cascading failures and blackouts.

Evolution. In early control centers measurements of important power flows, voltages, and currents

were displayed on banks of meters. Circuit-breaker and switch positions were sometimes monitored and displayed by the switching of lights, placed appropriately on a wall map of the power system. Hard-wired control switches were used to open or close breakers remotely from the control centers. These supervisory control and data acquisition (SCADA) functions became more versatile and flexible with the computerization of the control center.

Automatic generation control (AGC), which adjusts the generator outputs to follow the load changes, became absolutely necessary when utilities started to interconnect, and customized controllers were initially built for the purpose. The use of first analog and then digital computers greatly improved the capabilities of automatic generation control. For example, the problem of meeting the load with the cheapest combination of generating sources requires a multistep calculated solution, for which digital computers are ideal.

Digital computer-based SCADA-AGC control centers became common in the 1960s. They could handle more information, provide the information to the operator more efficiently, and provide monitoring and control of the power system more reliably. Also, the availability of unprecedented computational power made it possible to develop new functions for better control. Security analysis and scheduling functions based on on-line computer programs are now available to the control-center operator, providing more economic and secure control.

Configuration. The control-center computers are connected through communication channels, usually microwave links, although fiber optics are increasingly used, to remote terminal units. A remote terminal unit is placed at each substation and generating station, and it can gather all the measurements at that station and execute the control commands sent from the control center.

At the control center itself, many computers must be linked together to accomplish the various functions. For example, one workstation may handle the SCADA functions, another may be dedicated to security analysis, while several more may handle the graphics interface needed by the operator. Such distributed processing configurations (**Fig. 4**) are common, and open-system standards are used for computer operating systems and communications, making it much easier to expand and modify existing hardware configurations. *See* DISTRIBUTED SYSTEMS (COMPUTERS); DISTRIBUTED SYSTEMS (CONTROL SYSTEMS).

Redundancy has to be carefully designed into the control-center configuration so that the probability of losing the critical functions is very small. In addition to redundant computers, backup is usually provided for communication channels, remote terminal unit circuits, operator consoles, and other equipment. *See* MULTIPROCESSING.

Human-machine interface. The banks of control panels have been replaced by operator consoles which typically consist of color display screens, a keyboard with extra function buttons, logging

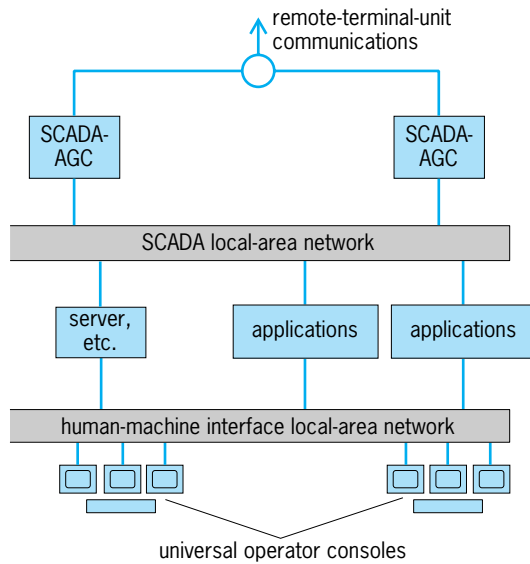


Fig. 4. Typical distributed-processing computer configuration for an energy control center.

printers, and telephones. All the measured and calculated data can be viewed at the display screens, and all control can be initiated and monitored from these consoles. A control center will usually have several consoles (Fig. 5), each dedicated to a particular aspect of system operation, like interchange monitoring or transmission control. Because the consoles are computer-controlled, their responsibilities can be changed at any time to accommodate changing procedures or availability of operators. The use of wall-mounted maps of the power system is quite common. Increasingly, these maps are very large liquid-crystal displays (LCDs) or plasma displays and are driven by the computers to denote breaker positions, line overloading, or out-of-limit voltages. See CATHODE-RAY TUBE; ELECTRONIC DISPLAY.

The control-center computers must handle thousands of displays, many of which are continually accessed and updated. This imposes a high demand on the computers, especially during emergencies, when the operator has to absorb and react to rapidly changing power-system conditions.

Control centers use display generator programs that allow the on-line creation of displays by the users, without computer-programming changes. Both the hardware and software techniques for graphics have improved dramatically over the years, and the quality of the human-machine interface is very high.

Supervisory control and data acquisition. All measurement data are scanned periodically by the computers, usually every 2 to 4 s. These data and others calculated from them are available to the operators through the displays. The data are checked for limit violations, and alarms are generated if necessary. Computerization and good database management programs have made it feasible for a control center to handle data from tens of thousands of information points. See DATABASE MANAGEMENT SYSTEM.

The supervisory control for switching of circuit breakers is usually done by using the displays and the function keys on the keyboard, with automatic checking and verification by the computers.

Automatic generation control. To match the load at all times by controlling generation, an automatic feedback control called load frequency control is used. The errors in frequency and scheduled power interchanges with neighboring utilities are continually checked, and control signals are sent to the generating units to adjust their generation levels. The units that can be controlled are those that the generation owners have agreed to put under control and those that have the ability to ramp their outputs up and down. The generators are compensated for this service at a contracted or market rate. (Fig. 6).



Fig. 5. Operator consoles at the energy control center of the New England Power Exchange (NEPEX) in Holyoke, Massachusetts. (ESCA Corporation, © Richard G. Shaw)

The generation control functions also include reserve monitoring to ensure adequate generation reserve, and interchange scheduling to automatically schedule contracted sales and purchases of energy between companies. If generation reserve is not adequate, the operator has to buy reserve capacity, and if interchange scheduling cannot be maintained at the contracted levels because of reliability considerations, the operator has to follow procedures to change the scheduled interchanges.

Security analysis functions. This set of functions (Fig. 7) is the direct result of the availability of powerful digital computers at control centers. A mathematical model of the power system can be stored in the computer and automatically updated in real time by using the measurements that are being continually monitored by the data-acquisition system. A program called the state estimator can track the real-time conditions of the power system with the stored model.

This real-time model can be used to provide the operator with answers to security questions. The contingency analysis program can analyze the effect of all probable contingencies on this model. If equipment failures at particular locations are determined to cause unacceptable operating conditions, the operator is alerted and can change the operating conditions so that these contingencies are no longer threatening to the system. This contingency analysis can be a static analysis (power flow), but for those systems that have stability problems a dynamic analysis is also used.

By using a study power flow program, the operator can also study this model for other purposes. If changes to the operating condition are contemplated by the operator, they can be first studied on this model to determine if they produce the desired results. An enhanced version, called the optimal power flow, can even determine the best operating strategies that should be tried to get the desired results.

Scheduling functions. For generation owners, optimizing costs is important. Steam-generating units have limits on how fast they can be started up or shut down. Hydroelectric units have limited available energy depending on the storage capacity of their reservoirs. These and other constraints affect the commitment and decommitment of generating units to meet the daily and weekly cycles of load. Unit commitment and hydrothermal coordination computer programs are available to determine the most economic scheduling pattern. With the availability of large computers these programs can be integrated with operations data so that the scheduler can update the results conveniently for changing conditions. These operations-planning functions, especially those that plan over many weeks, are often performed by engineering personnel distinct from system operators, on a separate computer system that shares data with the control-center computers.

In regions where generation owners are deregulated, the scheduling of generation depends on long-term (weeks) bilateral contracts and short-term

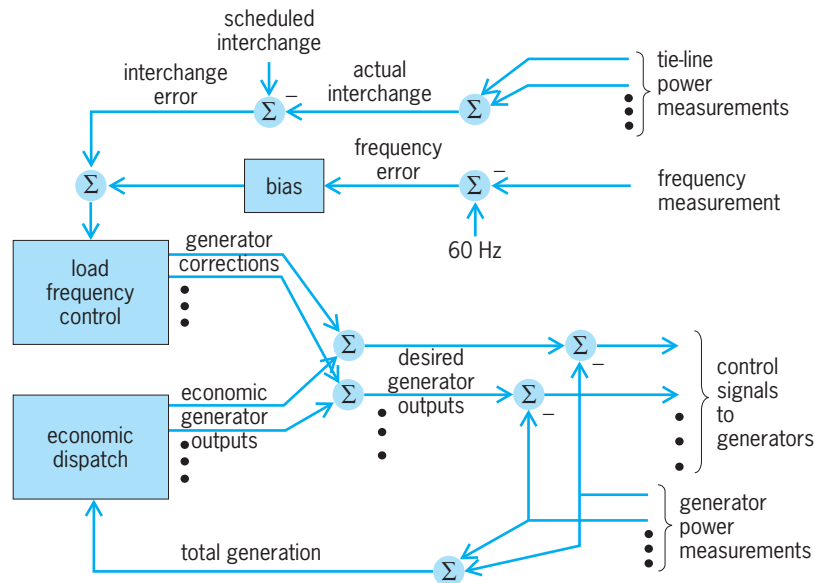


Fig. 6. Flow chart of the automatic generation control functions of an energy control center.

(24 hours or less) market transactions. These markets are computerized and use some sort of auction method to pick the bids and clear the market. Although the basic energy market is straightforward, the markets for reserve capacities and other ancillary services can be rather complicated. Moreover, even after the markets have cleared, the buy-sell transactions have to be checked by the system operator for their feasibility. This is because the markets clear on bid prices without regard to whether the transactions cause congestion (transmission overloads). Thus the operator has to check for congestion and manage the congestion by altering the transactions. This can be done iteratively between the market computers and the reliability studies or in one step by using the optimal power flow.

Logging and accounting. Computers have largely automated the logging and accounting functions, resulting in more comprehensive record keeping of the actual power flows and reduced labor costs. Daily, weekly, and monthly reports are produced automatically.

Computer programs are used to determine the cost of electric energy using the contract or market prices and the actual flows which are always

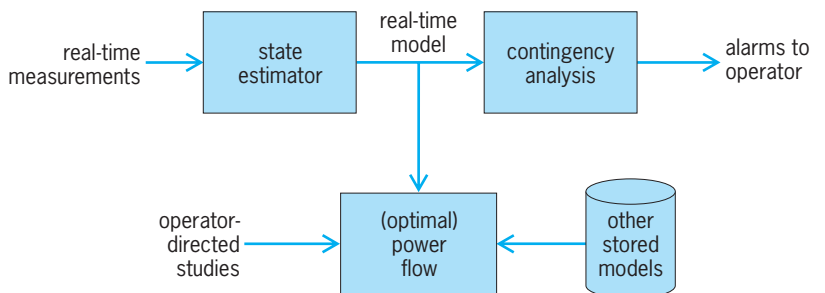


Fig. 7. Flow chart of the security analysis functions of an energy control center.

different from the contracted flows. The purpose is to automatically calculate the billing for energy that was bought and sold between generating companies and the distribution companies or loads. The same data can be used to compensate the transmission companies for their transmission service.

Operator training simulator. As the functions conducted at the control center have become more sophisticated and operating the power system more complex, the training of the operator has also become more important. The use of simulation training, like that used to train airline pilots and air-traffic controllers, has become common because there are no good alternatives for learning to handle emergency situations. The capability of computers to model the power system in real time has made this simulation training possible.

Impact of deregulation. The worldwide trend toward deregulating the electric power industry has had a significant impact on how the electric grid is operated. The presence of more organizations that has resulted from the separation of generating, transmission, and distribution companies has increased the need for new operational procedures to ensure the reliability and security of electricity supply. Here also, superior communication and computation technology play a vital role; in fact, such deregulation would not be feasible without them. Some of the trends in separating the economic functions from the reliability functions are described above, including some of the new tools like market systems, but changes are happening at different rates in different parts of the world and best practices for operating these deregulated systems are still emerging.

Anjan Bose

Power-System Planning

The two traditional objectives of power-system planning are (1) to provide electricity as cheaply as possible while (2) satisfying standards of reliability, safety, and quality. However, since about 1980, social, economic, and technical changes have altered the orientation of power-system planning everywhere. Power system planning, once done by vertically integrated utilities in isolation, has become a cooperative and even a competitive process.

Load forecasting. Planning begins with load forecasts. For decades, demand in the United States grew at about 7% per year and was correlated to the gross domestic product, but since about 1970 the rate of load growth has averaged about 2% per year. Overall demand is forecast in three ways: extrapolation of historic growth, econometric modeling, and end-use forecasting. A fourth technique, based on historic experience, is used to forecast local-area load growth. For day-to-day or hour-to-hour operations forecasting, weather and social patterns dominate.

Extrapolation. If load grows by the same percentage each year, the logarithm of load grows at a linear rate. A straight line passing through and beyond a plot of the logarithm of historical data gives an extrapolation forecast of the logarithm of load; load itself is forecast by taking the antilog. This method fails if the rate of

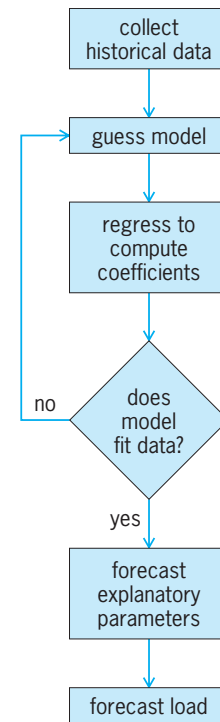


Fig. 8. Econometric load forecasting.

load growth changes. See ALGEBRA; EXTRAPOLATION; LOGARITHM.

Econometric modeling. Mathematical models (formulas) can be developed that relate changes in load to changes in explanatory variables like population, gross domestic product, prices, and so forth (Fig. 8). Historical data are collected on loads and on possible explanatory factors. Often, the first formula or model tried is linear, like the equation below. Regression is

$$\text{LOAD} = k_1 + k_2 \times \text{POPULATION} + k_3 \\ \times \text{GROSSDOMESTICPRODUCT} + \dots$$

a mathematical procedure for computing the values of the coefficients k_1 , k_2 , and so forth, to make the model consistent with the historical data. Statistical tests measure just how closely the model matches the data. See STATISTICS.

Once the model adequately relates the historical explanatory data to past loads, the explanatory factors are forecast. Finally, the formula or model is used to compute the load forecast from the forecasts of explanatory factors.

Econometric modeling often works better than extrapolation. But both methods infer future load growth from historical data—a reasonable approach only if major changes in the economy do not occur.

End-use forecasting. Residential load forecasts can be built up from energy-use forecasts. The number of homes is forecast first. Major appliances are identified and their market penetrations are projected. The annual energy used by each type of appliance is also projected. This is multiplied by the appliance penetration and by the number of homes to give an energy forecast per appliance. The residential forecast is the sum of the appliance forecasts. This

method takes more effort than the other two and does not work well for commercial or industrial loads, but it depends less on the assumption that the past foretells the future.

Local-area forecasting. Forecasting loads in local areas is important for subtransmission and distribution system planning. Load often develops in an S-shaped pattern, with low historic loads succeeded by rapid growth as new residential or commercial development occurs, followed by low growth again when development saturates. Forecasting is based on experience with prior development patterns.

Generation planning. Centralized generation planning is done by utilities or other organizations having a statutory or implied "obligation to serve" and guaranteed recovery of costs. Generation options include building new utility-owned power plants, buying power from independent power producers or other entities, retiring old generating units or extending their lives, and repowering older units to increase their capacity or efficiency.

Which options are most attractive varies from country to country, and among regions in larger countries, depending on local conditions. For example, there are still vast unexploited hydroelectric power resources in Latin America, Asia, and Africa, while in much of Europe and the United States most of the economically attractive hydroelectric power has been developed.

Plans traditionally are evaluated in terms of two main criteria: reliability and cost (variable and fixed). They must also satisfy environmental and other requirements. Planners gage reliability in terms of reserves (installed generation capacity minus peak demand) or any of several probabilistic measures of the likelihood that enough generators will be available to meet the load.

Variable costs are roughly proportional to how heavily a plant is operated; fuel costs are an example. Variable costs include not only the cost of a particular option but also how it will affect the costs of the rest of the system. Variable costs often are calculated by using a production-cost computer program that simulates the operation of the region's plants.

Fixed costs are independent of how a particular plant is run. They include interest, depreciation, profits, taxes, and so forth. In the United States, fixed and variable costs are often added and expressed as revenue requirements, the amount that customers will have to pay if the utility selects a particular option. To choose among several options, the present worths of their revenue requirements are compared.

Planners also assess the effects of load forecast uncertainty, generator forced outages, possible variations in fuel price, uncertainty in performance of new or rebuilt units, construction cost overruns, and, particularly in less-developed countries, uncertainty in construction time.

While an independent power producer (IPP) can be a generation-planning option from the perspective of a utility, it adds uncertainty and competes

with utility-owned generators. It is evaluated much like a utility-owned plant in terms of its effect on cost of electricity, system reliability, and so forth. If the IPP's cost or performance is uncertain, the utility must ensure that it will be allowed to pass cost variations to its customers through tariffs. If there is uncertainty with regard to whether the IPP will be built, the utility must have contingency plans to replace it. If the IPP is viewed as a competitor, the utility attempts to design and operate its own generators at a cost lower than the price the IPP can afford to charge.

Planning is done from an entirely different perspective by IPPs or other entities without guaranteed cost recovery and without an obligation to serve. Decisions are based on market analyses, where load growth is only one of the critical drivers. Other drivers include the costs of existing and possible future generating units with which the IPP will compete. The IPP also considers regulations and market rules, which can affect its operation and its profitability. This type of planning is different in detail but not in concept from the development planning that any commercial organization must make. Planning is complicated by the high initial cost and long economic lives of generating units. Utilities are happy to recover their investments over 40 years. Most investors want less than 10 years.

Transmission planning. A transmission system consists of current-carrying hardware (lines, transformers, structures, and so forth), control and communication equipment (computers, fuses, relays, and so forth), and practices and procedures. The transmission system is coupled to the generation system; together, they form an integrated energy conversion machine. The transmission system is not a transportation system like a gas pipeline, and it cannot be analyzed independent of the generation system.

One transmission planning option is construction of new alternating-current or direct-current lines, as discussed above. Other options include upgrading existing lines, installing transformers or compensation devices (capacitors and reactors), and introducing new operating procedures. Advances in semiconductor devices have greatly increased the variety of control techniques available. Collectively called FACTS (flexible alternating-current transmission system), these control devices allow utilities to utilize more fully the capability of the transmission system and to improve reliability. *See* STATIC VAR COMPENSATOR.

The final decision between two technically equivalent options is based on what each would cost the ratepayer, but the system must meet reliability criteria.

Reliability is assessed with computer programs that simulate network flows by using complex mathematical models of the generation and transmission system. The most sophisticated of these models simultaneously solve tens of thousands of nonlinear equations, with tens of thousands of unknowns; solution techniques are so powerful that this takes a few seconds on a personal computer.

Engineers often plan the transmission system to withstand the sudden loss of any single generating unit or network element under peak-load conditions; this is called the first-contingency criterion. Since this criterion requires the system to withstand the most severe contingency, the system is presumably more than adequate for less serious failures, and can generally withstand more than one of these.

Strategic and integrated resource planning. Strategic planning is concerned with long-range decisions affecting the fundamental nature of an enterprise. Options include acquisitions, divestitures, mergers, and diversification.

Least-cost planning or integrated resource planning (the terms are synonymous) is strategic planning for energy. It differs from tactical generation and transmission planning in three ways. First, integrated resource planning considers more options, with unbiased treatment of supply (generation) options and demand options such as conservation and load shifting. Second, integrated resource planning seeks to reconcile the objectives of a variety of stakeholders. Finally, integrated resource planning deals in a more fundamental way with risk and uncertainty.

Options. Numerous demand options are available, as discussed more fully below. Some are active; the utility can switch off loads when system conditions require it. Others are passive, relying on price signals to influence customers to reduce or shift demand. Others are both active and passive: in demand subscription service, for example, customers are paid to join a program requiring them to reduce load to their subscription level if the utility requests it.

Objectives. While generation and transmission planning concentrate on two main objectives (reliability and ratepayer cost), integrated resource planning also recognizes costs and benefits perceived by the investors, the regional or national economy, neighbors concerned with the environment, utility employees, and so forth. These costs and benefits cannot all be reduced to monetary terms, and they represent objectives that often conflict. For instance, a clean-burning premium fuel is good for the neighbors but expensive for the ratepayer.

The traditional approach for resolving conflicting objectives is to create a “utility function” by simply adding all of the objectives, each multiplied by an appropriate weighting factor. The utility function is then optimized. An extensive body of theory has been developed to support this approach. The basic difficulty is that objectives often are truly incommensurate. Even something as obviously economic as reliability cannot be expressed precisely in monetary terms.

A newer approach is to express all objectives in their natural units and to seek solutions that give the best trade-offs or compromises among these conflicting objectives.

Uncertainties and risk. An uncertainty is a parameter whose value is unknown. Uncertainties can be modeled probabilistically or as unknown-but-bounded variables, depending on the circumstances. Both types of models have their place. A common error is to insist on using a probability model, because it

appears to be more elegant mathematically, in circumstances where it is not valid.

For hundreds of years mathematicians and businessmen have studied the problem of risk. Some business risks are quantified in economic terms.

Nonetheless, many risks cannot be so quantified and managed. It is not reasonable to cover risks of terrorism, for example, by simply buying an insurance policy. The risks are not purely economic. Similarly, the risks of a blackout cannot be monetarized.

The approach to risk described here does not require that risks be reduced to dollar terms.

Risk is the hazard to which an entity is exposed because of uncertainties. Measures of risk include the likelihood or conditions under which a particular decision will be regrettable, and the magnitude of the possible regret. A decision that is not likely to be regretted, or where the possible regret is low, is said to be robust.

Integrated resource planners assess risk and develop hedges to manage it. Hedges include building small plants rather than big ones, buying futures contracts to protect against price volatility, doing pilot demand-management studies before selecting a demand-side option, and designing expansion so that it can be staged.

For example, one utility considered two major planning options: building a large coal-fired plant near the load, and building a subarctic hydroelectric plant with a long transmission line. The life-cycle costs were about the same, with the hydroelectric plant perhaps a bit cheaper, but the utility opted for the coal plant to avoid the risks associated with subarctic construction and a long line.

Regulatory changes. The regulatory compact is an unwritten agreement that permits electric utilities to enjoy monopoly protection and guaranteed cost recovery and profits. In exchange they accept an obligation to serve all comers and eschew monopoly profits. This century-old compact is being changed in a variety of ways in much of the world. Some feel that the industry is now mature and does not need special protection, and that a competitive market can do better than a regulated monopoly to meet customer needs at the lowest price.

Changes in the regulatory structure of the industry create important technical options and uncertainties. One key issue concerns the handling of transmission access and power wheeling, the provision of transmission services to energy buyers and sellers who are not directly interconnected. There are difficult engineering and equity problems associated with this issue, and they have not been resolved completely.

Planning in developing countries. Electric power is less available in the developing world because (1) its cost is beyond the means of many families and businesses, and (2) capital is not available for the massive investments needed to build generation and distribution facilities. These barriers can be overcome in part by planning and building the power system differently. Reliability and other quality-of-service standards in Addis Ababa should not be the same as in Tokyo. Furthermore, economic analyses

must use high discount rates (cost of capital) to reflect unavailability of capital to utilities and consumers. This means that minimizing first cost is more important than minimizing life-cycle costs. Equipment and development should be designed to be cheap initially, even if they must be replaced sooner rather than later.

Hyde M. Merrill

Power System Reliability

An electric power system serves the basic function of supplying customers, both large and small, with electrical energy as economically and as reliably as possible. The reliability associated with a power system is a measure of its ability to provide an adequate supply of electrical energy for the period of time intended under the operating conditions encountered. Modern society, because of its pattern of social and working habits, has come to expect the power supply to be continuously available on demand. This availability, however, is not physically possible in reality due to random system failures which are generally outside the control of power system engineers, operators, and planners.

The probability of customers getting disconnected can be reduced by increased investment during the planning phase, the operating phase, or both. Overinvestment can lead to excessive operating costs, which must be reflected in the tariff structure. Consequently, the economic constraints can be violated even though the system may be highly reliable. On the other hand, underinvestment can lead to the opposite situation. Design, planning, and operating criteria and techniques have been developed over the past three or four decades in an attempt to resolve and satisfy the dilemma between the economic and reliability constraints.

The criteria and techniques first used in practical applications were basically deterministic (rule-of-thumb) ones; for instance, installed generating capacity equals the expected maximum load demand plus a fixed percentage of the expected maximum demand. The essential weakness of these methods is that they do not account for the probabilistic and stochastic nature of system behavior, customer load demands, and component failures. Such aspects can be considered only through probabilistic criteria.

It is important to conjecture at this point what can be done regarding reliability assessment and why it is necessary. Failures of components, plant, and systems occur randomly; the frequency, duration, and impact of failures vary from one year to the next. Generally, all utilities record details of the events as they occur, and produce a set of performance measures, such as system availability, estimated unsupplied energy, number of incidents, and number of hours of interruption. These measures are valuable since they identify weak areas needing reinforcements and modifications, establish chronological trends in reliability performance, establish existing indices which serve as a guide for acceptable values in future reliability assessments, enable previous predictions to be compared with actual operating experience, and monitor the response to system design changes.

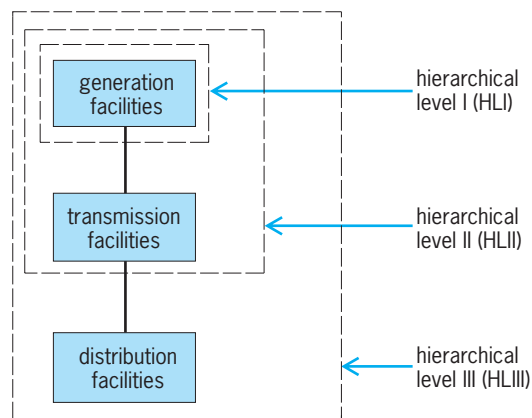


Fig. 9. Hierarchical levels in reliability evaluation.

Hierarchical levels. A modern power system is complex, highly integrated, and usually very large. Analyzing the entire system as a whole, apart from being extremely difficult, may provide results that are so vast that meaningful interpretation will be difficult, if not impossible. The power system can be categorized, for purposes of planning, organization, control, monitoring, operation, maintenance, and so forth, into three distinct zones of generation, transmission, and distribution. While this subdivision of the system may seem somewhat simplistic, most power utilities are divided into these zones or segments, and perhaps some smaller subzones such as subtransmission. The three zones perform quite different functions but must act in unison to supply the customer with reliable electric energy. These zones can also be combined to form three distinct hierarchical levels (HL) [Fig. 9], and reliability evaluation can be performed at one or more of these levels. There are two main approaches that can be used to evaluate the reliability of power systems regardless of the hierarchical level: analytical and simulation.

Adequacy and security. The concept of power system reliability, the overall ability of the system to satisfy the customer load requirements economically and reliably, is extremely broad. For the sake of simplicity, it can be divided into the two basic aspects of system adequacy and system security. Adequacy relates to the existence of sufficient facilities within the system to satisfy customer load demands. These include the facilities to generate power, and the associated transmission and distribution facilities required to transport the generated energy to the load points. Adequacy therefore relates to static system conditions. Security pertains to the response of the system to the perturbations or disturbances to which it is subjected. These may include conditions associated with local and widespread disturbances and loss of major generation or transmission. Most of the techniques presently available are in the domain of adequacy assessment.

Reliability cost and reliability worth. Due to the complex and integrated nature of a power system, failures in any part of the system can cause interruptions which range from inconveniencing a small number of local residents to a major and widespread

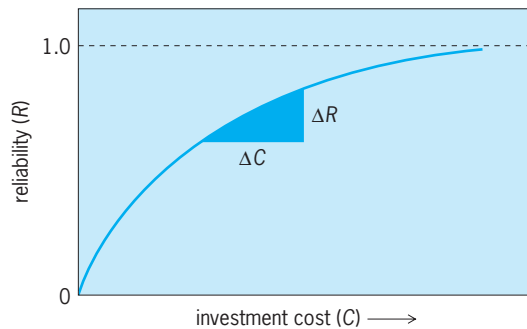


Fig. 10. Incremental cost of reliability.

catastrophic disruption of supply. The economic impact of these outages is not necessarily restricted to loss of revenue by the utility or loss of energy utilization by the customer but, in order to estimate the true costs, should also include indirect costs imposed on customers, society, and the environment due to the outage. For instance, 84% of the total costs of the 1977 New York blackout were attributed to indirect costs. In order to reduce the frequency and duration of these events, it is necessary to invest in the design phase, the operating phase, or both. This involves answering a few difficult questions—how much should be spent? is it worth spending any money? should the reliability be increased, maintained at existing levels, or allowed to degrade? who should decide, the utility, a regulator, the customer? on what basis should the decision be made?

The first step in answering the above questions is illustrated in Fig. 10, which shows how the reliability of a product or system is related to the investment cost; that is, increased investment is required in order to improve reliability. This figure clearly shows the general trend that the incremental cost ΔC to achieve a given increase in reliability ΔR increases as the reliability level increases. Alternatively, a given increase in investment produces a decreasing increment in reliability as the reliability is increased. In either case, high reliability is expensive to achieve.

HLI evaluation. The total problem can be divided into two conceptually different areas designated as static and operating capacity requirements. The static capacity area relates to the long-term evaluation of this overall system requirement. The operating capacity area relates to the short-term evaluation of the actual capacity required to meet a given load level.

One practice that has developed among utilities over many years is to measure the adequacy of both the planned and installed capacity in terms of the percentage reserve, which tends to compare the relative adequacy of capacity requirements provided for totally different power systems on the basis of peak loads experienced for each system over the same time period, and also does not attach any penalty to a power-generating unit because of size, unless the unit size exceeds the total capacity reserve.

During planning, it is necessary to determine how much capacity needs to be installed in order to satisfy the expected demand at some point in the future,

and to provide sufficient reserve to perform corrective and preventive maintenance. The ability to move the generated energy to bulk supply points and customers is not considered at HLI.

Generating capacity reliability (HLI reliability) is defined in terms of the adequacy of the installed generating capacity to meet the system load demand. Outages of generating units in excess of the estimates could result in “loss of load,” that is, the available capacity (installed capacity minus capacity on outage) being inadequate to supply the load. In general, this condition requires emergency assistance from neighboring systems and emergency operating measures such as system voltage reduction and voluntary load curtailment. Depending on the shortage of the available capacity, load shedding may be initiated as the final measure after the emergency actions.

Probabilistic criteria and indices used in HLI studies are:

1. Loss of load probability (LOLP), the probability that the load will exceed the available generation. Its weakness is that it defines the likelihood of encountering trouble (loss of load) but not the severity.
2. Loss of load expectation (LOLE), defined as the average number of days (or hours) per period on which the load is expected to exceed the available capacity.
3. Loss of energy expectation (LOEE), the expected energy that will not be supplied (in kilowatt-hours or megawatt-hours per period) due to those occasions when the load exceeds the available generation. It is presently less used than LOLE but is a more appealing index since it encompasses severity of the deficiencies as well as their likelihood.
4. Frequency and duration (F&D) indices, which measure how often a power outage occurs (frequency, in occurrences per period), and how long the outage lasts (duration, in hours per occurrence) each time it occurs.

The adequacy of the generating capacity in a power system is normally improved by interconnecting the system to another power system. The actual interconnection benefits depend on the installed capacity in each system, the total tie-line capacity, the load levels, and the type of agreement in existence between the systems. A loss-of-load situation is considered to arise in a system when the available assistance through the interconnections cannot offset a capacity deficiency arising due to capacity outages and load demands in that system.

HLII evaluation. An important element in the planning process, which is not considered in HLI studies, is the development of a suitable transmission network to transport the generated energy to the customer load points. The transmission network can be divided into the two general areas of bulk transmission and distribution facilities. The distinction between these two areas cannot be made strictly on a voltage basis but must include the function of the facility within the system.

In addition to providing the means to move the generated energy to the terminal stations, the bulk transmission facilities must be capable of

maintaining adequate voltage levels and loadings within the thermal limits of individual circuits, and also of maintaining system stability limits. The total problem of assessing the adequacy of the generation and bulk power transmission systems in regard to providing a dependable and suitable supply at the terminal stations can be designated as composite system reliability evaluation. Quantitative assessment of the adequacy of a composite system can be performed using a contingency evaluation approach. The basic procedure involves the selection and evaluation of contingencies (component failures), the classification of each contingency according to related failure criteria, and the accumulation of adequacy indices. Two types of indices exist: load point (bus) indices and system indices. Load point indices reflect the adequacy of individual bulk supply points, while the system indices reflect the overall adequacy of the entire system.

HLIII evaluation. HLIII includes all three functional zones in an actual assessment of customer load point reliability. The HLIII reliability performance evaluation is extremely complex. It can, however, be performed (under certain specified conditions) by utilizing the major load bus indices at HLII as input to the distribution networks.

Over the past few decades, distribution systems have received considerably less attention in terms of reliability modeling and evaluation as compared to the generating systems. The main reasons are that the generating stations are individually very capital intensive and that the generation inadequacy can have widespread catastrophic consequences for both society and its environment. Consequently, greater emphasis has been placed on ensuring the adequacy and meeting the needs of this part of the power system. A distribution system, on the other hand, is relatively cheap and outages have very localized effect. Therefore less effort has been devoted to quantitative assessment of various alternative designs and reinforcements. Analysis of customer failure statistics, however, indicates that for many utilities the distribution system makes the greatest contribution to the unavailability of supply to the customer.

The three basic reliability indices of average failure rate (number of failures per period of time), average outage time (duration of failure in hours per occurrence), and average annual outage time (hours per year of power interruptions) are calculated in HLIII studies. Although the three primary indices are fundamentally important, they do not always give a complete representation of the system behavior and response. For example, the same indices will be evaluated irrespective of whether 1 or 100 customers were connected to the load point or whether the average load at a load point was 10 kW or 100 MW. In order to reflect the severity of a system outage, additional reliability indices can be, and frequently are, evaluated: System average interruption frequency index SAIFI, customer average interruption frequency index (CAIFI), system average interruption duration index (SAIDI), customer average interruption duration index (CAIDI), aver-

age service availability index (ASAD), and average service unavailability index (ASUI). (ASAI indicates the percentage ratio of customer-hours of service provided to customer-hours requested.) ASUI indicates the percentage ratio of customer-hours of service not provided to customer-hours requested. *See* RELIABILITY, AVAILABILITY, AND MAINTAINABILITY. Lalit Goel

Demand-Side Management

Electricity is produced and distributed in direct response to the demands of residences, business, and industry. These demands come from the billions of electrical devices and appliances which are switched on and off. The demand which each device puts on the electrical system varies by its electrical characteristics, by the way it is used, and by the time of day. Power-system designs are based on plans in which power-system engineers try to forecast the future pattern and amount of the demand for electricity. That prediction is based on a complex set of formulas which consider a number of critical components. These considerations include the electrical demand from the appliances and devices that consume electricity, including predictions of how their population and efficiency may change over time; how new devices may become popular in the future; how people's behavior, or the way people use the devices, may change over time; how economic activity will allow business and industry to flourish and permit enough disposable income to cause consumers to purchase electricity; and the price of electricity and of competing energy forms (such as natural gas).

The role of electric power has grown steadily in both scope and importance over the past century. Developments in key technologies—including electric lighting, motors, computers, and telecommunications—have continuously reshaped daily life and increased the productivity of its commercial and industrial foundations. But the nature of end-use devices is changing dramatically. New uses are emerging from the convergence of developments in energy, telecommunications, transportation, the Internet, and electronic commerce, thus affecting the overall use of electricity.

These new devices require advanced, digital-quality service capabilities. In the United States, the electricity supply system was designed decades ago to nominally maintain 99.99% (5 sigma) reliability, while effectively ignoring outages lasting less than 5 minutes. This was generally acceptable in the twentieth century's analog, electromechanical economy. However, this traditional reliability standard is no longer relevant in today's digital world, where a power interruption or distortion of as little as a quarter of a cycle (0.004 second) can shut down a microprocessor-controlled end-use electrical device or process. This sensitivity infers a systemic (combined supply network plus end-use device) reliability requirement of the order of 10 sigma plus. By comparison, the current system performance results in an electricity unreliability cost to the U.S. economy that has grown to about a 50 surcharge on every dollar of electricity purchased, with no

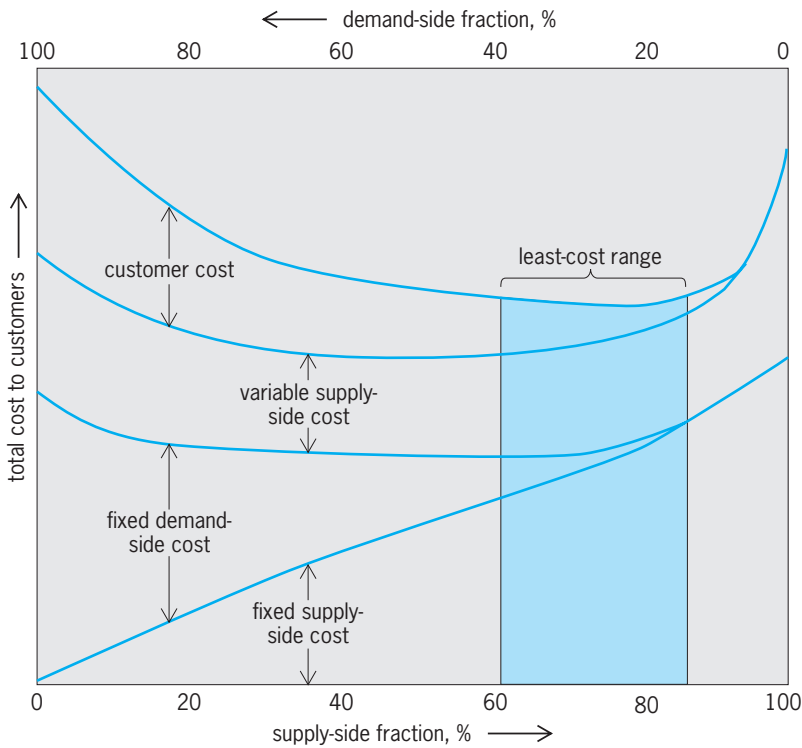


Fig. 11. Balance between supply-side and demand-side management activities.

end in sight until the system is modernized and its reliability transformed. This required performance improvement can most effectively be met by a “systems” approach involving a combination of electricity supply and end-use performance improvements.

Meeting the energy requirements of an increasingly digital society will require applying a combination of advanced technologies—including power generation (for example, central and distributed energy resources), interface devices, and less sensitive end-use equipment and circuits. In addition, the electricity infrastructure must be modernized to incorporate real-time communications, sensors, and computational capability.

The future demand for electricity, estimated based on assumptions about energy-consuming devices and their use, has traditionally been treated as a predetermined quantity by utility planners. Their responsibility has been to forecast that quantity and plan the addition of new generating units accordingly. By changing the efficiency of an appliance, or by changing the way it is used, a utility can alter the amount of future demand or alter the pattern of use in order to change the kind of new electrical generation required. However, traditionally, electrical utilities were not involved in customer programs to accomplish these objectives.

However, new technology is needed if society is to leverage the ever-expanding opportunities of networks such as the Internet and if electric utilities’ natural connectivity to consumers is to be effectively utilized to revolutionize both the role of the rapidly changing electricity service industry and the way

consumers may be connected to electricity markets in the future. This electricity energy service transformation will also open the door for significant efficiency improvements throughout the energy supply chain while facilitating greater use of relatively clean renewable and distributed energy resources.

In order to increase efficiency, control cost, and meet increasing demand for quality, electrical utilities have begun to address the increasingly complex needs of consumers and to shift from supply-side management, which considers only generating capacity, toward demand-side management. The latter promotes the efficient and lowest-cost utilization of facilities and augments the choices available to utilities and their customers. A compromise between supply-side and demand-side management can produce a range of low-cost resources called the least-cost range (Fig. 11).

Definition. Demand-side management is the planning and implementation of those utility activities designed to influence customer use of electricity in ways that will produce desired changes in the utility’s load shape, that is, in the time pattern and magnitude of a utility’s load. Utility programs in the category of demand-side management include load management, new uses, energy efficiency, electrification, customer generation, and adjustments in market share. Although this broad definition of demand-side management has many merits, in practice the term is generally used more narrowly to encompass increased efficiency, load management, and conservation, all with an emphasis on reducing the need for electrical energy or generation capacity.

Although there is an infinite combination of load-shape-changing possibilities, six changes in two categories can illustrate the range of possibilities: (1) load management including peak clipping, valley filling, load shifting, strategic load growth, and flexible load shape; and (2) energy efficiency. These changes are not mutually exclusive and may frequently be employed in combinations.

Load management. The first step in the shift toward demand-side management was the introduction and evolution of load management. Load-management programs focus on reducing customer use strategically at the time of high utility system loads. The goal is to avoid construction of generation facilities that would be operated for relatively few hours per year or costly power purchases when customer loads can be shifted or displaced at a lesser cost. There are several important components of load-management activities, including direct load control, pricing-based options (interruptible or curtailable rates), and thermal energy storage. The main costs to the utility for load-management programs are rate reductions or bill credits that are used to promote customer participation.

Energy efficiency. Efficiency-type demand-side management programs include various means by which a utility encourages the installation and use of high-efficiency appliances and devices. Programs often involve rebates for buying efficient appliances, bill credits, special rates, advertising, and promotion.

Programs have been focused on efficient, building insulation, heating and air-conditioning equipment installation, water heater wraps, high-efficiency lighting, and efficient refrigerators. Often appliance dealers, wholesalers, electricians, and others are involved in the utility efforts.

Consumer portal. Developments in electric power systems are stimulating the integration of a communications system throughout much of the power system. Once that system is present, connectivity to the ultimate consumers and facilitation of demand-side management can be enhanced with communications. This enhancement will allow three new areas of functionality: one that relates directly to electricity services (for example, added billing information or real-time pricing and demand-side management), one that involves services related to electricity (for example, home security or appliance monitoring); and a third that involves what are more generally thought of as communications services (for example, data services).

The integration of electric energy services and communications with consumer facilities and equipment opens a wide variety of opportunities to enhance consumer energy services, and to provide the consumer with access to a variety of new value-added services while enhancing the performance and operation of the energy system. Since consumer actions create the demand for energy that drives the operation of the energy system, communication with consumer equipment has the potential to profoundly impact the future of the energy system as well as to improve energy efficiency.

A consumer portal (CP) represents the technology, data standards, and protocols that enable the two-way communication of electronic messages with consumer-owned networks and intelligent equipment, and provide a pathway to the full development and implementation of a variety of demand-side management and energy-related services.

Energy saving through electricity. Because of the energy conversion which occurs in producing electricity, it has been argued that a unit of electricity saved is equivalent to preserving three units of primary energy. While this is often the case, concerns about electricity efficiency have sometimes extended to an assumption that electricity use is inefficient and adversely impacts the environment. However, a full understanding of energy efficiency requires a focus on total resource requirements. Electricity offers unique precision and control in its application, enhancing its efficiency. In addition, electricity offers benefits in comfort and environmental advantages. Because of these attributes, in almost every case new electric appliances and devices require less total resources than comparable natural gas or oil-fired systems. Increasing applications of electricity or electrification have been fundamental in increasing national productivity.

As originally defined, demand-side management not only is measured in terms of power-systems savings but also includes full-cycle energy efficiency (total resource use), environmental benefits, and pro-

ductivity. As such demand-side programs can include not only load management and energy efficiency but new technologies which can expand electricity use.

Full-cycle energy efficiency. Reduced electricity consumption is not necessarily the goal of policy aimed at energy efficiency. A more comprehensive goal is to reduce total consumption of energy resources. For example, the total resource efficiency of an advanced high-efficiency air-source electric heat pump may be compared with the total resource efficiency of the best available gas furnace. By tapping solar energy, the heat pump delivers 3.4 or more units of heat for every unit of electricity produced. An older production and delivery system supplying the heat pump might lose 68% of the fuel in power-plant generation (0.32 factor) and 8% in transmission and distribution (0.92 factor), still yielding 100% efficiency. By contrast, gas has only an 86% efficiency, without considering the requirement for an electric fan. See HEAT PUMP.

Other examples where electricity appears to have an advantage in terms of total resource efficiency include ground-source heat pumps, plug-in hybrid electric vehicles, and a variety of advanced industrial processes. In these cases, utility policy that is aimed at energy efficiency may imply increased electricity use.

Environmental benefits. The key environmental benefits stem from efficiency-related demand-side management. By obtaining the same end-use service with less electricity, there is a reduction in power-plant operation and fuel usage, implying a reduction in central-station emissions. These environmental benefits are an important figure of merit in measuring the success of demand-side management. However, most electric appliances and devices are far cleaner at the point of end use. A comprehensive demand-side management strategy must examine the comparative environmental benefits of alternative devices and systems.

Productivity. Application of automation and electric technologies can have important productivity implications. The benefits come in form of faster processing, more precise process control, reduced material and process waste, and higher product quality. The net impact is often reduced process cost, implying reduced product cost and increased competitiveness. Enhanced productivity is an important measure of success, even if these productivity gains require greater electricity usage.

Demand-side management practice. System expansion can be delayed or eliminated and the use of critical resources reduced by significantly reducing energy and peak consumption through demand-side management. Indeed, demand-side management can help achieve a broad range of operational objectives by merely changing the systems load shape. Numerous industries have found that changing the pattern of the demand for their product can be profitable. For example, telephone utilities have long offered reduced evening rates to shift demand and to encourage off-peak usage, airlines offer night coach fares, and movie theaters offer reduced matinee

prices—all examples of deliberate attempts to change the demand pattern for a product or service.

In an electric utility, the physical plant is designed to serve the projected demand for electricity in the least-cost manner, given a specified level of desired quality and reliability. If the load shape is not fixed but may be altered, the cost of serving the load can be reduced still further. Cost reduction due to the changes in the load shape arise primarily from three attributes: reduction in the requirements for new assets or energy, higher utilization of facilities, and more efficient operation of facilities.

Higher utilization of existing and planned facilities can also be achieved through a program of load growth. While such programs increase total costs due to higher fuel and other operating expenses, they reduce unit costs by spreading the fixed costs (debt service and dividends) over more units of energy sold.

Clark W. Gellings

Bibliography. P. M. Anderson and A. A. Fouad, *Power System Control and Stability*, 2d ed., 2003; R. Billinton and R. N. Allan, *Reliability Evaluation of Power Systems*, 2d ed., 1996; Edison Electric Institute, *Statistical Yearbook of the Electrical Utility Industry*, annually; D. G. Fink and H. W. Beaty (eds.), *Standard Handbook for Electrical Engineers*, 14th ed., 2001; J. D. Glover and M. S. Sarma, *Power System Analysis and Design*, 3d ed., 2001; T. Gonen, *Modern Power System Analysis*, 1988; J. J. Grainger and W. D. Stevenson, *Power System Analysis*, 1994; E. Kahn, *Electric Utility Planning and Regulation*, 2d ed., 1991; P. Kundur, *Power System Stability and Control*, 1994; R. H. Miller and J. H. Malinowski, *Power System Operation*, 3d ed., 1994; North American Electric Reliability Council, *Electricity Supply and Demand* (10-year forecast), annually; North American Electric Reliability Council, *Reliability Assessment* (10-year forecast), annually; F. Saccomanno, *Electric Power System: Analysis and Control*, 2003; *Transmission Line Reference Book, 200 kV and Above*, Rep. EL-2500, Electric Power Research Institute, 3d ed., 2005; A. J. Wood and B. F. Wollenberg, *Power Generation, Operation, and Control*, 2d ed., 1996.

Electric power transmission

The transport of generator-produced electric energy to loads. An electric power transmission system interconnects generators and loads and generally provides multiple paths among them. Multiple paths increase system reliability because the failure of one line does not cause a system failure. Most transmission lines operate with three-phase alternating current (ac). The standard frequency in North America is 60 Hz; in Europe, 50 Hz. The three-phase system has three sets of phase conductors, one to four conductors per phase, and system voltage is defined as the root-mean-square voltage between any two of the phases. Long-distance energy transmission occasionally uses high-voltage direct-current (dc) lines. *See*

ALTERNATING CURRENT; DIRECT CURRENT; DIRECT-CURRENT TRANSMISSION.

Overhead transmission lines are used in open areas such as interconnections between cities. However, environmentally preferable underground cables are used in congested areas within cities. The transmission lines generate electric and magnetic fields. The electric fields produce harmless but annoying sparks. For example, a worker standing on a wooden ladder may experience mild sparks when a coworker, standing on the ground, hands over a tool. The magnetic fields have become a source of concern, as discussed below. The use of underground cables eliminates the electric fields completely and reduces the magnetic fields in most cases. The sight of a large transmission line is not pleasing for most people. A further problem is the possibility of electrocution due to traffic accidents or accidental contact. The use of cable eliminates both problems.

The electric power system can be divided into the distribution, subtransmission, and transmission systems. With operating voltages less than 34.5 kV, the distribution system carries energy from the local substation to individual households, uses both overhead and underground lines, and is generally divided into high-voltage distribution (above 1000 V) and low-voltage distribution (110–480 V). *See* ELECTRIC DISTRIBUTION SYSTEMS.

With operating voltages of 69–138 kV, the subtransmission system distributes energy within an entire district and regularly uses overhead lines. For example, it may interconnect large substations located outside a town with the distribution substations inside the town. *See* ELECTRIC POWER SUBSTATION.

With operating voltage exceeding 230 kV, the transmission system interconnects generating stations and large substations located close to load centers by using overhead lines; high-voltage underground cables are used for special purposes only. Examples are the supply of large urban substations or underwater cable connections.

Requirements

The energy transportation system must be reliable and highly efficient. Both transmission-line capacity and supplied-power quality must meet variable load requirements.

Reliability. Transmission lines must transport energy continuously to customers in spite of being subject to the environment, that is, pollution, rain, storms, and lightning. System operation can itself produce disturbances which may interrupt line operation. Present utility practices require that the number of transmission line outages be less than one per year per 100 mi (160 km). The most frequent cause of line outages is the single-phase short circuit caused by lightning strikes, pollution-initiated flashover of the insulator, or trees touching the line during a storm.

Safety. Human activity near or under transmission lines may result in short circuits or accidents. A frequent cause of accidents is contact between the line and an operating crane. Transmission lines must

meet strict safety standards. Most important are strict clearances between the conductor and ground, and between the conductor and supporting structures. The electric and magnetic fields which are an inevitable by-product of system operation have also caused concern, as discussed below.

Capacity. The maximum power that can be transmitted through a line depends on the size of the conductor, voltage drop, and system stability. The current-carrying capacity of a conductor limits the maximum power for short lines, while voltage drop and system stability are the determining factors for long lines.

Economy. Economic operation of the line is determined by system losses, operation costs, maintenance costs, and investment expenses. The dominant losses are the resistance (I^2R) losses in the conductors.

Power quality. Transmission-line voltage depends on the load. Load current produces voltage drop on the line impedance, reducing line voltage, and line capacitance increases voltage at the line end. Voltage at the receiving end of the line is usually low at heavy load during peak hours, and high at light load during night hours. Proper operation of the electrical system requires maintenance of the line voltage within narrow limits. Typical line voltage is kept between $\pm 5\%$, and regulation is achieved by switched shunt capacitors and shunt reactances.

Overhead Alternating-Current Transmission

Overhead transmission lines distribute the majority of the electric energy in the transmission system. A typical high-voltage line has three (Fig. 1) phase conductors to carry the current and transport the energy, and two grounded shield conductors to protect the line from direct lightning strikes. The usually bare conductors are insulated from the supporting towers by insulators attached to grounded towers or poles. The normal distance between the supporting towers is a few hundred feet.

Routes. Most low- and medium-voltage transmission lines are located along highways, roads, or streets. High-voltage lines are placed in transmission-line corridors to reduce interference with the public. Several lines usually share the same corridor.

Hardware. The major components of a transmission line are conductors, insulators, and towers or poles.

Conductors. Transmission lines use ACSR (aluminum cable, steel reinforced) and ACAR (aluminum cable, alloy reinforced) conductors. In an ACSR conductor, a stranded steel core carries the mechanical load, and layers of stranded aluminum surrounding the core carry the current. An ACAR conductor is a stranded cable made of an aluminum alloy with low resistance and high mechanical strength. ACSR conductors are usually used for high-voltage lines, and ACAR conductors for subtransmission and distribution lines. Ultrahigh-voltage (UHV) and extrahigh-voltage (EHV) lines use bundle conductors. Each phase of the line is built with two, three, or four conductors connected in parallel and separated by about 1.5 ft (0.5 m). Bun-

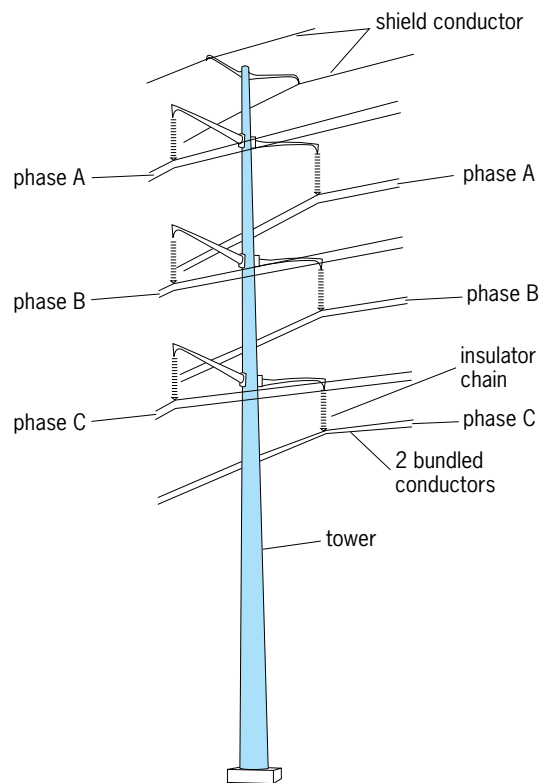


Fig. 1. Double-circuit, high-voltage transmission line.

dle conductors reduce corona discharge. See CONDUCTOR (ELECTRICITY).

Insulators. Lower-voltage lines use post insulators, while the high-voltage lines are built with insulator chains or long-rod composite insulators. A post insulator is made of either porcelain or composite material and bolted to the cross arm of the tower. The conductor is tied to the top of the insulator. A typical composite post insulator (Fig. 2a) consists of a fiberglass core covered by a silicone or ethylene propylene diene monomer (EPDM) rubber weathershed. The fiberglass provides high mechanical strength, and the rubber assures long leakage distance and good performance in bad weather.

An insulator chain is typically built with ball-and-socket glazed porcelain insulators (Fig. 2b). Under typical operating conditions, voltage is around 8–12 kV per unit. Wind swings the vertical chains and reduces clearances to the tower, a movement that is eliminated by a V arrangement.

In a suspension-type composite insulator (Fig. 2c), high-strength glass fibers bonded by epoxy or polyester resin form a fiberglass rod that is protected by an injection-molded weathershed made of silicone or EPDM rubber elastomer with hydrated aluminum filler. A corrosion-resistant iron or steel end fitting is glued with epoxy and wedged, crimped, or compressed to the fiberglass rod.

Industrial pollution covers insulators with a contamination layer that becomes conductive when wet by fog or dew, reducing insulator flashover voltage. Pollution-caused flashover results in many service interruptions. In polluted environments, composite

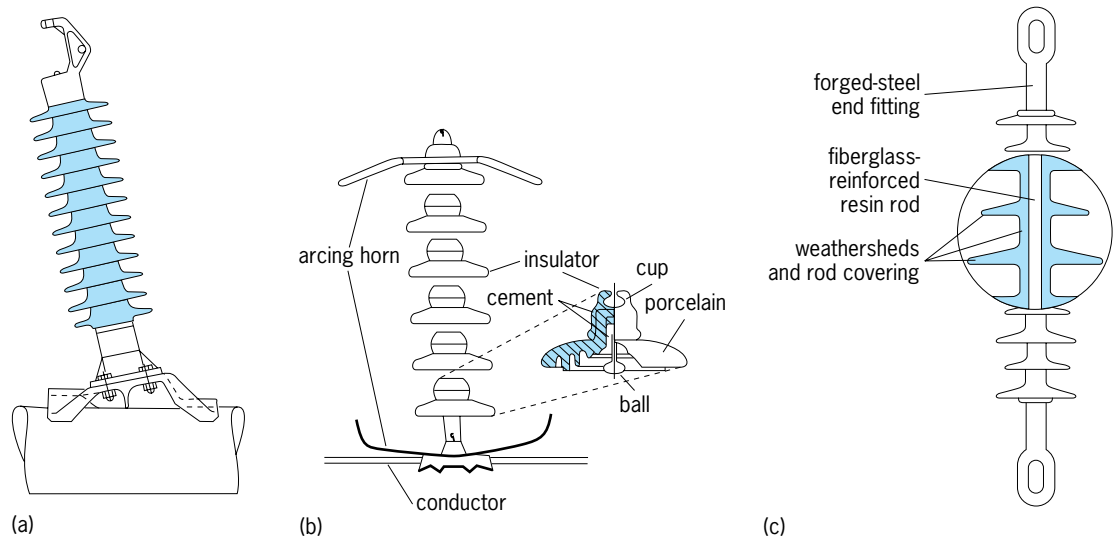


Fig. 2. Transmission-line insulators. (a) Composite post insulator. (b) Insulator chain. (c) Suspension-type composite insulator. (Sediver, Inc.)

insulators usually perform better than porcelain insulators. The performance of porcelain insulators can be improved by covering the insulators with silicone grease or a thin composite (silicone rubber) layer. In extremely severe conditions, insulators are washed regularly. Fog-type porcelain insulators are built with deep corrugation on the lower side to increase leakage distance. See ELECTRIC INSULATOR; ELECTRICAL INSULATION.

Towers and poles. A transmission line is supported by poles or towers. Distribution lines use wooden poles treated with chemicals (Figs. 3a, b). In some areas, wooden poles are replaced by concrete or tapered metal structures, whose higher mechanical strength permits longer span length, reducing the number of towers required for a line. These towers also have longer life spans and better appearance.

Some double-circuit high-voltage lines use lattice

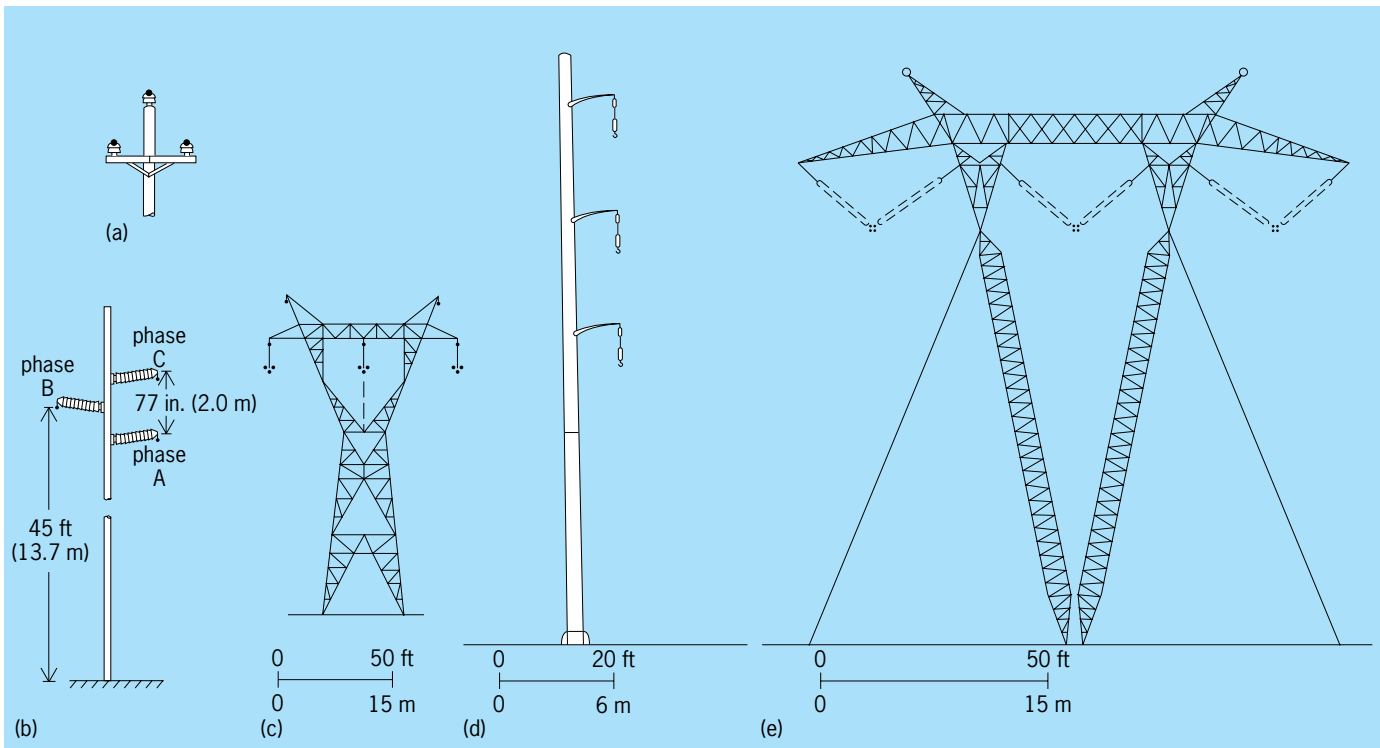


Fig. 3. Transmission-line towers and poles. (a) Distribution wood pole with porcelain post insulators. (b) Distribution wood pole with composite post insulators. (c) Steel tower with lattice structure and insulator chains. (d) Tapered, steel-tube tower with insulator chains. (e) Guyed V aluminum tower. (After D. G. Fink and H. W. Beaty, eds., *Standard Handbook for Electrical Engineers*, 13th ed., Mc Graw-Hill, 1993)

towers (Fig. 3c). Other high-voltage lines built with tubular steel poles (Fig. 3d) require less right-of-way and have a more pleasing appearance. Some extrahigh-voltage lines use guided towers in a V arrangement, supported by steel wires (Fig. 3e).

A transmission line is divided into sections, each terminated by a stronger dead-end tower with nearly horizontal insulators. From 10 to 25 suspension or tangent towers with vertical insulators are installed between the two dead-end towers, which have only transverse and vertical strength. Dividing the line assures that mechanical failures such as the breakage of conductors will affect only a small section of the line between dead ends. Stronger angle towers are used when lines change direction.

Mechanical performance. Transmission lines are subject to environmental adversities, including wide variations of temperature, high winds, and ice and snow deposits. Typically designed to withstand environmental stresses occurring once every 50–100 years, lines are intended to operate safely in adverse conditions.

Tension and sag. Conductor temperature varies with air temperature and power load, affecting sag and tension. Increased temperature reduces conductor tension and increases sag, reducing ground clearance (Fig. 4), while lower temperature increases tension of transmission lines.

Wind and ice loading. Extreme winter weather may coat conductors with ice, which increases conductor weight, as well as tension and sag. High wind produces a horizontal force on the conductor that increases tension and swings the conductor to an inclined position, reducing clearance between the conductor and the tower. Fortunately, high wind, ice loading, and minimum temperature do not occur simultaneously. Combined mechanical loads should not produce tension exceeding the conductor's tensile yield strength, and clearance should be greater than the minimum safety clearance.

Vibration. Nonuniform wind pressure produces conductor vibration. This is a low-amplitude, audio-frequency mechanical oscillation that bends the conductor at the point where it is clamped to the insulator. Stress from repeated bending may cause fatigue breakage of individual wires in the stranded conductor. Ultimately, complete breakage of the conductor may result in service interruptions. Dampers with weights are installed a few feet from the insulators to eliminate vibration-caused conductor breakage.

Galloping. The sudden drop of ice loading in a span causes galloping, a low-frequency, large-amplitude, mechanical oscillation that generates high mechanical stress and reduces the distance between phases in vertical conductors. The danger of galloping can be reduced by melting ice with high currents, using a low-voltage auxiliary transformer to drive the current through a short-circuited line.

Electrical performance. Variable weather affects line operation. Extreme weather reduces corona inception voltage, leading to an increase in audible noise, radio noise, and telephone interference. Load

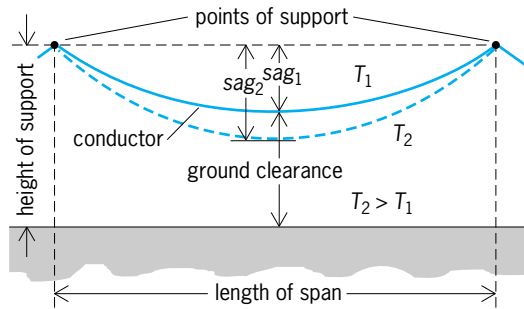


Fig. 4. Catenary curve showing the effect of temperature increase from T_1 to T_2 .

variation requires regulation of line voltage. A short circuit generates large currents, overheating conductors and producing permanent damage.

Line parameters. Transmission-line parameters are conductor resistance, inductance, and capacitance. Each phase has different inductance if the distance between conductors is not equal, and different capacitance if conductor heights are different. Longer lines are transposed by cyclic changes of conductor positions, permitting the use of average inductance and capacitance.

Line inductance and capacitance of a single-circuit, three-phase line can be calculated from the distances between the phases and the radius and geometry of the phase conductors. Conductor alternating-current resistances can be obtained from tables. Alternating-current resistance differs slightly from direct-current resistance because of the skin effect, and is also dependent on temperature. Use of a bundle conductor reduces both line reactance and resistance and increases shunt capacitance. See CAPACITANCE; ELECTRICAL RESISTANCE; INDUCTANCE; REACTANCE; SKIN EFFECT (ELECTRICITY).

Equivalent circuit. The parameters (inductance, resistance, and capacitance) are distributed along the line. For calculation of the relationship between sending- and receiving-end voltages and currents, distributed parameters are concentrated to an equivalent resistance, inductance, and capacitance (Fig. 5). Equivalent values for a short line are calculated by multiplication with the line length; for a long line, more complicated equations are used. An equivalent circuit represents only one phase of a three-phase system; thus, one-third of the three-phase power flows from the sending end (supply) to the receiving end (load). Voltages in this single-phase equivalent circuit are line-to-ground. The relationship between

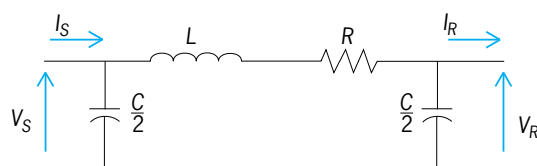


Fig. 5. Single-phase equivalent circuit for a transmission line. Line capacitance (C) is divided into two equal parts. V_S, I_S = voltage and current at sending end; V_R, I_R = voltage and current at receiving end.

sending- and receiving-end voltages and currents can be calculated by the Kirchhoff equations. See KIRCHHOFF'S LAWS OF ELECTRIC CIRCUITS.

Typically, in no-load or open-circuit conditions, receiving-end voltage is higher than sending-end voltage by as much as 5–25%, and can produce insulator flashover, voltage transformer fault, and so forth. The receiving-end voltage may be limited by shunt inductance, which loads the line. This inductance may be switched off when the load increases.

Increased load produces a voltage drop on the line and may reduce receiving-end voltage below sending-end voltage. This voltage drop can be compensated for by switched capacitors connected in series with the line. Capacitor reactance is negative and subtracted from line reactance, reducing the voltage drop and increasing receiving-end voltage for inductive load. Continuous control of series compensation can be achieved by shunting the series capacitor with a thyristor-controlled reactance. Delay of thyristor firing controls the time in which the inductance is connected in parallel with the capacitor in each half cycle, reducing the effect of capacitance. See SEMICONDUCTOR RECTIFIER.

Line power-transfer capacity. The power that a line can transport is limited by the line's electrical parameters. Voltage drop is the most important factor for distribution lines; where the line is supplied only from one end, the permitted voltage drop is about 5%.

The high-voltage line interconnects two buses in the electric system. Each of these buses is directly or indirectly connected to a voltage source. The amplitude and phase-angle difference between bus voltages regulates the direction and amount of power transferred through the line. The power transfer in a line with negligible losses is given by the equation below. Here, the bus voltages at the receiving and

$$P = \frac{V_S V_R}{X} \sin \delta$$

sending ends are V_R and V_S respectively; the phase angle between V_R and V_S is $\delta = \delta_S - \delta_R$, where δ_R and δ_S are the voltage phase angles at the receiving and sending ends; and the line reactance is X . Maximum power transfer occurs when $\delta = 90^\circ$. Stability of this operation requires that δ be not more than 40 – 60° . Maximum transferred power can be increased by reduction of series impedance by using series compensation. Voltage at both ends can be supported by shunt capacitors or static var compensators, which continuously generate and regulate reactive power. See STATIC VAR COMPENSATOR.

Thermal loading limit. Conductor temperature must be lower than the temperature which causes permanent elongation. A typical maximum steady-state value for ACSR is 212°F (100°C), but in an emergency temperatures 10–20% higher are allowed for a short period of time (10 min to 1 h). Steady-state temperature is calculated from the heat balance of the conductor. The sum of resistance (I^2R) losses and solar-radiation-generated heat is equal to the sum of convection- and radiation-produced heat losses. The line is typically designed for 104°F (40°C) ambient, 212°F (100°C)

conductor temperature, and 2 ft/s (0.6 m/s) wind speed. Maximum current is reduced by increased altitude and affected by weather conditions. The typical design value for a 500-kV line with two conductors in a bundle is 2400 A.

Corona. This discharge is generated when the electric field at the surface of the conductor becomes larger than the breakdown strength of the air. This high electric field accelerates free electrons to high speed, and the collision of these electrons with air molecules produces ionization. Ionization generates visible light, and fast ion movement produces losses and an audible hissing noise. The oscillatory nature of the discharge generates high-frequency, short-duration current pulses, the source of corona-generated radio and television interference. Because the electric field decreases rapidly with distance, the corona discharge occurs only in the vicinity of the conductor. Surface irregularities such as water droplets cause local field concentration, enhancing corona generation. Thus, during bad weather, corona discharge is more intense and losses are much greater. See ELECTRICAL INTERFERENCE.

Corona discharge generates audible noise with two components: a broad-band, high-frequency component, which produces crackling and hissing, and a 120-Hz pure tone. The loudest noise is generated during rain, but low-level noise during fog or wet conditions is equally disturbing; fair-weather noise is inconsequential. See CORONA DISCHARGE.

Electric field. Transmission-line conductors are surrounded by an electric field which decreases as distance from the line increases, and depends on line voltage and geometry. Near the conductors, this field generates corona discharge. At ground level, it induces current and voltage in grounded bodies, causes corona in grounded objects, and can induce fuel ignition. Current induced in the human body results in unpleasant sensations such as hair stimulation, and the field produces a secondary spark discharge when an ungrounded person touches a grounded object under the line. Utilities limit the electric field at the perimeter of right-of-ways to about 1000 V/m. See ELECTRIC FIELD.

Magnetic field. Load current produces an ac magnetic field around the transmission line that decreases with distance from the line. For example, the maximum magnetic field at ground level under a line which carries 2000 A is 176 milligauss (17.6 microtesla), and decreases to 9.6 milligauss (0.96 microtesla) at 330 ft (100 m) from the line. The field depends on line geometry and current, but is independent of line voltage. The Earth produces a dc magnetic field around 500 mG ($50 \mu\text{T}$).

Lightning protection. Lightning strikes produce high voltages and traveling waves on transmission lines, causing insulator flashovers and interruption of operation. Steel grounded shield conductors at the tops of the towers (Fig. 1) significantly reduce, but do not eliminate, the probability of direct lightning strikes to phase conductors. The ground resistance of the towers is important. A direct lightning strike to the shield wires or tower drives high current

to the ground. The ohmic voltage drop on the tower's ground resistance increases the tower's potential and may cause insulator flashover, referred to as backflash. Reduction of tower footing resistance reduces the probability of backflash. Proper line design results in 0.5–1 fault per year per 100 mi (160 km). See LIGHTNING AND SURGE PROTECTION.

Switching surges. The operation of circuit breakers causes switching surges that can result in interruption of inductive current, energization of lines with trapped charges, and single-phase ground fault. A switching surge is a unidirectional or oscillatory impulse with a 250–500-microsecond rise time and 2–10- μ s time to half value. These surges can have peak voltages 2–3.5 times the peak value of the 60-Hz voltage. A switching surge is most dangerous in extrahigh voltage lines (500 kV and above). Modern circuit breakers, operating in two steps, reduce switching surges to 1.5–2 times the 60-Hz voltage. First, the breaker inserts a resistance into the circuit. Then the breaker interrupts the already reduced current, or shorts the resistance and completes the line energization. See CIRCUIT BREAKER.

Telephone interference. Line current induces a disturbing voltage in telephone lines running parallel to transmission lines. Because the induced voltage depends on the mutual inductance between the two lines, disturbance can be reduced by increasing the distance between the lines and shielding the telephone lines. Telephone interference has also been reduced by the use of fiber optics and digital telephone techniques. A short circuit may induce dangerous voltages in neighboring telephone lines. See ELECTRICAL SHIELDING; INDUCTIVE COORDINATION.

Underground and Submarine Power Transmission

Most cities use underground cables to distribute electrical energy. These cables virtually eliminate negative environmental effects and reduce electrocution hazards. However, they entail significantly higher construction costs, typically 6–10 times those of transmission lines with similar capacity for medium-voltage distribution, and about 20 times higher for high-voltage transmission. High-voltage power cables are also used to transmit energy through wide deep rivers and bays where building transmission lines is difficult.

Cable construction. Underground cables are divided into two categories: distribution cables (less than 69 kV) and high-voltage power-transmission cables (69–500 kV).

Distribution cables. Extruded solid dielectric cables dominate in the 15–33-kV urban distribution system. In a typical arrangement (Fig. 6a), the stranded copper or aluminum conductor is shielded by a semiconductor layer, which reduces the electric stress on the conductor's surface. Cross-linked polyethylene (XLPE) or ethylene-propylene-rubber (EPR) is used most frequently for solid dielectric insulation. The insulation is protected by a semiconducting metallic tape shield which terminates the electric field. A poly(vinyl chloride) [PVC] jacket provides mechanical protection and prevents moisture penetration.

The major problem with these cables is trees generated by electric fields and water penetration. Trees are channels in the dielectric produced by electric discharges. In most cases, they start at the conductor and progress slowly through the insulation, ultimately producing dielectric breakdown and cable failure.

Oil-impregnated paper-insulated cables are used for higher voltages and in older installations (Fig. 6b). The stranded copper conductor is shielded by a metallized paper tape, which alleviates the unevenness caused by stranding. The conductors are insulated with oil-impregnated paper tape. Each insulated conductor is covered by metallized paper tape, which produces a concentric electric field. Oil-impregnated paper fillers are used to create a round cable assembly, held together by a metallized paper binder. The assembly is protected from moisture aggression by a lead or lead-alloy shield. Cables designed for insulation in a duct are protected by an extruded PVC sheet. Cables designed for direct burial are covered by a fibrous material (jute) impregnated by bituminous compound and double steel-tape armoring.

High-voltage power-transmission cables. Cable temperatures vary with load changes, and cyclic thermal expansion and contraction may produce voids in the cable. High voltage initiates corona in the voids, gradually destroying cable insulation. Low-pressure oil-filled cable construction reduces void formation. A single-phase concentric cable has a hollow conductor with a central oil channel. Three-phase cables have three oil channels located in the filler. Channels are supplied by low-pressure oil which penetrates the paper insulation and fills the voids

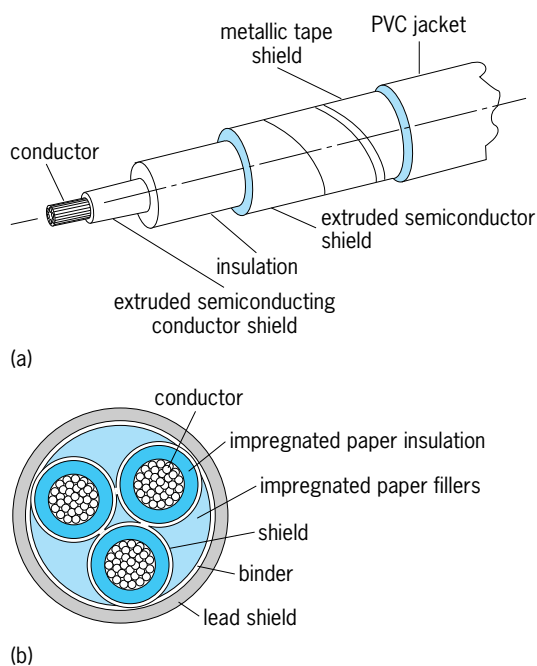


Fig. 6. Underground distribution cables. (a) Extruded solid dielectric cable. (b) Three-phase oil-impregnated paper-insulated cable. (After R. Bartnikas Srivastava, *Elements of Cable Engineering*, Sandford Educational Press, 1980)

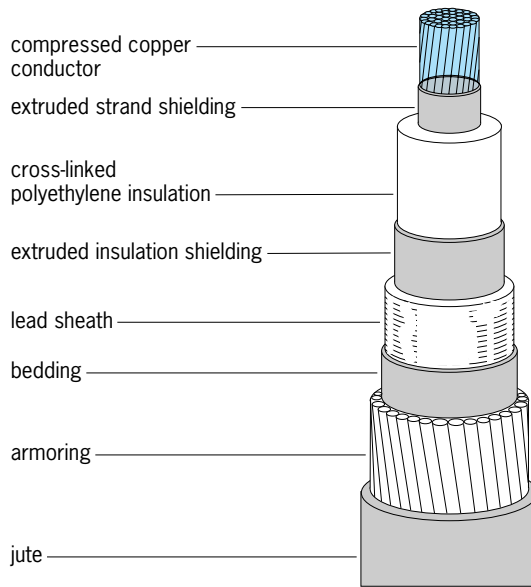


Fig. 7. Solid-dielectric submarine cable. (ASEA Kabel)

produced by thermal expansion. Oil is degassed before installation, and the system is hermetically sealed to avoid gas penetration during operation. These cables need complicated joints and oil tanks alongside them, which can cause difficulties. They are used between 69 and 500 kV, although 750-kV cable has been developed.

For extrahigh voltage, high-pressure oil-filled pipe cables are used. Three single-conductor shielded cables are installed in a steel pipe. The pipe is filled with high-pressure (200-lb/in.² or 1.4-megapascal) low-viscosity oil that penetrates the insulation and prevents corona generation. Oil pressure and temperature are maintained by automatically operated pumping stations.

Submarine cables. High-voltage cables are frequently used for crossing large bodies of water. Water provides natural cooling, and pressure reduces the possibility of void formation. A typical submarine cable (Fig. 7) has cross-linked polyethylene insulation, and corrosion-resistant aluminum alloy wire armoring that provides tensile strength and permits installation in deep water.

Cable parameters. Cables have significantly larger phase-to-ground capacitance and smaller inductances than overhead lines. Large capacitance results in high charging current, which reduces the power transfer capacity. At critical cable length, charging current is equal to load current. Typical critical length is in the range of 15 mi (25 km) at 500 kV and 45 mi (70 km) at 138 kV. Charging current can be compensated by shunt reactors or static var compensators. Because of the small inductance, reactive components of the voltage drop dominate in a cable system, and the voltage drop is significantly less than in an overhead line.

Cable power-transfer capability is dependent on the temperature of the dielectric, which depends on the heat generated by resistive power loss and on the thermal resistance. In a three-phase system,

the unbalanced load current induces a circulating current in the cable shield, which heats the cable and reduces power-transfer capacity.

Cable power-transfer capacity can be improved by externally cooling the cables. One method used in large substations is installing pipes cooled by circulating water near the cables.

Installation. In urban areas, cables are installed in underground concrete ducts requiring access holes every 1000–1500 ft (300–500 m). The ducts facilitate location of faulted cables and provide protection against accidental cutting of the cable during road construction.

In suburban areas, cables may be buried in trenches. Concrete slabs or bricks are often placed on top of the cable for protection and marking. The disadvantages of this method are difficulty in fault location and a lack of mechanical protection.

Cables are manufactured in sections 1000–1500 ft (300–500 m) long, with joints connecting the sections for urban cable. These joints, which are located in the access holes, frequently incur faults, and their preparation requires highly skilled workers. Bending and pulling may damage cables during installation, and proper installation requires expensive equipment and skilled labor. See ELECTRIC POWER SYSTEMS; TRANSMISSION LINES.

George G. Karady

Bibliography. Electric Power Research Institute, *Transmission Line Reference Book: 345 kV and Above*, EL 2500, 1987; D. G. Fink and H. W. Beaty (eds.), *Standard Handbook for Electrical Engineers*, 14th ed., 1999; J. A. Williams and P. L. Ostermann, *Underground Transmission Systems Reference Book*, Electrical Power Research Institute, 1992.

Electric protective devices

A particular type of equipment applied to electric power systems to detect abnormal and intolerable conditions and to initiate appropriate corrective actions. These devices include lightning arresters, surge protectors, fuses, and relays with associated circuit breakers, reclosers, and so forth.

From time to time, disturbances in the normal operation of a power system occur. These may be caused by natural phenomena, such as lightning, wind, or snow; by falling objects such as trees; by animal contacts or chewing; by accidental means traceable to reckless drivers, inadvertent acts by plant maintenance personnel, or other acts of humans; or by conditions produced in the system itself, such as switching surges, load swings, or equipment failures. Protective devices must therefore be installed on power systems to ensure continuity of electrical service, to limit injury to people, and to limit damage to equipment when problem situations develop.

Protective devices, like any type of insurance, are applied commensurately with the degree of protection desired or felt necessary for the particular system. Thus, application of these devices varies widely.

Protective relays. Protective relays are compact analog or digital networks, connected to various points of an electrical system, to detect abnormal conditions occurring within their assigned areas. They initiate disconnection of the trouble area by circuit breakers. These relays range from the simple overload unit on house circuit breakers to complex systems used to protect extrahigh-voltage power transmission lines. They operate on voltage, current, current direction, power factor, power, impedance, temperature, and so forth, as well as combinations of these. In all cases there must be a measurable difference between the normal or tolerable operation and the intolerable or unwanted condition. System faults for which the relays respond are generally short circuits between the phase conductors, or between the phases and grounds. Some relays operate on unbalances between the phases, such as an open or reversed phase. A fault in one part of the system affects all other parts. Therefore relays and fuses throughout the power system must be coordinated to ensure the best quality of service to the loads and to avoid operation in the nonfaulted areas unless the trouble is not adequately cleared in a specified time. *See* FUSE (ELECTRICITY); RELAY.

Zone protection. For the purpose of applying protection, the electric power system is divided into five major protection zones: generators; transformers; buses; transmission and distribution lines; and motors (Fig. 1). Each block represents a set of protective relays and associated equipment selected to initiate correction or isolation of that area for all anticipated intolerable conditions or trouble. The detection is done by protective relays with a circuit breaker used to physically disconnect the equipment. For other areas of protection *see* GROUNDING; UNINTERRUPTIBLE POWER SUPPLY.

Fault detection. Fault detection is accomplished by a number of techniques. Some of the common methods are the detection of changes in electric current or voltage levels, power direction, ratio of voltage to current, temperature, and comparison of the electrical quantities flowing into a protected area with the quantities flowing out. The last-mentioned is known as differential protection.

Differential protection. This is the most fundamental and widely used protection technique (Fig. 2). The inputs to the relays are currents from a current transformer. Current through the relay (nominally 5 or less amperes) is proportional to the high-power current of the main circuit. For load through the equipment or for faults either to the right or left of the current transformers, the secondary currents all sum to essentially zero (Fig. 2a). However, for internal trouble, they add up to flow through the relay (Fig. 2b).

A typical three-line diagram for a two-circuit differential system is shown in Fig. 3. For increased sensitivity, differential relays have restraint and operating windings. On through conditions (Fig. 2a) the restraint winding currents desensitize the operating winding, while full sensitivity is available for the internal faults (Fig. 2b). Thus for a fault

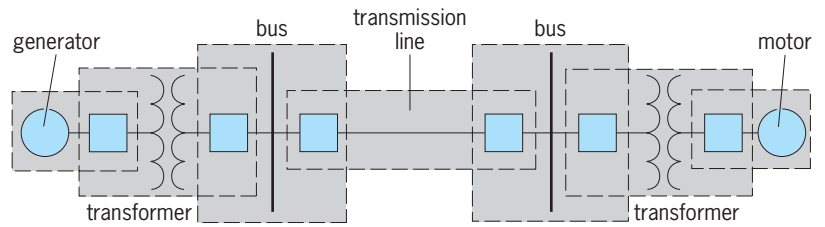
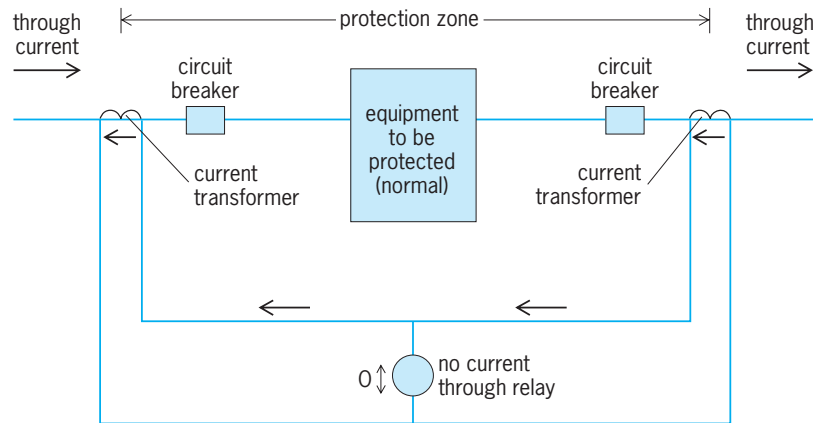


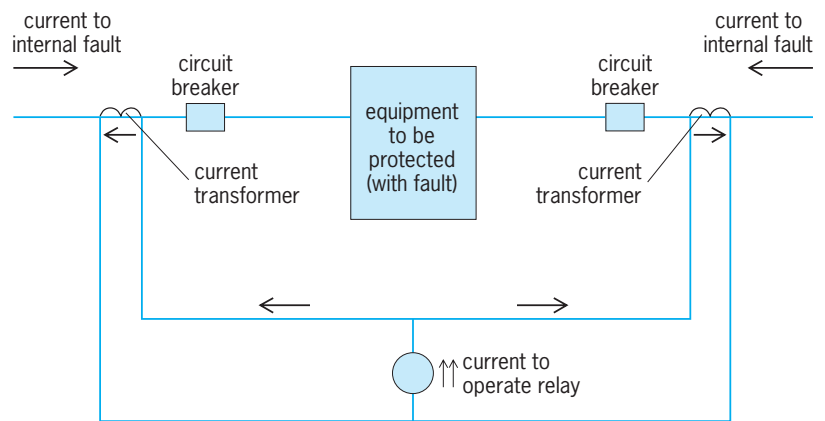
Fig. 1. Zones of protection on simple power system.

between the current transformers, secondary current flows through the restraining windings (with reversed phase in one of the coils) and in the operating coils to operate the relay, trip the two circuit breakers, and isolate the fault or damaged equipment from the rest of the power system.

The equipment to be protected can be a generator, transformer, bus, motor, or line. In the last case, with many miles between the circuit breakers, telephone-type low-energy pilot wires, optical fibers, radio frequency on the power line, or a microwave channel is used. The quantities at the line terminals are thus differentially compared via the channel, in what is known as pilot relaying.



(a)



(b)

Fig. 2. Simple differential protection scheme. The system compares currents to detect faults in the protection zone. Current transformers reduce the primary current to a small secondary value. (a) Normal operation. (b) With fault.

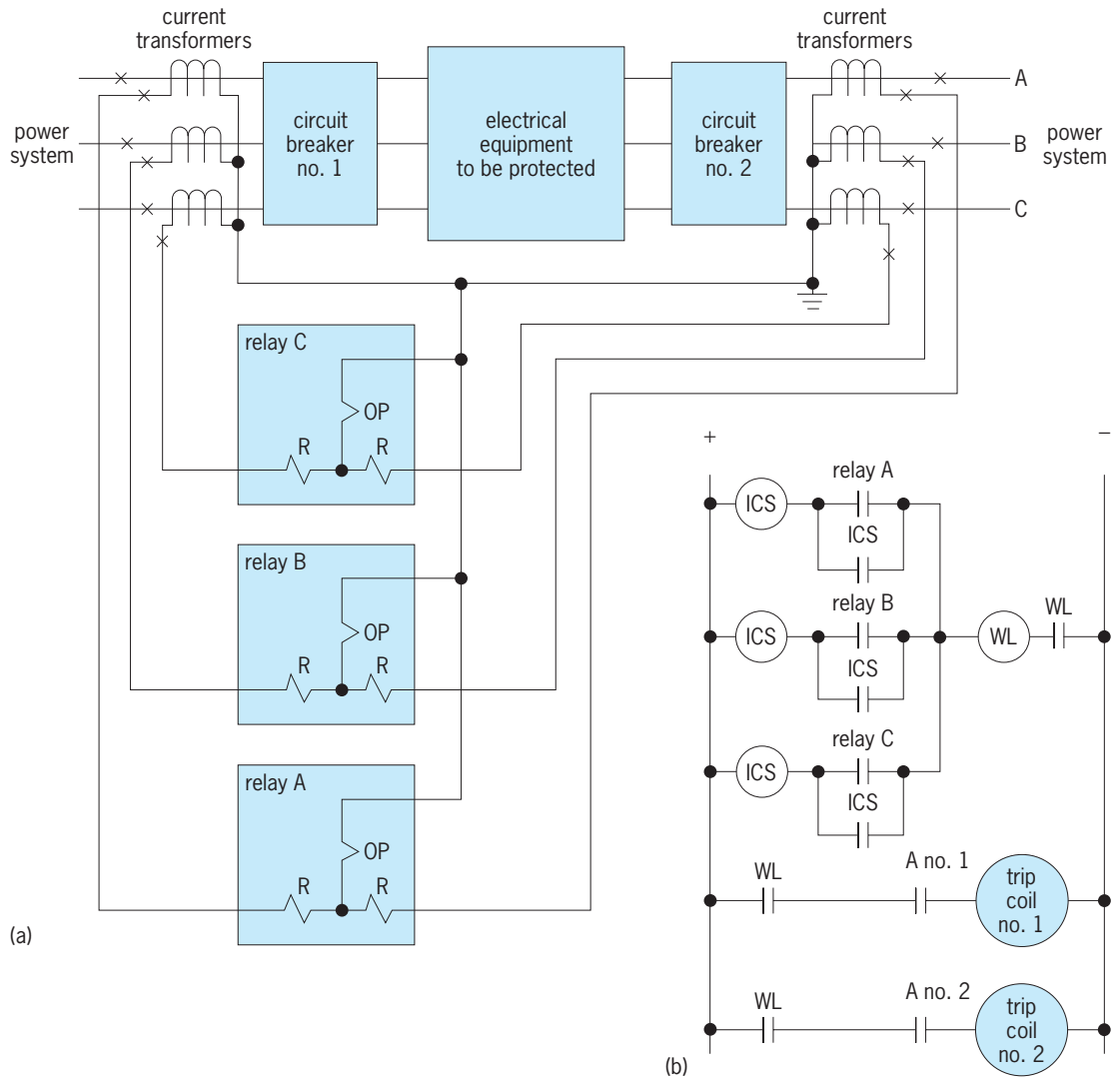


Fig. 3. Differential protection. (a) Typical arrangement for electrical equipment. Relays A, B, and C are differential relays. (b) Trip circuit. A = auxiliary out-off breaker; ICS = operation indicator and circuit seal-in auxiliary relay; OP = operating coil; R = restraining coil; WL = auxiliary hand-reset lock-out relay.

Overcurrent protection. This must be provided on all systems to prevent abnormally high currents from overheating and causing mechanical stress on equipment. Overcurrent in a power system usually indicates that current is being diverted from its normal path by a short circuit. In low-voltage, distribution-type circuits, such as those found in homes, adequate overcurrent protection can be provided by fuses that melt when current exceeds a predetermined value.

Small thermal-type circuit breakers also provide overcurrent protection for this class of circuit. As the size of circuits and systems increases, the problems associated with interruption of large fault currents dictate the use of power circuit breakers. Normally these breakers are not equipped with elements to sense fault conditions, and therefore overcurrent relays are applied to measure the current continuously. When the current has reached a predetermined value, the relay contacts close. This actuates the trip circuit of a particular breaker, causing it to open and thus isolate the fault. See CIRCUIT BREAKER.

Overcurrent relays are of two basic types: instantaneous and inverse-time. Instantaneous units have no intentional time delay and operate in the order of 10 to 50 ms. Inverse-time units operate quickly for high currents and with increasing time for smaller currents.

Overcurrent protection is applied in all parts of the power system. A typical application for the protection of a line between two stations R and S is shown in Fig. 4a. The overcurrent relays are connected to the current transformers to the left of breaker E and to the right of breaker F to operate for all faults in the primary protection zone between the two stations R and S.

If there is a significant magnitude difference between the close-in fault F_1 and the far bus fault F_2 , instantaneous overcurrent relays are set to operate for faults F_1 out to n_1 but not to operate for fault F_2 . The operating zone is indicated by the horizontal line from the relays at breaker E to n_1 , and for the relays at breaker F to n_2 . These lines represent the

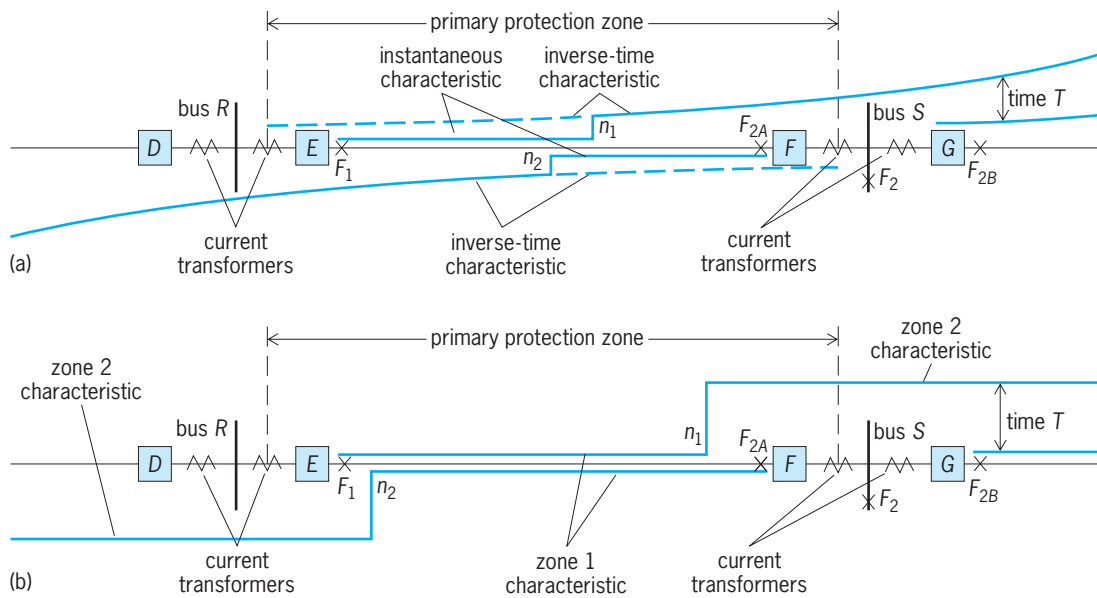


Fig. 4. Typical relay applications and settings. (a) Instantaneous and inverse-time overcurrent relays D , E , F , and G are circuit breakers. For the relay characteristics shown, operating times are shown vertically. Instantaneous relay time is instantaneous from the relay current transformers to points n_1 and n_2 . (b) Distance relay. Zone 1 relay characteristic is instantaneous. Zone 2 relay characteristic is instantaneous with operation delayed by time T .

maximum coverage. When the fault level decreases, less of the line will be protected.

The current magnitude through the relays at breaker E for a fault at either F_2 , F_{2A} , or F_{2B} is the same. Thus it is not possible for the relays at breaker E to determine if fault F_2 is at F_{2A} for which it should operate, or at F_{2B} where it should be cleared by the relays at breaker G . Inverse-time overcurrent relays with directional sensing into the line are applied at both breakers E and G . The relays at E are set to operate as fast as possible for faults at F_1 but to delay for a time interval T for the F_2 faults. Unfortunately this arrangement delays tripping for the internal fault F_{2A} , but it does provide time for the relays and breaker at G to clear the faults F_{2B} and faults to the right of G . If, however, either the relays or breaker G fails to operate for faults F_{2B} , then in time T breaker E will be opened to remove the fault current from the left side. This is backup protection by the relays and breaker at E for faults on the line to the right of breaker G . This process is known as relay coordination or selective settings.

Distance protection. Distance-type relays operate on the combination of reduced voltage and increased current occasioned by faults. They are widely applied for the protection of higher voltage lines. An application is illustrated in Fig. 4b. A major advantage is that the operating zone is determined by the line impedance and is almost completely independent of current magnitudes. Typical instantaneous settings n_1 and n_2 are 90% of the total line impedance.

The same difficulty of distinguishing internal faults F_{2A} from external faults F_{2B} at breaker E exists, and so a second distance unit is applied. It is set well into the adjacent lines as shown and is operated through a timer T . This permits the G relays and breaker an opportunity to clear faults F_{2B} and faults to the right

before E relays and breaker can operate as backup.

Overvoltage protection. Lightning in the area near the power lines can cause very short-time overvoltages in the system and possible breakdown of the insulation. Protection for these surges consists of lightning arresters connected between the lines and ground. Normally the insulation through these arresters prevents current flow, but they momentarily pass current during the high-voltage transient to limit overvoltage. Overvoltage protection is seldom applied elsewhere except at the generators, where it is part of the voltage regulator and control system. In the distribution system, overvoltage relays are used to control taps of tap-changing transformers or to switch shunt capacitors on and off the circuits. See LIGHTNING AND SURGE PROTECTION.

Undervoltage protection. This must be provided on circuits supplying power to motor loads. Low-voltage conditions cause motors to draw excessive currents, which can damage the motors. If a low-voltage condition develops while the motor is running, the relay senses this condition and removes the motor from service.

Undervoltage relays can also be used effectively prior to starting large induction or synchronous motors. These types of motors will not reach their rated speeds if started under a low-voltage condition. Relays can be applied to measure terminal voltage, and if it is below a predetermined value the relay prevents starting of the motor.

Underfrequency protection. A loss or deficiency in the generation supply, the transmission lines, or other components of the system, resulting primarily from faults, can leave the system with an excess of load. Solid-state and digital-type underfrequency relays are connected at various points in the system to detect this resulting decline in the normal system

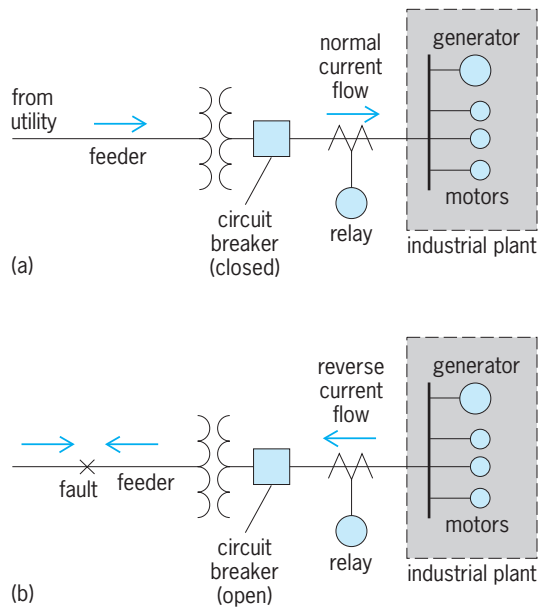


Fig. 5. Reverse-current protection. (a) Normal load conditions. (b) Internal fault condition; relay trips circuit breaker under reverse-current condition.

frequency. They operate to disconnect loads or to separate the system into areas so that the available generation equals the load until a balance is reestablished.

Reverse-current protection. This is provided when a change in the normal direction of current indicates an abnormal condition in the system. In an ac circuit, reverse current implies a phase shift of the current of nearly 180° from normal. This is actually a change in direction of power flow and can be directed by ac directional relays.

A common application of reverse-current protection is shown in **Fig. 5**. In this example, a utility supplies power to an industrial plant having some generation of its own. Under normal conditions, current flows from the utility to the plant (**Fig. 5a**). In the event of a fault occurring on the utility feeder (**Fig. 5b**) the current reverses direction and flows from the plant to the fault location. The relay operates and trips the circuit breaker, isolating the plant from the utility, thus preventing an excessive burden on the plant generator. Usually in these cases, the plant generator cannot carry the plant load, so that underfrequency relays are used to shed noncritical load. When the utility tie is restored, the shed loads then can be reconnected for full plant service.

Phase unbalance protection. This protection is used on feeders supplying motors where there is a possibility of one phase opening as a result of a fuse failure or a connector failure. One type of relay compares the current in one phase against the currents in the other phases. When the unbalance becomes too great, the relay operates. Another type monitors the three-phase bus voltages for unbalance. Reverse phases will operate this relay.

Reverse-phase-rotation protection. Where direction of rotation is important, electric motors must be protected against phase reversal. A reverse-phase-

rotation relay is applied to sense the phase rotation. This relay is a miniature three-phase motor with the same desired direction of rotation as the motor it is protecting. If the direction of rotation is correct, the relay will let the motor start. If incorrect, the sensing relay will prevent the motor starter from operating.

Thermal protection. Motors and generators are particularly subject to overheating due to overloading and mechanical friction. Excessive temperatures lead to deterioration of insulation and increased losses within the machine. Temperature-sensitive elements, located inside the machine, form part of a bridge circuit used to supply current to a relay. When a predetermined temperature is reached, the relay operates, initiating opening of a circuit breaker or sounding of an alarm.

J. Lewis Blackburn

Bibliography. J. L. Blackburn, *Protective Relaying: Principles and Applications*, 2d ed., 1997; W. A. Elmore (ed.), *Protective Relaying: Theories and Applications*, 1994; S. H. Horowitz, *Protective Relaying for Power Systems*, 2 vols., 1981, 1992; Institute of Electrical and Electronics Engineers, *IEEE Guides and Standards for Protective Relaying Systems*, rev. ed., 1991; National Association of Relay Manufacturers, *Engineers Relay Handbook*, 4th ed., 1992.

Electric rotating machinery

Any form of apparatus, having a rotating member, which generates, converts, transforms, or modifies electric power. Essentially all of the world's electric power is produced by rotating electrical generators, and about 70% of this energy is consumed in driving electric motors. Therefore, in a large sense, modern technology is built on the foundation of electric machinery. Electric machines are electromechanical energy converters; generators convert mechanical energy into electrical energy and motors convert electrical energy into mechanical energy. The basic laws of physics on which electric machinery is based were developed in the late nineteenth century.

Physical laws. It is possible to construct electric machines that are based on the energy stored in electric or magnetic fields. Those using electric fields are not efficient except in very small sizes, and as a result most machines are based on magnetic attraction or repulsion. The operation of electric machines can be explained in terms of two laws of physics: Ampère's force law and Faraday's law of induction. When a conductor is placed in a magnetic field, it experiences a force \mathbf{F} that is proportional to the current i and to the strength of the magnetic field \mathbf{B} . This force is in a direction at right angles to both the current and the field. This is known as Ampère's force law and can be expressed mathematically as Eq. (1), where \mathbf{l} is the length of the con-

$$\mathbf{F} = i\mathbf{l} \times \mathbf{B} \quad (1)$$

ductor. Ampère's law has numerous equivalent statements, from any of which the production of torque

can be explained. The second law is Faraday's law, which states that an electromotive force \mathbf{E} will be induced in a circuit when that circuit experiences a change in flux linkage λ . The magnitude of \mathbf{E} is proportional to the time rate of change of flux linkage. The electromotive force will be directed so that it produces currents that will oppose the change of flux linkage. Faraday's law can be written as Eq. (2).

$$\mathbf{E} = \frac{d\lambda}{dt} \quad (2)$$

See ELECTRIC FIELD; ELECTROMAGNETIC INDUCTION.

Principles of operation. An electric machine can be constructed on the principle that a magnet will attract a piece of permeable magnetic material such as iron or magnetic steel. In Fig. 1a, a pole structure is shown along with a magnetic block that is allowed to rotate. The magnetic block will experience a torque tending to rotate it counterclockwise to the vertical direction. This torque, called a reluctance or saliency torque, will be in the direction to minimize the reluctance of the magnetic circuit. In Fig. 1b, a winding is added to the rotor (the part which is allowed to rotate). In this case there is an additional torque on the rotor in the counterclockwise direction produced by the attraction of opposite poles. This torque will be approximately proportional to the sine of the angle θ . While the magnets in Fig. 1 are electromagnets, permanent magnets could be used with the same effect. See MAGNET; WINDINGS IN ELECTRIC MACHINERY.

In these examples, if the rotor were allowed to move under the influence of the magnetic forces, it would eventually come to rest at an equilibrium position, $\theta = 0$. Since most applications require continuous motion and constant torque, it is necessary to keep the angle between the rotor magnetic field and the stator magnetic field constant. Thus, in the above examples, the stator magnetic field must rotate ahead of the rotor. In a polyphase machine this is accomplished by distributing the winding so as to produce a rotating magnetic field. In the case of a three-phase machine the windings are distributed

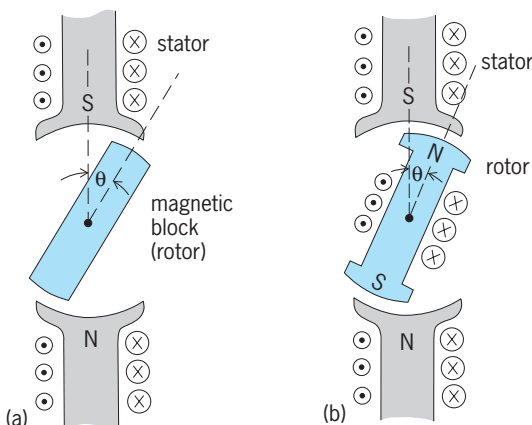


Fig. 1. Devices illustrating principles of electric machines. (a) Permeable rotor and stator with magnetic pole structure. (b) Device with magnetic pole structures on both stator and rotor.

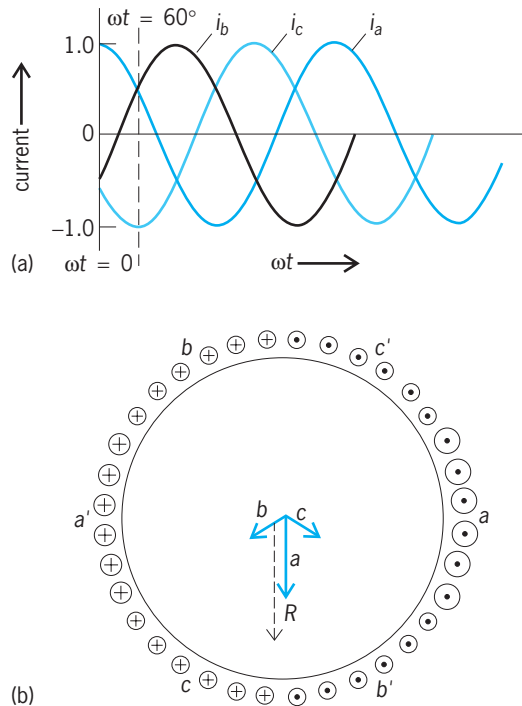


Fig. 2. Three-phase machine. (a) Three-phase balanced currents with angular frequency ω . (b) Distribution of currents and fields at $\omega t = 60^\circ$.

in space so that their magnetic axes lie 120° apart. During normal operation the currents in the three windings (i_a , i_b and i_c in phases a, b, and c; Fig. 2a) have the same magnitude but are 120° out of phase. (Such currents are called balanced currents.) At time $t = 0$ the current in phase a is maximum, and the currents in phases b and c are $-1/2$ of their peak values. The magnetic field R due to the combined three-phase system lies in the same direction that phase a current alone would produce (Fig. 2b). At a later time, $\omega t = 60^\circ$, where ω is the angular frequency of the currents, the current in phase c is at its negative peak, while the currents in phases a and b are $1/2$ of their peak value. The resultant field is along the (negative) magnetic axis of phase c, the total distribution having rotated 60° in space. It can be demonstrated that the fundamental spatial component of the field has constant magnitude and rotates uniformly at angular frequency ω . In this way it is possible to keep a constant angle between the stator and rotor magnetic fields and achieve constant torque. See ALTERNATING CURRENT.

In single-phase machines, there is no starting torque, and a device such as a capacitor or a short-circuiting ring is used to simulate the effect of a rotating distribution at starting. The single-phase winding with sinusoidal current produces two rotating pole distributions, one rotating in the forward direction and one in the backward direction. In single-phase alternating-current machines, once the rotor is rotating (in either direction) there is a net magnetic torque in the direction of rotation for motor operation and opposite the direction of rotation for generator operation.

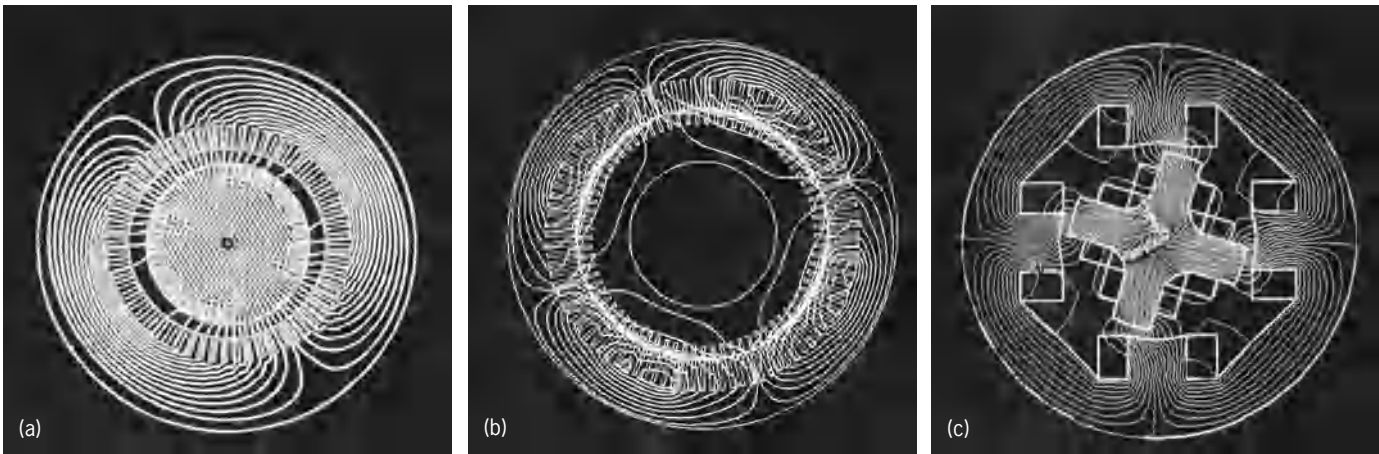


Fig. 3. Magnetic flux lines in rotating electric machines operating under load. (a) Three-phase, round-rotor synchronous generator (Rensselaer Polytechnic Institute). (b) Three-phase, squirrel-cage induction motor (Rodolpho Palma, Rensselaer Polytechnic Institute). (c) Direct-current motor (MAGSOFT Corp.).

Types. Although there are many variations, the three basic machine types are synchronous, induction, and direct-current machines. These machines may be used as motors or as generators, but the basic principles of operation remain the same.

Synchronous machines. The synchronous machine runs at a constant speed determined by the line frequency. There is an alternating-current winding (normally on the stator) and a direct-current winding (normally on the rotor). **Figure 3a** shows the magnetic field lines for a round-rotor synchronous generator operating under load. See ALTERNATING-CURRENT GENERATOR; ALTERNATING-CURRENT MOTOR; SYNCHRONOUS MOTOR.

Induction machines. The induction machine is another alternating-current machine which runs close to synchronous speed. The alternating-current winding of the stator is similar to that of the synchronous machine. The rotor may have an insulated winding (wound rotor) but more commonly consists of uninsulated bars embedded in a laminated structure and short-circuited at the end (squirrel cage). There is normally no voltage applied to the rotor. The voltages are produced by means of Faraday's law of induction. In an induction motor the stator-produced flux-density wave rotates slightly faster than the rotor during normal operation, and the flux linkages on the rotor therefore vary at low frequency. The rotor currents induced by these time-varying flux linkages produce a magnetic field distribution that rotates at the same speed as the stator-produced flux wave. The flux distribution of a squirrel-cage induction motor under load is shown in Fig. 3b. See INDUCTION MOTOR.

Direct-current machines. In a direct-current motor, direct current is applied to both the rotor and the stator. The stationary poles on the stator produce a stationary magnetic field distribution. Since the angle between the stator-produced poles and rotor-produced poles must remain constant, the direct-current machine uses a device known as a commutator which switches the current from one

rotor circuit to another so that the resulting field is stationary. A direct-current motor operating under load is depicted in Fig. 3c. See DIRECT-CURRENT GENERATOR; DIRECT-CURRENT MOTOR.

Additional systems. Besides the magnetic system that produces the energy conversion in rotating electric machines, a number of other systems operate in these devices.

Electric circuits. The windings of the machines may be any combination of series or parallel circuits, possibly with multiple turns in each circuit using one or many phases. The geometrical layout of these circuits is an important consideration in the electrical design because it affects the efficiency and performance of the device. See CIRCUIT (ELECTRICITY).

Thermal system. Heat is generated in the windings because of Joule heating and in the magnetic circuit by eddy currents and hysteresis. When the rotor is turning, there are also frictional and aerodynamic losses. The heat associated with all of these losses must be removed. Ventilation of the machine to limit the temperature to an acceptable level is often the limiting factor in machine design. The most commonly used cooling medium is air. Large machines may use high-pressure hydrogen or purified water as a coolant. See AIR COOLING; CORE LOSS; EDDY CURRENT; ELECTRICAL RESISTANCE.

Insulation system. Machines vary in voltage levels from a fraction of a volt to many thousands of volts. A dielectric material is used to coat the windings and insulate them from each other and from other parts of the machine. Deterioration of insulation with age and with increased temperature is one factor that limits the lifetime of machines. Mechanical forces on insulation within windings and between windings and other parts are also a limit on the design. See DIELECTRIC MATERIALS; ELECTRICAL INSULATION.

Mechanical system. The rotor on a rotating machine must be free to turn. This means that a system of bearings (often with a lubricating system) must be provided. The rotor windings must be supported

against the centrifugal forces due to rotation. The electromagnetic torque that is applied to the rotor is also applied, in an opposite direction, to the stator, which must be suitably supported. If vibrations are not limited to acceptable levels, excessive audible noise or fatigue failures may occur. See ACOUSTIC NOISE; ANTI-FRICTION BEARING; MECHANICAL VIBRATION; METAL, MECHANICAL PROPERTIES OF.

With these different modes of operation and design considerations, thousands of variations of motors and generators have been developed for industrial and domestic use. See GENERATOR; MOTOR.

Sheppard Salon

Bibliography. B. Chalmers and A. Williamson, *A. C. Machines, Electromagnetics and Design*, 1991; A. E. Fitzgerald, C. Kingsley, and S. D. Umans, *Electric Machinery*, 6th ed., 2003; I. L. Kosow, *Electric Machinery and Transformers*, 2d ed., 1991; S. A. Nasar and I. Boldea, *Electric Machines: Dynamics and Control*, 1992; M. G. Say, *Alternating Current Machines*, 5th ed., 1983; M. G. Say, *Direct Current Machines*, 2d ed., 1986.

Electric spark

A transient form of gaseous conduction. This type of discharge is difficult to define, and no universally accepted definition exists. It can perhaps best be thought of as the transition between two more or less stable forms of gaseous conduction. For example, the transitional breakdown which occurs in the transition from a glow to an arc discharge may be thought of as a spark.

Electric sparks play an important part in many physical effects. Usually these are harmful and undesirable effects, ranging from the gradual destruction of contacts in a conventional electrical switch to the large-scale havoc resulting from lightning discharges. Sometimes, however, the spark may be very useful. Examples are its function in the ignition system of an automobile, its use as an intense short-duration illumination source in high-speed photography, and its use as a source of excitation in spectroscopy. In the second case the spark may actually perform the function of the camera shutter, because its extinction renders the camera insensitive. See SPECTROSCOPY; STROBOSCOPIC PHOTOGRAPHY.

Mechanisms. The spark is probably the most complicated of all forms of gaseous conduction. It is exceedingly difficult to study, because it is a transient and because there are so many variables in the system. Some of these variables are the components of the gaseous medium, the gas pressure, the chemical form of the electrodes, the physical shape of the electrodes, the microscopic physical surface structure, the surface temperature, the electrode separation, the functional dependence of potential drop on time, and the presence or absence of external ionizing agents. One or more of these conditions may change from one spark to the next. Because

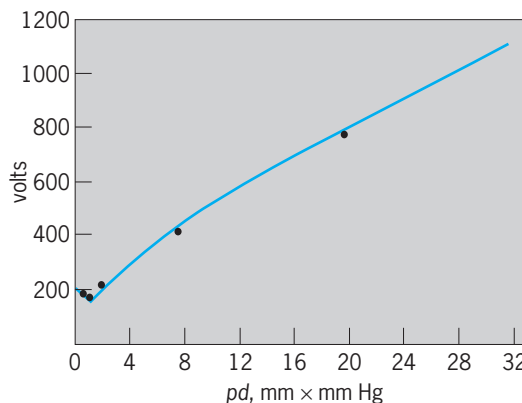


Fig. 1. Dependence of sparking potential on pd for a sodium cathode in hydrogen gas. (After L. B. Loeb and J. M. Meek, *The Mechanism of the Electric Spark*, Stanford University Press, 1941)

of the great complexity, it will be impossible to do more than touch on some of the main features in this article.

The dependence of breakdown, or sparking, potential on pressure p and electrode separation d is considered first. It was shown experimentally by F. Paschen and theoretically by J. S. Townsend that the sparking potential is a function of the product pd and not of p or d separately (Fig. 1). Further, there is a value of pd for which the sparking potential is a minimum. Thus, if it is desired to prevent sparking between two electrodes, the region may be evacuated or raised to a high pressure. The latter method is used in accelerators of the electrostatic generator variety. Here the entire apparatus is placed in a pressurized tank.

Qualitatively, one of the aspects of a spark is that the entire path between electrodes is ionized. It is the photon emission from recombination and decay of excited states which gives rise to the light from the spark. Further, if the spark leads to a stable conduction state, the cathode must be capable of supplying the needed secondary electrons, and the conduction state produced must permit the discharge of the interelectrode capacitance at the very minimum. See ARC DISCHARGE; ELECTRICAL CONDUCTION IN GASES; GLOW DISCHARGE.

In a consideration of the mechanism involved in the spark, the time required for the breakdown of the gas in a gap is an important element. L. B. Loeb pointed out that this time is often less than that required for an electron to traverse the gap completely. This implies that there must be some means of ionization present other than electron impact and that the velocity of propagation of this ionizing agent or mechanism must be much greater than the electron velocity. It seems definitely established that this additional method must be photoionization. In the intense electric field which is necessary for the spark, the initial electron will produce a heavy avalanche of cumulative ionization. Light resulting from the decay processes will produce ionization throughout the gas and electrons at the surfaces by the photoelectric

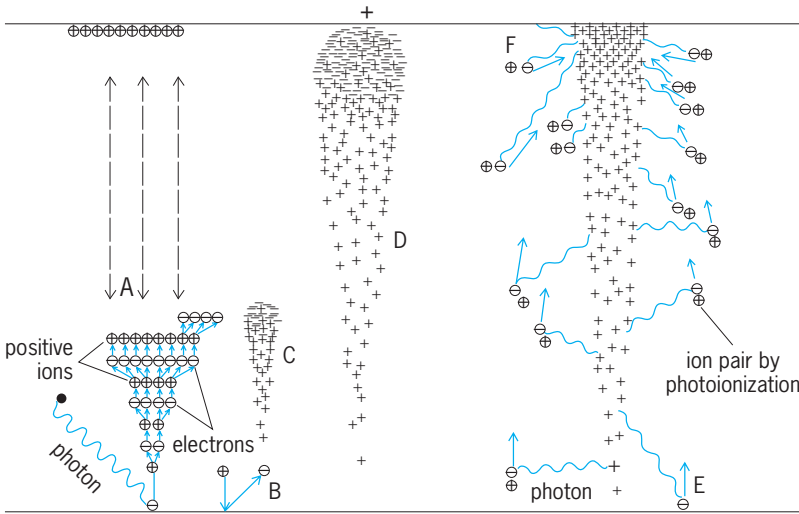


Fig. 2. Electron multiplication and avalanching during an electric spark discharge. (After L. B. Loeb and J. M. Meek, *The Mechanism of the Electric Spark*, Stanford University Press, 1941)

effect (Fig. 2). The electrons resulting from this will in turn produce further avalanches through the entire region, so that in a time of the order of 10^{-8} s the entire path becomes conducting. If the pressure is approximately atmospheric, the spark will be confined to a relatively narrow region, so that the conducting path, while not straight, will be a well-defined line. If the external circuit can supply the necessary current, the spark will result in an arc discharge. At lower pressure the path becomes more diffuse, and the discharge takes on either a glow or arc characteristic.

Figure 2 shows A, the electron multiplication of electrons by the cumulative ionization of a single electron liberated from the cathode by a photon; B, a secondary electron emitted from the cathode by a positively charged ion; C, the development and structure of an avalanche, with positively charged ions behind electrons at the tip; D, the avalanche crossing the gap and spreading by diffusion; and E, an older avalanche when electrons have disappeared into the anode. A positive space-charge boss appears on the cathode at F. Ion pairs out from the trail indicate the appearance of photoelectric ion pairs in the gas produced by photons from the avalanche. E shows a photoelectron from the surface of the cathode produced by the avalanche.

Theory. Mathematically, the theory of Townsend predicts that the current in a self-sustained discharge of the glow variety will follow Eq. (1), where I is

$$I = I_0 \frac{e^{\alpha x}}{1 - \gamma e^{\alpha x}} \quad (1)$$

the current with a given plate separation x , I_0 is the current when x approaches zero, and α and γ are constants associated with the Townsend coefficients. This equation represents the case where the electrode separation is varied while the ratio of electric field to pressure is held constant. The condi-

tion for a spark is that the denominator approach zero, which may be stated as in Eq. (2). Loeb indi-

$$\gamma = e^{-\alpha x} \quad (2)$$

cated that this criterion must be handled carefully. Townsend's equation really represents a steady-state situation, and it is here being used to explain a transient effect. If the processes which are involved are examined more carefully, it appears that there should be a dependence on I_0 as well.

Glenn H. Miller
Bibliography. M. N. Hirsh and H. J. Oskam (eds.), *Gaseous Electronics*, 1978; Y. P. Raizer and J. E. Allen, *Gas Discharge Physics*, 1997.

Electric susceptibility

A dimensionless parameter measuring the ease of polarization of a dielectric. In a vacuum, the electric flux density \mathbf{D} (measured in coulombs/m²) and electric field strength \mathbf{E} (volts/m) are related by Eq. (1) where ϵ_0 is the permittivity of free space,

$$\mathbf{D} = \epsilon_0 \mathbf{E} \quad (1)$$

having the approximate value 8.854×10^{-12} farad/m. In a dielectric material, polarization occurs, and the flux density has the value given by Eq. (2), where

$$\mathbf{D} = \epsilon_0 \epsilon_r \mathbf{E} = \epsilon_0 \mathbf{E} + \mathbf{P} \quad (2)$$

ϵ_r is the relative permittivity of the material and \mathbf{P} is the polarization flux density. This can be written as Eq. (3), where $\chi_e = \epsilon_r - 1$ is known as the electric

$$\mathbf{P} = \epsilon_0 (\epsilon_r - 1) \mathbf{E} = \epsilon_0 \chi_e \mathbf{E} \quad (3)$$

susceptibility of the dielectric material. It is a measure of that part of the relative permittivity which is due to the material itself.

The electric susceptibility can be related to the polarizability α by expressing the polarization in terms of molecular parameters. Thus Eqs. (4) hold, where

$$\begin{aligned} \mathbf{P} &= N \langle \mu \rangle_{\text{avg}} = N \alpha \mathbf{E}_L \\ \chi_e &= \frac{N \alpha \mathbf{E}_L}{\epsilon_0 \mathbf{E}} \end{aligned} \quad (4)$$

N is the number of molecules per unit volume, $\langle \mu \rangle_{\text{avg}}$ is their average dipole moment, and \mathbf{E}_L is the local electric field strength at a molecular site. At low concentrations, \mathbf{E}_L approaches \mathbf{E} , and the susceptibility is proportional to the concentration N . For a discussion of the properties and measurement of electric susceptibility see PERMITTIVITY; POLARIZATION OF DIELECTRICS. Robert D. Waldron; A. Earle Bailey

Electric switch

A device that makes, breaks, or changes the course of an electric circuit. Basically, an electric switch consists of two or more contacts mounted on an

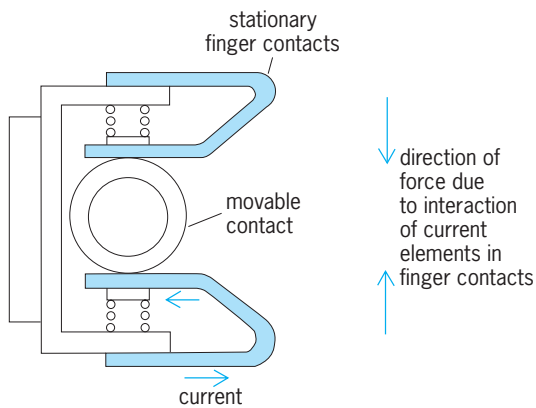


Fig. 1. With paired fingers for fixed contacts, short-circuit forces tend to hold contacts closed.

insulating structure and arranged so that they can be moved into and out of contact with each other by a suitable operating mechanism. See CIRCUIT (ELECTRICITY).

The term switch is usually used to denote only those devices intended to function when the circuit is energized or deenergized under normal manual operating conditions, as contrasted with circuit breakers, which have as one of their primary functions the interruption of short-circuit currents. Although there are hundreds of types of electric switches, their application can be broadly classified into two major categories: power and signal. See CIRCUIT BREAKER.

Power switches. In power applications, switches function to energize or deenergize an electric load. On the low end of the power scale, wall switches are used in homes and offices for turning lights on and off; dial and push-button switches control power to electric ranges, washing machines, and dishwashers. On the high end of the scale are load-break switches and disconnecting switches in power systems at the highest voltages (several hundred thousand volts).

For power applications, when closed, switches are required to carry a certain amount of continuous current without overheating, and in the open position they must provide enough insulation to isolate the circuit electrically. The latter function is particularly important in high-voltage circuits because it is the practice in the electrical industry to forbid people from working on electrical equipment unless it is isolated from the electrical supply system by a visible break in air. As an added precaution, a grounding switch is often used to connect the equipment to the ground before permitting any direct contact by a worker.

Load-break switches are required also to have the capability of interrupting the load circuit. Although this requirement is easily met in low-voltage and low-current applications, for high-voltage and high-current circuits, arc interrupters, similar to those used in circuit breakers, are needed.

In medium-voltage applications the most popular interrupter is the air magnetic type, in which

the arc is driven into an arc chute by the magnetic field produced by the load current in a blowout coil. The chute may be made either of ceramic materials or of organic materials such as plexiglass. For voltages higher than 15,000 V, gaseous interrupters are sometimes used. Special gases such as sulfur hexafluoride (SF_6) have been found to have good interrupting capability for load currents without a high-speed flow. For short-circuit currents, even with SF_6 , a high-speed flow is needed. Vacuum interrupters have been produced with high-voltage load-breaking duties. See BLOWOUT COIL; SHORT CIRCUIT.

Some load-break switches may also be required to have the capability of holding the contacts in the closed position during short-circuit conditions so that the contacts will be blown open by electromagnetic forces when the circuit breaker in the system interrupts the short-circuit current. For most contact configurations, the interaction between the short-circuit current and the magnetic field tends to force the contacts to part. Special configurations such as finger contacts are sometimes used to overcome this problem (Fig. 1).

Signal switches. For signal applications, switches are used to detect a specified situation that calls for some predetermined action in the electrical circuit. For example, thermostats detect temperature; when a certain limit is reached, contacts in the thermostat energize or deenergize another electrical switching device to control power flow. Centrifugal switches prevent overspeeds of motors. Limit switches prevent cranes or elevators from moving beyond preset positions. Many different types of switches perform such signaling purposes in communication systems, computers, and control systems for industrial processes. See ALTERNATING-CURRENT MOTOR; SWITCHING SYSTEMS (COMMUNICATIONS); THERMOSTAT.

Switches for signaling purposes are often required to have long life, high speed, and high reliability. Contaminants and dust must be prevented from interfering with the operation of the switch. For this purpose, switches are usually enclosed and are sometimes hermetically sealed.

Among the many different arrangements for contacts, the knife switch (Fig. 2), because of its

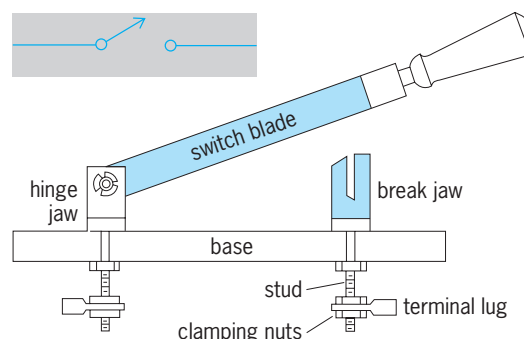


Fig. 2. Early knife switch with schematic symbol of single switch shown at top.

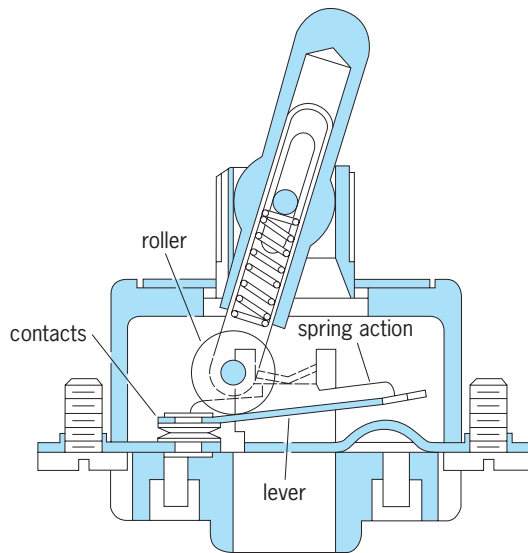


Fig. 3. Toggle mechanism actuates butt-contact switch.

early widespread use, is the basis for the elemental switch symbol. In the knife switch, a metal blade is hinged to a stationary jaw at one end and contacts a similar jaw at the other end.

In the leaf-spring switch, parallel strips of spring metal are sandwiched between blocks of insulation at one end and pushed into or out of contact at the other end. The sliding contact switch takes the form of a dial or drum with metal segments engaging contact fingers that wipe over the dial or drum surface. A butt-contact switch (Fig. 3) consists of one movable contact in the form of a disk or bar and one or two stationary contacts of similar form. In the mercury switch, a pool of mercury is made either to bridge the contacts or to separate from them (Fig. 4).

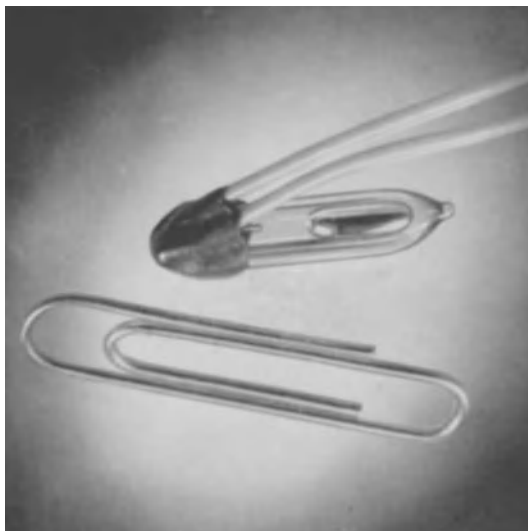


Fig. 4. Miniature mercury element may be used with a variety of tilting mechanisms to cause switch action. (Minneapolis-Honeywell Regulator Co.)

Switches frequently are composed of many single circuit elements, known as poles, all operated simultaneously or in a predetermined sequence by the same mechanism. Switches used in complex machines, such as computers, may have a large number of poles. Switches used in power circuits usually have from one to four poles, depending on the kind of circuit. Switches are often typed by the number of poles and referred to as single-pole or double-pole switches, and so on. It is common to express the number of possible switch positions per pole, such as a single-throw or double-throw switch.

Thomas H. Lee

Bibliography. D. G. Fink and H. W. Beaty (eds.), *Standard Handbook for Electrical Engineers*, 14th ed., 1999.

Electric transient

A temporary component of current and voltage in an electric circuit which has been disturbed. In ordinary circuit problems, a stabilized condition of the circuit is assumed and steady-state values of current and voltage are sufficient. However, it often becomes important to know what occurs during the transition period following a circuit disturbance until the steady-state condition is reached. Transients occur only in circuits containing inductance or capacitance. In general, transients accompany any change in the amount or form of energy stored in the circuit. After the following Introduction, transients in direct-current (dc) and alternating-current (ac) circuits with lumped elements are discussed. Transients of a more complex nature occur on distributed-parameter circuits, such as transmission lines. Such phenomena are elaborated upon in a separate section on transients in power systems. Finally, the modeling of transients in digital simulation is explained.

Introduction

An electric circuit or system under steady-state conditions of constant, or cyclic, applied voltages or currents is in a state of equilibrium. However, the circuit conditions of voltage, current, or frequency may change or be disturbed. Also, circuit elements may be switched in or out of the circuit. Any change of circuit condition or circuit elements causes a transient readjustment of voltages and currents from the initial state of equilibrium to the final state of equilibrium. In a sense the transient may be regarded as superimposed on the final steady state, so that Eq. (1) applies.

$$\left(\begin{array}{c} \text{Instantaneous} \\ \text{condition} \end{array} \right) = \left(\begin{array}{c} \text{final} \\ \text{condition} \end{array} \right) + \left(\begin{array}{c} \text{transient} \\ \text{terms} \end{array} \right) \quad (1)$$

Furthermore, since the instantaneous condition at the first instant of disturbance (time zero) must be

the initial condition, it may be described by Eq. (2).

$$\left(\begin{array}{c} \text{Initial} \\ \text{condition} \end{array} \right) = \left(\begin{array}{c} \text{final} \\ \text{condition} \end{array} \right) + \left(\begin{array}{c} \text{transient terms} \\ \text{at time zero} \end{array} \right) \quad (2)$$

A great deal of information may be obtained from these two “word equations” without recourse to mathematics. For example, if the weight on the end of a vertical spring is suddenly increased, its final displacement can be determined. Since the spring-weight combination is known to be an oscillating system, the amplitude of the transient oscillation follows from Eq. (2) as the difference between the initial and final displacements.

The nature of an electric transient is determined by three things: (1) the circuit or network itself—the interconnections of its elements and the circuit parameters (resistances R , inductances L , capacitances C , mutual inductances M); (2) the initial conditions of voltages, currents, charges, and flux linkages at the start of the transient; and (3) the nature of the disturbance which initiated the transient.

The circuit, or network, is usually defined by a diagram of connections showing all of the interconnections, junctions, meshes, circuit parameters, voltage and current sources and their polarities, and switches. Corresponding to the network a differential (or integral-differential) equation, or a set of such equations, may be written. These equations may also be written in terms of operational calculus, or as Laplace transforms, or in any other suitable mathematical equivalents. They are established in accordance with Kirchhoff’s laws: (I) The sum of the currents i at a junction is equal to zero: $\Sigma i = 0$. (II) The sum of the voltages e around a mesh is equal to zero: $\Sigma e = 0$. See KIRCHHOFF’S LAWS OF ELECTRIC CIRCUITS.

The voltage drops associated with a resistance, inductance, mutual inductance, and capacitance, respectively, are given by Eqs. (3). In applying

$$\begin{aligned} e_R &= Ri & e_L &= L \frac{di}{dt} \\ e_i &= M_{12} \frac{di_2}{dt} & e_c &= \frac{1}{C} \int i dt \end{aligned} \quad (3)$$

Kirchhoff’s law to a closed mesh, it is merely necessary (with due regard for signs and polarities) to equate the sum of all the voltage sources (such as generators and batteries) to the sum of all the voltage drops in the mesh, as in Eq. (4), where m refers

$$\sum_m e_m = \sum_m \sum_n \left(R_{mn} i_n + L_{mn} \frac{di_n}{dt} + \frac{1}{C_{mn}} \int i_n dt \right) \quad (4)$$

to any branch of the mesh in question, and n to any current causing a voltage drop in that branch (a branch may carry currents from meshes other than the mesh in question, or be mutually coupled with

other branches). Equation (4) formulates the differential equation of the circuit in terms of the voltages and currents. It is sometimes convenient to make use of charges q and fluxes ϕ by the substitutions of Eqs. (5).

$$i = \frac{dq}{dt} \quad \text{or} \quad q = \int i dt \quad \text{and} \quad N\phi = Li \quad (5)$$

Since Eq. (4) contains an integral, it is convenient to eliminate it either by differentiating once or by substituting Eq. (5). Then Eqs. (6) follow.

$$\begin{aligned} \sum_m \frac{de_m}{dt} &= \sum_m \sum_n \left(R_{mn} \frac{di_n}{dt} + L_{mn} \frac{d^2 i_n}{dt^2} + \frac{i_n}{C_{mn}} \right) \\ \sum_m e_m &= \sum_m \sum_n \left(R_{mn} \frac{dq_n}{dt} + L_{mn} \frac{d^2 q_n}{dt^2} + \frac{q_n}{C_{mn}} \right) \end{aligned} \quad (6)$$

See COUPLED CIRCUITS.

It is customary to factor out the current from Eq. (4) and write the equation in terms of a generalized impedance operator Z_{mn} defined by Eq. (7). Since an equation of this type can be

$$\begin{aligned} \sum_m e_m &= \sum_m \sum_n \left(R_{mn} + L_{mn} \frac{d}{dt} + \frac{1}{C_{mn}} \int dt \right) i_n \\ &= \sum_m \sum_n Z_{mn} i_n \end{aligned} \quad (7)$$

written for each of the N -meshes of the network, the totality of the differential equations for the entire network constitutes a system of simultaneous differential equations in the form of Eqs. (8).

$$\begin{aligned} e_1 &= Z_{11} i_1 + Z_{12} i_2 + \cdots + Z_{1N} i_N \\ e_2 &= Z_{21} i_1 + Z_{22} i_2 + \cdots + Z_{2N} i_N \\ &\dots\dots\dots \\ e_N &= Z_{N1} i_1 + Z_{N2} i_2 + \cdots + Z_{NN} i_N \end{aligned} \quad (8)$$

In general, the process of solution of such a set of equations for any current leads to a differential equation of order $2N$; there will be $2N$ integration constants associated with it, and these integration constants must be determined from the initial conditions at the first instant of the disturbance.

Equation (9) is the solution of an ordinary

$$\left(\begin{array}{c} \text{Complete} \\ \text{solution} \end{array} \right) = \left(\begin{array}{c} \text{particular} \\ \text{integral} \end{array} \right) + \left(\begin{array}{c} \text{complementary} \\ \text{solution} \end{array} \right) \quad (9)$$

differential equation with constant coefficients and is in two parts.

The particular integral depends on the form of the applied voltage, and represents the final steady-state solution.

The complementary solution is independent of the form of the applied voltages, but depends on the initial conditions in the circuit, and represents the transient terms. Thus Eq. (9) is the mathematical equivalent of Eq. (1).

The opening of a switch may be simulated by superimposing an equal and opposite current and thereby canceling the current through the switch. Likewise, the closing of a switch may be simulated by canceling the voltage across its terminals by an equal and opposite voltage. The effect of a cancellation current (or voltage) may be superimposed on the currents and voltages already existing in the system, as in Eq. (10).

$$\begin{pmatrix} \text{Resultant} \\ \text{voltages} \\ \text{and} \\ \text{currents} \end{pmatrix} = \begin{pmatrix} \text{those which} \\ \text{would exist if} \\ \text{switches were} \\ \text{not operated} \end{pmatrix} + \begin{pmatrix} \text{those due to} \\ \text{"cancellation"} \\ \text{voltages or} \\ \text{currents} \end{pmatrix} \quad (10)$$

In many practical cases transients are one of three types:

1. Single-energy transients, in which only one form of energy storage (electromagnetic or electrostatic) is present; the transient exhibits simple exponential decay from the initial to the final conditions.
2. Double-energy transients, in which both forms of energy storage are present; the transient is either aperiodic or a damped sinusoid.
3. Combination of 1 and 2.

DC Circuit with Lumped Elements

Transients are initiated in dc circuits either when a switch is closed on a dc voltage, or when a switch is opened and the detached circuit is permitted to discharge its energy sources. The assumption is made that the suddenly applied dc voltage is sustained at a constant value for an indefinite period after the switching operation, and time is counted from the instant, $t = 0$, at which the switch is operated. It is further usually assumed, where possible, that all circuit parameters (resistances R , inductances L , and capacitances C) are constant, so that the resulting differential equations are linear.

The energy sources in dc systems comprise dynamoelectric machines, mercury-arc rectifiers, vacuum tubes, and electric batteries. Energy W is stored in a circuit in inductances ($W = 1/2Li^2$) and capacitances ($W = 1/2Ce^2$). Energy is dissipated in circuit resistances ($W = Ri^2$).

Resistance and inductance in series. A lumped parameter circuit comprising a dc voltage source E , a resistance R , and an inductance L is shown in Fig. 1. When the switch S is closed, the voltage E is suddenly impressed on the series circuit. By Kirchhoff's second law the differential equation is given by Eq. (11).

$$E = Ri + L \frac{di}{dt} \quad (11)$$

Equating the right-hand side of Eq. (11) to zero, the complementary solution is found by assuming that Eq. (12), in which I and a are constants to be

$$i = Ie^{at} \quad (12)$$

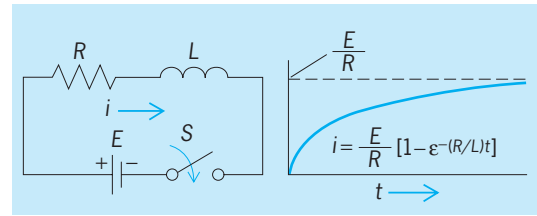


Fig. 1. Resistance and inductance in series.

determined, applies. Upon substitution in Eq. (11), Eq. (13) results, from which follows $a = -R/L$.

$$0 = RIe^{at} + aLIe^{at} = (R + aL)i \quad (13)$$

The particular integral is the final steady-state solution or $i(\infty) = E/R$. The complete solution, Eq. (14), is the sum of the complementary solu-

$$i = \frac{E}{R} + Ie^{-(R/L)t} \quad (14)$$

tion and the particular integral. But this expression contains the unknown integration constant I which must be found from the initial conditions. At $t = 0$, the instant at which the switch is closed, there can be no current in the circuit, since it is impossible to store energy in an inductance instantaneously. Therefore $i(0) = 0$ and when this is put in Eq. (14), Eq. (15) is obtained.

$$i(0) = 0 = \frac{E}{R} + I \text{ or } I = -\frac{E}{R} \quad (15)$$

The complete solution by Eq. (14) therefore is given by Eq. (16), which is in the form of Eq. (1).

$$i = \frac{E}{R} - \frac{E}{R}e^{-(R/L)t} \quad (16)$$

Here E/R is the final (steady-state) condition and $-(E/R)e^{-(R/L)t}$ is the transient term. Note also that at $t = 0$ this equation gives $i(0) = 0$, agreeing with Eq. (2). The graph of the transient is shown in Fig. 1, where the current starts at zero, rises exponentially, and approaches its steady-state value of E/R as t approaches infinity. See ELECTRICAL RESISTANCE; INDUCTANCE.

Resistance and capacitance in series. A resistance-capacitance series circuit is shown in Fig. 2. A residual charge q_0 is on the capacitor just prior to closing the switch S .

Equation (17) is the differential equation of the

$$E = Ri + \frac{1}{C} \int i dt \quad (17)$$

circuit. This equation may be solved as it is, or

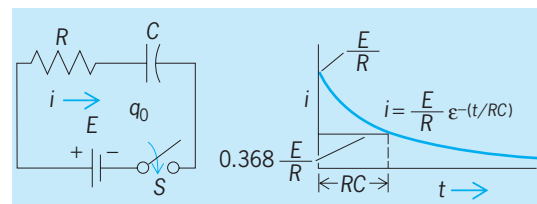


Fig. 2. Resistance and capacitance in series.

converted to a linear differential equation, Eq. (18),

$$0 = R \frac{di}{dt} + \frac{i}{C} \quad (18)$$

by differentiation or restated in terms of the instantaneous charge on the capacitor by substituting $i = dq/dt$, giving Eq. (19).

$$E = R \frac{dq}{dt} + \frac{q}{C} \quad (19)$$

The solution to this equation, following precisely the steps of the previous paragraph, is given by Eq. (20), in which Q is the integration constant. But

$$q = Qe^{-(t/RC)} + CE \quad (20)$$

initially, at $t = 0$, there was a residual charge q_0 on the capacitance; therefore the relationship given by Eq. (21) pertains, and upon the substitution of Eq. (21) in Eq. (20), Eq. (22) is obtained. This

$$q(0) = q_0 = Q + CE \text{ or } Q = q_0 - CE \quad (21)$$

$$q = CE - (CE - q_0)e^{-(t/RC)} \quad (22)$$

equation is in the form of Eq. (1), in which CE is the final (steady-state) condition and the other term is the transient. Initially, at $t = 0$, $q = q_0$ in agreement with Eq. (2). The current is given by Eq. (23).

$$\begin{aligned} i &= \frac{dq}{dt} = \frac{CE - q_0}{RC} e^{-(t/RC)} \\ &= \left(\frac{E}{R} - \frac{q_0}{RC} \right) e^{-(t/RC)} \end{aligned} \quad (23)$$

Equations (22) and (23) have been plotted in Fig. 2 for the case $q_0 = 0$. See CAPACITANCE.

Resistance, inductance, and capacitance in series.

This circuit is shown in Fig. 3. The capacitor is assumed to have an initial charge q_0 , or an initial voltage $V = q_0/C$. The differential equation of the circuit is given by Eq. (24).

$$E = Ri + L \frac{di}{dt} + \frac{1}{C} \int i dt \quad (24)$$

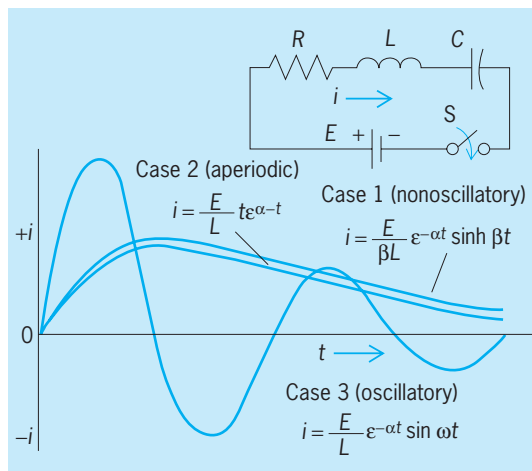


Fig. 3. Resistance, inductance, and capacitance in series for three cases.

Differentiating once to clear the integral, Eq. (25)

$$0 = R \frac{di}{dt} + L \frac{d^2i}{dt^2} + \frac{i}{C} \quad (25)$$

is obtained. Equation (25) is a second-order linear differential equation with constant coefficients; since it is equal to zero, it will possess only a complementary solution. Equation (26) is assumed, which, upon substitution in Eq. (25), yields Eq. (27).

$$i = Ae^{at} \quad (26)$$

$$0 = aRAe^{at} + a^2LAe^{at} + (A/C)e^{at} \quad (27)$$

Canceling the constant A and the exponential, there results a quadratic in a whose solution gives the two possible values shown in Eqs. (28).

$$\begin{aligned} a_1 &= \frac{-RC + \sqrt{R^2C^2 - 4LC}}{2LC} \\ a_2 &= \frac{-RC - \sqrt{R^2C^2 - 4LC}}{2LC} \end{aligned} \quad (28)$$

Associating the integration constant A_1 with a_1 and A_2 with a_2 , the solution takes the form of Eq. (29).

$$i = A_1e^{a_1t} + A_2e^{a_2t} \quad (29)$$

The voltage across the capacitor, Eq. (30), is obtained from Eq. (24).

$$\begin{aligned} e_c &= \frac{1}{C} \int i dt = E - Ri - L \frac{di}{dt} \\ &= E - (R + a_1L)A_1e^{a_1t} - (R + a_2L)A_2e^{a_2t} \end{aligned} \quad (30)$$

Initially, at $t = 0$, the current must be zero because of the inductance, and the capacitor voltage is $e_c = V$. By the first of these conditions, $i(0) = 0$, in Eq. (29) it is seen that $A_2 = -A_1$. And by the second condition, $e_c(0) = V$, in Eq. (30), Eq. (31) is obtained.

$$\begin{aligned} V &= E - (R + a_1L)A_1 - (R + a_2L)A_2 \\ &= E - (a_1 - a_2)LA_1 \end{aligned} \quad (31)$$

$$A_1 = \frac{E - V}{(a_1 - a_2)L} = \frac{C(E - V)}{\sqrt{R^2C^2 - 4LC}}$$

Then Eq. (32) is the complete solution.

$$i = \frac{C(E - V)}{\sqrt{R^2C^2 - 4LC}} (e^{a_1t} - e^{a_2t}) \quad (32)$$

There are three special cases of this solution, depending on the nature of the radical in Eq. (32).

Nonoscillatory case. $R^2C > 4L$. In this case the radical is positive, and the exponents a_1 and a_2 are real and negative, $a_1 = -\alpha + \beta$, $a_2 = -\alpha - \beta$, where $\alpha = R/2L$ and $\beta = \sqrt{R^2C^2 - 4LC}/2LC$. Then i is given by Eq. (33), which is shown in Fig. 3, where q_0 is

$$\begin{aligned} i &= \frac{C(E - V)}{\sqrt{R^2C^2 - 4LC}} e^{-\alpha t} (e^{+\beta t} - e^{-\beta t}) \\ &= \frac{2C(E - V)}{\sqrt{R^2C^2 - 4LC}} e^{-\alpha t} \sinh \beta t \end{aligned} \quad (33)$$

assumed to be 0 and $V = 0$.

Aperiodic case. $R^2C = 4L$. In this case $a_1 = a_2 = -R/2L$, and the radical in the denominator of Eq. (32) is zero, as is the numerator. The indeterminate is easily evaluated from Eq. (33) upon letting $\beta \rightarrow 0$, thus giving Eq. (34). This is the aperiodic or

$$\begin{aligned} i &= \frac{2C(E - V)}{2LC\beta} e^{-\alpha t} \sinh \beta t \Big|_{\beta \rightarrow 0} \\ &= \frac{E - V}{L} e^{-\alpha t} \end{aligned} \quad (34)$$

critical case, and is illustrated in Fig. 3.

Oscillatory case. $R^2C < 4L$. In this case the radical in Eq. (28) becomes imaginary and the exponents take the form of Eqs. (35), and Eq. (32) becomes Eq. (36). This is a damped oscillation and is illustrated in Fig. 3.

$$\begin{aligned} a_1 &= \frac{-R}{2L} + j \frac{\sqrt{4LC - R^2C^2}}{2LC} = -\alpha + j\omega \\ a_2 &= \frac{-R}{2L} - j \frac{\sqrt{4LC - R^2C^2}}{2LC} = -\alpha - j\omega \\ i &= \frac{C(E - V)}{2j\omega LC} e^{-\alpha t} (\epsilon^{j\omega t} - \epsilon^{-j\omega t}) \\ &= \frac{E - V}{\omega L} \epsilon^{-\alpha t} \sin \omega t \end{aligned} \quad (35)$$

AC Circuit with Lumped Elements

Alternating-current transients differ from direct-current transients in two important respects: (1) The final condition, or steady state, is an alternating or cyclic one, and (2) the amplitudes of the transient terms depend on the point on the ac applied voltage wave at which the transient is initiated, and can therefore have many different values, or even change sign.

In the section on dc transients, solutions were carried out for RL , RC , and RLC circuits switched onto a voltage source. In the present section, the solution will be carried out in detail for the RLC circuit only; the others will be regarded as special cases of this general solution by putting $C = \infty$ and $L = 0$, respectively.

Resistance, inductance, and capacitance in series.

Consider the circuit of Fig. 4 in which an RLC circuit is suddenly switched onto an ac voltage $e = E \sin(\omega t + \gamma)$ at an electrical angle γ displaced from $\omega t = 0$. Assume that a current I is flowing in the circuit and a voltage V is across the capacitance at the instant the switch is closed on the alternator. The voltage equation and the initial conditions then are given by Eqs. (37) and (38).

$$E \sin(\omega t + \gamma) = Ri + L \frac{di}{dt} + \frac{1}{C} \int i dt \quad (37)$$

$$i = I \quad \text{and} \quad e_c = V \quad \text{at} \quad t = 0 \quad (38)$$

Differentiating Eq. (37) once to clear it of the integral gives Eq. (39). The complementary solution

$$\omega E \cos(\omega t + \gamma) = L \frac{d^2i}{dt^2} + R \frac{di}{dt} + \frac{i}{C} \quad (39)$$

of this equation is the same as in the dc case of

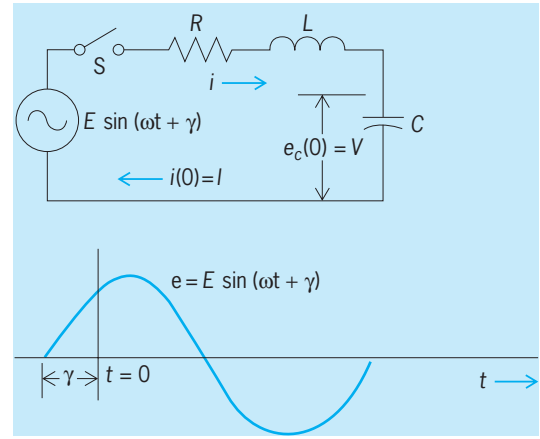


Fig. 4. RLC circuit with alternating current source.

Eq. (25), that is, Eq. (29), where a_1 and a_2 are given in Eq. (28).

The particular integral of Eq. (39) is its final steady-state current. This is most easily obtained as Eq. (40), the ordinary ac solution for the current, to which Eqs. (41) and (42) apply.

$$i_{ac} = \frac{E}{Z} \sin(\omega t + \gamma - \theta) \quad (40)$$

$$Z = \sqrt{R^2 + (\omega L - 1/\omega C)^2} \quad (41)$$

$$\tan \theta = \frac{\omega L - 1/\omega C}{R} = \frac{\omega^2 LC - 1}{\omega CR} \quad (42)$$

The complete solution then is given by Eq. (43),

$$i = \frac{E}{Z} \sin(\omega t + \gamma - \theta) + A_1 \epsilon^{a_1 t} + A_2 \epsilon^{a_2 t} \quad (43)$$

in which A_1 and A_2 are integration constants and a_1 and a_2 are given in Eq. (28).

The capacitor voltage is given by Eq. (44).

$$\begin{aligned} e_c &= \frac{1}{C} \int i dt = -\frac{E}{\omega CZ} \cos(\omega t + \gamma - \theta) \\ &\quad + \frac{A_1}{Ca_1} \epsilon^{a_1 t} + \frac{A_2}{Ca_2} \epsilon^{a_2 t} \end{aligned} \quad (44)$$

Now applying the initial conditions of Eq. (38) at $t = 0$, there results from Eqs. (43) and (44), respectively, Eqs. (45) and (46).

$$i(0) = I = \frac{E}{Z} \sin(\gamma - \theta) + A_1 + A_2 \quad (45)$$

$$\begin{aligned} e_c(0) &= V \\ &= \frac{-E}{\omega CZ} \cos(\gamma - \theta) + \frac{A_1}{Ca_1} + \frac{A_2}{Ca_2} \end{aligned} \quad (46)$$

Solving Eqs. (45) and (46) simultaneously for A_1 and A_2 and using Eq. (28), Eqs. (47)-(49) result.

$$\begin{aligned} A_1 &= \frac{1}{\sqrt{R^2 C^2 - 4LC}} \left\{ \frac{I}{a_2} - CV \right. \\ &\quad \left. - \frac{E}{Z} \left[\frac{1}{\omega} \cos(\gamma - \theta) + \frac{1}{a_2} \sin(\gamma - \theta) \right] \right\} \end{aligned} \quad (47)$$

$$A_2 = \frac{-1}{\sqrt{R^2C^2 - 4LC}} \left\{ \frac{I}{a_1} - CV - \frac{E}{Z} \left[\frac{1}{\omega} \cos(\gamma - \theta) + \frac{1}{a_1} \sin(\gamma - \theta) \right] \right\} \quad (48)$$

$$a_1 = \frac{-R}{2L} + \sqrt{\left(\frac{R}{2L}\right)^2 - \frac{1}{LC}} \quad (49)$$

$$a_2 = \frac{-R}{2L} - \sqrt{\left(\frac{R}{2L}\right)^2 - \frac{1}{LC}}$$

These values, together with Eqs. (41) and (42), inserted in Eqs. (43) and (44) give the solution in a form suitable for the critically damped case. However, if $R^2C < 4L$, then the radicals in Eqs. (47)-(49) become imaginary and these expressions become the complex numbers given in Eqs. (50), whereupon

$$\begin{aligned} A_1 &= \frac{1}{2}(M + jN) & A_2 &= \frac{1}{2}(M - jN) \\ a_1 &= -\alpha + j\omega & a_2 &= -\alpha - j\omega \end{aligned} \quad (50)$$

Eq. (43) takes the damped oscillatory form given in Eq. (51).

$$i = \frac{E}{Z} \sin(\omega t + \gamma - \theta) + \epsilon^{-\alpha t} (M \cos \omega t - N \sin \omega t) \quad (51)$$

It is evident from Eq. (43) or Eq. (51) that the final steady-state value, after the transient has died out, is the alternating current of Eq. (40). It is also clear from Eqs. (47) and (48) that the amplitudes of the transient terms depend upon the angle γ on the ac applied voltage wave at $t = 0$.

Resistance and inductance in series. This may be regarded as a special case of Eq. (43), in which $C = \infty$ (the capacitance short-circuited) and $V = 0$. Consequently, the solution is given by Eq. (52). Thus

$$i = \frac{E}{Z} \sin(\omega t + \gamma - \theta) + \left[I - \frac{E}{Z} \sin(\gamma - \theta) \right] \epsilon^{-(R/L)t} \quad (52)$$

the transient starting from an initial value I decays exponentially to its final ac steady-state value.

There will be no transient if the switch is closed at the angle γ on the voltage wave such that Eq. (53) applies.

$$I = \frac{E}{Z} \sin(\gamma - \theta) \quad (53)$$

Resistance and capacitance in series. This may be regarded as a special case of Eq. (43), in which $L = 0$ and $I = 0$. Under these conditions the solution is given by Eq. (54) and the capacitor voltage is,

$$i = \frac{E}{Z} \sin(\omega t + \gamma - \theta) - \left[\frac{V}{R} + \frac{E}{\omega CRZ} \cos(\gamma - \theta) \right] \epsilon^{-t/RC} \quad (54)$$

by Eq. (44), defined by Eq. (55). Thus the transient,

$$e_c = \frac{-E}{\omega CZ} \cos(\omega t + \gamma - \theta) + \left[V + \frac{E}{\omega CZ} \cos(\gamma - \theta) \right] \epsilon^{-t/RC} \quad (55)$$

starting with a voltage V , decays exponentially to its final ac steady-state value. If the switch is closed at an angle γ on the voltage wave such that Eq. (56)

$$V = -\frac{E}{\omega CZ} \cos(\gamma - \theta) \quad (56)$$

applies, there will be no transient. Loyal V. Bewley

Power Systems with Lumped and Distributed Elements

In an electric power system, a transient is a temporary disturbance leading to changes in the amounts of energy stored in the capacitances and inductances, causing variations of power system voltage and current waveforms. Such disturbances may also be referred to as electromagnetic transients to indicate that oscillatory variations are caused through the alternating conversion of electric energy of capacitors into magnetic energy of inductors. The frequency range typically extends from several hertz to several megahertz.

Classification. A classification of electric transients by cause is given in Fig. 5. Lightning strikes induce currents that cause electric transients in the form of traveling waves to propagate in forward and backward directions on lines.

Equipment switching events change the network topology, that is, the network configuration, and force redistribution of stored energy. Switching can also result in traveling waves on lines or cause energization and deenergization of lumped capacitances and inductances.

Conceptual analysis. The following study of the transient recovery voltage (TRV) serves as introduction to conceptual analysis of electric transients involving both harmonic oscillations and traveling waves. Figure 6a shows a line-to-ground short circuit on a transmission line. The fault current is supplied via a transformer and needs to be interrupted to clear the fault. The opening of a circuit breaker triggers electric transients due to a change in the

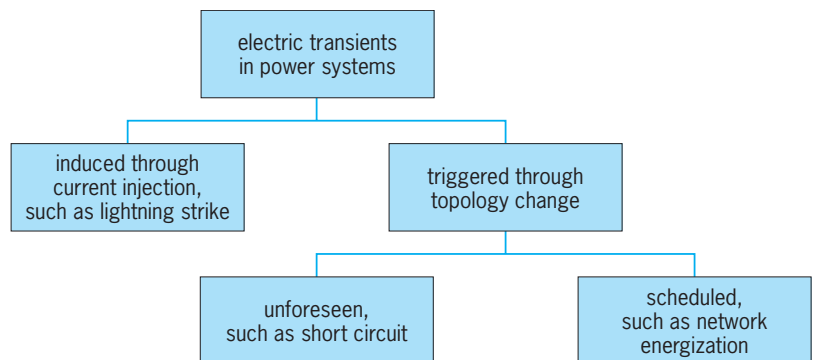


Fig. 5. Causes of electric transients in power systems.

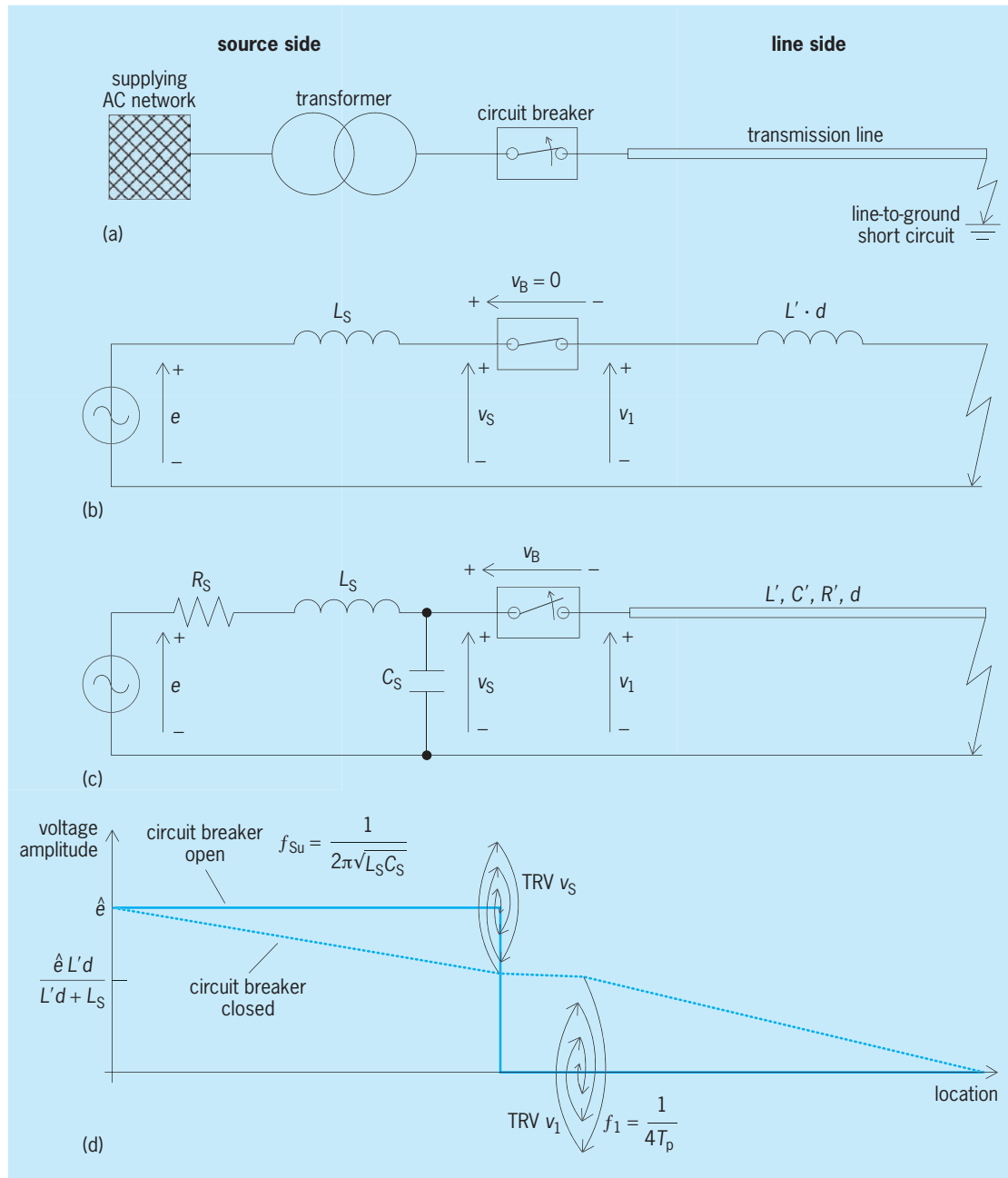


Fig. 6. Analysis of electric transients involving both harmonic oscillation and traveling waves. (a) Configuration of power system for study of transient recovery voltage (TRV). (b) Model for evaluating steady state before breaker opening. (c) Model for evaluating electric transients upon breaker opening. (d) Qualitative assessment of transients described by voltages.

topology of the system. In order to focus on the behavior of transients, only single-pole networks will be considered. (Multipole networks can be represented as single-pole equivalents.) See CIRCUIT BREAKER; SHORT CIRCUIT; TRANSMISSION LINES.

Before the circuit breaker opens, the voltage v_B across it is zero, as indicated in Fig. 6b, and hence $v_S = v_1$. The circuit of Fig. 6b models the configuration of Fig. 6a at source frequency f_c , which is 50 or 60 Hz in utility power systems depending on regional preferences. The purpose of this model is to estimate the voltage profile while the fault is still in place and the circuit breaker is closed. This allows the determination of the steady-state voltages

giving initial conditions from which the transients will start. The line has inductance per unit length L' , and the distance to the fault is d . The source side is modeled through an equivalent ac voltage source e behind a lumped inductive element L_S , which represents inductances of the transformer and the supplying ac network. Since the circuit breaker is to open at the instant when the fault current crosses zero to reduce stress, the voltages will reach their peak values at this instant. This is due to the inductive behavior of the line and the resulting phase shift of 90° between the fault current and the voltages v_S and v_1 . The voltage profile with the circuit breaker closed is shown by a line in Fig. 6d, which gives the

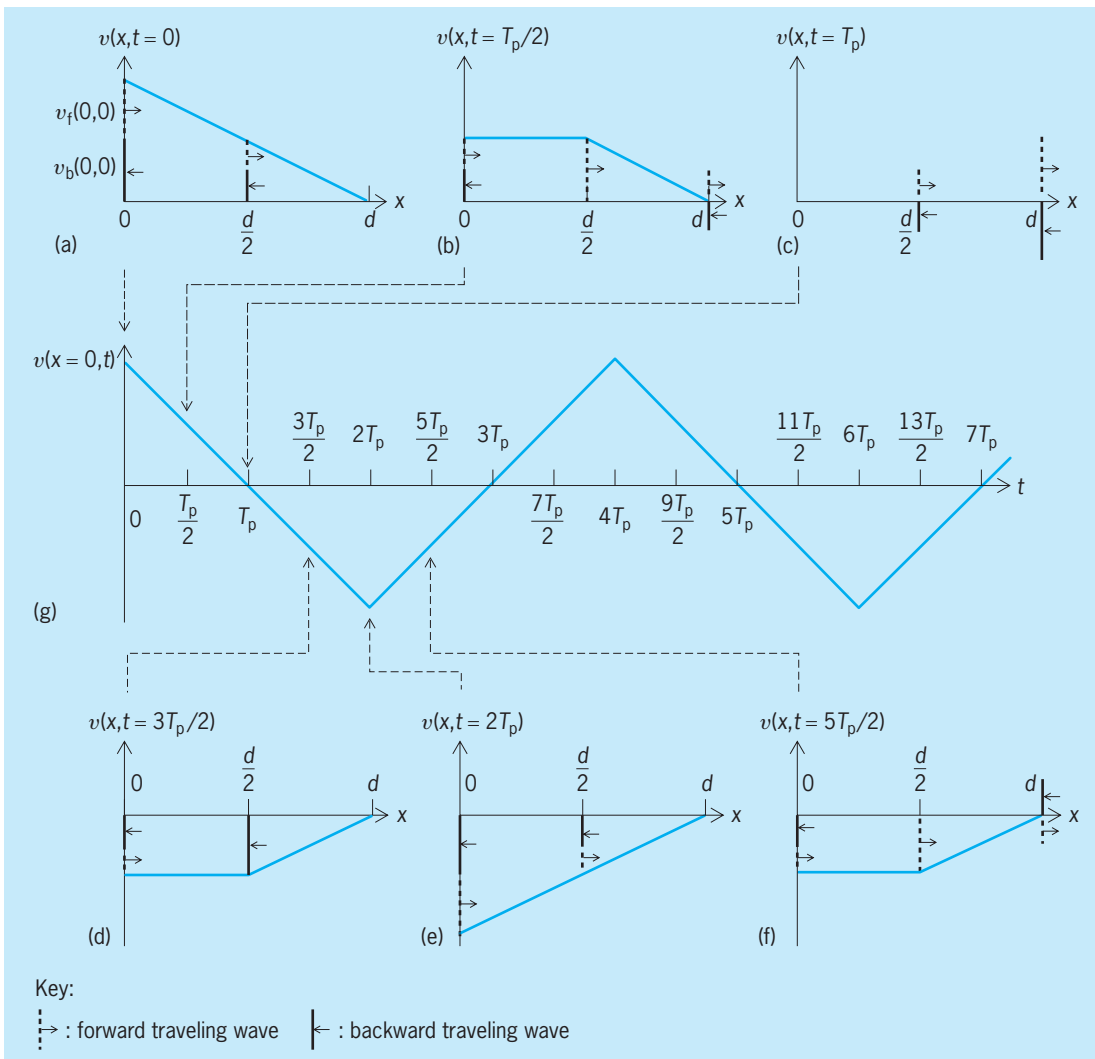


Fig. 7. Graphical interpretation of electromagnetic transients on transmission line. (a–f) Voltage profiles along line. (g) Derived voltage at beginning of line as a function of time.

amplitude as a function of location. The amplitude of the voltage source is denoted by \hat{e} . At the location of the short circuit, the voltage is zero. The amplitudes of v_s and v_l are obtained with the voltage divider expression given on the ordinate. Figure 6d also shows the steady-state voltage profile that is reached when the circuit breaker is open and the fault current is zero.

To assess the electric transients between the two steady states, the steady-state model of Fig. 6b is not adequate. To represent the redistribution of energy, both inductors and capacitors need to be modeled as indicated in Fig. 6c. The lumped capacitive element C_s represents the capacitance of transformer and bus bars on the source side, while C' gives the capacitance per unit length of the line. Losses are accounted for through R_s and R' on the source and line sides, respectively.

As the circuit breaker opens, the source and line sides become decoupled. The TRV v_s must first increase, and describes a damped harmonic oscillation around the new steady state (Fig. 6d). If R_s were neglected, the undamped frequency of this harmonic

oscillation on the source side would be given by Eq. (57), but because of the presence of R_s the ob-

$$f_{su} = 1/(2\pi\sqrt{L_s C_s}) \quad (57)$$

served value is slightly lower. The exact trajectory of the source-side TRV can be calculated with the techniques developed for circuits with lumped parameters. At the entrance to the line, the TRV v_l will first decrease.

The waveform of the TRV v_l is influenced by traveling waves which are investigated by means of Fig. 7 for the case of a lossless line. In Fig. 7a the initial voltage profile along the line is redrawn from Fig. 6d. The entrance to the line is at $x = 0$, and the circuit breaker is opened at $t = 0$. The waves propagate in all possible directions. Since the line is one-dimensional, there exist a forward-traveling wave marked with subscript f and arrow in direction of increasing x , and a backward-traveling wave marked with subscript b and arrow in direction of decreasing x . At any point on the line, the voltage is given as the sum of forward and backward traveling waves.

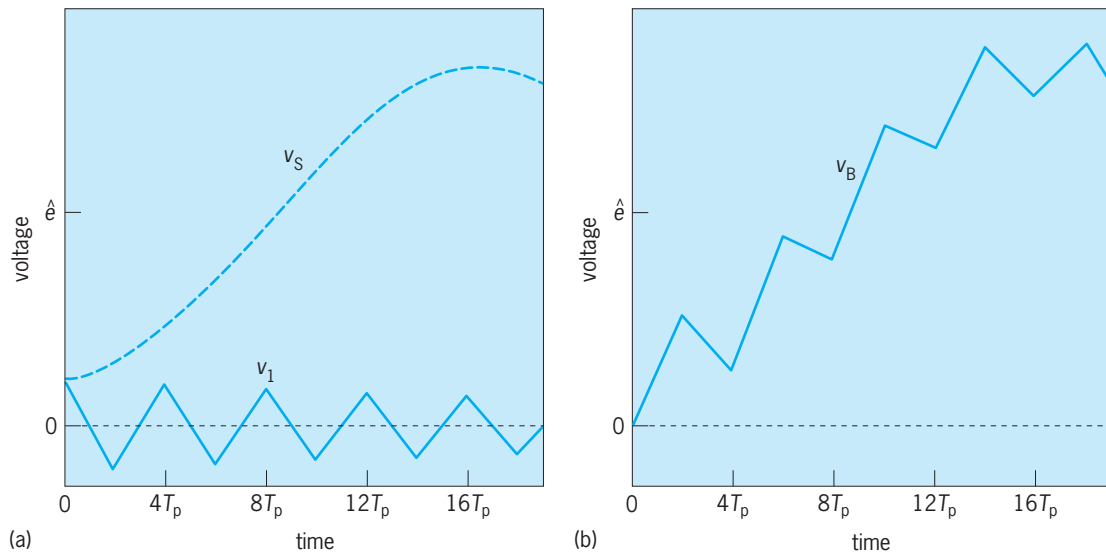


Fig. 8. Transient recovery voltages (TRVs). (a) TRVs on source side (v_s) and line side (v_l). (b) TRV across circuit breaker (v_b).

The propagation time that the waves need to travel from one end of the line to the other end is denoted by T_p .

Figure 7b illustrates how the original voltage profile of $t = 0$ has changed as $t = T_p/2$ is reached. The forward traveling wave of $x = 0$ at $t = 0$ traveled to $x = d/2$. The forward traveling wave of $x = d/2$ at $t = 0$ has moved by the same amount from $x = d/2$ to $x = d$. Likewise, the backward traveling wave of $x = d/2$ at $t = 0$ moved in opposite direction to reach $x = 0$. The waves are reflected at both ends. Since there is a short circuit at $x = d$, the voltage at $x = d$ must be zero at any time, that is, $v(d, t) = 0$, and therefore $v_f(d, t) = -v_b(d, t)$. At $x = 0$, there is an open circuit where the voltage reaches a maximum and $v_f(0, t) = v_b(0, t)$. Therefore, at the entrance of the line $v(0, t) = 2 v_b(0, t)$.

Figure 7c to 7f shows the further evolution of the voltage profile in increments of $T_p/2$. Figure 7g shows the voltage $v_l = v(0, t)$ at the entrance to the line as a function of time derived from the changing voltage profiles. The voltage is of sawtooth shape with a period of $4T_p$, that is, the frequency of oscillation of v_l , is $f_l = 1/(4T_p)$. On overhead lines, where the insulating material is air, the waves propagate at the speed of light $c_0 = 3 \times 10^8$ m/s (1.68×10^5 mi/s). At a distance to the fault of $d = 1$ km (0.62 mi), this leads to $f_l = c_0/(4d) = 75$ kHz. Such a fault is also known as a kilometric fault or short-line fault. For longer distances, the frequency is lower, but the voltage peak is higher as shown by the voltage divider in Fig. 6d. Thus, the stress on the circuit breaker due to the rate of voltage rise, which is directly related to the frequency, decreases, while the stress due to the reached voltage amplitude increases. For shorter distances, the opposite is the case. In the case of the kilometric fault, the combination of the stress due to rate of rise and amplitude of the voltage has been found to be particularly severe. Appropriate sizing of power system equip-

ment depends on knowledge of the properties of electric transients.

In Fig. 8a, the TRV waveforms on source and line sides are drawn in the same diagram. The source-side TRV (v_s) describes a harmonic oscillation in the form of a damped sinusoid which, as indicated in Fig. 6d, starts from its lowest point given by the voltage divider. In the steady state, after the transient has died out, v_s approximates the ac voltage source e . Since the frequency of the source-side TRV is much lower than the frequency of the line-side TRV, only the start of the source-side TRV is shown. On the line side, damping is taken into account (which is not done in Fig. 7), as evidenced by the decaying peak values of the sawtooth. The resulting TRV $v_b = v_s - v_l$ across the breaker is shown in Fig. 8b.

Digital Simulation

The conceptual analysis of the preceding section allows understanding of causes and consequences of electric transients. To the study specific cases and in particular those that involve complex circuits and systems, the use of digital simulation is required. Among the popular simulators are the Electromagnetic Transients Program (EMTP) and others of its type. Key to the success of EMTP is the companion modeling of network elements, where the differential equations characterizing a network element are approximated by a companion model consisting of a resistive circuit. A library of companion element models can so be established as shown in the table and derived hereafter.

Companion modeling for lumped-parameter elements. In the continuous time domain, the lumped inductance is described through differential equation (58).

$$v(t) = L \frac{di(t)}{dt} \quad (58)$$

For digital simulation, Eq. (58) must be discretized. Using the trapezoidal integration yields

Element type	Representation in continuous-time domain	Companion element model in discrete-time domain
Lumped resistance	$v(t) = i(t)R$	$i(k) = G_R v(k)$ $G_R = 1/R$
Lumped inductance	$v(t) = L \frac{di(t)}{dt}$	$i(k) = G_L v(k) + \eta_L(k)$ $G_L = \frac{\tau}{2L}$ $\eta_L(k) = i(k-1) + G_L v(k-1)$
Lumped capacitance	$i(t) = C \frac{dv(t)}{dt}$	$i(k) = G_C v(k) + \eta_C(k)$ $G_C = \frac{2C}{\tau}$ $\eta_C(k) = -i(k-1) - G_C v(k-1)$
Distributed parameter, lossless line	$\frac{\partial v(x,t)}{\partial x} = -L' \frac{\partial i(x,t)}{\partial t}$ $\frac{\partial i(x,t)}{\partial x} = -C' \frac{\partial v(x,t)}{\partial t}$	$i_1(k) = G_1 v_1(k) + \eta_{11}(k)$ $i_2(k) = G_1 v_2(k) + \eta_{12}(k)$ $\eta_{11}(k) = -i_2(k-\kappa) - G_1 v_2(k-\kappa)$ $\eta_{12}(k) = -i_1(k-\kappa) - G_1 v_1(k-\kappa)$ $G_1 = \sqrt{C'/L'}$ $\tau = \frac{T_p}{\tau}$ $T_p = l\sqrt{L'C'}$

Eq. (59), where τ is the time step size separ-

$$\frac{1}{2}[v(k) + v(k-1)] = L \frac{[i(k) - i(k-1)]}{\tau} \quad (59)$$

ating successive discrete time points at which solutions are calculated and k is the time step counter. The higher the frequencies of the electric transients to be studied, the smaller must be the time step size to allow for accurate tracking of the waveforms. Rearranging Eq. (59) yields the companion model of Eq. (60) [also shown

$$i(k) = G_L v(k) + \eta_L(k) \quad (60)$$

in the table], where G_L and $\eta_L(k)$ are given by

Eqs. (61) and (62).

$$G_L = \frac{\tau}{2L} \quad (61)$$

$$\eta_L(k) = i(k-1) + \frac{\tau}{2L} v(k-1) \quad (62)$$

The companion model simulates the inductance through a conductance and a current source, which provides the initialization for the solution at time step k based on the results of prior time step $k-1$. The same procedure is applied to derive the companion model for the capacitance. In the case of the resistance, no discretization is necessary.

Companion modeling for distributed-parameter elements. A lossless transmission line of length l can be

considered as being composed of an infinite number of incremental segments with incremental length dx for which Eqs. (63) and (64) apply, as indicated in the table.

$$\frac{\partial v(x, t)}{\partial x} = -L \frac{\partial i(x, t)}{\partial t} \quad (63)$$

$$\frac{\partial i(x, t)}{\partial x} = -C' \frac{\partial v(x, t)}{\partial t} \quad (64)$$

As already established, the voltage at any point is given by the sum of forward and backward traveling voltage waves: $v(x, t) = v_f(x - ct) + v_b(x + ct)$ with c , given by Eq. (65), being equal to the speed of light

$$c = 1/\sqrt{L'C'} \quad (65)$$

in the case of overhead lines. A forward traveling current wave $i_f(x - ct) = v_f(x - ct)/Z_1$ is associated with $v_f(x, t)$, and a backward traveling current wave $i_b(x + ct) = v_b(x + ct)/Z_1$ is associated with $v_b(x, t)$, where the characteristic impedance Z_1 , given by Eq. (66),

$$Z_1 = \sqrt{L'/C'} \quad (66)$$

relates the associated waves. Using the sign conventions given in the table, the voltages and currents at each point of the line are given by Eqs. (67) and (68).

$$v(x, t) = Z_1 i_f(x - ct) + Z_1 i_b(x + ct) \quad (67)$$

$$i(x, t) = i_f(x - ct) - i_b(x + ct) \quad (68)$$

From these two equations follow Eqs. (69) and (70).

$$i_f(x - ct) = \frac{1}{2} \left(i(x, t) + \frac{1}{Z_1} v(x, t) \right) \quad (69)$$

$$i_b(x - ct) = \frac{1}{2} \left(-i(x, t) + \frac{1}{Z_1} v(x, t) \right) \quad (70)$$

The argument $(x - ct)$ of i_f is zero for $x = 0, t = 0$ and $x = d, t = T_p$. Furthermore, at terminal 1, $x = 0, i_1(t) = i(0, t)$, and at terminal 2, $x = l, i_2(t) = -i(l, t)$. Therefore, from Eqs. (69) and (70) follows Eq. (71), which can be rearranged to Eq. (72),

$$i_1(t - T_p) + \frac{1}{Z_1} v(t - T_p) = -i_2(t) + \frac{1}{Z_1} v_2(t) \quad (71)$$

$$i_2(t) = \frac{1}{Z_1} v_2(t) - \left(i_1(t - T_p) + \frac{1}{Z_1} v_1(t - T_p) \right) \quad (72)$$

Similarly, Eq. (73) holds.

$$i_1(t) = \frac{1}{Z_1} v_1(t) - \left(i_2(t - T_p) + \frac{1}{Z_1} v_2(t - T_p) \right) \quad (73)$$

Discretization readily leads to the companion element model of the lossless line shown in the table.

Companion modeling for the network. The companion element models are connected to the overall companion network model in accordance with the topology of the power system under study. For the companion network model, node equations are formulated to obtain the nodal equation system. See NETWORK THEORY.

Solution procedure. Once the nodal equation system is known, the time step counter k is set to zero and a simulation loop is entered. The solution to the nodal equation system is calculated, and the quantities of interest are provided as output. Then, the time step counter is incremented, the loop is reentered, and the operations are repeated. The simulation stops once a termination condition is met. Simulation results can be displayed graphically, as shown in Fig. 8. See NUMERICAL ANALYSIS; SIMULATION. Kai Strunz

Bibliography. J. B. Aidala and L. Katz, *Transients in Electric Circuits*, Prentice Hall, 1980; H. W. Dommel, *EMTP Theory Book*, Microtran Power System Analysis Corp., Vancouver, 1992; A. P. S. Meliopoulos, *Power System Grounding and Transients*, Marcel Dekker, New York, 1988; K. Strunz, *Numerical Methods for Real Time Simulation of Electromagnetics in AC/DC Network Systems*, VDI-Verlag, Düsseldorf, 2002; L. van der Sluis, *Transients in Power Systems*, Wiley, Chichester, 2001; N. Watson and J. Arrillaga, *Power Systems Electromagnetic Transients Simulation*, IEE, London, 2003.

Electric vehicle

A ground vehicle propelled by a motor that is powered by electrical energy from rechargeable batteries or other source onboard the vehicle, or from an external source in, on, or above the roadway. Examples are the golf cart, industrial truck and tractor, automobile, delivery van and other on-highway truck, and trolley bus. In common usage, electric vehicle refers to an automotive vehicle in which the propulsion system converts electrical energy stored chemically in a battery into mechanical energy to move the vehicle. This is classed as a battery-only-powered electric vehicle. The other major class is the hybrid-electric vehicle, which has more than one power source. See AUTOMOBILE; BUS; TRUCK.

History. Construction of the first electric vehicle is credited to the French inventor and electrical engineer M. Gustave Trouvé, who demonstrated a motorized tricycle powered by lead-acid batteries in 1881. In the United States, Andrew L. Riker is credited with building the first electric vehicle (also a tricycle) in 1890, and by 1891 William Morrison had built the first electric four-wheeler. In France in 1899, a four-wheel electric vehicle driven by Camille Jenatzy became the first car to break 60 mi/h (96 km/h). By then, production of battery-powered vehicles for use as personal transportation, commercial trucks, and buses had already begun.

Electric vehicles, with their instant starting, quiet running, and ease of operation, peaked in their challenge to steam- and gasoline-powered cars in 1912. The limited performance, range, and speed of electric vehicles, plus the need for frequent battery charging, restricted their usefulness and dampened their popularity. By the 1920s, the piston-type internal combustion engine had prevailed as the dominant automotive powerplant. Most production and development work on electric vehicles ended during

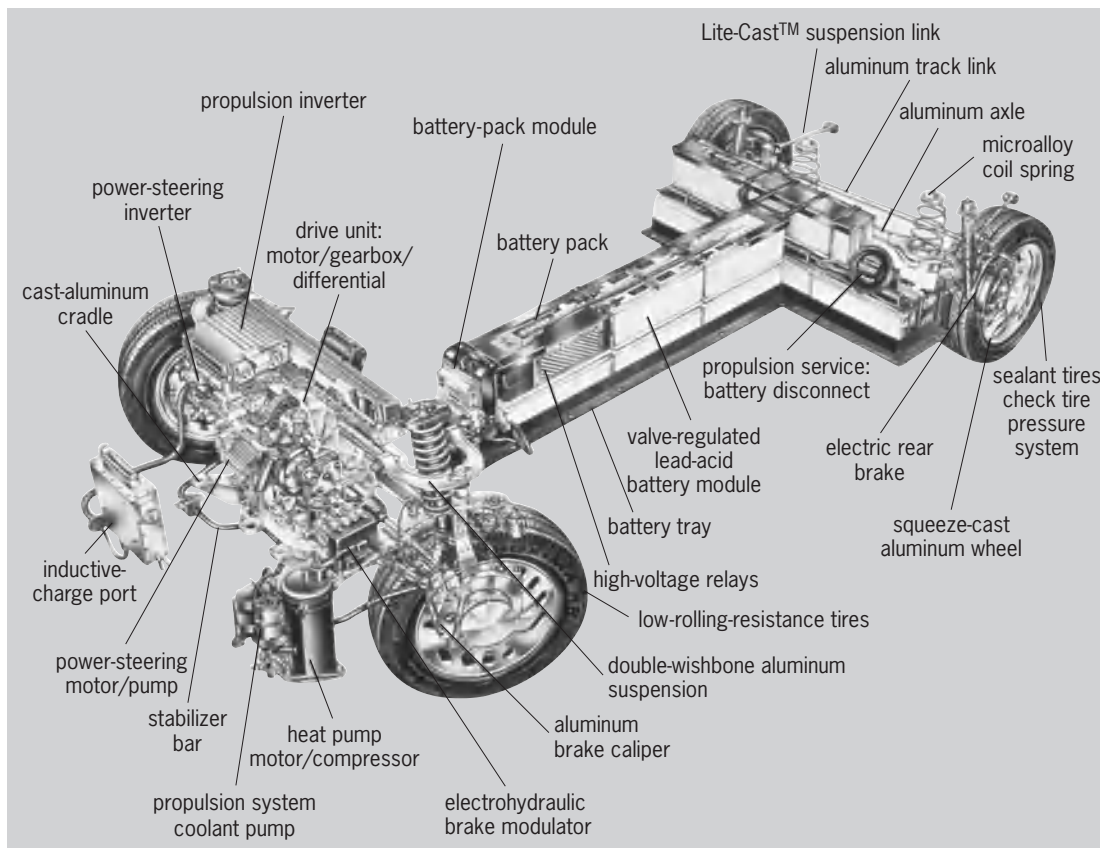


Fig. 1. EV1, a two-seat electric vehicle powered by lead-acid batteries. (General Motors Corp.)

the 1930s. See AUTOMOTIVE ENGINE; ENGINE; INTERNAL COMBUSTION ENGINE.

In the 1960s, interest revived in electric vehicles as a result of concern with diminishing petroleum reserves, rising cost of crude oil production, and air pollution from the automotive engine that burned gasoline which was refined from crude oil. Over the years, a few electric vehicles had been constructed, usually by converting small light cars and trucks into electric vehicles by removing the engine and fuel tank and installing an electric motor, controls, and batteries. However, during that time no major automotive manufacturer brought out an electric vehicle. See AIR POLLUTION; GASOLINE; PETROLEUM.

The Clean Air Act of 1963 and its amendments established limits on emissions from new vehicles sold in the United States. In 1990, the California Air Resources Board decided to further reduce air pollution by mandating (but later rescinding) that 2% of each automaker's sales must have zero emissions in the 1998 model year. This demand for a zero-emission vehicle (ZEV) could be met only by the electric vehicle, which typically was powered by lead-acid batteries. Used in an electric vehicle, lead-acid batteries have two major weaknesses: relatively high weight for the amount of energy stored, and reduced capacity in cold weather.

To help develop a better battery for electric vehicles, the U.S. Advanced Battery Consortium was formed in 1991. The purpose of this partnership among United States automakers and the electric

utility industry was to develop advanced batteries capable of providing future generations of electric vehicles with significantly increased range and performance.

In 1996 General Motors began limited marketing of the electric vehicle EV1 (Fig. 1). The EV1 was the first specifically designed electric car produced by a major automaker since before World War II. Other automakers also have developed and tested electric vehicles and vehicle conversions powered by lead-acid or advanced batteries.

Battery-only power. The General Motors EV1 is a battery-only-powered vehicle containing 26 lead-acid batteries which are assembled into a T-shaped pack that provides a nominal voltage of 312 volts (Fig. 2). These batteries differ from 12-V automotive batteries primarily in duty cycle. In the electric vehicle, the batteries must supply all the energy needs. Except for regenerative braking, there is no onboard charging. The batteries provide the power to propel the vehicle, and to power the lights and all accessories such as air conditioning and radio. As a result, electric-vehicle batteries go through much deeper discharge cycles than an automotive battery, which seldom is discharged more than 5% of its rated capacity. See BATTERY.

Because of the increased duty cycle, electric-vehicle batteries can deliver 85% of their charge without damaging the batteries or shortening their useful life. Ideally, this could provide the EV1 and similar electric vehicles with a useful range of 70 mi

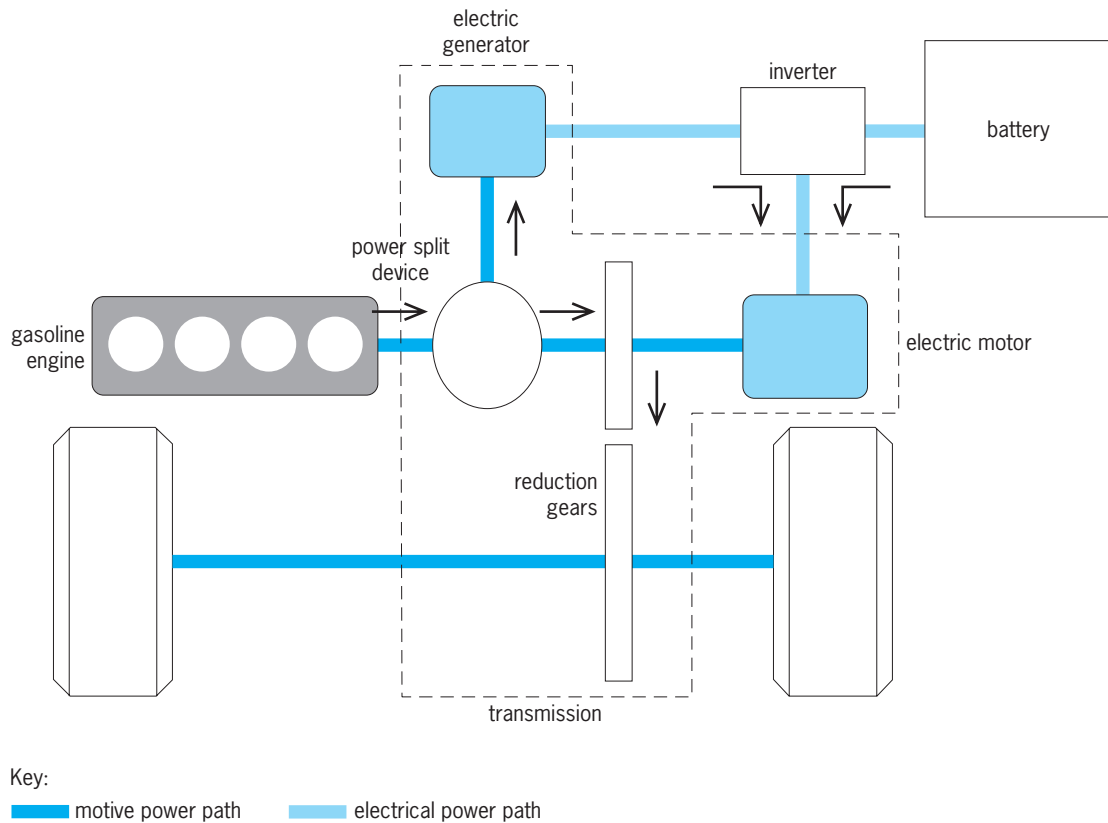


Fig. 2. Power flows through a hybrid vehicle which uses a gasoline engine to propel the vehicle and drive a generator that can operate the electric motor directly or charge the batteries, as necessary. (Toyota Motor Sales, U.S.A., Inc.)

(113 km) of city driving or 90 mi (145 km) of highway driving. A top speed of 80 mi/h (130 km/h) may be possible, but at a sacrifice in range, which also is shortened by hilly terrain and use of any electrical equipment on the vehicle.

In addition to lead-acid batteries, other batteries are used in electric vehicles. These include the newer nickel-metal hydride battery and the lithium-ion battery.

Hybrid power. A hybrid electric vehicle has more than one source of power. These sources can be different types of energy storage devices, power converters, and inverters. The first hybrid vehicle is credited to an Italian, Count Felix Carli. In 1894, Carli constructed an electric-powered tricycle that had a system of rubber springs which could release a short burst of additional power when needed.

Although power losses occur each time that energy is converted from one form to another, hybrid drive can be more efficient than a conventional automotive engine. However, having two or more power sources can increase the complexity, cost, and weight of a hybrid vehicle, as well as its manufacturing, safety, emissions, maintenance, and service problems. See AUTOMOTIVE ENGINE.

Internal combustion engine. Since the 1890s, major hybrid-vehicle research and development work has focused on adding a small internal combustion engine to an electric vehicle. Typically, the battery-powered motor drives the wheels, while the engine,

usually running at constant speed, drives a generator that charges the batteries. Operating the engine at constant speed reduces fuel consumption and produces cleaner exhaust gas than if the engine were larger, operating at variable speed, and providing the sole source of vehicle power.

In some hybrid vehicles, engine power is split (Fig. 2). Part of the engine power propels the vehicle, while part drives the generator which can operate the motor directly or charge the batteries, as necessary. Reportedly, engine exhaust emissions of hydrocarbons, carbon monoxide, and oxides of nitrogen are about 10% that of a conventional gasoline-engine vehicle. In addition, fuel efficiency is doubled.

Fuel cell. In this electrochemical device, the reaction between a fuel, such as hydrogen, and an oxidant, such as oxygen or air, converts the chemical energy of the fuel directly into electrical energy. The fuel cell is not a battery and does not store energy, although the fuel cell also has two electrodes separated by an electrolyte. As long as fuel is supplied to one electrode of the fuel cell and oxygen or air to the other, a voltage is produced between the electrodes. When an external circuit connects the electrodes, electrons will flow through the external circuit. Since fuel-cell voltage is less than 1 V, stacks of fuel cells are connected together to provide the needed electrical energy.

When fuel cells are the primary power source in a

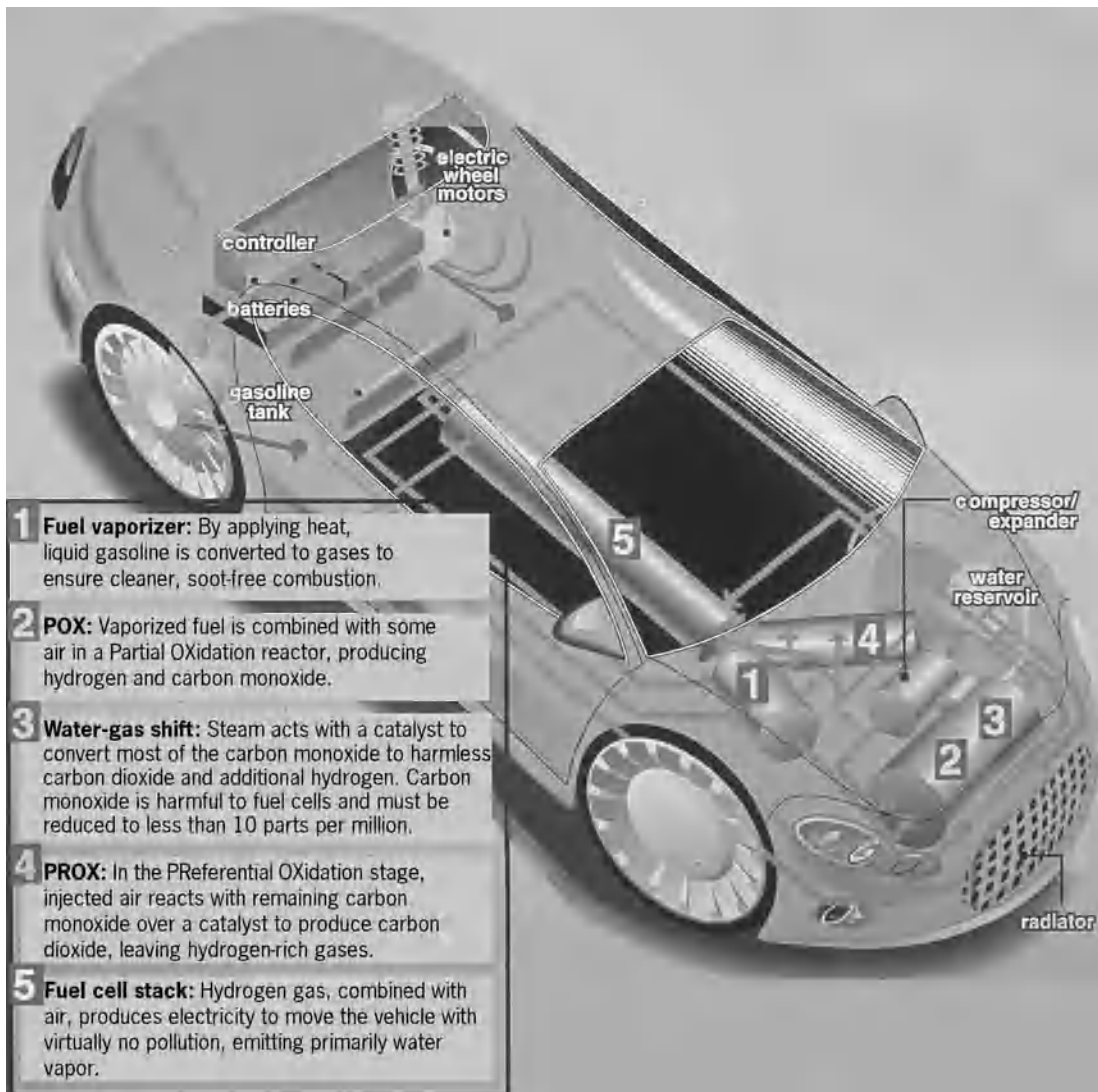


Fig. 3. Layout and five-step process that produces electricity from gasoline, which powers the fuel cells. (Chrysler Corp.)

hybrid vehicle, batteries provide secondary power. Fuel cells do not provide immediate output during a cold start. Until the fuel cells reach operating temperature, which may take about 5 min, a battery pack supplies the power for initial startup and vehicle movement. See FUEL CELL; HYDROGEN.

Three types of fuel cell under development for electric vehicles are the hydrogen-fueled, the methanol-fueled, and the gasoline-fueled. (1) A hydrogen-fueled cell runs on hydrogen gas stored in pressure tanks carried by the vehicle. Range of the vehicle is determined by the amount of compressed hydrogen that the tanks can hold. When hydrogen is used as the fuel, fuel cell operation produces no significant amounts of unwanted emissions. Water vapor and electricity are the only products. However, widespread use of hydrogen as a near-term vehicle fuel is unlikely because there exists no infrastructure of hydrogen refueling stations which are in place and accessible to the public. (2) To provide the fuel cell with hydrogen gas while avoiding the prob-

lems of hydrogen refueling, a methanol-to-hydrogen reformer onboard the vehicle can produce hydrogen gas from liquid methanol. Installed on production models, this would allow motorists to refuel vehicles in the conventional manner at existing service stations through any pump which dispensed methanol. However, use of a reformer to obtain hydrogen lowers vehicle efficiency and creates some emissions of carbon dioxide. (3) A gasoline-to-hydrogen reformer on the vehicle can extract hydrogen gas from gasoline. The hydrogen is then delivered to the fuel cell stack (Fig. 3). Use of gasoline could move fuel cell technology years closer to production in automotive vehicles, while reportedly improving fuel efficiency by 50% and emissions by 90%. See METHANOL.

Donald L. Anglin

Bibliography. Robert Bosch GmbH, *Automotive Handbook*, 6th ed., 2005; Society of Automotive Engineers, *SAE Handbook*, 3 vols., annually; M. H. Westbrook, *The Electric Car: Development and Future of Battery, Hybrid and Fuel-Cell Cars*, 2001.

Electrical breakdown

A large, usually abrupt rise in electric current in the presence of a small increase in electric voltage. Breakdown may be intentional and controlled or it may be accidental. Lightning is the most familiar example of breakdown.

Breakdown in gases and solids. In a gas, such as the atmosphere, the potential gradient may become high enough to accelerate the naturally present ions to velocities that cause further ionization upon collision with atoms. If the region of ionization does not extend between oppositely charged electrodes, the process is corona discharge. If the region of ionization bridges the gap between electrodes, thereby breaking down the insulation provided by the gas, the process is ionization discharge. When controlled by the ballast of a fluorescent lamp, for example, the process converts electric power to light. In a gas tube the process provides controlled rectification. *See* ARC DISCHARGE; BREAKDOWN POTENTIAL; CORONA DISCHARGE; ELECTRIC SPARK; ELECTRICAL CONDUCTION IN GASES; FLUORESCENT LAMP; GLOW DISCHARGE; IONIZATION.

In a solid, such as an insulator, when the electric field gradient exceeds 10^6 V/cm, valence bonds between atoms are ruptured and current flows. Such a disruptive current heats the solid abruptly; the rate of local temperature rise may fracture the insulator, the high temperature may carbonize or otherwise decompose the insulation, or as occasionally happens when lightning strikes a tree, the heat may ignite the insulator. *See* ELECTRIC INSULATOR; ELECTRICAL INSULATION.

In a semiconductor if the applied backward or reverse potential across a junction reaches a critical level, current increases rapidly with further rise in voltage. This avalanche characteristic is used for voltage regulation in the Zener diode. In a transistor the breakdown sets limits to the maximum instantaneous voltage that can safely be applied between collector and emitter. When the internal space charge extends from collector junction through the base to the emitter junction, even a voltage below the avalanche level can produce a short circuit, in which case the phenomenon is termed punch-through. *See* SEMICONDUCTOR; TRANSISTOR; ZENER DIODE.

Frank H. Rockett

Laser-induced breakdown in gases. For sufficiently high intensities, laser radiation can ionize (or breakdown) the gas atoms or molecules through which it is propagating. The laser radiation intensity required to produce ionization is dependent on the wavelength of the light and on the duration of the laser pulses as well as the gas species and pressure. In general, the breakdown by laser radiation can occur by two different mechanisms, a cascade- or collision-induced ionization, or a multiple-photon absorption process. The most often encountered mechanism is the cascade process, but multiple-photon ionization can be important under certain conditions.

Multiple-photon ionization. Multiple-photon ionization occurs when an atom absorbs directly from the laser

radiation field a sufficient number of photons to ionize the atom. The number of photons required is simply the ionization potential of the atom divided by the energy of a photon. For example, the photon energy of the ruby laser at 0.69-micrometer wavelength is 1.78 eV. Therefore, the multiple-photon ionization of argon having an ionization potential of 15.8 eV requires the simultaneous absorption of nine photons. Because of the nature of this process, the ionization is strongly dependent on the laser intensity and wavelength. Experimental measurements of multiple-photon ionization have shown that the threshold for air at pressures below 100 torr (2 psi or 13 kilopascals) is 10^{18} W/m² (10^{17} W/ft²) for a 50-picosecond-duration-pulse, 0.69- μ m-wavelength laser. Multiple-photon ionization can be important for laser-induced breakdown at short wavelengths, short pulse durations, and low gas pressures. *See* IONIZATION POTENTIAL; PHOTON.

Cascade breakdown. For the cascade breakdown process, a free electron is accelerated in the radiation field of the laser, but only gains an increment of energy upon collision with atoms. This energy absorption continues until the electron has sufficient energy to ionize an atom by an inelastic collision, producing a second electron. The two electrons then gain additional energy from the field, producing four electrons and so forth. Under breakdown conditions this process continues until the gas is fully ionized. An electron gaining energy from the laser radiation field can also lose energy by processes such as electronic excitation of atoms, or electrons can be lost from the breakdown process by diffusion out of the laser beam or by recombination or by attachment to form negative ions. The rate of ionization necessary to overcome these loss processes then defines the breakdown threshold intensity.

Breakdown threshold. Since high-intensity lasers are usually pulsed devices, an important parameter affecting the breakdown threshold is the duration of the laser pulse. If the laser pulse is long compared to the characteristic electron-loss time, breakdown occurs when the rate of electron production just exceeds the losses; breakdown is determined by the laser intensity and is independent of the laser pulse duration. If the laser pulse is short relative to the characteristic loss times, breakdown occurs when there is sufficient fluence (energy per unit area) in the pulse to generate full ionization of the gas. For air at atmospheric pressure and large-diameter beams (such that electron diffusion is not important), the threshold is fluence-dependent for pulses less than approximately 10^{-7} s and is intensity-dependent for pulses having a longer duration. As an example, for air at atmospheric pressure and laser beams larger in diameter than 4×10^{-3} in. (10^{-2} cm), the threshold fluence for the 10.6- μ m carbon dioxide (CO₂) laser wavelength is 10^5 joules/m² (10^4 joules/ft²) for pulse durations less than 10^{-7} s, and for longer pulses the intensity threshold is 3×10^{13} W/m² (3×10^{12} W/ft²). These theoretically predicted values have been experimentally verified.

Conditions for occurrence. Two important factors affecting the occurrence of breakdown are the source of the initiating electron and the presence of aerosol particles. The electron required for initiating the cascade process must be provided by some mechanism other than the cascade, such as photoionization from natural causes or a spark or discharge. For visible-wavelength laser-induced breakdown, the initial electron can be the result of multiple-photon absorption ionization of impurities, but for long-wavelength radiation, certainly $10.6\ \mu\text{m}$, an external source must provide it. Particulate matter, such as naturally occurring aerosols in the atmosphere, can have a dramatic effect on the breakdown threshold of gases. For particles of $0.1\ \mu\text{m}$ and larger, the interaction of the laser radiation with the particle vaporizes the particle, ionizes the vapor, and subsequently ionizes the surrounding gas. The threshold for this process can be several orders of magnitude lower than the threshold for gas breakdown in the absence of aerosols. See ATMOSPHERIC CHEMISTRY.

Pressure dependence. The laser-induced gas breakdown has an interesting pressure dependence. The rate of laser energy absorption by free electrons increases with pressure because of the increase in the electron-atom collision frequency. This energy absorption increase continues as gas pressure increases until the collision frequency equals the oscillatory frequency of the laser radiation field. For pressures above this level, the rate of energy absorption decreases with increasing pressure, and the breakdown approaches that induced by a constant electric field in the high pressure limit. Therefore, there is a minimum in a plot of the threshold as a function of pressure that occurs when the collision frequency equals the oscillatory frequency. As an example, for atmospheric air and $10.6\text{-}\mu\text{m}$ -wavelength radiation this minimum occurs at approximately 44 atm ($650\ \text{lb/in.}^2$ or 4.5 megapascals).

For gas pressures lower than the minimum, the gas breakdown threshold varies as the inverse square of the wavelength; that is, the shorter the wavelength, the more difficult it is to break down the gas. This behavior has been verified for wavelengths from 0.69 to $10.6\ \mu\text{m}$. The threshold for gas breakdown at pressures above the minimum is independent of the laser wavelength.

Applications. Laser-induced gas breakdown is important as a limitation for propagating high-intensity laser radiation, and as a source for the production of a high-density, high-temperature plasma for a number of applications. Gas breakdown represents the first step in research leading to the study of laser-generated controlled fusion. The plasma produced by laser-induced breakdown can be used as a convenient source for soft x-rays and an intense source of ultraviolet or visible radiation for a fast-rise-time light source. See LASER; NUCLEAR FUSION; PLASMA (PHYSICS).

David C. Smith

Bibliography. G. Bekef, *Principles of Laser Plasmas*, 1976; L. A. Dossado and C. J. Fothergill (eds.), *Electrical Breakdown and Degradation in Poly-*

mers, 1992; E. E. Kunhardt and L. H. Luessen (eds.), *Electrical Breakdown and Discharges in Gases*, 1983; Y. P. Raizer, *Laser-Induced Discharge Phenomena*, 1977; Y. P. Raizer and J. E. Allen, *Gas Discharge Physics*, 1997 reprint.

Electrical codes

Systematic bodies of rules governing the practical application, installation, and interconnection of electrically operated equipment, devices, and electrical wiring systems.

National Electrical Code. The basic code used throughout the United States, and a number of other countries throughout the world, for indoor and outdoor electrical installations to supply lighting, motors, appliances, and machines is the National Electrical Code, prepared under the direction of the National Fire Protection Association (NFPA). It is approved by the American National Standards Institute, and the 2005 edition of the code is known as NFPA 70-2005 (ANSI).

The National Electrical Code was originally drawn in 1897 as a result of the united efforts of various insurance, electrical, architectural, and allied interests. The original code was prepared by the National Conference on Standard Electrical Rules, composed of delegates from various interested national associations.

In 1911 the National Conference of Standard Electrical Rules was disbanded, and since then the NFPA has acted as the code's sponsor. Beginning with the 1920 edition, the National Electrical Code has been under the further auspices of the American National Standards Institute with the NFPA continuing in its role as administrative sponsor. Since that date, the NFPA committee that produces the code has been identified as ANSI Standards Committee C1 (formally USAS C1 or ASA C1).

The provisions of the National Electrical Code are under constant review by a number of panels whose members are selected to provide broad representation of electrical, industrial, and public interests. The code is amended in its periodic republication every 3 years or by tentative interim amendments (TIAs) which are announced by bulletins and through the technical press.

The National Electrical Code is purely advisory as far as the National Fire Protection Association is concerned, but it is very widely used for legal regulatory purposes. The code is administered by various local inspection agencies, whose decisions govern its application to individual installations. Local inspectors are largely members of the International Association of Electrical Inspectors. This organization, the National Electrical Manufacturers Association, the National Electrical Contractors Association, the Institute of Electrical and Electronic Engineers, the Edison Electric Institute, the Underwriters' Laboratories, Inc., the International Brotherhood of Electrical Workers, governmental groups, and independent experts all contribute to the

development and application of the National Electrical Code.

Compliance with the provisions of the code can effectively minimize fire and accident hazards in any electrical design. It sets forth requirements that constitute a minimum standard for the framework of electrical design. As stated in its introduction, the code is concerned with the “practical safeguarding of persons and property from hazards arising from the use of electricity.” The National Electrical Code is recognized as a legal criterion of safe electrical design and installation in many jurisdictions throughout the country and the world. In some large municipalities within the United States, the National Electrical Code either is accepted with certain modifications or serves as the basis for a local electrical code. Where accepted by the local governing body, the National Electrical Code is used in court litigation and by insurance companies as a basis for insuring buildings.

Other standards. In addition to the National Electrical Code volume itself, other standards and recommended practices are made available in pamphlet form by the National Fire Protection Association, which oversees the promulgation and periodic revision of over 270 standards. These cover such special subjects as hospital operating rooms, municipal fire alarm systems, garages, aircraft hangars, and other equipment with great potential hazards due to improper design.

The National Electrical Safety Code (to be distinguished from the National Electrical Code) is published by the Institute of Electrical and Electronic Engineers, Inc. Designed as ANSI C2, this code presents basic provisions for safeguarding persons from hazards arising from the installation, operation, or maintenance of (1) conductors and equipment in electric supply stations, and (2) overhead and underground electric supply and communications lines. Basically, this code applies to the outdoor circuits of electric utility companies and to similar systems or equipment on commercial and industrial premises.

Municipal codes. The National Electrical Code is incorporated bodily or by reference in many municipal building ordinances, often with additional provisions or restrictions applicable in the particular locality. Some large cities have independent electrical codes; however, the actual provisions in most such codes tend to be basically similar to the National Electrical Code.

Testing of electrical products. Standards on the construction and assembly of many types of electrical equipment, materials, and appliances are set forth in literature issued by the National Electrical Manufacturers Association (NEMA) in conjunction with Underwriters Laboratories, Inc. (UL). The Underwriters Laboratories examines, tests, and determines the suitability of materials and equipment to be used according to code regulations. Each year, it publishes three volumes listing commercially available electrical products which have been found acceptable with reference to fire and acci-

dent hazards and which conform with the application and installation requirements of the code. The three volumes are titled *Electrical Construction Materials Directory*, *Electrical Appliance and Utilization Equipment Directory*, and *Hazardous Location Equipment Directory*. The Underwriters' Laboratories publishes other literature such as *Gas and Oil Equipment Directory* and *Fire Protection Equipment Directory*, dealing with special equipment involving hazard to life or property.

Testing of electrical equipment and products is permitted to be performed by any “nationally recognized testing laboratory,” as stipulated by the Occupational Health and Safety Administration (OSHA). Underwriters Laboratories is one such company. Tests performed by other such companies are conducted in accordance with the applicable NEMA/UL standard.

Administration. Electrical codes are administered locally by inspectors. The inspectors, who may review plans and specifications, make periodic site visits during the construction process to visually examine electrical work during installation and after the work is completed to ensure compliance with applicable rules or ordinances.

Electrical inspection bureaus are maintained in many cities by the city itself. In other cities, a third-party inspection agency, such as the New York Board of Fire Underwriters, which provides electrical inspectors for any jurisdiction within New York State that wishes to avail itself of the Board's services, is used. Typically where third-party inspectors are used, the companies recognized to perform inspections within a given jurisdiction are identified by the local building department. In communities where codes are enforced by ordinance, inspections may be performed by municipal electrical inspectors. Utility inspectors examine the service entrance and metering installation for compliance with prevailing utility regulations.

Federal and state buildings are usually inspected by authorized government electrical inspectors. In these instances inspection includes both safety consideration and the requirements of the particular job specifications. Other specification (by underwriters or municipal inspectors) is often waived. See ELECTRONIC EQUIPMENT GROUNDING; GROUNDING; WIRING.

J. F. McPartland; Brian J. McPartland

Bibliography. J. F. McPartland, *Handbook of Practical Electrical Design*, 3d ed., 1999; J. F. McPartland, *National Electrical Code Handbook*, 25th ed., 2005; *National Electrical Safety Code*, NFPA 70-2005 (ANSI), 2005.

Electrical communications

That branch of electrical engineering dealing with the transmission and reception of information. Information can be transmitted over many different types of pathways, such as satellite channels, underwater acoustic channels, telephone cables, and

fiber-optic links. Characteristically, any communications link is noisy. The receiver never receives the information-bearing waveform as it was originally transmitted. Rather, what is received is, at best, the sum of what was transmitted and noise. In reality, what is more likely to be received is a distorted version of what was transmitted, with noise and perhaps interference. Consequently, the design and implementation of a communications link are dependent upon statistical signal-processing techniques in order to provide the most efficient extraction of the desired information from the received waveform. *See* COMMUNICATIONS CABLE; COMMUNICATIONS SATELLITE; DISTORTION (ELECTRONIC CIRCUITS); ELECTRICAL INTERFERENCE; ELECTRICAL NOISE; MICROWAVE; OPTICAL COMMUNICATIONS; RADIO-WAVE PROPAGATION; TELEPHONE SERVICE; TRANSMISSION LINES; UNDERWATER SOUND; WAVEGUIDE.

Broadly speaking, there are two basic classes of communication waveforms, those involving analog modulation and those involving digital modulation. The former type implies that the modulation process allows the actual information signal to modulate a high-frequency carrier for efficient transmission over a channel; this is achieved by using the continuum of amplitude values of an analog waveform. Examples of analog modulation systems include amplitude-modulation (AM) and frequency-modulation (FM) systems, as well as a variety of others such as single-sideband (SSB), double-sideband (DSB), and vestigial-sideband (VSB) systems. In digital modulation systems, the initial information-bearing waveform, assuming it is in analog form (such as voice or video), is first sampled, then quantized, and finally encoded in a digital format for carrier modulation and transmission over the channel. *See* MODULATION.

Analog modulation. Analog communication systems are often characterized as linear or nonlinear modulation systems. Linear systems include AM and single-, double-, and vestigial-sideband systems, whereas nonlinear systems include FM and phase-modulation (PM). Among the linear analog systems, AM is probably the most widely used, because its receiver structure can be implemented in the form of a simple envelope detector. However, AM is very inefficient from a power point of view, because its signal structure results in part of its total transmitted power being wasted in an unmodulated carrier instead of all its power going to the information-bearing portion of the waveform. *See* AMPLITUDE MODULATION; FREQUENCY MODULATION; PHASE MODULATION; SINGLE SIDEBAND.

This problem is remedied with either single- or double-sideband systems, but their receivers are more complex than the AM receiver. Furthermore, the single-sideband system, which has the most complex transmitter, requires only half the bandwidth of the other two, thereby giving it yet another advantage over AM. The vestigial-sideband system is essentially a compromise between the double sideband and the single sideband, for although implementation with it is somewhat easier than with a single sideband, its bandwidth is greater than that of the

single sideband but less than that of the double sideband.

The disadvantage of any linear analog modulation system is that its ultimate performance, as measured by the ratio of the signal power of the final demodulated waveform to the noise power of the demodulated waveform, called the signal-to-noise ratio (SNR), is limited by this ratio at the input to the receiver. In nonlinear analog communications, this is not necessarily the case. There the output signal-to-noise ratio can be made significantly higher than that at the input, and consequently significantly higher than the output signal-to-noise ratio of a linear modulation scheme. However, the price for this improvement in performance is increased transmission bandwidth. *See* SIGNAL-TO-NOISE RATIO.

Digital modulation. In digital communications, it is necessary to convert the information into a digital format. At times, the information might originate in an appropriate digital format, as in the transfer of computer files from one computer to another (so-called computer communications). However, in many other instances, the information originates in analog form. This is the case for both voice communications and video transmission. Therefore, some means of digitizing the analog waveform is required.

The classical technique for accomplishing this is probably the use of pulse-code modulation (PCM). There the analog signal is first sampled as a minimum at what is known as the Nyquist frequency (that is, at a rate corresponding to twice the bandwidth of the waveform). The samples are quantized and encoded into a binary word. The digits of this word are used to modulate a carrier for transmission over the communications channel. *See* PULSE MODULATION.

Pulse-code modulation is one example of the general technique of source encoding. Other schemes attempt to remove as much redundancy as possible from the message in order to minimize the number of binary digits needed to reproduce it. Examples of practical source-coding schemes are delta modulation and differential pulse-code modulation. In both of these techniques, an estimate of the current sample is made, the error between the sample and its estimate is ascertained, and this error signal is encoded into binary digits to be modulated on the carrier.

Once the waveform is in digital format, many modulation techniques can be used. These consist of both binary schemes, whereby one of two possible signals is transmitted across the channel, and M -ary schemes with $M > 2$, whereby one of M signals is used to convey the information in any given signaling interval. Typically, $M = 2^n$, for some integer n , and if all M signals are equally likely (that is, the a priori probability of transmitting any one of them is $1/M$), then each transmitted symbol conveys n bits of information. Thus, a binary system conveys one bit of information each time a symbol is transmitted, a quaternary system conveys two bits of information each time a symbol is sent, and so on.

For binary systems, the most widely used modulations are phase-shift keying and frequency-shift

keying. In the former scheme, a radio-frequency carrier has its phase changed to 0 and π radians each time a binary 1 or a binary 0 is to be transmitted, respectively. In frequency-shift keying, the information is transmitted by sending one of two possible frequencies during each signaling interval.

Many modulations can be used with M -ary schemes. Perhaps conceptually the simplest are the M -ary versions of binary phase-shift or frequency-shift keying, for example, the transmission of a carrier which uses one of M possible phases during each signaling interval. However, other schemes make use of a combination of both amplitude and phase modulation; these techniques are termed quadrature amplitude modulation (QAM).

Communication theory. The theoretical foundations of digital communications stem from the results of statistical decision theory and information theory. By using concepts developed in the former discipline, it is possible, in principle, to design an optimal receiver for a given transmitted waveform and a given noise environment. For example, in the detection of one of M known signals, all of which have equal energy and equal probability to be transmitted a priori, in a noise environment that contains only radio-frequency receiver noise, one form of an optimum receiver consists of a parallel filter bank of M matched filters with one filter matched to each of the M possible transmitted waveforms. A matched filter has an impulse response which is the time inverse of the signal to which it is matched. *See* DECISION THEORY; ELECTRIC FILTER.

Information theory has as its key results the ultimate limitations and performance capabilities of any communication system. The most important of these is embodied in the so-called noisy channel-coding theorem, which states that, under a fairly broad set of conditions, it is possible to achieve perfectly reliable communication over a noisy (unreliable) channel. This is accomplished with the help of error-correction coding, a technique of adding a controlled amount of redundancy to the information sequence in order to correct the errors that are made when the data are transmitted over the noisy channel. This seemingly astounding result is not especially practical, since it leads to system characteristics which are not realistic, such as infinite delay before the message can be decoded. However, its tremendous significance is due to the fact that it shows the power of error-correction coding in the design of communication systems. *See* INFORMATION THEORY.

Multiple-access techniques. In many communication systems, multiple sources and multiple destinations are present, and the manner in which accessing the channel is achieved becomes important. Perhaps the most common examples of such multiple-access links are those used by commercial radio and television stations. These systems operate by assigning to each transmitted waveform a distinct frequency band which is adjacent to but nominally not overlapping with its neighboring bands. In this way, there is approximately no interference among the various users. This type of operation, known as

frequency-division multiple accessing (FDMA), can be used with either analog or digital modulation formats. *See* RADIO; TELEVISION.

Because of certain problems that arise with frequency-division multiple accessing, such as the generation of intermodulation products when such a waveform is passed through a nonlinear amplifier, other types of multiple-accessing schemes have been devised. (An intermodulation product corresponds to a frequency generated at the output of a nonlinear device due to the nonlinear interaction of two or more frequencies at the input of the device.) In particular, when digital modulation is being employed, both time-division multiple access (TDMA) and code-division multiple access (CDMA) are possible alternatives.

In time-division multiple access, all users occupy the full bandwidth of the channel, but they transmit sequentially in time, and only one user has access to the channel at any given time. The system operates as follows. Each user is assigned a certain amount of channel time, and when the time slot of one user ends, the next user begins. After all the users have had an opportunity to transmit, the channel returns to the first user and the process is repeated. The time for the channel to return to a given user is referred to as the frame time. A predetermined pattern of bits, known as a frame synchronization pattern, is transmitted at the beginning of each frame. This pattern is known by each user, and it allows the users to synchronize to it so that they will know to transmit their data only during the allocated time interval, thus avoiding interference between the users. The advantage of the system over frequency-division multiple accessing is that the intermodulation problem disappears since, even if the channel is nonlinear, only one signal is on it at any time.

The technique of code-division multiple access operates by allowing multiple users to individually occupy the full bandwidth of the channel while simultaneously transmitting. This is accomplished by superimposing on the data of each user a wide-bandwidth digital waveform which identifies that user. Typically, a different identification sequence is assigned to every user. At the receiver this sequence for a given user allows the receiver to demodulate the data of that user while rejecting the waveforms of the other users.

Code-division multiple access is an application of a more general technique known as spread spectrum communication, in which the transmission bandwidth is purposely made much larger than the information bandwidth. This has the advantage that the resulting waveform is much less vulnerable to interference than a conventional waveform. Since the interference could be due to any one of a variety of sources, such as multipath or intentional jamming, this system has an inherent natural interference immunity which both frequency-division and time-division multiple access lack. *See* SPREAD SPECTRUM COMMUNICATION.

Other types of multiple-accessing schemes are designed to operate over a specific channel. One

example is space-division multiple access, which is used over a satellite communications link. In this scheme, spot beams on the satellite are used to illuminate different areas on the earth, so that the same frequency band can be used simultaneously by each of the spot beams without causing any interference. See MULTIPLEXING AND MULTIPLE ACCESS.

Random access. In many cases, the number of potential users is much greater than the number which can be simultaneously accommodated on the channel. However, if the percentage of time employed by each user is statistically very small, the users can compete for access to the channel. When a given user has a message ready for transmission, most of the time a vacant slot on the channel will be found, allowing data transmission. However, at times, all available slots on the channel are taken, and the user has to delay sending the message.

Systems that operate in this manner are referred to as random-access systems, and are typical of computer communication networks. Depending on the geographical size of such networks, they have come to have their own specialized names. For example, they are known as local-area networks (LANs) if they are concentrated over an area roughly the size of a few blocks or less. If the terminals which are furthest apart are within a few miles of one another, they are referred to as metropolitan-area networks (MANs). Networks that span still larger geographical distances are called wide-area networks (WANs). See DATA COMMUNICATIONS; LOCAL-AREA NETWORKS; WIDE-AREA NETWORKS. Laurence B. Milstein

Bibliography. M. Schwartz, *Information Transmission, Modulation and Noise*, 4th ed., 1990; H. Taub, *Principles of Communication Systems*, 2d ed., 1986; R. E. Ziemer and W. H. Tranter, *Principles of Communications*, 3d ed., 1990.

Electrical conduction in gases

The process by means of which a net charge is transported through a gaseous medium. It encompasses a variety of effects and modes of conduction, ranging from the Townsend discharge at one extreme to the arc discharge at the other. The current in these two cases ranges from a fraction of 1 microampere in the first to thousands of amperes in the second. It covers a pressure range from less than 10^{-4} atm (10 pascals) to greater than 1 atm (10^5 pascals). See ARC DISCHARGE; TOWNSEND DISCHARGE.

In general, the feature which distinguishes gaseous conduction from conduction in a solid or liquid is the active part which the medium plays in the process. Not only does the gas permit the drift of free charges from one electrode to the other, but the gas itself may be ionized to produce other charges which can interact with the electrodes to liberate additional charges. Quite apparently, the current voltage characteristic may be nonlinear and multivalued. See ELECTRICAL CONDUCTIVITY OF METALS; ELECTROLYTIC CONDUCTANCE; SEMICONDUCTOR.

The applications of the effects encountered in

this area are of significant commercial and scientific value. A few commercial applications are thyratrons, gaseous rectifiers, ignitrons, glow tubes, and gas-filled phototubes. These tubes are used in power supplies, control circuits, pulse production, voltage regulators, and heavy-duty applications such as welders. In addition, there are gaseous conduction devices widely used in research problems. Some of these are ion sources for mass spectrometers and nuclear accelerators, ionization vacuum gages, radiation detection and measurement instruments, and thermonuclear devices for the production of power.

The discussion of this complicated process will be divided into two parts. The first will deal with the basic effects involved, including production and loss of charges within the region and the motion of charges in the gas. The second part will deal with the mechanism of conduction.

Basic effects. To produce gaseous conduction, two conditions must obtain. First, there must be a source of free charges. Second, there must be an electric field to produce a directed motion of these charges. Considering the first of these, one finds that the free charge concentration is a result of a number of processes which produce and remove charges.

Sources of free charges. In many gaseous devices, a thermionic-emission cathode is included. The process of electron emission from a heated electrode is well known. Closely related to this as a source of electrons is field emission. Here, a strong positive field at a metallic surface lowers the barrier for electron emission. Thus, the electron current from a surface at a given temperature may be significantly increased. See THERMIONIC EMISSION.

Both of these effects result in electron production. Another effect, photon absorption, may result only in electrons if the absorber is a solid. However, if the photon interacts with a gas molecule or atom, ionization may result and both an electron and a positive ion can be obtained. The photon may come from some external source or it may be a secondary effect of the gaseous conduction. It may have a wavelength in the visible, ultraviolet, or x-ray region.

Conduction in flames is largely a result of the thermal production of ionization. This is a specialized field which has long been of interest in chemistry and combustion studies. To produce appreciable thermal ionization, the temperature must be high, as in a flame. If the effective temperature is known, the ionization concentration may be determined from statistical mechanics. Thermal ionization is also of tremendous importance in devices for production of power by thermonuclear processes, and in ion-propulsion equipment. A special form of this is surface ionization, in which a hot surface may cause ionization of a gas atom that comes in contact with it.

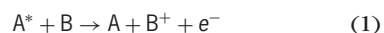
Another source of ionization is particle radiation due to cosmic rays, radioactive material in the walls or in the gas, or particles produced from an external source. These particles may then produce an ionized track in the gas. An example is the ionization produced by an α -particle from a polonium source in an ionization chamber.

In most of these methods of charge production, the sources are primary. That is, the presence of other free charges is not important in the production of ionization or electrons. Other processes are secondary in origin, although they may be of prime importance. It was pointed out that photons could originate either externally or as a secondary effect. Field emission could be a secondary effect, also. Other methods of ionization, however, are generally thought of as being secondary in origin. Ionization of the gas by electron impact is such a case. Here free electrons may gain enough energy in an electric field to interact with an atomic or molecular electron to produce an ion pair.

Cumulative ionization is an extension of ionization by impact. If the original electron and its offspring gain enough energy so each may produce another electron, and if this process is repeated over and over, the result is called an avalanche, and the ionization thus produced is referred to as cumulative ionization. This is the basis for particle detection in some ionization devices. See GEIGER-MÜLLER COUNTER; IONIZATION CHAMBER.

Another secondary source is electron emission from either electron or positive ion bombardment of a surface. This should not be confused with thermionic emission resulting from heating under bombardment. See SECONDARY EMISSION.

Other sources which may be important are atomic collisions, sputtering, and collisions of the second kind. In the first case, an atom or heavy ion may collide with an atom to produce ionization. This is quite unlikely until an energy of many times the ionization energy is obtained. The second is somewhat analogous to secondary electron emission. Here the positive ions strike a surface and knock out atoms or groups of atoms. Some of these come off as ions. In the third case, an excited atom may interact with an atom or molecule which is chemically different and has an ionization potential lower than the excitation potential of the excited atom. The result may be the decay of the excited state with ionization of the struck molecule or atom. Symbolically, this is shown by reaction (1), where A^* is the excited atom, B the



struck atom or molecule, and e^- an electron. See EXCITATION POTENTIAL; IONIZATION POTENTIAL; SPUTTERING.

Free-charge removal. The net free-charge concentration is a balance between charge-production and charge-removal processes. Recombination is one such process. Here, an electron or heavy negative ion and a positive ion may recombine. The energy transition may appear as electromagnetic radiation or may be carried off by a third body, if one is present. There are a wide variety of conditions which may lead to recombination. Where the temperature and electric field are high, the recombination will occur predominantly at the walls.

The method of charge removal is important from the aspect of conduction, however. If the charges

move to the appropriate electrodes under the influence of the field and there recombine, then they contribute to the current. If they simply diffuse to the walls and recombine there with ions of the opposite sign, or if they recombine in the gas volume, they may not appear as part of the external current.

Motion of the charges. The motion of the charges within the gas will be largely influenced by the potential function. For the usual regular geometries, this could be calculated in principle if there were no charges present. However, in a gas with free charges distributed throughout, the problem is quite different. The charges modify the charge-free potential, but the potential itself determines how the charge will move. The motion of the charges further modifies the potential and so on. Although the situation can be described physically by Poisson's equation, it is generally impractical to carry out an analysis. As a practical result, the potential function must be determined by measurements which are made with probes. This requires careful procedures to obtain significant results.

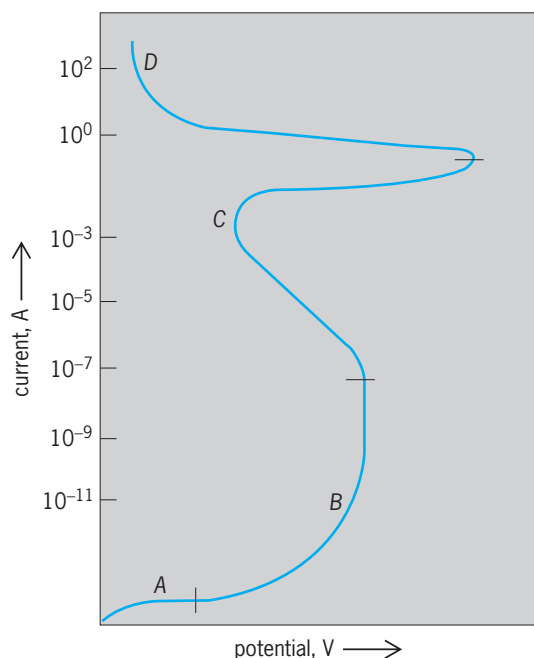
Diffusion of ions. Ion diffusion is a type of random motion which is always present and is the result of thermal or agitation energy. The randomness of the motion is brought about by the many collisions with molecules and other ions. A great difference exists in the motions of electrons and heavy ions. Because of low mass, electrons are easily deflected, so they move erratically. They diffuse badly, and follow field lines only generally. Again because the mass is small, an electron can give up appreciable energy only in an inelastic collision, in which excitation or ionization takes place. Hence, electron agitation and diffusion will be much greater in a pure inert gas than in a gas having many low-energy molecular states. Conversely, heavy ions exchange energy effectively at every collision. Diffusion is much less, so that they follow the electric field lines more closely than do electrons.

Mechanism of conduction. The ionic mobility μ relates drift velocity v to electric field X by Eq. (2). For

$$v = \mu X \quad (2)$$

electrons, the mobility is high, and a drift velocity of 10^6 cm/s or greater may be obtained. The electronic mobility is not a true constant, but varies with field, pressure, temperature, and gas composition. For heavy ions, the mobility is much more nearly constant, but is still dependent on these quantities to some extent. Drift velocities are usually of the order of 10^3 - 10^4 cm/s. Thus, in a typical conduction device, an electron may move from one electrode to the other before a heavy ion is displaced appreciably.

It would appear from the foregoing that if accurate information about the important processes existed, one could predict the characteristics of the conduction process under given conditions. Unfortunately this is not the case. Generally, the situation is so complicated that the theory can yield only qualitative



Current-potential characteristics for a two-electrode device with constant pressure.

predictions. Accordingly, most of the information concerning the various forms of gaseous conduction is empirical. In the present description, it will be possible to mention the main features of a few of these modes.

The **illustration** shows a sample voltage-current characteristic for a two-electrode device with constant pressure. It is assumed that there is a constant source of ionization which could be any of the primary sources previously discussed. In region *A*, the current first rises and then over a limited range is relatively constant as the voltage across the electrodes is increased. The initial rise is the result of the collection of charges which were either recombining or diffusing to the walls. The nearly constant current region is the result of the collection of almost all of the charges.

In region *B*, further increase in voltage produces an increase in current. Here, ionization by electron impact is occurring. The situation is described by specifying that each free electron makes α additional ion pairs in traveling 1 cm in the direction of the field. The number of ion pairs produced per second in 1 cm at a distance x from the cathode (assuming parallel plate electrodes) is given by Eq. (3),

$$n = n_0 e^{\alpha x} \quad (3)$$

where n_0 is a constant depending on the initial number of electrons. This is a form of the Townsend equation, and α is the first Townsend coefficient. In the region *B* in the illustration, the increase in current represents an increase in α . Near the end of this region, the current increases more rapidly with applied field. Here, additional effects are taking place, such as the photoelectric process and secondary emission. This situation is described by

Eq. (4), where β is the second Townsend coefficient,

$$i = i_0 \frac{(\alpha - \beta) e^{(\alpha - \beta)x}}{\alpha - \beta e^{(\alpha - \beta)x}} \quad (4)$$

cient, i_0 is the initial electron current at the cathode, and i is the anode current as a function of plate separation x ; β is also a function of electric field.

At the end of the region, the slope becomes infinite, and if the external resistance is not too large, the current will jump in a discontinuous fashion. The transition is referred to as a spark, and the potential at which it occurs is the breakdown or sparking potential. The region *B* is called a Townsend discharge and is not self-sustained. Thus, if the source of primary ionization were removed, the discharge would cease. See BREAKDOWN POTENTIAL; ELECTRIC SPARK.

As the potential reaches the sparking potential, a transition occurs to the region *C*. This is the self-sustained glow discharge region. Over an extensive current range, the voltage drop remains substantially constant. During the current increase, a glow occurs at the cathode, and at the upper end of the range, the cathode is completely covered. At this point, a further current increase can be achieved only if the potential drop across the discharge is increased. This portion of the characteristic is known as the abnormal glow. Throughout this portion of the discharge characteristic curve, secondary effects are quite important. Particularly vital are the effects of cumulative ionization and secondary emission at the cathode. See GLOW DISCHARGE.

Further increase in current leads to another mode of discharge, the arc. This is shown as region *D* in the illustration. Characteristic of this mode is the low cathode potential fall and the very high current density. It is generally felt that the predominant effect in the production of the large number of electrons at the cathode necessary for the arc is thermionic emission. This is consistent with the very high temperatures known to exist either generally or locally on the cathode. Although the arc type of discharge has very great commercial value, the mechanism of its operation is not very well understood.

In addition to these general types of conduction, there are very special cases of considerable interest. Some of these are the corona discharge, radiofrequency or electrodeless discharge, hot-cathode discharge, and discharges in the presence of a magnetic field.

Glenn H. Miller

Bibliography. M. N. Hirsh and H. J. Oskan (eds.), *Gaseous Electronics*, vol. 1: *Electrical Discharges*, 1978; E. E. Kunhardt and L. H. Luessen (eds.), *Electrical Breakdown by Discharges in Gases*, 2 vols., 1983; Y. P. Raizer and J. E. Allen, *Gas Discharge Physics*, 1997.

Electrical conductivity of metals

The property of a metal that measures its ability to conduct electricity, following Ohm's law. Electrical conductivity σ (measured in $\Omega^{-1} \cdot \text{m}^{-1}$) is the

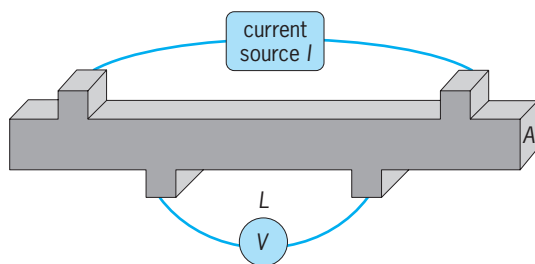


Fig. 1. Schematic of a four-probe measurement of electrical conductivity. Current I is fed through the outer leads, and voltage drop V is measured on the inner leads. In this way the contact potential drop experienced by the applied current is localized at the junctions with the outer leads. Measuring V with minimal current through the voltmeter minimizes the contact potential contribution to the measured conductance $G = I/V$.

reciprocal of the resistivity ρ [in $\Omega \cdot \text{m}$, where resistance R in Ohms (Ω) is voltage drop V in volts (V) divided by current I in amperes (A)]. Ohm's law, $V = IR$, may also be written in the form of Eq. (1), where j is current density (current per

$$j = \sigma E \quad \text{or} \quad E = \rho j \quad (1)$$

unit area, I/A , measured in A/m^2) and $E = V/L$ is electric field or electrical potential gradient, measured in V/m , where V is the voltage drop measured across a length L of material (**Fig. 1**). The conductivity σ is an intrinsic property of a pure material, related to the measured conductance $G = I/V$ via $G = \sigma A/L$, just as resistance $R = V/I$ relates to resistivity ρ via $R = \rho L/A$. Thus a large-gauge conductor [large area (A)] of short length has a high conductance, but the physical dimensions do not affect the conductivity σ or the resistivity ρ . Positive current flows from higher to lower voltage, and σ is never negative. See ELECTRICAL RESISTANCE; ELECTRICAL RESISTIVITY; OHM'S LAW.

Metals versus insulators. Metals have large conductivity (typically greater than $10^5 \Omega^{-1} \cdot \text{m}^{-1}$ but never greater at room temperature than that of silver, $0.66 \times 10^8 \Omega^{-1} \cdot \text{m}^{-1}$.) As temperature T is lowered, conductivity almost always increases in metals (**Fig. 2**). These properties contrast with semiconductors and insulators, which have smaller conductivity and different temperature dependence. The reason concerns the spectrum of quantum excited states in matter (**Fig. 3**). Outer electrons (occupying frontier orbitals in chemical language) have a continuous range of energy levels in solids, called a band of quantum states. In metals the highest-energy occupied state and the least-energy empty state differ only infinitesimally in energy, whereas in semiconductors or insulators they are separated by a significant energy gap. In both semiconductors and metals the states that correspond to the frontier orbitals are mobile, meaning that they permit electrons to travel from atom to atom. In semiconductors the energy gap is moderate (typically 0.5–3.0 eV). If electrons are introduced into empty states above the gap, a semiconductor behaves much like a metal. In an insulator the energy gap is larger, 5 eV or more, and elec-

trons typically get trapped on atoms or defects, rendering them comparatively immobile. See BAND THEORY OF SOLIDS; ELECTRIC INSULATOR; HOLE STATES IN SOLIDS; SEMICONDUCTOR.

Propagation of electrical current. To outside observers, electrical current flow in a wire resembles water flow in a pipe, but the analogy does not explain what is going on internally. In water flow, molecules have a net drift velocity, which is preserved by intermolecular collisions. Locally there is Galilean invariance, which means that an internal observer moving with the hydrodynamic current and observing individual molecules would not be able to detect the current. In a metal the electrons disengage from their parent positive ions. The ions vibrate around fixed positions, not participating in the current. When electrical current j flows, the electrons have a net drift velocity which can be detected by a local observer watching ions as well as electrons. The system does not have Galilean invariance. Collisions of electrons do not conserve j , and act to restore electrons to the $j = 0$ state of thermal equilibrium, not the moving local equilibrium of hydrodynamics. This is the source of resistance in metals. Electrical current is less “collective” than hydrodynamic current. However, at very low temperature many metals become superconductors, seen in the data of **Fig. 2**. This is a different state of electronic matter in which electrons behave more collectively than they do in the “normal” or higher-temperature resistive state of metallic electronic matter. The collective (and quantum-mechanically phase-coherent) superconducting state allows $j \neq 0$ with $E = 0$, thus making $\sigma = \infty$. See SUPERCONDUCTIVITY.

Electrons in quantum theory are waves as well as particles, and can diffract around the crystalline

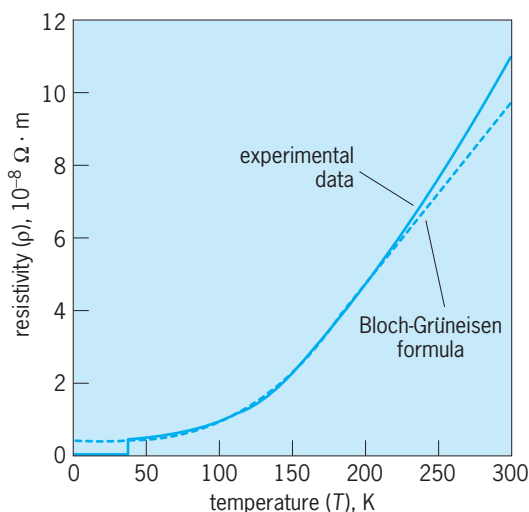


Fig. 2. Electrical conductivity of magnesium diboride (MgB_2), plotted, as is customary for metals, as resistivity (ρ) versus temperature (T). The unit $10^{-8} \Omega \cdot \text{m}$ is often written as 1 micro-ohm cm or $1 \mu\Omega\text{-cm}$. Below $T = 40 \text{ K} = -233^\circ\text{C}$, MgB_2 is a superconductor with $\rho = 0$. A fit using Bloch-Grüneisen theory, explained in the text, is also shown. (Data from Z. X. Ye et al., *Electron scattering dependence of dendritic magnetic instability in superconducting MgB_2 films*, *Appl. Phys. Lett.*, 86:5284–5286, 2004)

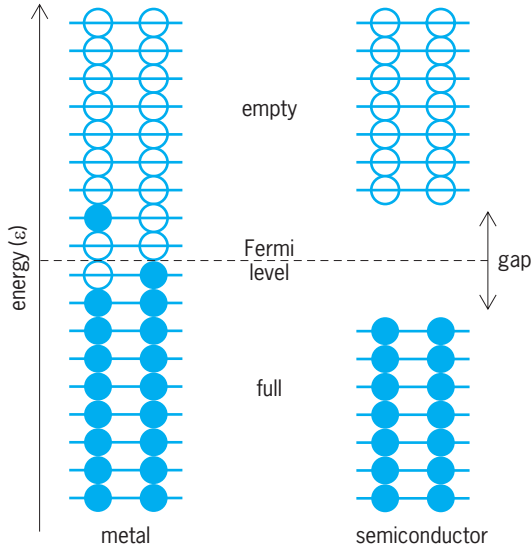


Fig. 3. Schematic energy spectrum in a metal versus a semiconductor. Filled and empty circles represent full and empty quantum states, two per level corresponding to the two spin orientations, just as in an atom. Adjacent levels are infinitesimally close except at the Fermi level of a semiconductor, which lies in the energy gap where there are no quantum states. The metal states are shown with a filled state above the Fermi level and a vacant state or hole below the Fermi level, as would happen by ordinary thermal excitation at room temperature. Semiconductors are likely to have few thermal excitations. In such a state the conductivity σ would be close to zero.

lattice of ions. In perfect crystals, if ions were not vibrating, the electron waves would remain coherent over the whole sample size, and conductivity would be infinite. The actual finite value occurs because electrons scatter from defects in the lattice, from atomic vibrations, and from each other. The last source of scattering is small, and the other sources are minimized in pure crystals at low temperature. Unlike molecules in water, which collide frequently, electrons in good metals collide infrequently and have long “mean free paths” or distances of travel between collisions. The collision frequency for water molecules at room temperature (300 K), denoted $1/\tau$, is about $6 \times 10^{12}/s$. This comes from dividing the thermal velocity, given by Eq. (2), by the distance $\sim 10^{-10}$ m between colli-

$$v_{th} = \sqrt{3k_B T/M} \approx 600 \text{ m/s} \quad (2)$$

sions. For electrons in metals the frequency of collisions with atomic vibrations at $T = 300$ K is given by Eq. (3). In these formulas, k_B is Boltzmann's

$$1/\tau \approx (1 - 10)k_B T/\hbar \sim 10^{14}/s \quad (3)$$

constant, M is the mass of a water molecule, and \hbar is Planck's constant divided by 2π . Compared to atoms, electrons have much higher mean velocities, $v_{el} \sim (2-10) \times 10^5$ m/s, and v_{el} is independent of temperature. This is a quantum effect caused by the Pauli principle. In a free-electron model, where all diffraction from atoms is ignored, the value is called the

Fermi velocity and is given Eq. (4), where n is the

$$v_F = \hbar(3\pi^2 n)^{1/3}/m \quad (4)$$

number of electrons per unit volume, of the order of 10^{29} m^{-3} , and m is the electron mass. As a consequence of the larger velocity, the mean free path for electrons, estimated as velocity divided by collision rate, is typically 10 to 20 interatomic spacings, 50 times longer than for water molecules. See CRYSTAL DEFECTS; EXCLUSION PRINCIPLE; FREE-ELECTRON THEORY OF METALS; LATTICE VIBRATIONS; QUANTUM THEORY OF MATTER.

Bloch-Boltzmann theory. The long mean free path of electrons in metals is the key to formulating a theory. In a perfect crystal the quantum description of electrons as particle-wave entities is similar to the idea of electromagnetic waves confined in a resonant cavity. The wavelengths have to “fit” the dimensions of the cavity, which causes quantization of wavelength λ . The individual quantum states or orbitals are labeled by their wavevectors $k = 2\pi/\lambda$. The orbital labeled by k has an energy $\varepsilon(k)$, and is pictured as a wave oscillating with a space and time dependence given approximately by $\cos [kx - \varepsilon(k)t/\hbar]$, like a traveling electromagnetic wave with frequency $\omega(k) = \varepsilon(k)/\hbar$. A wave packet built from such orbitals moves with velocity given by Eq. (5). There is no sim-

$$v(k) = [d\varepsilon(k)/dk]/\hbar \quad (5)$$

ple formula for the function $\varepsilon(k)$. It has to be measured by photoemission or computed using band theory, and varies from metal to metal. The wave packets in real metals do not remain phase-coherent over the whole sample, only over the distance between collisions. Provided this distance is longer than a wavelength, the wavelength and wavevector remain meaningful quantities, which can be used to build a theory, as was done by Felix Bloch in 1928.

Bloch's theory starts from the recognition of Bloch orbitals with wavevectors k . It then adopts from Boltzmann the notion of the distribution $F(k)$ which gives the probability that the state k is occupied. Because of the Pauli principle, this probability cannot exceed 2 (one for spin “up” and one for spin “down”). For states of low energy [$\varepsilon(k)$ far below the Fermi level ε_F], the occupancy is very close to 2, and for states of high energy [$\varepsilon(k)$ far above the Fermi level ε_F], the occupancy is very close to 0. In thermal equilibrium the occupancy $F(k)$ is equal to the Fermi-Dirac function, given by Eq. (6). The con-

$$f(k) = \frac{2}{\exp\left[\frac{\varepsilon(k) - \varepsilon_F}{k_B T}\right] + 1} \quad (6)$$

ductivity is calculated by summing the velocities of the occupied states, yielding Eq. (7).

$$\sigma = j/E = \frac{-e}{E \cdot Vol} \sum_k v(k)F(k) \quad (7)$$

When the field E is zero, the distribution $F(k)$ is $f(k)$ in Eq. (6), which has $j = 0$ because states with

positive and negative velocities $v(k)$ are equally populated. The effect of the field E is to cause the k -vectors of the occupied states to change according to Bloch's law, Eq. (8). This is the quantum version

$$\hbar dk/dt = -eE \quad (8)$$

of Newtonian acceleration by the electric force $-eE$, using de Broglie's momentum $\hbar k$ of a quantum wave. The acceleration is continuously resisted by scattering. Therefore the total shift of the typical k -vector is given by Eq. (9), where τ is the time between col-

$$\Delta k \approx -eE\tau/\hbar \quad (9)$$

lisions. Then, using $F(k) \approx f(k - \Delta k)$, we can write the answer for the conductivity as Eq. (10), where $(n/m)_{\text{eff}}$ is given by Eq. (11).

$$\sigma = \left(\frac{n}{m}\right)_{\text{eff}} e^2 \tau \quad (10)$$

$$\begin{aligned} \left(\frac{n}{m}\right)_{\text{eff}} &= \frac{-1}{\hbar \text{Vol}} \sum_k v(k) \frac{df(k)}{dk} \\ &= \frac{1}{\text{Vol}} \sum_k \left[\frac{d^2 \varepsilon}{\hbar^2 dk^2} \right] f(k) \end{aligned} \quad (11)$$

See NONRELATIVISTIC QUANTUM THEORY.

The last part of Eq. (11) is derived by integration by parts after substituting the expression given by Eq. (5) for $v(k)$, and can be interpreted as electron density, given by Eq. (12), times reciprocal effective mass, given by Eq. (13). The Bloch theory recovers a

$$n = \left(\sum f(k) \right) / \text{Vol} \quad (12)$$

$$1/m^* = d^2 \varepsilon / d(\hbar k)^2 \quad (13)$$

result proposed on classical grounds by Paul Drude shortly after the discovery of the electron, namely $\sigma = ne^2\tau/m$. However, the interpretation is quite different. A proper quantum theory for the scattering rate τ is needed, and was provided by Bloch by making a quantum generalization of the Boltzmann scattering operator. See BOLTZMANN TRANSPORT EQUATION.

An approximate solution of Bloch's theory, known as the Bloch-Grüneisen formula, is often very successful in fitting data, as shown in Fig. 2 for magnesium diboride (MgB_2). In this example, theory deviates from experiment at higher temperatures because high-frequency lattice vibrations, involving boron atoms, are not well represented by the form of the approximate theory. The extrapolated $T = 0$ resistivity value, $\rho_0 = 0.4 \times 10^{-8} \Omega \cdot \text{m}$, is called the residual resistivity, and would be zero in a perfect crystal.

Bloch's theory has failures, but they are outnumbered by the successes. One failure is that collective behavior can sometimes be seen at low temperature, the outstanding example being superconductivity. Other examples are more controversial, and include so-called Luttinger liquid behavior, which definitely happens in ideal one-dimensional theoretical models and probably is seen in careful experiments on

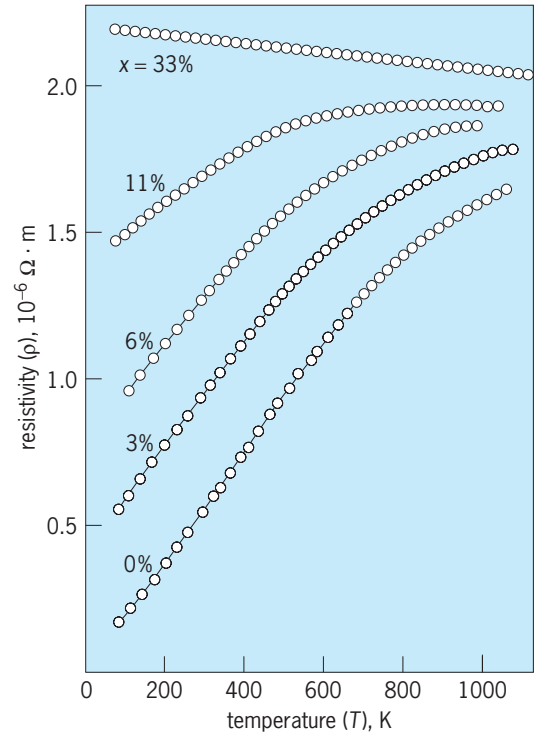


Fig. 4. Resistivity ρ of $\text{Ti}_{1-x}\text{Al}_x$ alloys as function of temperature (T) for various aluminum concentrations x . (From J. H. Mooij, *Electrical conduction in concentrated disordered transition-metal alloys*, *Phys. Stat. Sol. (a)*, 17: 521-530, 1973)

such systems as carbon nanotubes. Another case is transport by so-called sliding charge-density waves, perhaps seen in quasi-one-dimensional metals like NbSe_3 . A more common failure is that in metals which are very impure or which have very strong scattering, the mean free path between scattering events may become so short that the wavelength and wavevector of electron quantum states is no longer definable. This undermines the basis of Bloch's theory. See CHARGE-DENSITY WAVE.

Such a failure of the theory (technically, a failure of the quasiparticle picture) is illustrated in Fig. 4, showing the resistivity of $\text{Ti}_{1-x}\text{Al}_x$ alloys. At the lower aluminum concentrations, $x = 3\%$ and 6% , the $\rho(T)$ curves shift upward, obeying Matthiessen's rule, given by Eq. (14), where the residual resistivity

$$\rho = \rho_0 + \rho_{\text{pure}}(T) \quad (14)$$

ρ_0 is proportional to impurity concentration x . For the larger concentrations, a too short mean free path causes the quasiparticle picture underlying Bloch-Boltzmann theory to fail. This in turn causes failure of Matthiessen's rule. The quasiparticle picture also fails at high temperature for all concentrations x . See MATTHIESSEN'S RULE.

A more abstract formulation of nonequilibrium effects like conductivity exists under the name of Kubo formulas. These provide a route to improved theories.

Philip B. Allen

Bibliography. N. W. Ashcroft and N. D. Mermin, *Solid State Physics*, Saunders, Philadelphia, 1976; J.

Bass, Electrical resistivity of pure metals and dilute alloys in *Landolt-Börnstein Tables 15a*, Springer, Berlin, 1980; O. Gunnarsson, M. Calandra, and J. E. Han, Colloquium: Saturation of electrical resistivity, *Rev. Mod. Phys.*, 75:1085–1099 2003; C. Kittel, *Introduction to Solid State Physics*, 8th ed., Wiley, New York, 2005; J. M. Ziman, *Electrons and Phonons*, Oxford, 1960, reprint 2001.

Electrical connector

A device that joins electrical conductors mechanically and electrically to other conductors and to the terminals of apparatus and equipment. The term covers a wide range of devices designed, for example, to connect small conductors employed in communication circuits, or at the other extreme, large cables and bus-bars.

Electrical connectors are applied to conductors in a variety of ways. Soldered connectors have a tube or hole of approximately the same diameter as the conductor. The conductor and connector are heated, the conductor inserted, and solder flowed into the joint until it is filled. Solderless connectors are applied by clamping the conductor or conductors in a bolted assembly or by staking or crimping under great mechanical force, usually by means of special tools designed for the purpose.

Industrial and power types. Connectors for electric power systems are commonly cast and assembled with bolts. The material most commonly used is copper alloy, which may vary in composition and properties depending upon the intended use of the connector. Where heavy currents flow through the body of the connector, an alloy of relatively high conductivity is used. Where the connector serves principally to clamp conductors together, an alloy of higher strength and lower conductivity may be used. An important consideration governing the choice of alloy is that it must have a coefficient of thermal expansion very close to that of the conductor material itself (copper or aluminum), so that the connection will remain secure through wide temperature changes.

Connectors designed for severe outdoor service and heavy current and mechanical loading are called power connectors and are widely used in power substations for cable-to-stud connections at equipment, cable-to-bus, bus-to-bus, and bus-to-line connections.

For industrial application, more compact constructions and greater versatility may be required, and connectors with such features are called industrial connectors. In practice, both classes are employed in power systems and industrial plants, the choice depending upon the actual service conditions.

Typical connector types are in-line splice couplers, T-tap connectors, terminal lugs, and stud connectors. Couplers join conductors end to end. T-tap connectors join a through conductor to another conductor at right angles to it (Fig. 1a). Terminal lugs join the conductor to a drilled tongue for bolting to the terminals of equipment (Fig. 1b). Stud connectors join

the conductor to equipment studs; the stud clamp is threaded or smooth to match the stud. Many variations of these types are made.

Split-bolt connectors are a compact construction widely used for splices and taps in building wiring. The bolt-shape casting has a wide and deep slot lengthwise. The conductors are inserted in the slot and the nut is drawn up, clamping the conductors together inside the bolt (Fig. 1c). Newer types of connectors are made to serve this same function.

Expansion connectors or flexible connectors allow some limited motion between the connected conductors. The clamp portions of the connector are joined by short lengths of flexible copper braid and may also be held in alignment by a telescoping guide.

Heavily tinned copper-alloy bolted connectors or aluminum-alloy-body connectors may be used to connect aluminum conductors. They are applied with an oxide-inhibiting compound. Where copper and aluminum conductors are to be clamped together, they are separated by bimetal or tinned washers to prevent galvanic corrosion.

In electronic equipment assembly, involving hookup of wires to terminals and of other components onto printed circuit boards, connections are made by soldering the wire or component to the terminal or printed circuit board. Although soldered terminals are by far the most common type of wiring connections in electronic chassis applications and in manufacture of electronic and electrical equipment, the pressure-wrap connection is a tool-applied method used in the field that speeds up connections where very many have to be made, as in wiring telephone and signal and control equipment. The tool that applies the wire wrap to the terminal stud assures an extremely tight, low-resistance, high-reliability connection by slightly elongating the wire as it is wrapped, putting a tight constrictive force on the terminal.

Separable types. These consist of matched plugs and receptacles, which may be readily separated to disconnect a conductor or group of conductors from the circuit or system. Separable connectors are commonly used for the connection of portable appliances and equipment to an electric wiring system. In electronics and communication they are used for connecting various components so that a single component can be readily removed for servicing without disturbing the circuitry.

For connections to an electric power supply, separable connectors consist of receptacles, usually fixed and permanently connected to the wiring system,

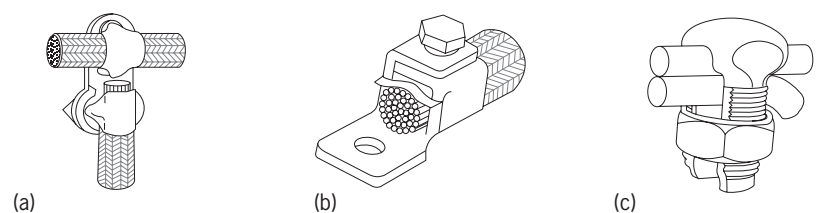


Fig. 1. Types of connectors. (a) T-tap connector. (b) Terminal lug. (c) Split-bolt connector.

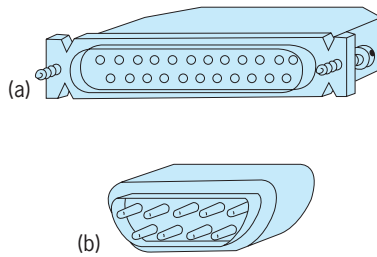


Fig. 2. Typical subminiature computer communications connectors. (a) 25-pin RS-232 plug. (b) 9-pin RS-422 jack.

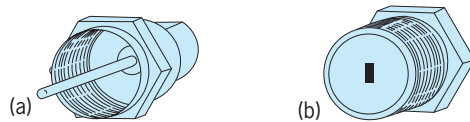


Fig. 3. Typical coaxial cable connection: 5/8 in., 24 threads per inch. (a) Jack. (b) Plug.

and attachment plugs, usually connected to a cord or cable extending to the appliance or equipment served. Other constructions provide cable-to-cable and cable-to-appliance connections.

Connectors for power systems are typically two-wire, three-wire, or four-wire, depending upon the circuit requirements. They are made in a wide range of types and sizes for applications involving various voltage, current, mechanical, and environmental requirements. They are so constructed that a particular receptacle will accept only a matching plug of identical characteristics, rating, and circuit configuration. Two-wire connectors that can be coupled in only one position are said to be polarized. Multiconductor connectors are designed, by a nonuniform arrangement of contacts or by a key in the assembly housing, so that they may be coupled in only one position.

Locking types. These are designed so that, when coupled, they may not be separated by a strain or pull on the cord or cable. In a typical construction the plug is inserted and twisted through a small arc, locking it securely in place. For hazardous areas, where the spark caused by separating a connector under load may be dangerous, connectors are designed so that the circuit is made or broken within the receptacle only in the fully engaged and locked position. Receptacles may also be interlocked mechanically with a switch, allowing insertion or removal of the plug only when the switch is in the open position.

Multewire connectors used in electronic, communication, and control circuits contain more connections of much smaller current capacity than power connectors (Fig. 2). Connector constructions are designed primarily for mechanical and electrical contact reliability.

Coaxial connectors provide for connecting and disconnecting a coaxial cable from an equipment circuit (Fig. 3). They are so constructed as to maintain the conductor configuration, through the separable connection, and the characteristic impedance of the coaxial cable. Coaxial connectors are made of copper or aluminum with contact and current-carrying

surfaces plated with silver or gold. They are made in a variety of sizes adaptable to the cables with which they are used. They are designed for cable-to-chassis and cable-to-cable connections in plug, jack, receptacle, feed-through, right-angle, and T-adaptor designs. See COAXIAL CABLE.

Common receptacles. Plug receptacles, sometimes called convenience outlets, are a type of wiring device distributed throughout buildings for the attachment of cord plug caps from portable lamps and appliances. In residences at least one such outlet must be provided for every 12 linear feet (3.66 m) or major fraction of wall perimeter. The receptacle is conventionally two-wire, 120-V, 15-A, and supplied by a 15- or 20-A branch circuit. The slots which accept the attachment plug are parallel. The device may consist of a single receptacle, a duplex receptacle, or a triplex receptacle (for one, two, or three pairs of blade slots for plug attachment, respectively) in a single outlet box. More attachment plugs can be provided by ganging two or more devices in a common outlet box under a common wall plate. Receptacles may also be combined with other devices (as switches and pilot lights).

In split-circuit duplex receptacles, each half of the device is supplied independently through a separate terminal, with typically one part connected directly to the branch circuit conductors and the other connected through a switch. Grounding receptacles have an additional contact that accepts the third round or U-shaped prong of a grounding attachment plug (Fig. 4). Grounding receptacles will accept either grounding or nongrounding attachment plugs,

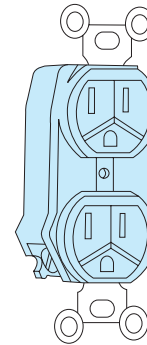


Fig. 4. Grounding-type 15-A receptacle.

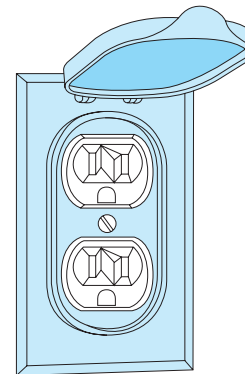


Fig. 5. Weatherproof receptacle with hinged cover.

but grounding attachment plugs can be used only in grounding receptacles. The grounding contact of the receptacle is connected internally to a supporting strap providing a ground both through the outlet box and the grounding conductor, armor, or raceway of the wiring system. See GROUNDING.

Multioutlet assemblies are a type of surface raceway designed to hold conductors and plug receptacles, which are usually spaced uniformly along the length of assembly.

Weatherproof receptacles (Fig. 5) are equipped with a weatherproof plate, gasketed around the edges, with a hinged or threaded cover over the device. See WIRING.

J. F. McPartland; Brian J. McPartland

Bibliography. T. Croft and W. I. Summers (eds.), *American Electricians Handbook*, 13th ed., 1996; IEEE, *National Electrical Safety Code 1997*, 1996; J. F. McPartland, *The Handbook of Practical Electrical Design*, 3d ed., 1999; J. F. McPartland, *National Electrical Code Handbook*, 23d ed., 1999.

Electrical degree

A time interval equal to 1/360 of the time required for one complete cycle of alternating current. Mechanical rotation is often measured in degrees, 360° constituting one complete revolution. In describing alternating voltages and currents, the time for one complete cycle is considered to be equivalent to 360 electrical degrees (360°) or 2π electrical radians. For example, if the frequency *f* is 60 cycles per second (60 Hz), 360° corresponds to 1/60 second and 1 electrical degree to 1/12,600 second.

There is a definite relationship between electrical and mechanical degrees in rotating electric generators and motors. Figure 1 shows typical coil and angular relationships in a two-pole alternator. As the magnetic field in the machine moves relative to the coils in the armature winding, the coils are linked se-

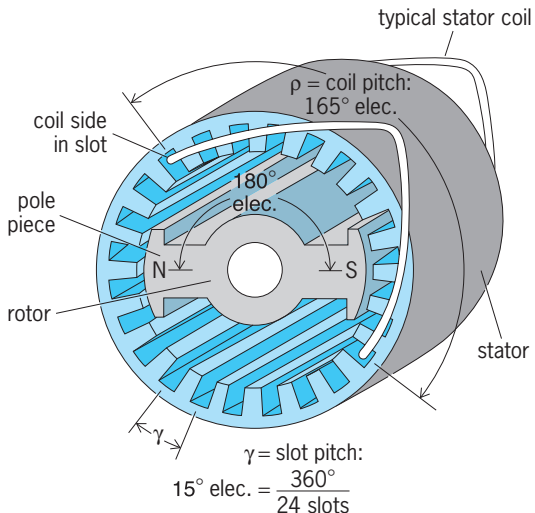


Fig. 1. Coil and angular relationships in a two-pole alternator.

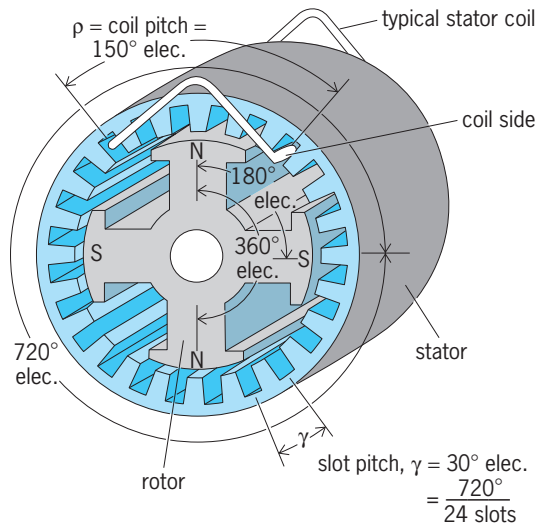


Fig. 2. Coil and angular relationships in a four-pole alternator.

quentially by the fluxes of north and south magnetic poles; two flux reversals induce one cycle of voltage in a given coil. Thus, in a two-pole machine such as in Fig. 1, 360° of electrical cycle corresponds to 360° of mechanical rotation, and an angle measured in mechanical degrees has the same value in electrical degrees. However, in a machine with more than two poles, one electrical cycle is generated per pair of poles per revolution. For example, a six-pole machine generates three cycles of voltage in each armature coil per revolution. In this case, each mechanical degree is equivalent to 3 electrical degrees. The four-pole case is illustrated in Fig. 2. In general, the relationship below is valid, where *p* is the number

$$\text{Number of electrical degrees in a given angle} = \frac{p}{2} \cdot \left(\text{number of mechanical degrees in that angle} \right)$$

of magnetic poles of either the rotor or the stator. It follows that the electrical angle between the centers of succeeding poles of opposite polarity is always 180 electrical degrees.

The concept of electrical degrees simplifies the analysis of multipolar machines by allowing them to be analyzed on a two-pole basis. Furthermore, it permits trigonometry to be used in solving alternating-current problems. See ALTERNATING CURRENT; ELECTRIC ROTATING MACHINERY; GENERATOR; MOTOR; WINDINGS IN ELECTRIC MACHINERY.

George McPherson, Jr.

Electrical energy measurement

The measurement of the integral, with respect to time, of the power in an electric circuit. The absolute unit of measurement of electrical energy is the joule, or the charge in coulombs times the potential difference in volts. The joule, however, is too small (1 watt-second) for commercial use, and the more commonly used unit is the watt-hour (3.6 × 10³ joules). The most common measurement

application is in the utility field. See ELECTRICAL UNITS AND STANDARDS.

Electrical energy is one of the most accurately measured commodities sold to the general public. Many methods of measurement, with different degrees of accuracy, are possible. The choice depends on the requirements and complexities of the measurement problems. Basically, measurements of electrical energy may be classified into two categories, direct-current power and alternating-current power. The fundamental concepts of measurement are, however, the same for both.

Methods of measurement. There are two types of methods of measuring electrical energy: electric instruments and timing means, and electricity meters.

Electric instruments and timing means. These make use of conventional procedures for measuring electric power and time. The required accuracy of measurement dictates the type and quality of the measuring devices used (for example, portable instruments, laboratory instruments, potentiometers, stopwatches, chronographs, and electronic timers). Typical methods are listed below. See ELECTRIC POWER MEASUREMENT.

1. Measurement of energy on a direct-current circuit by reading the line voltage and load current at regular intervals over a measured period of time. The frequency of reading selected (such as one per second, one per 10 s) depends upon the steadiness of the load being measured, the time duration of the test, and the accuracy of measurement desired. The first dc electric energy meter was an electroplating cell in series with the load. It deposited a mass of metal on an electrode exactly proportional to the total charge transported to the load. The electrode was weighed periodically to determine the energy used. Errors were introduced if line voltage was not constant. It was replaced by more convenient instruments. See CURRENT MEASUREMENT; VOLTAGE MEASUREMENT.

In electrical energy measurements, the losses in the instruments must be considered. Unless negligible from a practical standpoint, they should be deducted from the total energy measured. If the voltmeter losses are included in the total energy measured, then watt-hours = $(VI - V^2R)t/3600$, where V is the average line voltage (volts), I is the average line current (amperes), R is the voltmeter resistance (ohms), and t is the time (seconds).

2. Measurement of energy on a direct-current circuit by controlling the voltage and current at constant predetermined values for a predetermined time interval. This method is common for controlling the energy used for a scientific experiment or for determining the accuracy of a watt-hour meter. For best accuracy, potentiometers and electronic timers are desirable.

3. Measurement of energy on an alternating-current circuit by reading the watts input to the load at regular intervals over a measured period of time. This method is similar to the first, except that the power input is measured by a wattmeter.

4. Measurement of energy on an alternating-current circuit by controlling the voltage, current, and watts input to the load at constant predetermined values. This method is similar to the second, except that the power input is measured by a wattmeter. A common application of this method is to determine the standard of measurement of electrical energy, the watt-hour.

5. Measurement of energy by recording the watts input to the load on a linear chart progressed uniformly with time. This method makes use of a conventional power record produced by a recording wattmeter. The area under the load record over a period of time is the energy measurement.

Electricity meters. These are the most common devices for measuring the vast quantities of electrical energy used by industry and the general public. The same fundamentals of measurement apply as for electric power measurement, but in addition the electric meter provides the time-integrating means necessary to electric energy measurement.

A single meter is sometimes used to measure the energy consumed in two or more circuits. However, multistator meters are generally required for this purpose. Totalization is also accomplished with fair accuracy, if the power is from the same voltage source, by paralleling secondaries of instrument current transformers of the same ratio at the meter. Errors can result through unbalanced loading or use of transformers with dissimilar characteristics.

Watt-hour meters are generally connected to measure the losses of their respective current circuits. These losses are extremely small compared to the total energy being measured and are present only under load conditions.

Other errors result from the ratio and phase angle errors in instrument transformers. With modern transformers these errors can generally be neglected for commercial metering. If considered of sufficient importance, they can usually be compensated for in adjusting the calibration of the watt-hour meter. For particularly accurate measurements of energy over short periods of time, portable standard watt-hour meters may be used. Errors may also arise due to integral-cycle control. This is a well-established method of power control in which a silicon-controlled rectifier circuit acts to turn the power on for several cycles and off for a different number of cycles. The main source of error in the induction meter is the difference between the mechanical time-constants for two on-off intervals, making the meter read high.

Watt-hour meters used for the billing of residential, commercial, and industrial loads are highly developed devices. Over the last years many significant improvements have been made, including improvements in bearings, insulating materials, mechanical construction, and sealing techniques which exclude dust and other foreign material. See WATT-HOUR METER.

Automatic remote reading. Various aspects of the energy crisis spurred active development of automatic meter-reading systems with the functional

capability of providing meter data for proposed new rate structures, for example, time-of-day pricing, and initiating control of residential loads, such as electric hot-water heaters.

Automatic meter-reading systems under development generally consist of a utility-operated, minicomputer-controlled reading center, which initiates and transmits commands over a communication system to a terminal at each residential meter. The terminal carries out the command, sending the meter reading back to the reading center, or activating the control of a residential load.

Several communication media have been proposed and tested by system developers, including radio, CATV, use of the existing subscriber phone lines, and communication over the electric distribution system itself.

Quantities other than watt-hours. Included in the field of electrical energy measurement are demand, var hours, and volt-ampere hours.

Demand. The American National Standards Institute defines the demand for an installation or system as "the load which is drawn from the source of supply at the receiving terminals, averaged over a suitable and specified interval of time. Demand is expressed in kilowatts, kilovolt-amperes, amperes, kilovars and other suitable units" (ANSI C12.1-1988).

This measurement provides the user with information as to the loading pattern or the maximum loading of equipments rather than the average loading recorded by the watt-hour meter. It is used by the utilities as a rate structure tool.

Var hour. ANSI defines the var hour (reactive volt-ampere hour) as the "unit for expressing the integral, of reactive power in vars over an interval of time expressed in hours" (ANSI C12.1-1988).

This measurement is generally made by using reactors or phase-shifting transformers to supply to conventional meters a voltage equal to, but in quadrature with, the line voltage.

Volt-ampere hour. This is the unit for expressing the integral of apparent power in volt-amperes over an interval of time expressed in hours. Measurement of this unit is more complicated than for active or reactive energy and requires greater compromises in power-factor range, accuracy, or both. Typical methods include: (1) Conventional watt-hour meters with reactors or phase-shifting transformers tapped to provide an in-phase line voltage and current relationship applied to the meter at the mean of the expected range of power-factor variation. (2) A combination of a watt-hour and a var-hour meter mechanically acting on a rotatable sphere to add vectorially watt-hours and var-hours to obtain volt-ampere hours, volt-ampere demand, or both.

Measurement of volt-ampere hours is sometimes preferred over var-hours because it is a more direct measurement and possibly gives a more accurate picture of the average system power factor. This would not necessarily be true, however, where simultaneous active and reactive demand are measured and recorded. *See* ELECTRICAL MEASUREMENTS.

William H. Hartwig

Bibliography. D. A. Bell, *Electronic Instrumentation and Measurement*, 2d ed., 1994; *Code for Electrical Metering*, ANSI C12.11988, 8th ed., 1988; D. Fink and H. W. Beaty (eds.), *Standard Handbook for Electrical Engineers*, 14th ed., 1999; Instrument Society of America, *ISA Standards and Practices for Instrumentation*, 12th ed., 1995; L. Schnell (ed.), *Technology of Electrical Measurements*, 1993; L. M. Thompson, *Electrical Measurements and Calibration: Fundamentals and Applications*, 2d ed., 1994.

Electrical engineering

The branch of engineering that deals with electric and magnetic forces and their effects. *See* ELECTRICITY; ENGINEERING; MAGNETISM.

Scope. Electrical engineers design computers and incorporate them into devices and systems. They design two-way communications systems such as telephones and fiber-optic systems, and one-way communications systems such as radio and television, including satellite systems. They design control systems, such as aircraft collision-avoidance systems, and a variety of systems used in medical electronics. Electrical engineers are involved with generation, control, and delivery of electric power to homes, offices, and industry. Electric power lights, heats, and cools working and living space and operates the many devices used in homes and offices. Electrical engineers analyze and interpret computer-aided tomography data (CAT scans), seismic data from earthquakes and well drilling, and data from space probes, voice synthesizers, and handwriting recognition. They design systems that educate and entertain, such as computers and computer networks, compact-disk players, and multimedia systems. *See* CHARACTER RECOGNITION; COMMUNICATIONS SATELLITE; COMPACT DISK; COMPUTER; COMPUTERIZED TOMOGRAPHY; CONTROL SYSTEMS; DIGITAL COMPUTER; ELECTRIC HEATING; ELECTRIC POWER GENERATION; ELECTRIC POWER SYSTEMS; ELECTRIC POWER TRANSMISSION; ELECTRICAL COMMUNICATIONS; ILLUMINATION; MULTIMEDIA TECHNOLOGY; OPTICAL COMMUNICATIONS; OPTICAL FIBERS; RADIO; TELEPHONE SERVICE; TELEVISION; VOICE RESPONSE.

Early development. Early experimenters developed batteries, electromagnets, and electric motors. In 1832, M. Faraday showed that electric and magnetic effects are not separate, and combined the results of various observations into a single theory of electromagnetic induction. The first major commercial application of electromagnetism was the telegraph, demonstrated in 1844. *See* BATTERY; ELECTROMAGNET; ELECTROMAGNETIC INDUCTION; MOTOR; TELEGRAPHY.

The last quarter of the nineteenth century saw rapid developments in lighting, stemming from the invention of the incandescent lamp, and in electric power generation and delivery. During this period, lighting moved to homes, offices, and schools, and motors powered new forms of transportation and began to transform factories. A major technical

debate in this period was whether to generate and distribute electric energy as direct current or alternating current. After the invention of the transformer in 1883, alternating-current generation and delivery quickly became the dominant system, because of its greater efficiency and versatility. *See* ALTERNATING CURRENT; DIRECT CURRENT; ELECTRIC ROTATING MACHINERY; INCANDESCENT LAMP; TRANSFORMER.

The latter nineteenth and early twentieth centuries saw the development of the telephone and wireless telegraphy. Telegraph wires provided communications for railroads and nearly instantaneous message capability for business, industry, and personal data, while the telephone industry grew rapidly. Wireless telegraphy linked continents with high-powered spark transmitters.

Electronics. The audion, invented in 1906 and later called the triode, was a vacuum tube (or valve) with three electrical elements. The triode was used to design and build amplifiers, a basic element in modern electronics. Amplifiers enable a weak information signal to control a local power source, yielding a more powerful output. Amplification and other applications of tubes enabled the rapid development of the electronics industry. The junction transistor, whose invention in 1948 followed that of the point-contact transistor in the previous year, became the dominant electronic device by the 1960s, replacing vacuum tubes. *See* AMPLIFIER; TRANSISTOR; VACUUM TUBE.

The integrated circuit, invented in 1959, quickly revolutionized electronics. This device or process enables designers to put many transistors and associated components into a single package. The number of electronic devices in a single integrated-circuit chip now routinely exceeds 1 million. These units have made the personal computer a reality, and are installed in home appliances, cars, offices, and factories. *See* ELECTRONICS; EMBEDDED SYSTEMS; INTEGRATED CIRCUITS; MICROCOMPUTER; MICROPROCESSOR.

Feedback. Electronic devices can produce, along with amplification, distorted output signals. An invention of the 1920s, called negative feedback, reduces this problem. A fraction of the amplifier output signal is subtracted from the input signal. While this subtraction reduces gain and signal level, it also reduces distortion. The idea illustrates the trade-offs that must be made throughout the design process, as competing requirements are balanced. While positive feedback can sometimes find application, it is often a serious problem in electronics. Positive feedback is responsible for the screeching sound in a public address system when a person inadvertently points a microphone at a loudspeaker. *See* DISTORTION (ELECTRONIC CIRCUITS); FEEDBACK CIRCUIT.

Control systems. Control systems use feedback. Common examples are thermostats that control furnaces or air conditioners. Automobile cruise controls and automatic pilots for aircraft are applications. Other examples are found in chemical processes and steel-rolling mills. In control systems, electronic circuits enable desired results to be compared with

actual results, and the feedback enables the reduction of the error to nearly zero. Electrical engineers work with other engineers in the design and development of control systems. *See* AUTOPILOT; PROCESS CONTROL; THERMOSTAT.

Computers. Electrical engineers have a major role in the continuing development and application of computers, both large and small. They work with computer engineers, computer scientists, and other professionals. Electrical engineers design the data-storage systems, commonly called disks (floppy or hard), the central processing unit (CPU), and the display. The disks are magnetic storage devices, and proper operation of the disks requires a sophisticated control system. The CPU, which is the heart of the computer, includes very large scale integrated (VLSI) circuits for computations, for temporary storage during computation, and for control of printers, disk drives, and displays. The display may be a large vacuum tube with a phosphorescent screen, whose design requires knowledge of electromagnetic theory and materials, or it may be a liquid-crystal display, a technology that has developed rapidly. Computer capability continues to double every 18–24 months. *See* COMPUTER STORAGE TECHNOLOGY; ELECTRONIC DISPLAY.

The computer-communication system known as the internet or worldwide web is a linkage of computers around the world. The internet uses telephone lines, and in turn the telephone system uses satellites, microwave communications, wire communications, and radio. This system allows users to exchange data and ideas rapidly, reliably and at very low cost. *See* DATA COMMUNICATIONS.

Optics. In optical communications, another major contribution of electrical engineering, light pulses and waves carry data over very small strands of glass. The glass replaces copper wire, and the new systems are less expensive, more reliable, and in many respects safer to the user than those they replace. Fiber-optic networks provide two-way video and audio links for universities, colleges, public schools, libraries, hospitals, and court houses.

Challenges. The integration of communications equipment, control systems, computers, and other devices and processes into reliable, easily understood, and practical systems is a major challenge, which has given rise to the discipline of systems engineering. Electrical engineering must respond to numerous demands, including those for more efficient and effective lights and motors; better communications; faster and more reliable transfer of funds, orders, and inventory information in the business world; and the need of medical professionals for access to medical data and advice from all parts of the world. *See* INFORMATION SYSTEMS ENGINEERING; MEDICAL INFORMATION SYSTEMS; SYSTEMS ENGINEERING.

Edwin C. Jones, Jr.

Bibliography. J. E. Brittain (ed.), *Turning Points in American Electrical History*, 1977; D. G. Fink and H. W. Beaty (eds.), *Standard Handbook for Electrical Engineers*, 14th ed., 1999; T. P. Hughes, *Networks of Power*, 1983, reprint 1993; J. D. Irwin and

D. V. Kerns, Jr., *Introduction to Electrical Engineering*, 1995; C. R. Paul, S. A. Nasar, and L. E. Unnewehr, *Introduction to Electrical Engineering*, 2d ed., 1992; J. D. Ryder, *Engineers and Electrons*, 1984.

Electrical impedance

The measure of the opposition that an electrical circuit presents to the passage of a current when a voltage is applied. In quantitative terms, it is the complex ratio of the voltage to the current in an alternating-current (ac) circuit.

The **illustration** shows a generalized ac circuit, which may be composed of the interconnection of various types of circuit elements. If the voltage is $v = V_0 \sin(\omega t + \theta)$ and the current is $i = I_0 \sin(\omega t + \beta)$, then the circuit can be analyzed by considering a complex voltage $V = V_0 e^{j\theta}$ and a complex current $I = I_0 e^{j\beta}$. The impedance of the circuit is then given by Eq. (1), and Z is a complex number given by Eq. (2). R , the real part of the impedance,

$$Z = \frac{V}{I} \quad (1)$$

$$Z = R + jX \quad (2)$$

is the resistance of the circuit, and X , the imaginary part of the impedance, is the reactance of the circuit. The units of impedance are ohms. See ELECTRICAL RESISTANCE; REACTANCE.

The modulus of the impedance is given by Eq. (3),

$$|Z| = \sqrt{R^2 + X^2} \quad (3)$$

and this is the ratio of the maximum voltage to the maximum current. The phase angle of the impedance is given by Eq. (4), and this is the angle by

$$\psi = \tan^{-1} \frac{X}{R} \quad (4)$$

which the current waveform lags the voltage waveform. The power factor of the circuit is $\cos \psi$ and is given by Eq. (5). The power dissipated in the circuit is given by Eq. (6).

$$\cos \psi = \frac{R}{\sqrt{R^2 + X^2}} \quad (5)$$

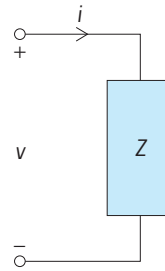
$$P = \frac{1}{2} |V| |I| \cos \psi \quad (6)$$

For a circuit composed of resistors, inductors, and capacitors, R is always positive, but X may be positive or negative. For example, in a circuit consisting of a resistor R and inductor L in series, the impedance at frequency ω is given by Eq. (7), so R is greater than

$$Z = R + jL\omega \quad (7)$$

0 and X is greater than 0. For a circuit consisting of a resistor R and capacitor C in series, the impedance at frequency ω is given by Eq. (8) so R is greater than

$$Z = R + \frac{1}{jC\omega} = R - j\frac{1}{C\omega} \quad (8)$$



Generalized alternating-current circuit.

0 and X is less than 0. See ALTERNATING-CURRENT CIRCUIT THEORY. J. O. Scanlan

Electrical instability

A condition that is generic to all physical systems that use the amplification of signals, power, or energy in various forms. Electrical circuits are perhaps the best-known example because the use of amplification is very common there. In order for instability to occur, a closed (circular) path must exist, and the net amplification around it must be sufficiently large that any random signal (for example, noise) that travels around the closed path grows uncontrollably. The condition is eventually limited by some physical aspect of the system such as nonlinearity, saturation, or power available. The energy to create the uncontrollable signal growth is obtained by conversion from some other form of energy. For example, a high-frequency electrical oscillation can draw its energy from a direct-current (dc) power supply. See ENERGY; LINEARITY; SATURATION.

There can be various degrees of stability and instability in an electrical circuit or system, and instability may be desirable or undesirable depending on the context. Desirable instability is displayed by a cross-coupled positive-feedback latch in a computer memory cell, or an oscillator that is used to create the high-frequency carrier wave of a radio transmitter. An undesirable instability might arise in the loss of control over a power transmission system due to surges on the generators, leading to a necessary decoupling of the system. Another example is the high-pitched whine heard when a radio-station telephone caller leaves the home radio receiver on and thus allows feedback from the radio receiver to the telephone, that is, a closed-loop signal path. See FEEDBACK CIRCUIT; OSCILLATOR.

Yet another example is a computer simulation program that fails because of lack of convergence when trying to solve for the signal levels in a circuit. This could be due to a numerical instability in the computer algorithm. See ALGORITHM; SIMULATION.

Electrical instability can occur at zero frequency (dc), for example, in the cross-coupled memory latch, but its accidental occurrences are generally at high frequencies. A configuration may be designed to be stable at low frequencies, but inadvertent phase shifts or time delays at higher frequencies may

create a positive-feedback situation. For this reason, the design of amplifiers, feedback control systems, and so forth necessarily involves much attention to the higher-frequency aspects even though the desired use may be primarily at lower frequencies. *See* AMPLIFIER; CONTROL SYSTEMS.

The intellectual basis for the design of electrical circuits and systems to ensure their stability is quite well developed, based on linear mathematics. Computer tools are extensively used in complex application areas. However, the situation is less satisfactory in the design of deliberately unstable systems, such as oscillators, because their behavior is finally explained by less tangible nonlinear limiting mechanisms. *See* CONTROL SYSTEM STABILITY; LINEAR SYSTEM ANALYSIS.

Miles A. Copeland

Electrical insulation

A nonconducting material that provides electric isolation of two parts at different voltages. To accomplish this, an insulator must meet two primary requirements: it must have an electrical resistivity and a dielectric strength sufficiently high for the given application. The secondary requirements relate to thermal and mechanical properties. Occasionally, tertiary requirements relating to dielectric loss and dielectric constant must also be observed. A complementary requirement is that the required properties not deteriorate in a given environment and desired lifetime. *See* CONDUCTOR (ELECTRICITY).

Electric insulation is generally a vital factor in both the technical and economic feasibility of complex power and electronic systems. The generation and transmission of electric power depend critically upon the performance of electric insulation, and insulation now plays an even more crucial role because of the energy shortage.

Insulation Requirements

The important requirements for good insulation will now be discussed.

Dielectric strength. The basic difference between a conductor and a dielectric is the availability of charge carriers with high mobility in a conductor, whereas free charge has little mobility in a dielectric. Dielectric strength is a measure of the electric field required to cause sudden and substantial motion of charge across the surface or through the bulk of a dielectric. Dielectric strength usually deteriorates with increased moisture content and with elevated temperature. For high-voltage (on the order of kilovolts) applications, dielectric strength is the most important single property of the insulation.

Steady-state strength. The voltage at which a sudden high-current discharge punctures an insulator is called the breakdown voltage. In a uniform field this voltage divided by the electrode gap is called the dielectric strength, the breakdown field strength, the electric strength, or the breakdown stress. In a non-uniform field it is important to differentiate between the maximum field and the average field at breakdown. *See* BREAKDOWN POTENTIAL; ELECTRICAL BREAKDOWN.

An insulator with an operational dielectric strength of 4–12 kV/mm (100–300 V/mil) is considered good. Many of the insulators listed with a dielectric strength more than an order of magnitude higher than this are ordinarily stressed at these lower values in practical operation. The reason is that imperfections in the form of voids, cracks, filamentary defects, and so forth occur whenever long lengths and thick sections of insulating material are manufactured.

Whenever two dielectrics are in series in an alternating electric field, the electric field E is higher in the medium of lower dielectric constant κ by the ratio of dielectric constants, $E_2/E_1 = \kappa_1/\kappa_2$. Therefore the electric field is higher in a cavity, which generally has a dielectric constant of 1, by a factor of 2–3, than it is in the insulator. Furthermore, the dielectric strength is much less in the defect than in the insulator (**Table 1**; **Fig. 1**). Thus, in addition to operating insulators at much lower than intrinsic dielectric strength, schemes such as oil impregnation are used

TABLE 1. Properties of various good insulators

Insulator	Resistivity, 10^{14} ohm-cm	Dielectric strength, 10^3 V/mil*	Power factor $\times 10^{-4}$	Dielectric constant	Tensile strength, 10^3 lb/in. ² †	Chemical resistance	Flammability
Lexan	10–1000	8–16	10–30	3	8	Good	Self-extinguishing
Kapton H	1000–10,000	3–8	20–50	3	22	High	Self-extinguishing
Kel-F	10–1000	2–6	20–40	2–3	5	High	No
Mylar	0.1–1	4–16	30–120	3	20	Good	Yes
Parylene	10–10,000	6–10	2	2–3	10	Very good	Yes
Polyethylene N	100–500	1–17	3–30	2	4–6	Good	Yes
Polytetrafluoroethylene	0.1–1000	1–7	1–4	2	4–5	High	Very low
Air (1 atm, 1 mm gap)	—	0.1	10^{-3}	1	0	Stable	No
Sulfur hexafluoride (1 atm, 1 mm gap)	—	0.2–0.3	10^{-3}	1	0	Stable	No
Vacuum ($<10^{-5}$ torr or 1.3×10^{-3} Pa, 1 mm gap)	—	2–3	0	1	0	Stable	No

* 10^3 V/mil = 4×10^5 V/cm.

† 10^3 lb/in.² = 6.9 (MPa).

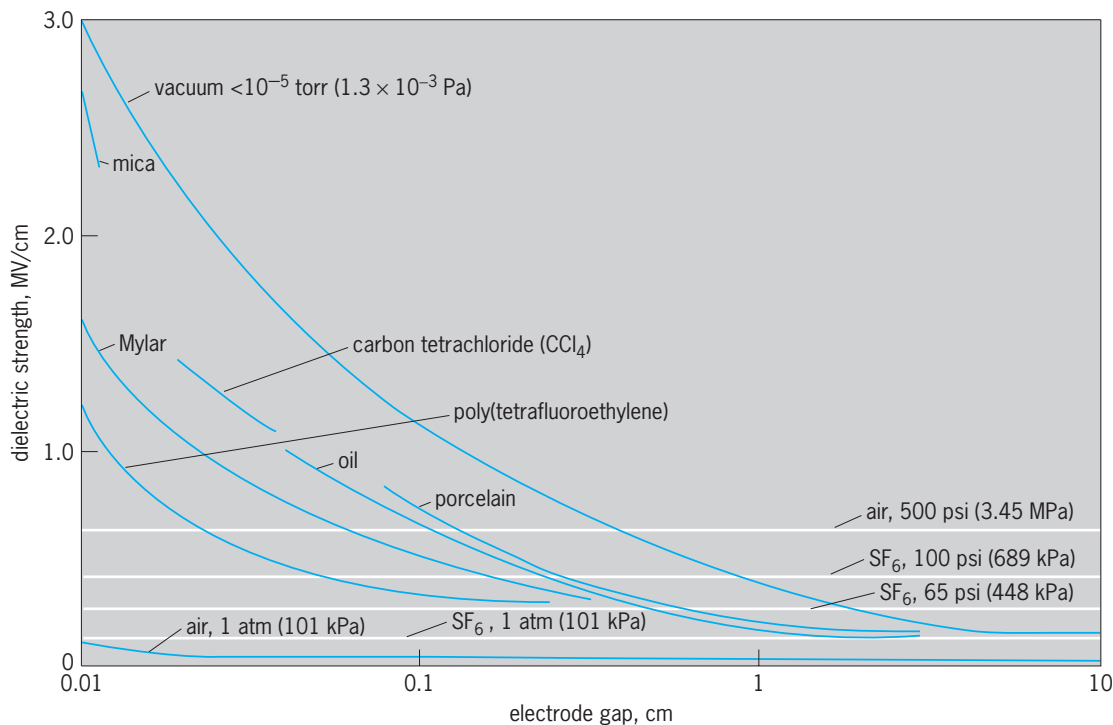


Fig. 1. Uniform field dielectric strength for solid, liquid, gas, and vacuum insulation.

when possible to fill the defects to increase both their dielectric constant and their dielectric strength.

Although solid extruded dielectric cable is now employed up to at least 420 kV system voltage and operates at average fields in the range of 15 kV/mm, oil-impregnated-paper or paper-polypropylene laminate (PPLP) wound-tape insulation can be used at even higher voltages. The probability that defects in the tape will line up in proximity is much less than in a solid piece of the same insulation material. The tape thickness also limits the size of the defect in the direction of the electric field. See DIELECTRIC MATERIALS.

Laboratory testing of insulation over a wide range of conditions is conducted to determine the effects of aging, high temperature, heat cycling, ultraviolet radiation, moisture, and so forth. A statistical analysis of the data based upon the weak-link Weibull statistics is then used to determine safe operating voltages and electric fields for the required lifetimes and environments. Operational withstand stress is sometimes arbitrarily defined as being at least two standard deviations below the average dielectric strength in the laboratory.

Impulse strength. Impulse strength is a measure of an insulator's ability to withstand a rapid rise and rapid fall (measured in microseconds) of electric stress. In addition to steady-state dielectric strength, an insulator must be able to meet a basic impulse level (BIL) and a basic switching level (BSL) surge (not as high a rate of voltage increase as for BIL) requirement. The impulse breakdown electric field strength is greatest for solids, followed by liquids, gases, and vacuum. For example, a development in tape insulation consisting of a combination of cellulose paper bound to

polypropylene film or fibers [known as polypropylene laminated paper (PPLP)] and impregnated with oil has an impulse strength of 2.7 MV/cm, whereas the impulse strength of sulfur hexafluoride from 100 to 400 kilopascals (15 to 60 lb/in.²) goes only from 0.2 to 0.45 MV/cm. Thus, in designing transmission cables, sulfur hexafluoride cables are limited by their BIL, whereas most solid insulation cables are limited only by their steady-state dielectric strength. This is why solid insulation cables are so much more compact than sulfur hexafluoride cables for the same ratings, in addition to the fact that they have higher steady-state dielectric strength. However, when a solid dielectric breaks down electrically, the punctured and surrounding regions must be replaced. After sulfur hexafluoride breaks down, the gas is self-healing and recovers dielectric strength. If no damage has been incurred by the conductors, operation may resume.

The impulse strength for sulfur hexafluoride increases with increasing pressure but tends to level off above 400 kPa (60 lb/in.²). The impulse strength of liquids increases less rapidly with increasing pressure, leveling off above 2 megapascals (300 lb/in.²). For example, carbon tetrachloride's impulse strength goes from 2.2 to 2.6 MV/cm in going from 100 kPa to 2 MPa (15 to 300 lb/in.²). This pressure dependence for liquids tends to disappear for rates of field increase greater than 4 MV/cm- μ s. This is experimental evidence that a "bubble" mechanism (see below) cannot be operative in such short times.

Resistivity. Resistivity is a measure of how much current will be drained away from the conductor through the bulk (bulk resistivity) or along the surface (surface resistivity) of the dielectric.

An insulator with resistivity equal to or greater than 10^{13} ohm-cm may be considered good. While resistivity measurements of conductors are straightforward, measurements on insulators are much more difficult and are subject to a much wider range of values for a given material. A dc measurement of current is made after application of a uniform electric field for a specified period of time. A time variation from 1 s to 1 min can lead to orders-of-magnitude variation in the measured resistivity, and the "true" resistivity can be nearly impossible to determine.

Since a resistivity measurement is nondestructive, it can be used as an on-line quality and uniformity check during manufacturing processes. During actual operation in the field, a sharp decrease in resistivity usually signals an ensuing service failure. See ELECTRICAL RESISTIVITY; INSULATION RESISTANCE TESTING; NONDESTRUCTIVE EVALUATION.

Dielectric loss. When a dielectric is subjected to an alternating field, a time-varying polarization of the atoms and molecules in the dielectric is produced. The response time of polarization phenomena varies from nearly instantaneous for electronic polarization to seconds or even hours for some forms of orientational polarization. If the time constant for polarization results in a lag between the applied field and polarization, a loss results, of which the power factor is a measure. The importance of the power factor varies with application and tends to become more important at high temperatures, where the loss tends to increase and can cause thermal runaway, and in cryogenic systems, where the energy required to remove the heat that results from dielectric loss can be great.

Dielectric loss is a function of both temperature and frequency and generally decreases as these two variables decrease. However, it is not a monotonic function of these parameters, and at various discrete values of temperature and frequency quasi-resonances occur where the loss can be large. A sharp rise in dielectric loss, usually accompanied by a rapid temperature increase, is an indication of impending breakdown.

Neglecting these high-loss sites, **Table 2** indicates the power factor in the three temperature regions of general interest, 300 K or 80°F (conventional), 77 K or -321°F (cryoresistive), and 4 K or -452°F (superconducting). Values of less than 10^{-3} are good at 300 K, values less than 10^{-4} are good at 77 K, and values less than 10^{-5} are good at 4 K. With the increas-

ing need to deliver large blocks of power (greater than 3000 MW) for long distances (greater than 80 km or 48 mi) from remote energy sources, advanced power transmission must have insulation capable of meeting cryogenic requirements. A power factor of less than 10^{-5} at the operating cryogenic temperature is generally desirable, since the overall power loss is amplified by the refrigeration inefficiency, which may require as much as 500 W per watt of dissipation. See SUPERCONDUCTING DEVICES.

Ideally the power factor should be independent of electric field stress. However, owing to processes unrelated to polarization, such as electron emission and ionization, the dielectric loss tends to increase with increasing electric field, and this is reflected in an increasing power factor.

The power factor is the cosine of the phase angle ϕ between the voltage and current, or equivalently the ratio of resistance to impedance for the medium. Occasionally confusion arises regarding terminology when the dielectric loss is related to the complementary angle δ , which satisfies Eq. (1) for small δ . Thus

$$\text{Power factor} = \cos \phi = \sin \delta \doteq \tan \delta \doteq \delta \quad (1)$$

the power factor is often used interchangeably with the terms loss tangent and loss angle.

Dielectric constant. The dielectric constant, also known as the relative permittivity or specific inductive capacity (SIC), is a measure of the ability of the dielectric to become polarized, taken as the ratio of the charge required to bring the system to the same voltage level relative to the charge required if the dielectric were vacuum. It is thus a pure number, but is in fact not a constant, and may vary with temperature, frequency, and electric-field intensity. For highly polar materials such as water, the variation with temperature and frequency is dramatic, going from 80 at 300 K (80°F) to 3 at 77 K (-321°F), or from 80 at 60 Hz to 20 at 10 kHz.

In addition to the problem of intensification of the electric field in regions of relatively lower dielectric constant, a low-dielectric-constant insulation is desirable for two more reasons. In ac transmission cables, the lower the dielectric constant, the more the current and the voltage will be in phase. Thus, more usable power will be delivered, without the need for reactive compensation. Furthermore, in reducing the charging current (which is proportional to the dielectric constant), concomitant power losses (related to the square of the charging current) are also reduced. A high dielectric constant is desirable in capacitors, since the capacitance is proportional to it. See CAPACITANCE; DIELECTRIC MEASUREMENTS; PERMITTIVITY.

Dielectric absorption. When high voltage is applied to a solid dielectric, charge tends to be absorbed into the dielectric. When the voltage is removed, some of this charge comes back out of the dielectric over time. This property is known as dielectric absorption. For example, if high voltage is applied to a polymer film capacitor for several minutes, the

TABLE 2. Power factor of insulators at three temperature

Insulator	Power factor $\times 10^{-4}$		
	300 K (80°F)	77 K (-321°F)	4 K (-452°F)
Kapton H	20-50	5	0.5
Polyethylene (low density)	3-30	0.1	0.02
Polytetrafluoroethylene	1-4	0.9	0.02
Kraft paper	50	20	7
Vacuum ($<10^{-5}$ torr or 1.3×10^{-3} Pa)	0	0	0

capacitor is shorted for a few seconds, and then the capacitor is left open-circuit, a voltage will appear across the capacitor terminals as a result of the absorption charge which evolves from the capacitor film. A "perfect" capacitor, such as a vacuum capacitor, would have no dielectric absorption. Many capacitors used in electrical or electronic applications, including most high-voltage capacitors, are based on polymer films and do exhibit dielectric absorption. Poly(tetrafluoroethylene) [PTFE] and polystyrene capacitor films exhibit among the lowest dielectric absorption.

Insulator Properties

Table 1 summarizes the properties of various good insulators at or near 60 Hz and ambient temperature (about 300 K or 80°F). The range of values indicates not only differences between measurements of different investigators, but differences resulting from small variations in a parameter. For example, the dielectric strength depends on the thickness and area of the sample being tested. It generally decreases with increasing thickness and area. The voltage at which electrical breakdown occurs increases less than linearly with thickness. Thus the corresponding value of electric field decreases.

All insulators may be classified as either solid or fluid. Solid insulation is further divided into flexible and rigid types.

Flexible insulation. Flexible hydrocarbon insulation is generally either thermoplastic or thermosetting. Thermosets are initially soft, and can be extruded by using only pressure. Following heat treatment, when they return to ambient temperature, they are tougher and harder. After thermosetting, nonrubber thermosets are harder, stronger, and have more dimensional stability than the thermoplastics. Thermoplastics are softened by heating, and when cool become hard again. They are heat-extruded.

Thermoplastics. The first seven insulators listed in Table 1 are thermoplastics, as are nylon, polystyrene, and poly(vinyl chloride). Those listed in Table 1 are solid polymers that are used at ambient temperature and have good potential for cryogenic applications, and their monomeric structures are shown in Fig. 2. Polymers come in amorphous and ordered states which have different dielectric properties. The ordered state is frequently called "crystalline," but is not the same as the crystalline state of solids. See POLYMER.

Lexan is a trade name for a group of polyesters formed from carbonic acid, and generically called polycarbonate (PC). Polycarbonate has good electrical and mechanical properties, good dimensional stability, good resistance to creep, high distortion temperature, and ease of molding and extrusion. See POLYESTER RESINS.

Kapton H is a trade name for poly(pyromellitimide) (PPMI). PPMI has exceptional resistance to thermal degradation, because it is one of the polyimides that incorporates multiple bonds along the backbone of the chain. It has good mechanical strength and good

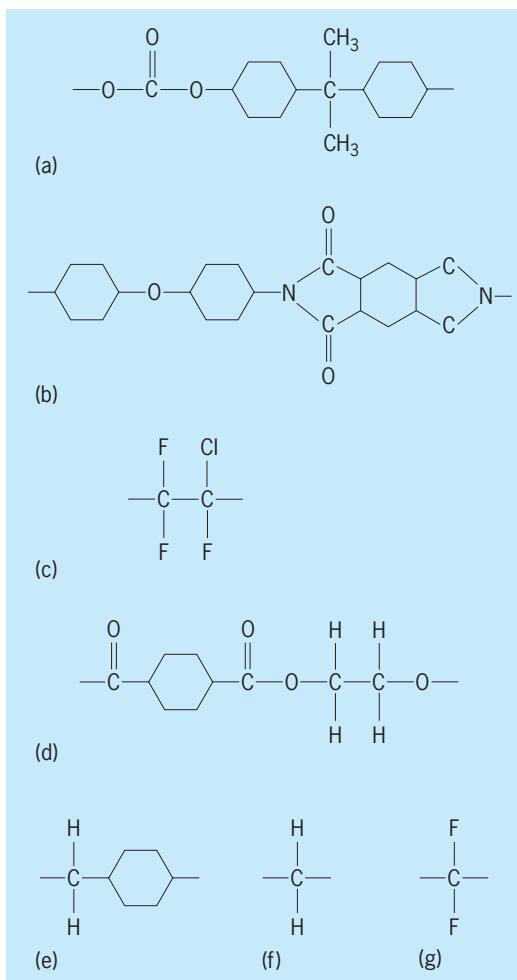


Fig. 2. Monomeric structures of polymers used as electrical insulators. (a) Lexan. (b) Kapton H. (c) Kel-F. (d) Mylar. (e) Parylene. (f) Polyethylene. (g) Polytetrafluoroethylene.

dimensional stability at both high and low temperatures, retaining ductility down to cryogenic temperatures. See HETEROCYCLIC POLYMER.

Kel-F is a trade name for poly(chlorotrifluoroethylene) [CTFE]. It is one of the best insulators in its ability to exclude moisture, in high-temperature strength, in chemical inertness, and in low coefficient of friction. See POLYFLUOROOLEFIN RESINS.

Mylar is a trade name for poly(ethylene terephthalate) [PET]. PET is a high-molecular-weight polyester with a stiff polymer chain and resilient interchain bonds that retain good mechanical properties beyond 150°C (300°F). Because of their toughness, PET films are used not only for electrical insulation, but also for magnetic tapes. See MAGNETIC RECORDING.

Parylene is a trade name for poly(*p*-xylylene) [PPX]. PPX is linear and highly regular, and was developed to be highly radiation-resistant. It has good dimensional stability down to cryogenic temperatures. It has not had widespread use owing to its cost. See POLY-PXYLYLENE RESINS.

Polyethylene (PE) is a widely used inexpensive electrical insulator with good mechanical and

excellent dielectric properties. Polyethylene can be cross-linked chemically or by irradiation to increase its mechanical strength and improve its high-temperature properties. Chemical cross linking leaves impurities which increase loss. Electron-beam cross linking results in lower loss. Cross-linked polyethylene (XLPE) power cables are manufactured for up to 420 kV system voltage using chemical cross linking, while XLPE high-frequency signal coaxial cables are made using electron-beam cross linking to reduce high-frequency loss. Although polyethylene has a higher dielectric strength than chemically cross-linked XLPE, its low melting point of 118°C (244°F) makes it less desirable for extruded dielectric cables than XLPE. Even after cross linking, melting of crystallites in XLPE results in a precipitous drop in yield stress above 90°C (194°F) and very large thermal expansion in the same temperature range. These characteristics make implementation of splices and terminations problematic for XLPE cables rated to 105°C (221°F). Ethylene propylene rubber (EPR) with inorganic filler is also be used in extruded solid dielectric power cables. Since the EPR polymer is nearly amorphous and the compound gains its physical properties from the inorganic filler, the mechanical and thermal properties of EPR cable dielectric are much more stable with temperature, both low and high, than XLPE, to the point that EPR cable can be rated to operate at 140°C (284°F). *See* POLYOLEFIN RESINS.

Polytetrafluoroethylene and its copolymer with hexafluoropropylene are highly ordered and orientable polymers made of chains of CF₂ with its strong C-F bond. This is why their structure involves little or no branching or cross linking. Thus this polymer has a low coefficient of friction, is stiff, inert, and insoluble, and has a high melting point of 327°C (621°F). Because of its high symmetry and tight bonding, it has one of the lowest known power factors and dielectric absorption for a solid, as well as other excellent dielectric properties. It may be processed by powder and sinter techniques as well as by extrusion or injection molding.

Polystyrene is a good insulator that is used in thin films or where it will not be stressed or bent, because it is stiff and brittle when it is thick. Fillers in the form of silica, alumina, mica, and so forth reduce the brittleness (probably by inhibiting crack propagation), improve the thermal conductivity, and reduce thermal expansion. Hydrated fillers, such as aluminum trihydrate, inhibit tracking, since the waters of hydration tend to be arc-quenching. Polystyrene capacitor film has among the lowest loss ($\sim 10^{-5}$) and dielectric absorption of all capacitor films. *See* POLYSTYRENE RESIN.

Thermosets. The various rubber compounds, neoprene, and epoxies are thermosets. Natural rubber, butyl rubber, buna S rubber, and silicone rubber are good insulators that are tough and flexible. Although chemically and mechanically tough, neoprene does not have good dielectric properties. The epoxies are not very flexible but are used at joints of flexible insulation because of their good bonding, electrical,

and mechanical properties. *See* POLYETHER RESINS; RUBBER.

Cellulose paper. Cellulose paper insulation is neither thermoplastic nor thermosetting. It is widely used in cables and rotating machinery in multilayers and impregnated with oil. It has a relatively high dielectric loss that hardly decreases with decreasing temperature (Table 2), which rules it out for cryogenic applications. Oil-impregnated cellulose paper has very high dielectric strength perpendicular to the paper (through the layers of oil-impregnated paper) but relatively low dielectric strength parallel to the paper surface. Because of its high dielectric strength, the high loss has not been a deterrent to its use in conventional ambient-temperature applications. However, the high dielectric strength deteriorates quickly if moisture permeates the paper. *See* CELLULOSE; PAPER.

Rigid insulation. This second major group of solid insulators includes glass, mica, epoxies, ceramoplastics, porcelain, alumina, other ceramics, and fiber-reinforced plastic (FRP). Rather than being used to insulate wires and cables except for mica, these materials are used in equipment terminations (potheads) and as support insulators (in tension or compression) for overhead lines whose primary dielectric is air. These rigid structures must be shock-resistant, be relatively water-impervious, and be able to endure corona discharges over their surfaces. *See* CORONA DISCHARGE.

For use as termination insulators, such as for transformers, underground transmission lines, reactors, and generators, both high bulk dielectric strength and high surface dielectric strength are important. For use in connection with equipment such as overhead transmission lines, the requirement for bulk dielectric strength is less severe, and the emphasis is on the surface. Surface dielectric strength of several kilovolts per inch is adequate, and the creepage paths are usually designed for approximately 400 V/cm (1 kV/in.).

In addition to normal electromagnetic operating forces and windage stress, these insulators must be able to withstand switching surges and shocks from short-circuit currents. These fault currents may be 10 times higher than the normal current, resulting in forces that are 100 times higher, since the force is proportional to the square of the current.

Porcelain is widely used for such applications. It is a hard, brittle material made by firing feldspar, quartz, clay, and other minerals at high temperature. Glass is a completely amorphous material containing 50–90% silica. Glass has a dielectric strength of about 2 MV/cm at 300 K (80°F). Mica is a mineral that occurs naturally in a laminated form. For thin sections of 0.002 cm thickness, mica has a dielectric strength as high as 10 MV/cm perpendicular to the direction of lamination. Fiber-reinforced plastics, often composites of glass fiber and epoxy, covered with silicone rubber are replacing glass and porcelain in many applications, including transmission-line insulators and high-voltage terminations. This structure has the advantages of much lower weight, improved

performance under polluted conditions as a result of the silicone rubber exterior insulation, and being explosion-proof when operated under pressure. See GLASS; MICA; PORCELAIN.

Fluid insulation. Liquids, gases, and vacuum fall in the category of fluid insulation. For all of these, the electrical structure must be such as to contain the fluid in the regions of high electric stress.

Liquids. The main types of insulating liquids are the mineral oils, silicones, chlorinated hydrocarbons, and the fluorocarbons with dielectric strengths on the order of megavolts per centimeter. Many other liquids also have good dielectric strength, such as carbon tetrachloride, toluene, hexane, benzene, chlorobenzene, alcohol, and even deionized water. However, special problems with liquids, such as chemical activity, high thermal expansion, thermal instability, low boiling point, and tracking (conducting carbon residue) after arcing, rule out many liquids. Even transformer oil in combination with the insulating paper in activated transformers produces water which must be filtered out.

Gases. Most gases have a dielectric constant of about 1, and low dielectric loss. **Table 3** compares the relative dielectric strengths of various gases. Those containing fluorine or chlorine are strongly electronegative. Of course, the temperature at which the gas condenses out of the gaseous state is an important consideration, as well as chemical stability, tracking, reactivity, and cost.

Air is used as a dielectric in a wide variety of applications, ranging from electronics to high-voltage (765-kV) and high-power (2000-MW) electric transmission lines. Dry air is a reasonably good insulator (Table 1). However, its dielectric strength decreases with increasing gap (Fig. 1).

Sulfur hexafluoride (SF₆; Table 1) has, at 1 atm (101 kPa), about three times the dielectric strength of air, carbon dioxide, or nitrogen, with this difference increasing with pressure and gap in uniform electric fields (Fig. 1). Sulfur hexafluoride is an electronegative gas which impedes electrical breakdown by capturing free electrons and forming negative ions. Not only does sulfur hexafluoride possess a

high dielectric strength, but also, as a result of very rapid recombination deionization after being decomposed in an arc, it recovers its dielectric strength after arcing much more rapidly than air. This property makes it 100 times as effective as air and other gases for use in power circuit breakers. It also is used frequently in electrical utility substations as an insulating medium.

Vacuum. Vacuum (that is, pressures of less than 10⁻³ Pa or 10⁻⁵ torr) has one of the highest dielectric strengths in gaps ranging from 0.1 to 1 mm (0.004 to 0.04 in.; Fig. 1; Table 1). However, as the gap increases, its dielectric strength decreases rapidly. A perfect vacuum should be a perfect insulator, since there would be no charge carriers present to contribute to electrical conductance. The limitation of vacuum as a dielectric results from the effects of a high electric field at the surface of electrodes in vacuum, which can cause field emission of electrons, rather than because a perfect vacuum is far from being realized in the laboratory.

The dielectric properties of vacuum can degenerate rapidly because vacuum offers no resistance to the motion of charge carriers once they are introduced into the vacuum region. The collision mean free path is of the order of meters at pressures less than 10⁻² Pa (10⁻⁴ torr), and this is below the pressure range for a gas discharge. However, avalanche processes occur at the electrodes, and vapor from the electrodes becomes the arcing medium.

Vacuum has found increasing applications as an insulator in electron tubes, photocells, high-frequency capacitors, electron microscopes, particle accelerators, distribution voltage (~15 kV) circuit breakers, and so forth. It is used as an arcing medium with exceptional rate of recovery of dielectric strength. The realization of its full potential is strongly linked to the constitution of the electrodes.

Nature of Electrical Breakdown

Electrical breakdown in gases first began to be understood in 1889, and breakdown in solids in 1935. Breakdown in liquids is less well understood. Breakdown in vacuum (in terms of an encompassing predictive theory) is least well understood, despite experimental investigations going back to 1897. See ELECTRICAL BREAKDOWN.

Gases. The Townsend avalanche criterion satisfactorily accounts for the threshold for electrical breakdown in nonelectronegative gases. This theory may be modified to include the process of electron capture by electronegative gases. When the space-charge field becomes large, other mechanisms enter in, such as streamers. At very high pressures, the electric field strength for breakdown is so high that field emission from the electrodes enters into the process.

The Townsend criterion for avalanche is given by Eq. (2). Here α , the first Townsend coefficient, is

$$\gamma(e^{\alpha\delta} - 1) = 1 \quad (2)$$

the number of ionizing collisions per unit length in

TABLE 3. Comparison of various gaseous dielectrics

Gas	Condensation temperature (at 1 atm or 101 kPa), °C (°F)	Dielectric strength relative to nitrogen
N ₂	-195.8 (-320.4)	1.0
CO ₂	-78 (-108)	0.9
SF ₆	-63.8 (-82.8)	2.5-3
CCl ₄	76 (169)	6.3
CFCl ₃	23.8 (74.8)	3-4.5
CF ₂ Cl ₂	-29.8 (-21.6)	2.4-2.5
CF ₂ ClCF ₂ Cl	3.6 (38.5)	2.8
CF ₃ CN	63 (-81)	4
C ₂ F ₅ CN	-30 (-22)	5
C ₃ F ₇ CN	1 (34)	6
C ₃ F ₈	-36.7 (-34.1)	2-2.9
C ₄ F ₆	-5 (23)	4
C ₄ F ₁₀	-2 (28)	2.5
C ₅ F ₈	25 (77)	6
C ₅ F ₁₀	22 (72)	>4
He	-269 (-452)	<0.2

the direction of the field made by an electron; and γ , the second Townsend coefficient, is the number of secondary electrons produced at the cathode per electron produced by impact ionization in the gap, δ . Thus the dielectric strength of a gas is a function of both the gas species and the cathode material.

Although impact ionization in the gas produces equal numbers of positive ions and electrons, a space charge would develop even with equal displacement of the negative and positive charges. Since the electrons have greater mobility, this further separates the charges. In those situations when enough positive-ion space charge is produced at the original avalanche head to cause the space-charge field to be almost as large as the applied field, streamer breakdown occurs.

In highly nonuniform fields, where the electric field is very high at one electrode and low at the other, localized breakdown (called corona) which does not bridge the gap can occur. Corona can reduce the electric field around the pointed electrode, in effect rounding it off. This results in a higher breakdown voltage than otherwise, a phenomenon known as corona stabilization.

Most systems in the field become contaminated with conducting particles which greatly reduce the gaseous dielectric strength. This problem is important enough to warrant incorporating particle traps in such equipment, particularly those using electronegative gases. *See* ELECTRICAL CONDUCTION IN GASES.

Solids. For ambient temperature and pressure, the densities of solids are 10^3 – 10^4 times those of gases. This results in dielectric strengths 10 – 10^2 times higher than gases. At these higher fields, quantum-mechanical processes such as electron tunneling enter in. Additionally, through excitation of lattice vibration a new mechanism for reducing the energy of the accelerated electrons is provided. *See* LATTICE VIBRATIONS; PHONON; QUANTUM MECHANICS.

For the purpose of considering electrical breakdown, solids may be put in the categories of polar crystals, nonpolar crystals, quasicrystalline (highly ordered), and amorphous. On the basis of only density, the dielectric strength of solids and liquids may be expected to be comparable, and this is so. The alkali halides such as potassium bromide, KBr, and sodium chloride, NaCl, are examples of polar crystals. Mica is an example of a nonpolar crystal. The polymeric solids such as polyethylene can vary between highly ordered and amorphous. Glass is totally amorphous. In their absence of structure, the amorphous materials more nearly resemble and behave like highly viscous liquids. *See* AMORPHOUS SOLID; POLAR MOLECULE.

Above 0 K (-459.67°F) some electrons in the high-energy tail of the electron energy distribution have enough energy to find themselves in the conduction band of an insulator. This is why insulators have a slight conductivity. A. R. Von Hippel assumed that above a critical electric field, these thermally excited electrons would gain energy faster in the conduction band than they would lose it. Thus they would gain

enough energy to ionize the solid and cause breakdown. The values of critical field, calculated from this theory, agreed reasonably well with the measured values of the breakdown fields for the alkali halides. *See* BAND THEORY OF SOLIDS.

A simplifying assumption in Von Hippel's treatment is that the electrons in the conduction band all have the same energy. A consequence of this assumption is that since there is no ionization below the critical field, the prebreakdown current can be associated only with these free electrons, implying that (1) the prebreakdown current should be extremely small, (2) the prebreakdown current should be independent of the electric field below the critical field, and (3) the ionization coefficient should be discontinuous—going from zero to a large value at the critical field. Since these predictions are at variance with experiment, new theories were developed. H. Fröhlich assumed an energy distribution for the electrons in the conduction band. C. M. Zener assumed tunneling of electrons into the conduction band because of the high electric field.

Both the Von Hippel and Fröhlich theories neglected interaction of the conduction electrons with the lattice vibrations (phonons). Introducing a further refinement, F. Seitz demonstrated that the phonon interaction is dominant in nonpolar crystals, and is nonnegligible in polar crystals.

Despite the success of the above approaches, there are solids in which other mechanisms also occur, such as thermal, electromechanical, gas-discharge, and electrochemical breakdown. In thermal breakdown, part of the solid reaches a critical temperature at the breakdown field, causing chemical deterioration or melting. Even when thermal breakdown is avoided by using thin specimens, massive electrodes, and pulsed voltages, a gas-discharge mechanism of breakdown may initiate in voids in the specimen, or between the electrodes and the solid. Such discharges may cause treeing (Lichtenberg figures) in the solid, and may chemically and mechanically degrade the solid to the point of breakdown.

Liquids. Strictly speaking, liquids do not have an intrinsic dielectric strength. It is necessary to specify the electrode composition, geometry, and the time duration of the applied field. For very long stress durations, the time is important even for solids. Thus, except for electrode composition, this distinction between liquids and solids may be only a matter of degree. In liquids, microprojections on the cathode and migration of suspended particles (due to the gradient in the field) to the high-field regions contribute to a decrease in dielectric strength. Dissociation of the liquid by field-emitted electrons, electrochemical processes, and thermal bubble formation further complicate electrical breakdown in liquids.

A conditioning effect occurs, also common to vacuum, in which the breakdown strength increases progressively with the first few breakdowns, provided that the breakdown current is limited. Other commonalities with vacuum are curvature and area effects. Contrary to expectation, up to a point the dielectric strength is higher for electrodes of smaller

radius of curvature, and higher for electrodes of smaller area.

The fact that the breakdown strength increases with pressure is indicative that a bubble is formed during the breakdown process. However, if the voltage is applied in the nanosecond range, the pressure dependence vanishes and another breakdown mechanism becomes dominant. A piece of evidence in support of the bubble mechanism of breakdown for longer stress duration is the direct relationship between the boiling point and the breakdown strength for the *n*-alkanes.

Vacuum. As pointed out above, breakdown in vacuum has a number of features in common with liquid breakdown. L. Cranberg suggested that breakdown is initiated when a charged clump of loosely adhering material is removed from one electrode surface by the electric field, strikes the opposite electrode, and thus causes a high enough temperature to produce local evaporation. A discharge then ensues in the metal vapor. (The clump here is analogous to particle-initiated gas or liquid breakdown.)

Cranberg's model predicts that $VE > K$, where V is the breakdown voltage, E is the macroscopic electric field at the electrode where the clump originates, and K is a constant characteristic of a given pair of electrodes. For a uniform field and gap d , this implies that V is proportional to $d^{1/2}$, which agrees reasonably well with experiment. However, it does not predict the area effect, nor the polarity effect where the breakdown voltage is much higher when the pointed electrode is the anode. Moreover, it incorrectly predicts the curvature effect.

L. B. Snoddy was the first to suggest anode vaporization as a result of electron bombardment to account for breakdown. A. J. Ahearn was the first to suggest local heating of the cathode as initiating breakdown. L. C. Van Atta, R. J. Van de Graaff, and H. A. Barton were the first to hypothesize a particle interchange multiplication process. The basic assumption here is that, at a critical voltage, a free charged particle upon striking an electrode produces an avalanche of charged particles by secondary emission, with photoemission also playing a role.

There have since been many variations on these models. However, none has yet correctly predicted (even qualitatively) the area and curvature effects, as well as the gap dependence of approximately $d^{1/2}$. M. Rabinowitz introduced a hypothesis that predicts at least qualitatively the known experimental results, without reference to any particular model in terms of processes that are assumed to occur. He observed that the initiation of vacuum breakdown and gap conduction processes occurs so fast (on the order of 10^{-9} s), compared with the time constants of most breakdown circuits, that only the capacitively stored energy of the electrodes and electrode supports discharges within this time. (Light travels less than 30 cm or 1 ft in 10^{-9} s). His assumption was simply that breakdown cannot occur until there is sufficient stored energy in the electric field in the gap, because only this energy (or a fraction of it) is available to break down the gap.

The energy in the field is the capacitively stored energy, $\frac{1}{2}CV^2$. The hypothesis leads to Eq. (3), where

$$W = \frac{1}{2}fCV^2 \quad (3)$$

W is the critical energy needed to initiate breakdown (characteristic of the electrodes), C is the capacitance of the electrode system, and f is the fraction of the capacitively stored energy available to produce breakdown, $0 < f < 1$. For a uniform field, this predicts that the breakdown voltage V is proportional to $d^{1/2}$. The curvature and area effects are also qualitatively predicted. Mario Rabinowitz; Steven A. Boggs

Bibliography. R. Bartnikas et al. (eds.), *Engineering Dielectrics*, vol. 1, 1979, vol. 2A, 1983, vol. 2B, 1987, vol. 3, 1993; A. Bradwell et al. (eds.), *Electrical Insulation*, 1983; L. A. Dissado and J. C. Fothergill, *Electrical Degradation and Breakdown in Polymers*, 1992; H. G. Erdman (ed.), *Electrical Insulating Oils*, 1988; D. Kind and H. Karner, *High Voltage Insulation Technology*, 1985; G. G. Raju, *Dielectrics in Electric Field*, 2003.

Electrical interference

Disturbance to the normal or expected operation of electrical or electronic devices, equipment, and systems. Electrical interference is sometimes called radio-frequency interference (RFI) or electromagnetic interference (EMI). Such disturbances may range from nuisances (for example, an electric shaver jamming a nearby broadcast radio receiver) to catastrophes (for example, midair aircraft collision due to EMI navigation errors in holding patterns during bad weather in terminal areas). Electrical noise is a broader term that includes those phenomena that are generally termed electrical interference, but also includes naturally occurring currents or voltages that are more or less continuous and cannot be completely eliminated.

Electrical interference originates from one or more of the following: transmitters such as those used for broadcast, communication, radar, and navigation; artificial incidental emission sources such as from sparking of motor brushes, automotive ignition, and fluorescent lamps; and natural phenomena such as lightning and electrostatic discharges. The interference emissions may be radiated through space or conducted along the paths of electrical wires in power lines and signal cables. See COMMUNICATIONS CABLE; ELECTROMAGNETIC RADIATION; TRANSMISSION LINES.

Electrical interference is usually controlled by national regulations, such as those recommended by the International Special Committee on Radio Interference (CISPR), of which there are 25 member nations, including the United States. In the United States, the Federal Communications Commission, through Part 15 (low-level radiators such as computers, wireless microphones, and radio receivers) and Part 18 (industrial, scientific, and medical devices) of its regulations, attempts to mitigate

interference through radiated and conducted emission control. The military sector uses its own standard, MIL-STD-461B, to designate required control of both emission and susceptibility on all purchased equipment and subsystems having electrical or electronic elements. The intelligence agencies regulate allowable emission levels from equipment carrying classified information (known as TEMPEST) through their specifications, NACSIM-5100A.

Excessive electrical disturbance to power mains and through radiation is a form of environmental pollution. Unlike other forms of pollution, however, electrical noise cannot usually be detected with the human senses. Instrumented antennas and conduction probes are used with receivers or spectrum analyzers to detect the presence of emissions and to determine if they might create (or to diagnose) an EMI problem. They are also used for specifications and regulation certification. See SPECTRUM ANALYZER.

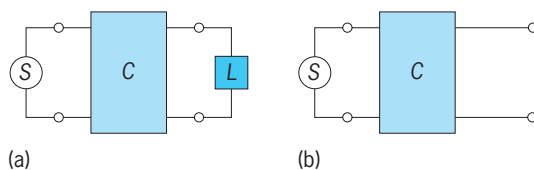
There are several means of mitigating electrical interference, including grounding, bonding, shielding, filtering, and transient control. There are two major approaches. One is to reduce radiation from boxes, cabinets, racks, and consoles by the design of multi-layer printed circuit boards and backplanes, or to use metal or metallized plastic cases with EMI-protected apertures. The other involves reduction of radiation to or from the interconnecting cable by using filter pin connectors or foil-wrap shields terminated with complete coverage by the connector backshell at the respective box bulkhead. See BONDING; ELECTRIC FILTER; ELECTRIC TRANSIENT; ELECTRICAL NOISE; ELECTRICAL SHIELDING; ELECTROMAGNETIC COMPATIBILITY.

Donald R. J. White

Bibliography. K. Javor, *Introduction to the Control of Electromagnetic Interference*, 1993; M. Mardiguian, *Electromagnetic Control in Components and Devices*, 1971; M. Mardiguian, *EMI Control Methodology and Procedures*, 1989; D. R. J. White, *EMI Control in the Design of Printed Circuit Boards*, vol. 1, 1996.

Electrical loading

The connection of a specific circuit (called the load) to another circuit. **Illustration a** shows a source of electrical energy (S), a coupling circuit (C), and a load (L). In power transmission, S is the generator of the power plant, C represents the connecting transmission lines and transformers, and L is, for instance, a household with its electrical appliances. The entire system operates to supply the load demand of L , that is, the electrical energy required by the load. In a high-fidelity sound system, S is the output of the



Systems (a) under load and (b) under no load.

phonograph cartridge or the tape playback head, C is the amplifier, and L is the coil of the loudspeaker. See ELECTRIC POWER SYSTEMS; SOUND-REPRODUCING SYSTEMS.

When the load is disconnected, the system operates under no-load conditions (illus. b).

Loading is used in certain wire-line transmissions in order to improve their characteristics. This is often done by connecting loading inductors along telephone lines, or along certain data transmission circuits. These inductors reduce the attenuation of the line, thus improving the transmission of signals. Waveguides, special antennas, and coaxial cables are often loaded with specific circuits to achieve specific purposes. See ANTENNA (ELECTROMAGNETISM); COAXIAL CABLE; INDUCTOR; TRANSMISSION LINES; WAVEGUIDE.

Shlomo Karni

Bibliography. D. J. Fink and H. W. Beaty (eds.), *Standard Handbook for Electrical Engineers*, 14th ed., 1999; M. Kaufman and A. H. Seidman, *Handbook for Electronics Engineering Technicians*, 2d ed., 1985; J. Markus and C. Weston, *Essential Circuits Reference Guide*, 1988.

Electrical measurements

Measurements of the many quantities by which the behavior of electricity is characterized. Such measurements are of great economic importance because of their impact upon industrial, commercial, and scientific activities. Measurements are required of all electrical quantities, possibly extending over a wide dynamic range and at frequencies ranging from 0 to 10^{12} Hz. The International System of Units (SI) is in universal use for all electrical measurements. Electrical measurements are ultimately based on comparisons with realizations, that is, reference standards, of the various SI units. These reference standards are maintained by the National Institute of Standards and Technology in the United States, and by the national standards laboratories of many other countries. Many electrical measurements are made to measure nonelectrical quantities. For example, platinum-resistance-thermometer measurements facilitate accurate temperature measurements, and eddy-current conductivity measurements are made to assess heat treatment of aluminum alloys. See ELECTRICAL UNITS AND STANDARDS.

Direct-current measurements. Direct-current (dc) measurements include measurements of resistance, voltage, and current in circuits in which a steady current is maintained.

Resistance measurement. Resistance is defined as the ratio of voltage to current. For many conductors this ratio is nearly constant, but depends to a varying extent on temperature, voltage, and other environmental conditions. The best standard resistors are made from wires of special alloys chosen for low dependence on temperature and for stability. They may have drift rates as low as 0.01 part per million (ppm) per year. See ELECTRICAL RESISTANCE; RESISTOR.

The SI unit of resistance, the ohm, is realized by means of a quantized Hall resistance standard. This is based upon the value of the ratio of fundamental constants h/e^2 , where h is Planck's constant and e is the charge of the electron, and does not vary with time. Wire-wound standard resistors can be compared with this standard with an uncertainty of about 1 part in 10^8 . See HALL EFFECT.

For resistance values below about 1 kilohm, the greatest accuracy is achieved by the use of four-terminal measurements, although this introduces considerable complication into resistive bridges.

Resistors in series produce a value equal to their sum, but resistors in parallel produce a resistance whose reciprocal, called its conductance, is equal to the sum of the conductances of the components. For this reason it is sometimes more convenient to use conductance in siemens, rather than the resistance (in ohms); for example, a resistance of 0.1Ω is equally well expressed as a conductance of 10 S . See CONDUCTANCE.

The principal instruments for accurate resistance measurement are bridges derived from the basic four-arm Wheatstone bridge, and resistance boxes. In the Wheatstone bridge the same current is passed through the unknown and reference resistors, and the resistance ratio is effectively deduced from the ratio of potential differences across them. Alternatively, by using a dc current comparator, the currents through the resistors are adjusted so as to produce equal potential differences across them, when the resistance ratio is given by the ratio of currents. The achievement of constant resistance ratios is much easier than that of constant resistance, so that if a bridge is used only to measure the ratio between a standard resistor and an unknown, it need not be constructed of resistors of the high and expensive quality of the standard. When used in this way, many commercial bridges can provide considerably better accuracy than is available by using their built-in reference resistors. See WHEATSTONE BRIDGE.

Many multirange digital electronic instruments measure resistance potentiometrically, that is, by measuring the voltage drop across the terminals to which the resistor is connected when a known current is passed through them. The current is then defined by the voltage drop across an internal reference resistor. For high values of resistance, above a megohm, an alternative technique is to measure the integrated current into a capacitor (over a suitably defined time interval) by measuring the final capacitor voltage. Both methods are capable of considerable refinement and extension. See OHMMETER; RESISTANCE MEASUREMENT.

Voltage measurement. The SI unit of voltage, the volt, is realized by using arrays of Josephson junctions. This standard is based on frequency and the ratio of fundamental constants e/h , so the accuracy is limited by the measurement of frequency. The secondary standards, which may be either saturated Weston mercury-cadmium standard cells or semiconductor Zener-diode-based electronic voltage standards, can be compared with the Josephson array with uncer-

tainties of about 1 part in 10^8 . Josephson arrays can produce voltages between $200 \mu\text{V}$ and 10 V .

Weston standard cells are sensitive to temperature, vibration, and the passage of current, and have to be kept in temperature-controlled enclosures stabilized to about 1 millikelvin. Electronic standards do not have these problems, but they may be sensitive to humidity and pressure, and their short-term noise and long-term drift in value are both slightly worse than those of standard cells. Voltage standards are usually compared back to back so that the resulting small voltage difference can be measured by a sensitive digital voltmeter which need not be calibrated to such a high accuracy as if the full voltage of the standard had to be measured.

At the highest levels of accuracy, higher voltages are measured potentiometrically, by using a null detector to compare the measured voltage against the voltage drop across a tapping of a resistive divider, which is standardized (in principle) against a standard cell. Dividers intended for this use up to 2 kV are usually called volt ratio boxes, and sometimes contain built-in self-calibration facilities. Resistive dividers for the accurate measurement of still higher voltages, up to about 500 kV , need increasing care to minimize the effects of the high electric fields, which may give errors due to current leakage through corona discharges. These can be avoided by the use of guarding techniques. See CORONA DISCHARGE; POTENTIOMETER.

The Zener diode reference standard is the basis for most commercial voltage measuring instruments, voltage standards, and voltage calibrators. The relative insensitivity to vibration and other environmental and transportation effects makes the diodes particularly useful as transfer standards. Under favorable conditions these devices are stable to a few parts per million per year. See ZENER DIODE.

Most dc digital voltmeters, which are the instruments in widest use for voltage measurement, are essentially analog-to-digital converters which are standardized by reference to their built-in reference diodes. Some of these instruments may be microprocessor-controlled, in which case they may be remotely or automatically operated, and may include built-in error correction for linearity over the indicated range. The basic range in most digital voltmeters is between 1 and 10 V , near the reference voltage. Other ranges are provided by means of resistive dividers, or amplifiers in which gain is stabilized by feedback resistance ratios. In this way these instruments provide measurements over the approximate range from 10 nanovolts to 10 kV . See DIGITAL-TO-ANALOG CONVERTER; MICROPROCESSOR; VOLTAGE MEASUREMENT; VOLTMETER.

Current measurement. The most accurate measurements of direct currents less than about 1 A are made by measuring the voltage across the potential terminals of a resistor when the current is passed through it. It is important that the resistor has been calibrated under the conditions of use, or that allowances are made for the differences. Higher currents, up to about 50 kA , are best measured by means

of a dc current comparator, which accurately provides the ratio of the high current to a much lower one which is measured as above. At lower accuracies, resistive shunts may be used up to about 5000 A, but the effective calibration of such shunts is a difficult process. *See* CURRENT COMPARATOR; CURRENT MEASUREMENT.

Power frequency measurements. The principal needs for accurate measurement at power frequencies (usually 60 or 50 Hz, but possibly up to 35 kHz for lighting systems) arise from the sale of electrical energy. There are also widespread needs for less accurate measurement for the monitoring and control of generators, distribution systems, and plants.

Alternating-current voltage and current measurements. Alternating-current (ac) voltages are established with reference to the dc voltage standards by the use of thermal converters. These are small devices, usually in an evacuated glass envelope, in which the temperature rise of a small heater is compared by means of a thermocouple when the heater is operated sequentially by an alternating voltage and by a reference (dc) voltage. Suitable series resistors, which have been independently established to be free from variation with frequency, permit direct measurement of power frequency voltages up to about 1 kV. Greater accuracy is provided by multijunction (thermocouple) thermal converters, although these are much more difficult and expensive to make. Thermal converters essentially measure the root-mean-square voltage; for a sinusoidal wave the peak value is $(2)^{1/2}$ times the root-mean-square value. Practical wave-forms are usually not sinusoidal, and the relationship between the mean and peak values will depend on the harmonic content or form factor of the waveform. Improvements in digital electronics have led to alternative approaches to ac measurement. For example, a line frequency waveform may be analyzed by using fast sample-and-hold circuits and, in principle, be calibrated relative to a dc reference standard. Also, electronic root-mean-square detectors may now be used instead of thermal converters as the basis of measuring instruments. *See* NONSINUSOIDAL WAVEFORM; THERMAL CONVERTERS.

Voltages above a few hundred volts are usually measured by means of a voltage transformer, which is an accurately wound transformer operating under lightly loaded conditions. These transformers are usually calibrated by comparison with a high-voltage capacitive divider, which can be independently evaluated. This permits the precise measurement of the voltage ratio and phase error of the transformer under its specified conditions of use.

The principal instrument for the comparison and generation of variable alternating voltages below about 1 kV is the inductive voltage divider. In principle, this resembles a series of tapped autotransformers connected to give a multidecade arrangement. Inductive voltage dividers are very accurate and stable devices, and at best provide an uncertainty of division of less than 1 in 10^8 . They are widely used

as the variable elements in bridges or measurement systems. *See* INDUCTIVE VOLTAGE DIVIDER.

Alternating currents of less than a few amperes are measured by the voltage drop across a resistor, whose phase angle has been established as adequately small by bridge methods. Higher currents are usually measured through the use of current transformers, which are carefully constructed (often toroidal) transformers operating under near-short-circuited conditions. The performance of a current transformer is established by calibration against an ac current comparator, which establishes precise current ratios by the injection of compensating currents to give an exact flux balance. *See* INSTRUMENT TRANSFORMER.

Commercial instruments for measurement of ac quantities are usually dc measuring instruments, giving a reading of the voltage obtained from some form of ac-dc transducer. This may be a thermal converter, or a series of diodes arranged to have a square-law response, in which case the indication is substantially the root-mean-square value. Some lower-grade instruments measure the value of the rectified signal, which is usually more nearly related to the peak value.

Power and energy measurements. The principal complication in power measurement is that the power is the vector product of the voltage V and current I , so a wattmeter is essentially a vector multiplier of these quantities. Although electromechanical instruments such as dynamometers still exist, electronic instruments are now cheaper and more accurate, and are much more widely used, especially where integration into an automated system is required. A long-established type of energy meter consists of a revolving conducting disk driven by the interaction between induced currents generated by voltage and current coils. Such meters are remarkable for their accuracy (particularly at low energy levels), stability, and low cost. Nevertheless, modern electronic energy meters that are competitive with the rotating disk types are available; they offer the advantages of digital indication and the possibility of remote interrogation via power line communications. *See* ELECTRICAL ENERGY MEASUREMENT; WATT-HOUR METER.

The great majority of wattmeters work at levels of current and voltage less than about 10 A and 250 V, single phase. High-power levels such as the output from power stations at several hundred megawatts per phase are measured through the use of current and voltage transformers.

Several wattmeters have appeared in which the multiplication is carried out electronically. In one method a pulse height is determined by the instantaneous value of one variable and the pulse width by the instantaneous value of the other. The integrated energy of the pulse string then gives the power. A second approach is direct simultaneous sampling of the voltage and current, which are then reduced to digital values. The product of each sample pair is stored, and the accumulated total of the series is a

measure of the power. Another approach to power measurement is the use of the quarter-squares principle, using the identity below. In this method the

$$(V + I \cos \theta)^2 - (V - I \cos \theta)^2 = 4 VI \cos \theta$$

outputs of two square-law devices are summed to give a direct power reading. All these approaches have been shown to be capable of producing accuracy comparable to that of the best wattmeters of more conventional design. *See* ELECTRIC POWER MEASUREMENT; WATTMETER.

Audio- and radio-frequency measurements. At the highest level of accuracy, audio-frequency (af) voltages are measured in the same way as at power frequencies, through the use of thermal converters, and employing inductive voltage dividers to provide variable ratios. Such methods provide uncertainties of about 5 ppm at 1 kHz, increasing to about 30 ppm at 1 MHz.

Both pointer and digital instruments for voltage measurement are similar to those in use at power frequencies. Current measurement is needed less often at these frequencies, and generally only small currents (less than about 1 A) are involved. These can readily be measured using a resistive transducer.

Since a typical af electrical component has both reactance and resistance, the ratio of af voltage to current depends upon the impedance of the component. A large number of bridges are available for the measurement of impedance at these frequencies, and most of these are especially suited to the measurement of particular parameters or ranges of values. The most accurate bridges are those based on transformers and inductive dividers, generally at frequencies below 100 kHz, and using coaxial systems throughout. Automatic network analyzers permit measurements to be made with great convenience within the range 10 kHz–40 GHz. Both vector and scalar instruments are available, with the latter being less expensive but not providing phase information. *See* CAPACITANCE MEASUREMENT; INDUCTANCE MEASUREMENT; RADIO-FREQUENCY IMPEDANCE MEASUREMENTS.

For power measurement, the use of dynamometer wattmeters is usually restricted to frequencies below 10 kHz. The range of electronic wattmeters, that is, those performing electronic vector multiplication, may often extend to 1 MHz, but in this frequency range the measurement of power by vector multiplication of voltage and current is superseded by the use of resistive bolometers or calorimeters. These devices provide a response when they are heated by a known fraction of the power to be measured; this approach may be used up to and beyond microwave frequencies. *See* BOLOMETER.

Microwave frequencies. At frequencies above about 300 MHz, where the circuit dimensions become an appreciable fraction of the wavelength, it becomes necessary to use totally enclosed circuits (waveguides), and the measurements of voltage, current, and lumped impedance are replaced by those

of wave impedance, reflection and transmission coefficients, voltage standing-wave ratio, attenuation, and power. The measurement techniques, other than for power measurement, are drastically changed. *See* ATTENUATION (ELECTRICITY); MICROWAVE IMPEDANCE MEASUREMENT; MICROWAVE NOISE STANDARDS; MICROWAVE POWER MEASUREMENT.

Fast pulse measurements. The measurement of electrical pulses with rise times of a few picoseconds is of growing importance in the fields of communications, computing, and digital electronics. The most commonly used instrument for measuring repetitive pulse trains is the sampling oscilloscope, and this is calibrated through electrooptical techniques in which short laser pulses are used to study electrical pulses as they travel along transmission lines on electrooptic crystals. In a variant of this, a small electrooptic crystal probe carrying fast laser pulses is brought close to the transmission line carrying the electric pulse, so that the electric field of the pulse affects the optical transmission of the probe and hence modifies the laser pulses. *See* ELECTROOPTICS; OPTICAL PULSES.

Automated measurement systems. There has been a noticeable trend toward the use of automated measurement systems for electrical measurements, facilitated by the readiness with which modern digital electronic instruments may be interfaced with computers. Many of these instruments have built-in microprocessors, which improve their convenience in use, accuracy, and reliability. A built-in reference standard may be used for self-calibration, and the results of many measurements may be stored, processed, and displayed in a variety of ways. A complex measurement system may comprise several such instruments operating under computer control and automatically producing measurement reports. The overall effect is that the performance of accurate measurements becomes less labor-intensive and more cost-effective. R. Gareth Jones; Owen C. Jones

Bibliography. D. A. Bell, *Electronic Instrumentation and Measurements*, 2d ed., 1994; *Code for Electrical Metering*, ANSI C12.1-1995, 1995; C. H. Dix and A. E. Bailey, Electrical standards of measurement, *Proceedings of the IEE*, vol. 122, no. 10R, October 1975; D. Fink and H. W. Beaty (eds.), *Standard Handbook for Electrical Engineers*, 14th ed., 2000; A. Hartland, Quantum standards for electrical units, *Contemp. Phys.*, 29:477–498, 1988; Institute of Measurement and Control, London, *A Guide to Measuring Direct and Alternating Current and Voltage below 1 MHz*, 2003; Institute of Measurement and Control, London, *A Guide to Measuring Resistance and Impedance below 1 MHz*, 1999; Instrument Society of America, *ISA Standards Library for Measurement and Control*, 1999; L. D. Jones and A. F. Chin, *Electronic Instruments and Measurements*, 2d ed., 1991; B. P. Kibble and G. H. Rayner, *Coaxial AC Bridges*, 1984; L. Schnell (ed.), *Technology of Electrical Measurements*, 1993; L. M. Thompson, *Electrical Measurements and Calibration: Fundamentals and Applications*, 2d ed., 1994.

Electrical model

A mathematical description or electrical equivalent circuit that represents the behavior of a device or system. Models for complex systems are often represented by networks of models for simpler electrical devices such as resistors, capacitors, transistors, and transformers. By using analogies between current or voltage and other physical parameters, equivalent circuits can also be used to analyze thermal, mechanical, magnetic, and acoustic systems. A recurring issue in development and in application of models is the use of simplifying assumptions to allow a compromise between accuracy and complexity. A hierarchy of models often exists, ranging from highly accurate but complex physics-based nonlinear time-domain computer models to linear equivalent-circuit models suitable for hand calculation. The best model for a particular application is the simplest one that predicts the relevant behavior with acceptable accuracy. See EQUIVALENT CIRCUIT; NETWORK THEORY.

Electrical models can be divided into nonlinear, large-signal and linear, small-signal models. Each of these can be further divided into time-invariant and time- or frequency-dependent categories.

Large-signal models. These are usually derived by applying physical laws to generalized devices. Since electromagnetic effects propagate across space at the speed of light, the most detailed and accurate version of these models needs to contain a description of device behavior across space and time. Such “distributed-circuit” models take the form of partial differential equations whose independent variables consist of three-dimensional spatial coordinates and time. Electromagnetic field models of this degree of complexity are required where the dimensions of the circuit are large compared with the wavelength at which the circuit operates, or when electromagnetic wave radiation is of interest. Unfortunately, these equations are difficult to solve, even numerically by computer simulation, and usually need to be simplified. See MAXWELL’S EQUATIONS.

At lower operating frequencies, where the physical dimensions of a circuit are a small fraction of a wavelength, spatial variations of electrical quantities can be neglected and a system of ordinary differential equations, with time as the sole independent variable, can model the circuit. In most electrical systems modeled at this “lumped-circuit” level, the variables of interest are the voltages at certain locations (nodes) in the circuit and the currents passing through these nodes. Other variables such as electric charge or magnetic flux also may be of interest. Physical models at this level of detail, without additional approximation, are often still highly complex and require lengthy computer solutions. See DIFFERENTIAL EQUATION.

An alternative approach is empirical modeling, in which a simple mathematical form is assumed for the model, with key physical parameters being determined by measurements on representative devices.

Such models further reduce predictive power for increased simplicity.

Successive simplifications are required to reduce circuit-model complexity to a point where pencil-and-paper calculations can produce an initial new system design. Once this preliminary design is completed, computer simulations using more detailed models can further refine physical circuit-performance predictions. Depending on the degree of simplification employed in the design process, it may be necessary to construct circuit prototypes to verify adequate circuit performance. If the agreement between circuit-model predicted performance and actual physical performance is inadequate, the simplified circuit models needs to be modified to account for these parasitic (unmodeled) effects and the entire design sequence needs to be repeated. In an attempt to avoid this expensive and time-consuming process, much effort has been expended in the continual development and enhancement of computer tools for circuit modeling. See COMPUTER-AIDED CIRCUIT DESIGN.

Many devices with lumped-circuit models can be represented with three controlling nodes or terminals. One of these can be selected as a reference to which the other voltages are compared. The result is a two-port network (Fig. 1a), with four variables: two port voltages and two terminal currents. Such a system can be characterized by using two independent equations. Ideally these equations should be expressible in closed-form: two of the variables are expressed as functions of the remaining two. For example, with voltages v_1 and v_2 in Fig. 1a as the independent variables, the closed-form model equations can be expressed as Eqs. (1) and (2). For devices

$$i_1(t) = f_1[v_1(t), v_2(t)] \tag{1}$$

$$i_2(t) = f_2[v_1(t), v_2(t)] \tag{2}$$

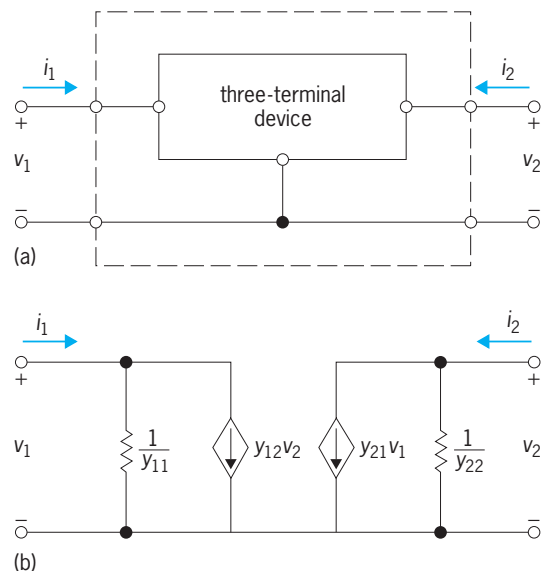


Fig. 1. Three-terminal device. (a) Two-port representation. (b) Y-parameter equivalent circuit. The diamond-shaped symbols represent voltage-controlled current sources.

controlled by more than three terminals, more port variables and more equations must be added to the model.

Once a suitable model is chosen, it can be used to determine the terminal currents and voltages that result when the device is connected in an electrical circuit. There are two main application areas for large-signal models. One is to find the response of a system to time-varying inputs. The other is to determine the voltages and currents of a circuit with zero input variations, as a preliminary to small-signal analysis.

Small-signal modeling. In many applications, current and voltage variations are small enough that the device's characteristics are approximately independent of signal size. This allows the device to be represented by a simplified linear equivalent circuit. To find the values of the equivalent circuit's elements, the input signal is replaced with its average value, and all currents and voltages are found by using the large-signal model. The results represent the quiescent operating point or Q point for the circuit. All currents and voltages can be approximated by series expansions around the Q point. For small enough signal variations, terms above first order can be neglected. For example (putting aside time derivatives for the moment), the currents in Eqs. (1) and (2) can be written as Eqs. (3) and (4), where the Δ 's denote

$$i_1 = I_1 + \Delta i_1 \cong I_1 + \frac{\partial i_1}{\partial v_1} \cdot \Delta v_1 + \frac{\partial i_1}{\partial v_2} \cdot \Delta v_2 \quad (3)$$

$$i_2 = I_2 + \Delta i_2 \cong I_2 + \frac{\partial i_2}{\partial v_1} \cdot \Delta v_1 + \frac{\partial i_2}{\partial v_2} \cdot \Delta v_2 \quad (4)$$

signal-dependent variations, I_1 and I_2 are the Q-point currents, and the derivatives are evaluated at the Q point.

An equivalent circuit can then be drawn in which only the variations of the signal are modeled; all signal-independent voltage and current sources are set to zero. The nonlinear device is replaced in the circuit by its linearized equivalent, whose element values are given by the derivatives in the series expansion. Equations (3) and (4) lead to the equivalent circuit shown in Fig. 1b, where the parameters of Eqs. (5) are the derivatives of the nonlinear equa-

$$\begin{aligned} y_{11} &= \frac{\partial i_1}{\partial v_1} y_{12} = \frac{\partial i_1}{\partial v_2} \\ y_{21} &= \frac{\partial i_2}{\partial v_1} y_{22} = \frac{\partial i_2}{\partial v_2} \end{aligned} \quad (5)$$

tions evaluated with v_1 and v_2 set equal to their Q-point values. These so-called y -parameters have units of conductance. Other small-signal representations result from different choices of dependent and independent variables.

Linearity has several benefits. One is the applicability of systematic solution methods for linear circuits that simplify both manual and computer solutions. Linearity also allows superposition: a linear system's response to the sum of two inputs equals the sum of its responses to the inputs applied separately. An important consequence of superposition is that

the system is completely specified by its response to sinusoids of arbitrary frequency. By using Fourier analysis, any input can be written as a unique sum of sinusoids whose magnitude and phase as functions of frequency make up the spectrum of the input. The linear system's response to each sinusoid is a sinusoid of the same frequency. The gain (amplitude increase) and phase shift vary with frequency but not with input amplitude. This frequency response determines the output spectrum for arbitrary inputs. Linear circuits can be analyzed entirely in the frequency domain. See FOURIER SERIES AND TRANSFORMS; GAIN; LINEARITY; RESPONSE; SUPERPOSITION THEOREM (ELECTRIC NETWORKS).

Neglecting time derivatives in the small-signal circuit derivation restricts the model's use to the midband region, where the gain is independent of frequency; in the midband equivalent circuit, all inductors and capacitors are treated either as short or open circuits.

Frequency-domain equivalent circuit. At high frequencies, parasitic effects inevitably reduce the gain. This frequency dependence is usually modeled by addition of capacitors or inductors to the model. The values of these elements depend (usually nonlinearly) on the Q point. The external circuit, of course, may also contain frequency-dependent elements. In the time domain, the current and voltage in an inductor or capacitor are related through time integrals or derivatives. These elements can be represented in the frequency domain with an imaginary impedance, directly or inversely proportional to frequency. The circuit can then be analyzed by using complex algebra, which is much easier than solving the time-domain differential equations. As an example of a frequency-domain model, Fig. 2 shows a high-frequency small-signal equivalent circuit for a bipolar transistor. See ALTERNATING-CURRENT CIRCUIT THEORY; TRANSISTOR.

Distortion. If the signal amplitudes are too large, the small-signal approximation is invalid, causing distortion: the generation of frequency components in the output that are not present in the input. Distortion is analyzed in the time domain by using large-signal models. See DISTORTION (ELECTRONIC CIRCUITS).

Noise models. Many electrical devices produce random current variations due to the discrete charge of the electron. This electrical noise limits the

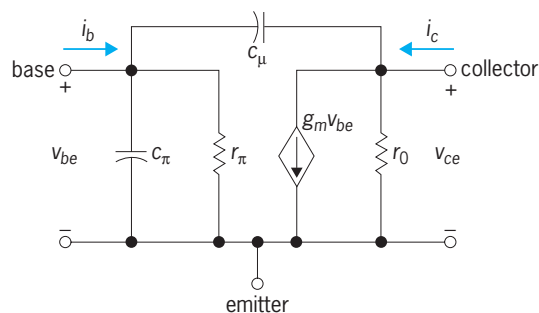


Fig. 2. High-frequency hybrid- π model for a bipolar transistor.

minimum usable signal amplitude. In the time domain, noise signals can be described only statistically. The statistical distributions can be used to represent noise in the frequency domain. The noise power at any frequency can be determined from the autocorrelation function, a measure of how closely two values of the signal measured at different times correspond as the time between the measurements is varied. Physical or empirical models are used to predict the spectrum's dependence on the Q point. See CIRCUIT (ELECTRONICS); ELECTRICAL NOISE.

Robert M. Fox; Philip V. Lopresti

Bibliography. G. Massobrio and P. Antognetti, *Semiconductor Device Modeling with SPICE*, 2d ed., 1993; A. S. Sedra and K. C. Smith, *Microelectronic Circuits*, 5th ed., 2004; Y. P. Tsividis, *Operation and Modeling of the MOS Transistor*, 3d ed., 2006; M. E. Van Valkenburg, *Network Analysis*, 3d ed., 1974.

Electrical noise

Interfering and unwanted currents or voltages in an electrical device or system. Electrical noise, or simply noise, has a significant effect on the design and operation of almost all electrical and optical systems which are used to communicate or process information. Noise is responsible for the familiar static observed on home radio receivers, the clicking sounds on frequency-modulation (FM) radios operating in fringe (near-threshold) areas, and the "snow"-type granularity on the viewing screen of a television receiver displaying a weak signal. In general, noise provides the fundamental limitation to the range over which radio or optical signals can be transmitted and received with integrity. Noise is, therefore, of great importance to the engineers who design and operate such systems.

It is convenient to differentiate between noise which results from human activity and that which is naturally occurring. Noise which results from human activity, such as that generated by an electrical appliance or an automotive ignition, can usually be eliminated or minimized by good design practice (shielding, filtering, equipment location, and so forth). Naturally occurring noise can be further subdivided into that which is irregular or erratic in nature and that which is more or less continuous. An example of noise which is irregular or erratic is that associated with an electrical storm. This type of noise is sometimes dealt with in the system design, but since it is only occasionally present, it does not ordinarily constitute a design limitation. On the other hand, naturally occurring noise which is essentially continuous in time is responsible for the fundamental limitation cited above. The remainder of the article therefore concentrates on this type of noise. See ATMOSPHERIC ELECTRICITY; CROSSTALK; ELECTRIC FILTER; ELECTRICAL INTERFERENCE; ELECTRICAL SHIELDING; ELECTROMAGNETIC COMPATIBILITY; ELECTRONIC EQUIPMENT GROUNDING; GROUNDING; STATIC.

Sources. Most noise generation is a consequence of the spontaneous fluctuations which occur within

matter at the microscopic level. In electrical circuits these fluctuations give rise to what are commonly referred to as thermal noise and shot noise. Thermal noise is generated by the random motion of free electrons in a resistor or any conductor with resistance. The random motion, and thus the noise generated, is proportional to the temperature of the medium. At absolute zero temperature on the Kelvin scale (-459.67°F), all motion ceases and no noise is generated. Shot noise is most commonly identified with the fluctuations in the current of a vacuum tube caused by the random emission of electrons from its heated cathode. Shot noise is also observed in semiconductor devices as random fluctuations in carrier density when an electric field is applied. There are other types of noise associated with electrical circuits, but shot noise and thermal noise are by far the most important. See FREE-ELECTRON THEORY OF METALS; KINETIC THEORY OF MATTER; SEMICONDUCTOR; VACUUM TUBE.

In a system in which signals are transmitted through the atmosphere [for example, amplitude-modulation (AM) or FM radio broadcast, or satellite communications], the receiving system will always receive noise as well as the desired signals. This noise is a result of thermal radiation from the Earth, planets, Sun, Moon, the galaxy (galactic noise), radio-emitting stars, and atmospheric gases. In addition, there is a small background level of thermal radiation, uniformly distributed, which is believed associated with the big bang origin of the universe. All of these noise sources, weighted by the directional characteristics of the receiving antenna, will contribute to the overall system noise. See COSMIC BACKGROUND RADIATION; HEAT RADIATION; TERRESTRIAL RADIATION.

In an optical communications system, a signal level is represented by a number of energy packets called photons. The mean arrival rate of the photons at the detector is proportional to the optical intensity or signal strength. At the detector (a photodiode), the photons are absorbed, each creating a hole-electron pair and thus a current in which the electrons are randomly positioned in time and in which the mean number of electrons is proportional to the optical intensity. The statistical nature of this process gives rise to fluctuations in the number of photons representative of a given level and, subsequently, the number of electrons generated to represent that level. If the detector has internal gain as in an avalanche photodiode, each hole-electron pair can create additional hole-electron pairs. This process, however, is statistical in nature, resulting in a mean value of gain but giving rise to additional fluctuations in the generated current. See MICROWAVE SOLID-STATE DEVICES; OPTICAL COMMUNICATIONS; OPTICAL DETECTORS; PHOTODIODE; PHOTON.

Mathematical analysis. Because of the statistical or random nature of noise, noise voltages and currents must be dealt with by using the branch of mathematics that deals with random variables, that is, probability theory. A random variable can be characterized by a probability density function $p(x)$. The definition of

$p(x)$ is such that the probability that the random variable x lies between x_1 and x_2 is given by Eq. (1). Ob-

$$\text{Prob}(x_1 < x < x_2) = \int_{x_1}^{x_2} p(x) dx \quad (1)$$

viously, $p(x)$ can be normalized as shown in Eq. (2).

$$\int_{-\infty}^{\infty} p(x) dx = 1 \quad (2)$$

Investigations have shown that fluctuation noise, in general, has a probability density function which is gaussian or normal, provided that the noise can be represented as a large number of independent overlapping samples. The gaussian distribution has the form given in Eq. (3), where σ is the standard de-

$$p(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-x^2/2\sigma^2} \quad (3)$$

viation, x is the random variable (noise voltage or current), and the average value is zero. See DISTRIBUTION (PROBABILITY); PROBABILITY.

If the average value is nonzero, the distribution will have the form illustrated by Eq. (4), where V is

$$p(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-(x-V)/2\sigma^2} \quad (4)$$

the average value. V might be a signal voltage, for example, with the total instantaneous voltage fluctuating about that level.

In the foregoing, it has been assumed that the statistical properties of the random variable are invariant with time. The statistical process associated with such a system is said to be a stationary process. In experimentally collecting data to determine the probability distribution of a random variable (for example, a noise voltage), two methods are available. First, a large number of samples can be collected from a system at a sequence of times. Second, and alternately, the outputs of a large number of identical systems can be sampled simultaneously. The latter method is called the ensemble method. If the statistics collected by either of these methods are invariant with time, the process is stationary. Furthermore, if both methods yield identical statistics, the process is ergodic. Most noise processes are both stationary and ergodic. For such processes the mean value [V in Eq. (4)] is just the time average of the noise, and the variance σ^2 is the average of the square of the fluctuations about the mean value. See PROBABILITY (PHYSICS); STATISTICAL MECHANICS.

Two additional functions that are useful in characterizing random noise are the autocorrelation function $R(\tau)$ and the power spectral density $G(w)$. $R(\tau)$ provides a measure of the dependence between any two values of the noise variable $x(t)$ and $x(t + \tau)$, as specified in Eq. (5). $G(w)$ gives the distribution of

$$R(\tau) = \lim_{T \rightarrow \infty} \frac{1}{2T} \int_{-T}^T x(t) x(t + \tau) dt \quad (5)$$

average power in noise as a function of frequency w . $G(w)$ and $R(\tau)$ are Fourier transform pairs; that is,

they satisfy Eqs. (6), where $j = \sqrt{-1}$. This is a use-

$$\begin{aligned} G(w) &= \int_{-\infty}^{\infty} e^{-jw\tau} R(\tau) d\tau \\ R(\tau) &= \int_{-\infty}^{\infty} e^{-jw\tau} G(w) dw \end{aligned} \quad (6)$$

ful relationship because it relates the autocorrelation function of a noise voltage at the output of a linear system to a function $G(w)$ which can be expressed in terms of the transfer function of the system. For example, if white noise (uniform power spectrum) is passed through a system with transfer function $H(w)$, the autocorrelation function at the output of the system is simply the Fourier transform of $H(w)^2$. See FOURIER SERIES AND TRANSFORMS.

Many network calculations involving shot noise and thermal noise can be handled simply by considering the noise sources as ordinary signal sources but with the distinction that mean-squared values are employed. Total voltages (or currents) can be determined by adding mean-squared values.

Noise figure and noise temperature. A quantity of considerable interest in describing the characteristics of an electrical or optical signal is the ratio of signal power to noise power (S/N) measured at various points in the system. In particular, the signal power-to-noise power ratio at the output of a network, $(S/N)_{\text{out}}$, can be related to the signal power-to-noise power ratio at the input of the network, $(S/N)_{\text{in}}$, by a term called the noise figure F , as expressed by Eq. (7). F is a quality factor of the net-

$$F = \frac{(S/N)_{\text{in}}}{(S/N)_{\text{out}}} \quad (7)$$

work because it provides a measure of the degradation introduced by the network to the signal-to-noise ratio. An ideal, lossless, noise-free network is one for which $F = 1$. Frequently F is expressed in decibels as shown by Eq. (8). Typical signal-to-noise ratio min-

$$F_{\text{dB}} = 10 \log F \quad (8)$$

imum requirements may range from 3 to 15 dB for many detection systems or may be as large as 50 dB for a high-quality voice communications circuit. See SIGNAL-TO-NOISE RATIO.

In radio astronomy and space communications and, in general, systems which employ very low-noise first-stage amplifiers in the receiving system (for example, masers and parametric amplifiers), it has become common practice to characterize the amplifier and, indeed, the system in terms of an equivalent noise temperature. The noise figure defined previously can be expressed in terms of this equivalent noise temperature, as shown in Eq. (9),

$$F = \frac{T_e}{T_o} + 1 \quad (9)$$

where T_e is the equivalent noise temperature and T_o is the ambient temperature (usually taken as 290 K or 62°F). Then equivalent noise temperature is

expressed in terms of noise figure by Eq. (10). Tem-

$$T_e = (F - 1)T_o \quad \text{K} \quad (10)$$

peratures are always expressed in kelvins. See AMPLIFIER; MASER; PARAMETRIC AMPLIFIER; SPACE COMMUNICATIONS.

If the system consists of a series of cascaded networks, each with an equivalent noise temperature T_i and gain G_i , it can be shown that the overall equivalent noise temperature referred to the input terminals of the first stage is given by Eq. (11). For

$$T_e = T_1 + \frac{T_2}{G_1} + \frac{T_3}{G_1 G_2} + \dots + \frac{T_n}{G_1 G_2 \dots G_{n-1}} \quad (11)$$

this system or any system in which the overall equivalent noise temperature is known, the noise power spectral density is simply given by kT_e watts/hertz, where k is Boltzmann's constant (1.38×10^{-23} joule/K) and T_e is in kelvins. See BOLTZMANN CONSTANT.

Thermal noise. Spontaneous fluctuations in a resistor at temperature T will give rise to a noise voltage at its terminals with mean-squared value given by Eq. (12), where k is Boltzmann's constant, R is the

$$\overline{V_N^2} = 4kTRB \quad (12)$$

resistance in ohms, and B is the bandwidth in hertz. Equation (12) is often referred to as the Nyquist formula and the thermal noise as Nyquist noise or Johnson noise in recognition of the original investigators.

Equation (12) is actually an approximate formula, but it is valid with useful accuracy over all radio and microwave frequencies of general interest. It is not, however, valid at optical frequencies. The Nyquist formula shows that the spectral density $4kTR$ is constant. Noise with a constant or uniform distribution with frequency is referred to as white noise.

In a communications system, thermal noise is evidenced throughout the receiving system: the antenna, the transmission lines, amplifiers, and so forth. In systems where very low noise is vital to operation (such as those used in space communications), cryogenic cooling of the first stage of amplification is frequently employed. Liquid helium and liquid nitrogen have commonly been used for this purpose. Receiver noise temperatures approaching 4 K (-452°F) have been demonstrated. In these systems, antenna noise temperature becomes very important.

Antenna noise temperature is defined formally in Eq. (13), where A_e is the effective receiving area of

$$T_A = \frac{A_e}{\lambda^2} \int \int T(\theta, \phi) P_N(\theta, \phi) d\Omega \quad (13)$$

the antenna in m^2 , λ is the wavelength in m, $P_N(\theta, \phi)$ is the normalized power pattern of the antenna, and $T(\theta, \phi)$ is the temperature distribution over all space in kelvins. $T(\theta, \phi)$ includes contributions from the

atmosphere, the galaxy, radio-emitting stars, and so forth, as well as thermal radiation from the Earth itself. In Eq. (13), Ω is a solid angle and the integration is carried out over 4π steradians. $P_N(\theta, \phi)$ will depend on the pointing direction of the antenna; consequently, the antenna temperature will depend on the pointing direction of the antenna. At microwave frequencies, the antenna temperature is usually a minimum when the antenna is pointed to zenith. See ANTENNA (ELECTROMAGNETISM).

Shot noise. Shot noise, like thermal noise, is thermal in origin. However, it differs from thermal noise in two significant ways: (1) for shot noise to be present in a circuit, a voltage (for example, cathode to anode) or electric field must be impressed; (2) shot noise is not characterized by a uniform, wide-band (white) spectral density, but rather has a high-frequency cutoff given approximately by the inverse transit time of an electron from cathode to anode.

The probability that N electrons will be emitted from a heated cathode in a time interval τ is given by the Poisson distribution, Eq. (14), where \bar{n} is the

$$P(N) = \frac{(\bar{n}\tau)^N e^{-\bar{n}\tau}}{N!} \quad (14)$$

average number of electrons emitted per second. If the electric field is sufficiently strong, all the emitted electrons are attracted to the anode, and the fluctuations in current are given by the fluctuations in emission. This is called the temperature-limited case, and the mean-squared noise current is given by the shot-effect formula of Eq. (15), where e is the

$$\overline{i_s^2} = 2e\bar{I}B \quad (15)$$

charge on an electron, \bar{I} is the average current, and B is the bandwidth in hertz.

If the applied electric field is not sufficiently strong to pull all the electrons to the anode, a negative space charge develops in the vicinity of the cathode which tends to inhibit shot-effect fluctuations. This is the space-charge-limited case, and the shot noise is then approximated by Eq. (16), where Γ^2 is a space-charge

$$\overline{i_s^2} = 2e\bar{I}\Gamma^2 B \quad (16)$$

reduction factor varying commonly between 0.01 and 1. The value of Γ^2 depends on physical geometry, cathode temperature, and applied voltage. See SPACE CHARGE.

Noise in optical detectors. The reception of an optical signal involves the use of a photodetector in which the optical signal is converted to an electric current. Thus, both optical noise and electronic noise can be expected to influence the signal-to-noise ratio at the receiver output.

Optical noise relates to the random arrival of photons at the detector and subsequent generation of hole-electron pairs. The number of hole-electron pairs generated during the period t to $t + \tau$ is a random variable having the Poisson distribution given

by Eq. (17), where Λ is given by Eq. (18), where η is

$$P(n) = \frac{\Lambda^n e^{-\Lambda}}{n!} \quad (17)$$

$$\Lambda = \int_t^{t+\tau} \frac{\eta}{hf} p(t) dt \quad (18)$$

the quantum efficiency (fraction of light absorbed), h is Planck's constant, f is the optical frequency, and $p(t)$ is the optical power.

For the case where $p(t)$ has a constant value p_0 , an expression for the signal-to-noise ratio is given as Eq. (19), where e is charge on the electron, P_N is the

$$S/N = \frac{\left(\frac{\eta e p_0}{hf}\right)^2}{\frac{2\eta e^2 p_0}{hf} B + P_N} \quad (19)$$

mean-square thermal (electronic) noise power associated with amplification of the generated current, and B is the bandwidth in hertz.

Many modern optical communications systems employ avalanche photodiodes. In the avalanche photodiode, a received photon creates a hole-electron pair, which in turn, by means of ionization collisions, creates additional hole-electron pairs. Thus the detector has internal gain. However, the process which gives rise to the gain is itself a random process, and therefore, additional fluctuation noise is generated. The signal-to-noise ratio for the optical receiver employing an avalanche photodetector is given by Eq. (20), where G represents the interval

$$S/N = \frac{\left(\frac{neGp_0}{hf}\right)^2}{\frac{2ne^2 G^2 F p_0 B}{hf} + P_N} \quad (20)$$

gain and F is the excess noise factor.

If, in Eq. (19) or (20), the thermal noise term is sufficiently small compared to the optically generated noise so as to be neglected, the resulting S/N is referred to as the quantum limit. In most present-day systems, however, thermal noise is found to be the dominant term.

Noise measurement. The most common apparatus used to measure noise is a total power receiver consisting of a wide-band amplifier, a square-law detector, an integrator, and an indicator. The sensitivity of the receiver (the minimum detectable noise fluctuation) is generally improved as the front-end bandwidth and the integration time is increased. Sensitive receivers of this type are called total power radiometers and are frequently employed in radio astronomy. See **RADIOMETRY**.

The noise temperature (or noise figure) of an amplifier can be determined by what is commonly called the Y -factor method. The measurement utilizes two calibrated noise sources and a power receiver of the type just described. One of the noise sources is connected to the input of the amplifier, and the output power is measured with the receiver. This is repeated with the second noise source. The

ratio of power received in the two cases is recorded and denoted Y . The amplifier noise temperature is then given by Eq. (21), where T_1 and T_2 (where T_2 is

$$T_{\text{amp}} = \frac{T_2 - Y T_1}{Y - 1} \text{ K} \quad (21)$$

greater than T_1) are the noise temperatures of the calibrated terminations. It is assumed that the amplifier has sufficiently high gain that the noise contribution from the receiver can be neglected.

The total system noise temperature can be measured by using the noise-adding method. In this method an increment of noise is injected into the system and the total received noise power is compared to the noise power without the added increment. Then the total system noise power is given by Eq. (22), where Y is given by Eq. (23), and ΔT is the

$$T_{\text{sys}} = \frac{\Delta T}{Y - 1} \text{ K} \quad (22)$$

$$Y = \frac{\text{system noise power} + \text{injected noise power}}{\text{system noise power}} \quad (23)$$

noise temperature of injected noise.

Modern test equipment is available which greatly simplifies the measurement of the noise performance of circuits. Automated instrumentation which provides output directly in terms of noise figures is, for example, readily available. See **ELECTRICAL NOISE GENERATOR**.

Randall W. Kreutel

Bibliography. M. S. Gupta, *Noise in Circuits and Systems*, 1988; J. D. Kraus, *Radio Astronomy*, 2d ed., 1986; J. E. Midwinter, *Optical Fibers for Transmission*, 1979, reprint 1991; R. Morrison, *Noise and Other Interfering Signals*, 1991; R. Pettai, *Noise in Receiving Systems*, 1984; M. Schwartz, *Information, Transmission, Modulation and Noise*, 4th ed., 1990.

Electrical noise generator

A device that produces electrical noise for use in electrical measurements. Electrical noise generators are commonly employed to measure the noise figure or noise temperature of radio receivers. They are also used in various other tests in radar and communications systems. Celestial noise sources are used to calibrate large antennas.

Some standard types of noise generators are hot-wire, diode, gas-discharge tube, hot and cold loads (terminations), and radio-star. A hot-wire noise source consists of the filament of a lamp heated by direct current. Thermal noise having spectral density $4kTR$, where k is Boltzmann's constant and T and R are the temperature and resistance of the filament respectively, is generated across the terminals of the filament. A diode noise generator utilizes the temperature-limited shot effect to generate noise. At frequencies less than the reciprocal transit time of the diode, the noise spectral density is $2e\bar{I}$, where e is the charge on the electron and \bar{I} is the average anode current. A gas-discharge noise generator, commonly

referred to as a noise tube, consists of a fluorescent light tube enclosed in a waveguide. Noise generation is essentially thermal. The noise tube is commonly employed at microwave frequencies. Hot and cold loads consist of well-matched terminations, either transmission line or waveguide, held at a given temperature by using an oven or by applying cryogenic refrigeration. Noise generation is thermal. Common temperatures for noise-generating terminations are nominally 80 and 300 K (−316 and 80°F)

Celestial radio sources (radio stars) are commonly employed as reference noise sources for evaluating the characteristics of very low-noise, high-gain space communications receiving antennas. There are a number of accurately calibrated sources available—the choice depending on system parameters (frequency, antenna gain, system noise temperature, elevation angle, and so forth) and physical location. The most common radio sources employed are Cassiopeia A, Taurus A, Cygnus A, and Orion A. The first three are classified as nonthermal sources in which radiation results from relativistic electrons interacting with an interstellar magnetic field. The electrons are rotated in a plane perpendicular to the magnetic field direction, and radiation is characterized by a component polarized parallel to that plane. The polarized component is small, however, and the major portion of the radiation is unpolarized. The nonthermal sources have flux densities which decrease with increasing frequency and, consequently, tend to have a cutoff frequency above which they are not usable. Orion A is a thermal source in which radiation occurs from a hot, ionized cloud. Orion A has a constant flux density at frequencies above 2 GHz. As a point of reference, a typical value of thermal noise received from Cassiopeia A by a 10-m-diameter (33-ft) antenna operating in the C band (4–6 GHz) would be about 150 K (−190°F). See POLARIZED LIGHT.

Other possible sources include the Sun, Moon, and the planets. The solar flux density is equal to that from blackbody radiation at 6000 K (10,300°F) at wavelengths less than 1 cm, but is greater than this at longer wavelengths. At radio wavelengths the radiation can be greatly enhanced during periods of sunspot activity. In general, at longer wavelengths the solar radio emission can be considered as the superposition of a slowly varying component dependent on the sunspot number and a rapidly varying component resulting from solar flares. The polarization of the former has a strong circularly polarized component, while that of the latter is generally random. See SUN.

At radio wavelengths, the Moon emits like a thermal blackbody with temperature in the range from 150 to 280 K (−190 to 44°F) depending on the lunar phase. At long wavelengths the lunar-phase dependence is minimal, with small (less than 10%) variations about a nominal temperature of approximately 190 K (−118°F).

It is important that variations in brightness over the source be corrected, that polarized components of flux be accounted for, and that attenuation charac-

teristics of the intervening medium, the atmosphere, be included. See ELECTRICAL NOISE; RADIO ASTRONOMY.

Randall W. Kreutel

Bibliography. M. S. Gupta (ed.), *Electrical Noise: Fundamentals and Sources*, 1977; M. S. Gupta, *Noise in Circuits and Systems*, 1988; J. D. Kraus, *Radio Astronomy*, 2d ed., 1986; R. Morrison, *Noise and Other Interfering Signals*, 1991; A. Van der Ziel, *Noise: Sources, Characterization and Measurement*, 1970.

Electrical resistance

Opposition of a circuit to the flow of electric current. Ohm's law states that the current I flowing in a circuit is proportional to the applied potential difference V . It is obeyed very accurately for metallic and many other conductors over a wide range of temperatures. The constant of proportionality is defined as the resistance R . Hence, Eq. (1) holds. If V and I

$$V = IR \quad (1)$$

are measured in volts and amperes, respectively, R is measured in ohms. Microscopically, resistance is associated with the impedance to flow of charge carriers offered by the material. For example, in a metallic conductor the charge carriers are electrons moving in a polycrystalline material in which their journey is impeded by collisions with imperfections in the local crystal lattice, such as impurity atoms, vacancies, and dislocations. In these collisions the carriers lose energy to the crystal lattice, and thus Joule heat is liberated in the conductor, which rises in temperature. The Joule heat P is given by Eq. (2).

$$P = I^2R = IV = \frac{V^2}{R} \quad (2)$$

See CRYSTAL DEFECTS; ELECTRICAL CONDUCTIVITY OF METALS; ELECTRICAL RESISTIVITY; JOULE'S LAW; OHM'S LAW.

Peter A. Schroeder

Electrical resistivity

Opposition of a material to the flow of electricity. For a conductor of uniform cross section A (in square meters), length l (in meters), and resistance R (in ohms, Ω), the resistivity is defined by Eq. (1).

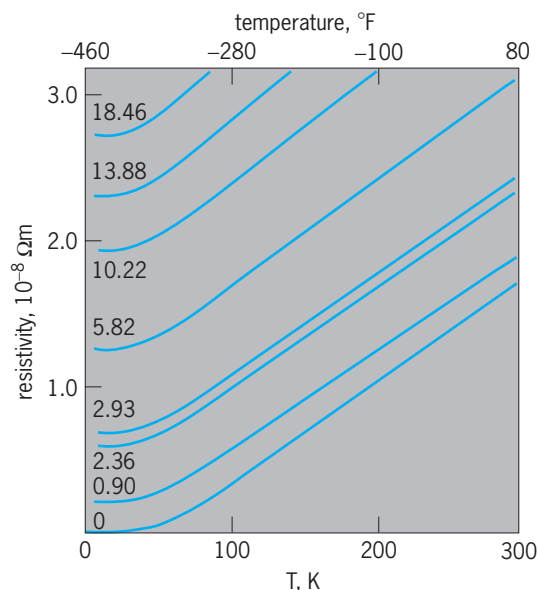
$$\rho = \frac{RA}{l} \quad (\Omega \cdot m) \quad (1)$$

Unlike the extensive quantity R , ρ is an intensive quantity, that is, an intrinsic property of a material, independent of its amount or shape. The range of resistivity values is extremely wide (see table). At the extremes are conductors such as copper with a room-temperature resistivity of $1.7 \times 10^{-8} \Omega \cdot m$ and insulators such as quartz with a resistivity greater than $10^{16} \Omega \cdot m$. The range is even wider when temperature dependence (a property useful in classifying materials) is taken into account. See ELECTRIC INSULATOR; ELECTRICAL RESISTANCE.

Resistivities at room temperature	
Material	$\Omega \cdot \text{m}$
Poly(tetrafluoroethylene)	10^{18}
Quartz	10^{16}
Diamond	10^{12}
Glass	10^{11}
Sodium chloride*	10^2
Silicon, pure	10^3
Germanium, pure	5
Metglass 2204 ($\text{Ti}_{50}\text{Be}_{40}\text{Zr}_{10}$)	3
Bismuth	1.2×10^{-6}
Antimony	4.2×10^{-7}
Nichrome	10^{-6}
Iron	10×10^{-8}
Potassium	7×10^{-8}
Copper	1.7×10^{-8}
Copper, pure†	10^{-12}

*At 1000 K (1340°F).
†At 4.2 K (-452°F).

Crystalline metals. For a pure metal such as copper (see *illus.*), the room-temperature resistivity is almost entirely caused by scattering of electrons by lattice vibrations, or, in quantum terms, phonons. As the temperature decreases, so does the phonon density, resulting in a decrease of resistivity. At the temperature of liquid helium (4.2 K or -452°F) this is so low that phonon scattering makes very little contribution to the resistivity, which is then determined by the concentration of imperfections in the crystal lattice. Thus, at 4.2 K and below, the resistivity may be $10^{-11} \Omega \cdot \text{m}$ or even lower. In an unstrained metal the imperfections responsible for the residual resistivity at 4.2 K is often taken as a measure of purity. See CRYSTAL DEFECTS; LATTICE VIBRATIONS; PHONON.



Resistivity of the disordered copper-zinc alloy system. Numbers on the curves give the concentration of zinc in atomic percent. (After W. E. Henry and P. A. Schroeder, *The low-temperature resistivities and thermopowers of α -phase copper-zinc alloys*, *Can. J. Phys.*, 41:1076-1093, 1963)

Many other scattering mechanisms contribute to resistance. Among these is mutual scattering of the charge carriers. Inasmuch as there are about 10^{23} electrons per cubic centimeter in a typical metal, a large contribution might be anticipated. However, the number of electrons that can interact with each other is severely restricted by energy considerations, and the T^2 variation of resistivity with absolute temperature T , typical of electron-electron scattering, is seen only at temperatures of a few kelvins where the phonon contribution to the resistivity becomes negligible. See ELECTRICAL CONDUCTIVITY OF METALS.

Disordered crystalline alloys. These are alloys in which the component ions are randomly distributed on the crystal lattice sites. For such an alloy the resistance can be written as in Eq. (2), where ρ_r is

$$\rho = \rho_r + \rho_i \quad (2)$$

the residual resistance produced by electron scattering by the impurity ions, and ρ_i is the resistivity contribution arising from the scattering of electrons by phonons. Matthiessen's rule states that ρ_i is the same for alloys as for the pure metal. For the more dilute alloys in systems such as copper-zinc, this is nearly true, as indicated by the parallelism of the graphs of resistivity versus temperature (see *illus.*), but Matthiessen's rule is at best an approximation which worsens as the concentration of the alloying constituent increases. Nichrome is an alloy with a high resistivity because of the scattering from a disordered array of the constituent ions. Because of its high resistivity and stability at high temperatures, it is frequently used as a heating element. The resistivity can become very high in an amorphous metal, such as Metglass 2204 (see table), in which all long-range order is lost. But the resistivity of a concentrated alloy may be considerably reduced if the component ions form an orderly array. This occurs in intermetallic compounds, which can be represented chemically as A_pB_q , where p and q are small integers, all the A atoms occupy one orderly array of sites, and all the B atoms occupy another. For example, the gold-gallium compound AuGa_2 has a residual resistivity smaller than that of many metals in their purest available form. See INTERMETALLIC COMPOUNDS; MATTHIESSEN'S RULE; METALLIC GLASSES.

Semimetals. Bismuth, antimony, arsenic, and graphite are classified as semimetals and have higher resistivities than iron, potassium, and copper (see table), because they have a much lower concentration of mobile electrons that can participate in the conduction process. Their resistivity, like that of metals, increases with temperature.

Semiconductors. The semiconductors silicon and germanium also have a higher resistivity than metals, again because the number of charge carriers is much reduced. They differ from semimetals in that, in their pure or intrinsic form, the number of charge carriers and the conductivity increase exponentially with absolute temperature, and consequently their resistivity decreases as the temperature increases. See SEMICONDUCTOR.

Ferromagnetic materials. Several transport properties of ferromagnetic alloys can be explained by assuming conduction in parallel by the spin-up and spin-down electrons. Such a model is assumed to explain the resistance of magnetic multilayers. For example, in a sample consisting of alternate layers of copper and cobalt with layer thicknesses up to 10 nanometers, the magnetizations of the cobalt layers lie in the plane of the layers and, for the appropriate thicknesses of the copper layers, the magnetizations of alternate cobalt layers are in opposite directions. When the magnetizations are made to line up by applying a magnetic field, the resistance is observed to drop so dramatically that this phenomenon is frequently referred to as giant magnetoresistance. The phenomenon is explained in terms of spin-dependent scattering. *See* MAGNETORESISTANCE.

Charge-transfer salts. A metallic conducting state with conductivity decreasing with temperature is also exhibited by organic charge-transfer salts. These consist of large organic donor molecules (cations) which can stack in various ways, and smaller inorganic acceptor molecules which arrange themselves between the stacks. TMTSF (tetramethyltetraselenafulvalene) and BEDT-TTF [bis(ethylenedithiolo)tetrathiafulvalene] are typical of complicated organic donor molecules. The TMTSF stacks in chains which give rise to quasi-one-dimensional conductivity. The BEDT-TTF salts are more two-dimensional. These materials are exciting in that the available choice of anions yield materials of low dimensionality with a multitude of interesting properties, including superconductivity. *See* ORGANIC CONDUCTOR.

Ionic conductors. The resistivity of ionic conductors such as sodium chloride (NaCl; see table) decreases with temperature. However, conduction in these materials is a thermally activated process, in which the conducting ions have to pass over the potential barriers separating them from an adjacent site. *See* IONIC CRYSTALS.

Resistance thermometry. The sensitivity of the resistivity to temperature has practical applications in thermometry. Platinum resistance thermometers are used between 14 and 330 K (−434 and 134°F); rhodium-iron alloys between 1 and 800 K (−458 and 980°F); and carbon and germanium at very low temperatures (0.01 to 100 K or −459.65 to 280°F) where their sensitivity arises from the rapid decrease in charge carrier density as the temperature is lowered. *See* LOW-TEMPERATURE THERMOMETRY; THERMOMETER.

Peter A. Schroeder

Bibliography. J. S. Dugdale, *The Electrical Properties of Metals and Alloys*, 1977; K. M. Hellwege and O. Madelung (eds.), *Landolt-Börnstein Numerical Data and Functional Relationships in Science and Technology, New Series, Group 3: Crystal and Solid State Physics*, vol. 15a: *Metals: Electronic Transport Phenomena*, 1982; C. Kittel, *Introduction to Solid State Physics*, 7th ed., 1996; K. Schröder, *Handbook of Electrical Resistivities of Binary Metallic Alloys*, 1983.

Electrical shielding

The imposition of a metal or composite barrier between one or more sources of electrical noise and their victims with the objective of reducing or eliminating electrical interference. Examples of the barrier are the case or housing of equipment; shields covering interconnecting cables between equipment; large cabinets, racks, or consoles; shielded (screen) rooms; and entire shielded buildings or vehicles.

Shielding effectiveness. The principal measure of a shield's performance is the shielding effectiveness. It is defined by the equation below, where SE_{dB} is

$$SE_{dB} = 20 \log_{10} \frac{F_b}{F_a}$$

the shielding effectiveness in decibels, F_b is the electric (or magnetic) field strength before imposition of the barrier, and F_a is the electric (or magnetic) field strength after imposition of the barrier. *See* DECIBEL.

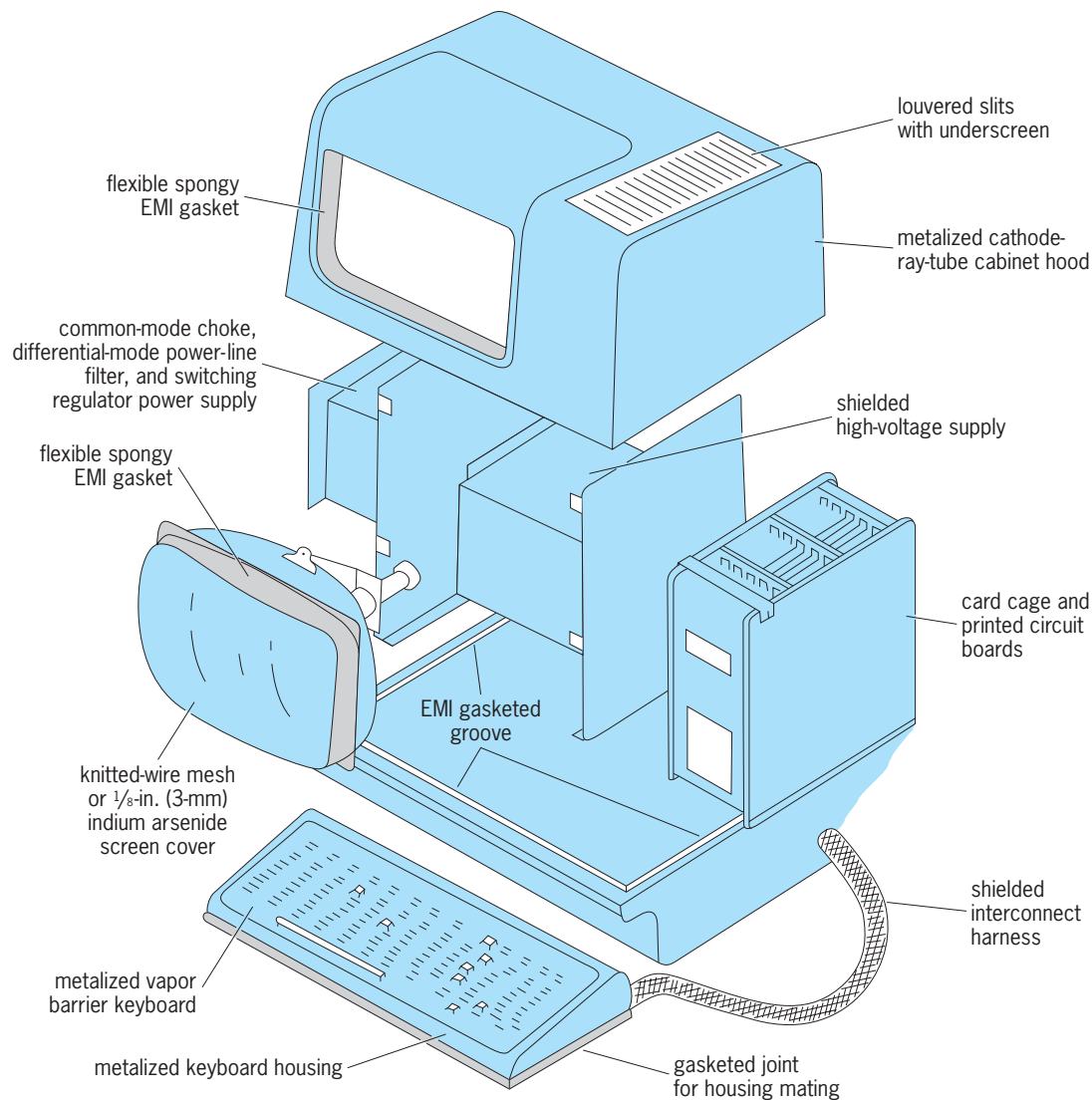
Shielding is obtained by the combination of reflection loss and absorption loss. The former is due to the impedance mismatch between the wave impedance of an oncoming wave and the surface impedance of the interposed barrier. Absorption loss corresponds to the attenuation due to skin effect at higher frequencies, and is dependent upon the frequency, conductivity, permeability, and thickness of the barrier. Shielding effectiveness is the sum of both losses. *See* ABSORPTION OF ELECTROMAGNETIC RADIATION; ATTENUATION (ELECTRICITY); RADAR-ABSORBING MATERIALS; REFLECTION AND TRANSMISSION COEFFICIENTS; REFLECTION OF ELECTROMAGNETIC RADIATION; SKIN EFFECT (ELECTRICITY).

Metal shields. Most intentional shields are made of metal to ensure high reflection losses. Even thin metals, such as household aluminum foils (with a thickness of about 1.5 mils or 0.038 mm), offer shielding effectiveness in excess of 100 dB. At low frequencies, these foils become electromagnetically transparent (that is, they do not attenuate magnetic fields). Thus, if there is a shielding problem due to the selection and makeup of the metal barrier, it is likely to occur only for low-frequency magnetic fields.

To obtain significant shielding to magnetic fields at low frequencies, the metal barrier must be very thick or composed of a highly permeable material such as Mumetal, Supermalloy, or Hypernom. Often such shields are fabricated in two or more layers, frequently laminated, to obtain good shielding properties per unit size or weight.

Leaky apertures. The major problem encountered in achieving adequate shielding at other than low frequencies is leaky apertures, not the metal choice. To provide the integrity required of the basic shield metal, all apertures must be protected or secured, or special electrical components and gaskets should be used (see *illus.*).

Shielded rooms. Shielded enclosures, ranging in size from small rooms to entire buildings, are



Some measures for controlling electromagnetic interference in video display terminals.

commercially available for performing as radiated emission and susceptibility test chambers, secure facility radiation containment, and protection for sensitive equipment from outside electromagnetic radiation. These enclosures typically provide in excess of 120 dB attenuation from 10 kHz to 1 GHz, and undergo reduced performance at lower and higher frequencies. Where extended performance is required, shielded rooms are usually made of solid seam-welded panels of galvanized steel. Special attention is given to leaky areas such as doors, and entrances where cables or wires such as power-line and telephone are brought in through filters. See ELECTRIC FILTER; ELECTROMAGNETIC PULSE (EMP).

Cable shields. Interconnecting cables between equipment often act as an "antenna farm" in which they behave as undesired pickup antennas and provide radiation escape from internal signals conducted along their length. Shielding such cables and harnesses is the dominant protection mechanism; other options include absorbers and filters. Cable shields may vary in complexity from

one cover braid, through multibraids and composite braids or foils, to extruded tubes or conduits. See COMMUNICATIONS CABLE; ELECTRICAL INTERFERENCE; ELECTRICAL NOISE; ELECTROMAGNETIC COMPATIBILITY.

Donald R. J. White

Bibliography. L. H. Hemming, *Architectural Electromagnetic Shielding Handbook*, 1992; M. Mardiguian and D. R. J. White, *Electromagnetic Shielding*, vol 3: *Interference Control Technologies*, 1988; D. R. J. White, *Handbook on Electromagnetic Shielding Materials and Performance*, 2d ed., 1980.

Electrical units and standards

The process of measurement consists in finding out how many times the quantity to be measured contains a fixed quantity of the same kind, called a unit. The definitions of the units often involve complex physical theory and do not lend themselves readily to practical realization. The concrete representations

of units are known as measurement standards. In practice, measurements are made by using an instrument calibrated against a local reference standard, which itself has been calibrated either directly or by several links in a traceability chain against the national standard held by the national standards laboratory.

In order that measurements of similar quantities made in different countries may be compared, it is important that the values of the standards held by these laboratories should be known in relation to one another. International harmonization of national standards is the concern of the International Bureau of Weights and Measures (Bureau International des Poids et Mesures, or BIPM) that was set up at Sèvres, near Paris, under the Convention of the Meter (1875). The General Conference on Weights and Measures (CGPM) has overall control of the organization at the political level and meets every few years. Scientific questions and the running of the bureau are the responsibility of the International Committee of Weights and Measures (CIPM), which meets annually and is assisted by a number of consultative committees, including the Consultative Committee for Electricity and Magnetism (CCEM). The CCEM encourages research on electrical standards and arranges international comparisons of national standards, including standards of radio frequency and microwave quantities. *See* PHYSICAL MEASUREMENT.

Electrical and Magnetic Units

With the growth of the electric telegraph and supply industries in the nineteenth century, the need was recognized for accurate measurements of electrical quantities, based on generally accepted units and standards. The scientific community played an important part in developing these. From the start, the electrical units had a logical scientific basis, rather than deriving from ancient arbitrary practice.

The cgs system. A proposal by W. E. Weber in 1851 led to the absolute cgs system in which all units of quantities to be measured could be derived from the base units of length, mass, and time—the centimeter, gram, and second. This system was widely adopted although it had three weaknesses: the size of the units was inconvenient for practical use; it was difficult to realize the units from their definitions; and there were separate sets of units for electrostatic and electromagnetic quantities, based respectively on the inverse-square laws of force between electric charges and between magnetic poles. These inverse-square laws involved the permittivity ϵ_0 and the permeability μ_0 of free space (also known as the electric and magnetic constants). Ratios of the electrostatic and electromagnetic units for the same quantity involved the product $\epsilon_0\mu_0$, which J. C. Maxwell showed to be equal to the reciprocal of the square of the speed of light, c_0^{-2} . This was the first example of a relationship between the electrical units and a fundamental physical constant.

International units. The first weakness was resolved by international agreement in 1881 to fix the

practical units—the volt, the ohm, and the ampere—at 10^8 , 10^9 , and 0.1 times the respective cgs electromagnetic units. The other weaknesses were avoided by the decisions of the 1908 International Congress in London, where realizations of these units in terms of easily reproducible standards were defined: the international ohm as the resistance of a column of mercury of specified dimensions; the international ampere as the current which would deposit silver in a silver voltameter at a given rate; and the international volt as their product, with an alternative definition in terms of the electromotive force of a Weston cell. The values were chosen to approximate as closely as the available information allowed to the exact multiples of the cgs units, but the international units were defined in their own right and their values were not changed when the results of later absolute determinations became available. They continued in use until officially abolished by the CGPM in 1948. *See* ELECTROMOTIVE FORCE (CELLS).

The mksa units. A more fundamental change resulted from a proposal by G. Giorgi in 1902. This led to the adoption of the mksa system of units, in which there are four base units: the meter, the kilogram, the second, and the ampere. Use of the meter and the kilogram instead of the centimeter and the gram gave units of a size more convenient for practical use, and use of the ampere as a base unit resolved the conflict between electrostatic and electromagnetic units while maintaining the magnitudes of the widely used practical units. This was a truly coherent system, in the sense that other units were derived from the base units without the need for factors of proportionality other than unity.

SI units. From the mksa system the present-day SI (Système Internationale) has developed, by the addition of further base units to include other fields of measurement. SI units were formally adopted by the CGPM in 1954.

The seven base units of SI are listed in **Table 1**. Of these, the first two, the kilogram and the second, are defined independently of the others. The kilogram is the mass of the international prototype of the kilogram held at the BIPM; the second is the duration of 9 192 631 770 periods of oscillation of the radiation corresponding to the transition between the two hyperfine levels of the ground state of the cesium-133 atom. The definitions of the other base units involve these two units. Until October 1983, the meter was also independently defined, in terms of the wavelength of a krypton lamp, but at that time

TABLE 1. SI base units

Quantity	Unit	Symbol
Mass	kilogram	kg
Time	second	s
Length	meter	m
Electric current	ampere	A
Thermodynamic temperature	kelvin	K
Luminous intensity	candela	cd
Amount of substance	mole	mol

TABLE 2. Some derived electrical units

Quantity	Unit and symbol	Derivation
Potential difference, emf	volt, V	$W \cdot A^{-1} = m^2 \cdot kg \cdot s^{-3} \cdot A^{-1}$
Resistance	ohm, Ω	$V \cdot A^{-1} = m^2 \cdot kg \cdot s^{-3} \cdot A^{-2}$
Electric charge	coulomb, C	$s \cdot A$
Capacitance	farad, F	$C \cdot V^{-1} = m^{-2} \cdot kg^{-1} \cdot s^4 \cdot A^2$
Conductance	siemens, S	$A \cdot V^{-1} = m^{-2} \cdot kg^{-1} \cdot s^3 \cdot A^2$
Magnetic flux	weber, Wb	$V \cdot s = m^2 \cdot kg \cdot s^{-2} \cdot A^{-1}$
Inductance	henry, H	$Wb \cdot A^{-1} = m^2 \cdot kg \cdot s^{-2} \cdot A^{-2}$
Magnetic flux density	tesla, T	$Wb \cdot m^{-2} = kg \cdot s^{-2} \cdot A^{-1}$
Magnetic field strength	ampere per meter	$m^{-1} \cdot A$
Current density	ampere per square meter	$m^{-2} \cdot A$
Electric field strength	volt per meter	$V \cdot m^{-1} = m \cdot kg \cdot s^{-3} \cdot A^{-1}$
Permittivity	farad per meter	$F \cdot m^{-1} = m^{-3} \cdot kg^{-1} \cdot s^4 \cdot A^2$
Permeability	henry per meter	$H \cdot m^{-1} = m \cdot kg \cdot s^{-2} \cdot A^{-2}$

the CGPM decided to take a fixed value for the speed of light and redefined the meter as the length of the path traveled by light in vacuum in 1/299 792 458 of a second. As a result, it is now possible to realize the meter with a lower uncertainty than previously. *See* LENGTH; LIGHT.

The units of other physical quantities (derived units) are derived from the base units by simple numerical relations, since SI is a coherent system. Examples are the unit of force, the newton ($N = m \cdot kg \cdot s^{-2}$); of energy, the joule ($J = N \cdot m$); and of power, the watt ($W = J \cdot s^{-1}$). Some of the more important derived electrical and magnetic units are listed in **Table 2**.

The SI base unit for electrical measurements is the ampere (A), the unit of electric current. It is defined in terms of a hypothetical experiment as that constant current which, if maintained in two straight parallel conductors of infinite length, of negligible circular cross section, and placed 1 meter apart in vacuum, would produce between these conductors a force equal to 2×10^{-7} newton per meter of length. According to electromagnetic theory, the force per unit length between two such conductors a distance d apart and carrying a current I is given by Eq. (1).

$$F = \frac{\mu_0 I^2}{2\pi d} \quad (1)$$

The definition of the ampere is thus equivalent to fixing the value of μ_0 as $4\pi \times 10^{-7}$ henry per meter.

The volt (V) is the unit of potential difference and of electromotive force. It is defined as the potential difference between two points of a conducting wire carrying a constant current of 1 ampere when the power dissipated between these points is equal to 1 watt. From the ampere and the volt, the ohm is derived by Ohm's law, and the other derived quantities follow in a similar manner by the application of known physical laws. *See* OHM'S LAW.

The following are the remaining units of electrical and magnetic quantities that have special names in the SI, together with their formal definitions, where appropriate.

Coulomb (C): The unit of electric charge, equal to 1 ampere-second. The coulomb is the quantity

of electricity carried in 1 second by a current of 1 ampere.

Farad (F): The unit of capacitance, equal to 1 coulomb per volt. The farad is the capacitance of a capacitor between the plates of which there appears a potential difference of 1 volt when it is charged by a quantity of electricity of 1 coulomb.

Henry (H): The unit of inductance, equal to 1 weber per ampere. The henry is the inductance of a closed circuit in which an electromotive force of 1 volt is produced when the electric current in the circuit varies uniformly at the rate of 1 ampere per second.

Ohm (Ω): The unit of electrical resistance, equal to 1 volt per ampere. The ohm is defined as the resistance between two points of a conductor when a constant potential difference of 1 volt, applied to these points, produces in the conductor a current of 1 ampere, the conductor not being the seat of any electromotive force.

Siemens (S): The unit of electrical conductance (the reciprocal of resistance), equal to 1 ampere per volt. It was formerly known as the mho.

Tesla (T): The unit of magnetic flux density, equal to 1 weber per square meter.

Weber (Wb): The unit of magnetic flux, equal to 1 volt-second. The weber is the magnetic flux which, linking a circuit of one turn, would produce in it an electromotive force of 1 volt if it were reduced to zero at a uniform rate in 1 second.

The mechanical units of frequency (hertz), energy or work (joule), and power (watt) are frequently involved in expressing electrical and magnetic quantities. The cgs units, such as the gauss, gilbert, maxwell, and oersted, formerly used, are not part of the SI and are now obsolete. *See* UNITS OF MEASUREMENT.

Electrical Standards

Realization of the values of the electrical and other units from their SI definitions involves great experimental difficulties. For this reason, it is customary for national standards laboratories to maintain stable primary standards of the units against which other reference standards can be compared. From time to time, absolute determinations of the values of these

primary standards were made in terms of their definitions, using the methods described below.

For a number of years the ampere and the ohm were determined using a current balance and a calculable mutual inductance. The latter were superseded by the calculable capacitor and the moving-coil balance until the late 1980s, when higher accuracies became necessary. Other methods were explored, but by then the Josephson effect and the quantum Hall effect had made possible the standardization of the volt and the ohm by relation to fundamental physical constants. The recommendation by the CCEM of the values to be adopted for these constants in 1990 led to a complete change in primary electrical standards and the method of handling them; and all the major national laboratories, and the BIPM, now use this method.

For many years the primary standards maintained by most laboratories were the volt, in terms of the mean electromotive force of a group of Weston cells, and the ohm, using a group of standard resistors. A range of reference standards of other quantities are derived from these, including direct-current (dc) voltage and resistance at a variety of levels; alternating-current (ac) voltage, resistance, and power; capacitance and inductance; radio-frequency (rf) and microwave quantities; magnetic quantities and properties of materials; dielectric properties; and other quantities. These secondary standards are used for day-to-day measurements and for the calibration of local reference standards of other users in the national measurement system. See CAPACITANCE MEASUREMENT; ELECTRIC POWER MEASUREMENT; INDUCTANCE MEASUREMENT; MAGNETIC MATERIALS; MICROWAVE MEASUREMENTS; PERMITTIVITY; RESISTANCE MEASUREMENT; VOLTAGE MEASUREMENT.

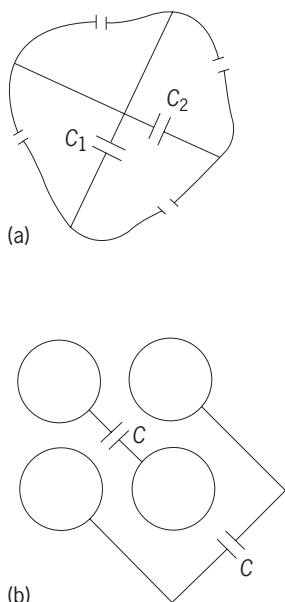


Fig. 1. Calculable capacitor. (a) Cross section of general cylindrical capacitor with four segments. (b) Cross section of symmetrical capacitor with cylindrical electrodes. (c) Calculable capacitor used in absolute determination of the ohm. (U.S. National Institute of Standards and Technology)

Absolute Determinations

It has long been the practice that national standards laboratories measure their primary electrical standards in terms of the kilogram and second, with the help of electromagnetic theory and the definition of the ampere. Two different measurements are needed to establish the values of the standards of the volt and the ohm.

For many years the current balance, which measures the force between pairs of cylindrical coils by weighing, has been a mainstay of absolute methods. The other traditional method was the determination of the ohm by comparison with a mutual inductance whose value was calculable from its dimensions. But the accuracy of both measurements has been limited to a few parts in 10^6 by difficulties in determining the dimensions of the coils involved in the experiments and the heating effects of the currents. This order of accuracy is not adequate for modern requirements, and these methods are now of historic interest only. See CURRENT BALANCE.

Calculable capacitor. In 1956, A. M. Thompson and D. G. Lampard published a new theorem in electrostatics which showed that if a conducting cylinder of any cross section is divided parallel to the axis into four parts (Fig. 1a) the cross capacitances per unit length, C_1 and C_2 , obey Eq. (2).

$$\exp \frac{-\pi C_1}{\epsilon_0} + \exp \frac{-\pi C_2}{\epsilon_0} = 1 \quad (2)$$

This theorem has been applied as the basis for a capacitor whose change in capacitance as a guard electrode is moved can be calculated exactly from the measurement of a single length. This can be done to a very high degree of accuracy by using a laser interferometer. The problems of making many complex dimensional measurements which limited the accuracy of the current balance and the calculable mutual inductor are thus no longer a limitation. If the capacitor is made symmetrical with four cylindrical electrodes (Fig. 1b), then $C_1 = C_2$ and the capacitance per unit length is given by Eq. (3). This is a

$$C = \frac{\epsilon_0 (\log_e 2)}{\pi} \quad \text{farads per meter} \quad (3)$$

little less than 2 picofarads per meter.

An early capacitor of this kind is shown in Fig. 1c. By means of a series of bridges, the value of a 1-ohm resistor can be measured in terms of the calculated reactance of the capacitor. In this way an absolute determination of the value of a standard ohm can be made with an uncertainty (1 standard deviation) of about 5 parts in 10^8 .

Electrometer. A number of laboratories have made determinations of the volt by an electrometer method, in which the electrostatic potential difference between two electrodes causes a force which can be measured mechanically. W. G. Clothier devised an electrometer in which the lower electrode is the surface of a pool of mercury. This rises against the force of gravity when the potential is applied,

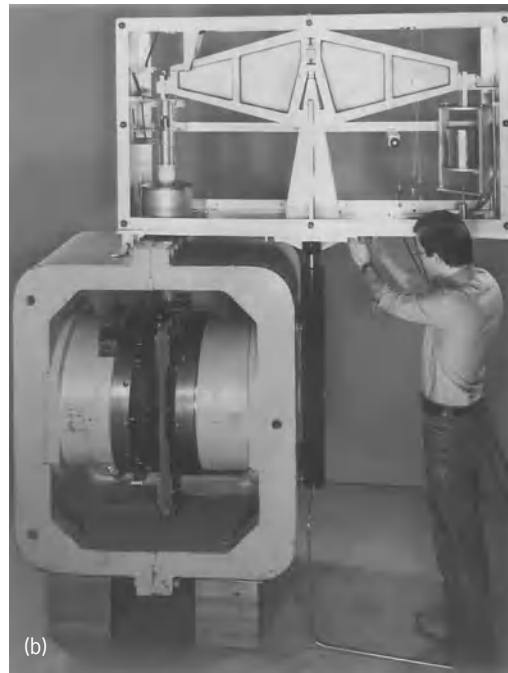
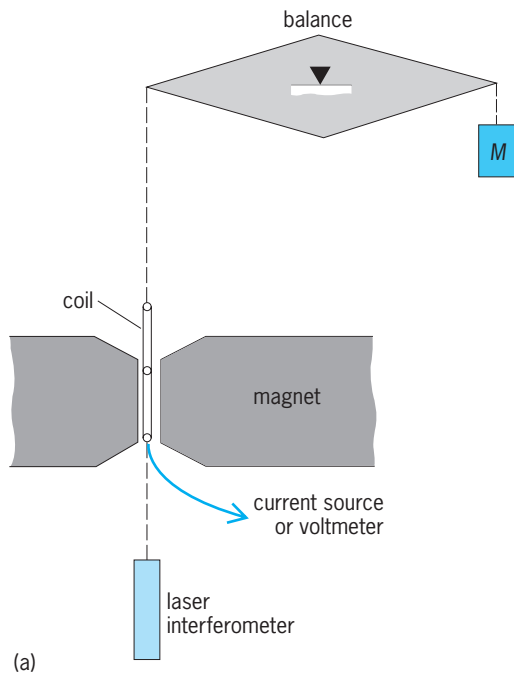


Fig. 2. Moving-coil balance. (a) Schematic diagram. (b) Equipment. (U.K. National Physical Laboratory)

and its movement can be measured very accurately by interferometric means. In this way an uncertainty (1 standard deviation) of about 3 parts in 10^7 has been achieved.

Moving-coil balance. B. P. Kibble developed equipment which in effect measures the electrical watt in terms of the mechanical watt (Fig. 2). From this the value of the volt can be determined absolutely. A coil is suspended from the arm of a balance, partly in the field of a powerful magnet. Its position in the vertical plane is measured by a laser interferometer. By using the apparatus successively in two different modes, it is not necessary to know the dimensions of the coil or the strength and distribution of the magnetic field; many errors cancel out, and the problems of accurate dimensional measurements are avoided.

In the first part of the experiment, a current I flows through the coil and the induced force is balanced by a mass M . In the second part, the coil is connected to a voltage-measuring system and made to move through the field with a velocity u , measured by the laser interferometer. The induced voltage V is measured in terms of the primary standard of the volt. Then Eq. (4) follows from the equality of the

$$VI = Mgu \quad (4)$$

electrical and mechanical units of power, where g is the acceleration due to gravity. If the current and the velocity are adjusted so that V is equal to the voltage drop when the current I flows through a resistor whose value R is known accurately in terms of the primary standard of resistance, then V is given by Eq. (5). With equipment based on this principle, an

$$V = (MguR)^{1/2} \quad (5)$$

uncertainty (1 standard deviation) of about 1 part in 10^7 has been achieved.

Josephson effect. If a Josephson junction, consisting of two superconductors separated by an insulating layer so thin that the wave functions are coupled, is irradiated with microwave energy at a frequency ν , it has a voltage-current characteristic that consists of a series of voltage steps (Fig. 3a). The height of these steps is given by Eq. (6), where K_J is known as the

$$\Delta V = \frac{\nu}{K_J} \quad (6)$$

Josephson constant; its value is independent of the materials used and of the experimental conditions. For a practical irradiation frequency the voltage step for a single junction is in the millivolt range, but by using many junctions in series (Fig. 3b) it is possible to use the effect to provide a voltage standard in the region of 1 volt and thus to monitor the voltage of a standard cell. See JOSEPHSON EFFECT.

Quantum Hall effect. If an electric current I is passed through a semiconductor placed in a magnetic field of induction B_z , at right angles to the current (Fig. 4a), a voltage V_H is developed across the semiconductor in the direction orthogonal to both the current and the field. This is known as the Hall effect, and the ratio R_H , given by Eq. (7), is known as the Hall resistance.

$$R_H = \frac{V_H}{I} \quad (7)$$

In 1980, K. von Klitzing demonstrated that the Hall resistance in certain semiconductors could be quantized in high magnetic fields at very low temperatures. If the semiconductor in Fig. 4a consists of a GaAs/AlGaAs heterostructure, in which a layer of

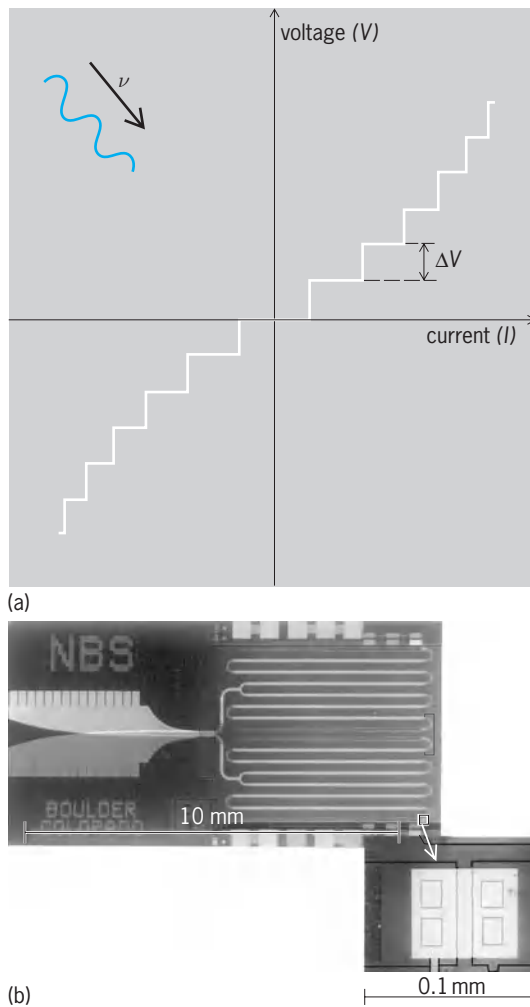


Fig. 3. Josephson effect. (a) Josephson junction characteristic. (b) Linear array of more than 2000 Josephson junctions in a folded microwave stripline, giving a total voltage step in the region of 1 volt. (U.S. National Institute of Standards and Technology)

aluminum gallium arsenide is deposited on a gallium arsenide substrate so that current flow is confined to the interface and is effectively two-dimensional, the Hall resistance takes the form shown in Fig. 4b and given by Eq. (8), where $i = 1, 2, 3, \dots$, and R_K is a

$$R_H = \frac{R_K}{i} \tag{8}$$

constant having a value in the region of 25 kilohms and known as the von Klitzing constant.

Again, the constant is independent of experimental conditions, and it is possible to set up a system working on a selected plateau of the resistance curve of Fig. 4b which acts as a resistance standard and can be used to monitor the values of the primary standards of resistance held by national standards laboratories. See HALL EFFECT.

International comparisons. For measurements made in one country to be valid in another, it is essential that national primary standards should be extremely stable and their values accurately known in relation to corresponding standards in other countries. The voltages of Weston cells and the resistances

of the best wire-wound resistors drift typically by a few parts in 10^7 per year. This is significantly greater than the precision, a few parts in 10^8 , with which the values of different standards can be compared. For many years the BIPM has maintained groups of standard cells and standard resistors, similar to those of the national laboratories. It has been the practice of these laboratories to compare their standards every 3 years with those of BIPM by sending traveling standards to Sèvres. In this way it has been possible to monitor variations in national standards over a period of several decades and to adjust their values when the drift has become excessive.

The BIPM also arranges the direct intercomparison of national standards of derived quantities, under the auspices of the CCEM. This has particular value for rf and microwave quantities.

The uncertainties associated with traveling standards and the drifts in value between comparisons have left little margin between national standards and the demands of users for improved accuracy. There are clear advantages if national standards can be referred with sufficient accuracy to some invariant natural physical constant, allowing standards in different countries to be related without transporting them. This has become possible for the volt and the ohm by using the Josephson and quantum Hall effects.

Values of constants. Theory indicates that the values of the Josephson and von Klitzing constants should be equal to $2e/h$ and h/e^2 , respectively, where e is the charge on the electron and h is Planck's constant. These are further examples of relations between electrical standards and fundamental physical constants. It might be thought that the effects could be used for absolute determinations of the volt and

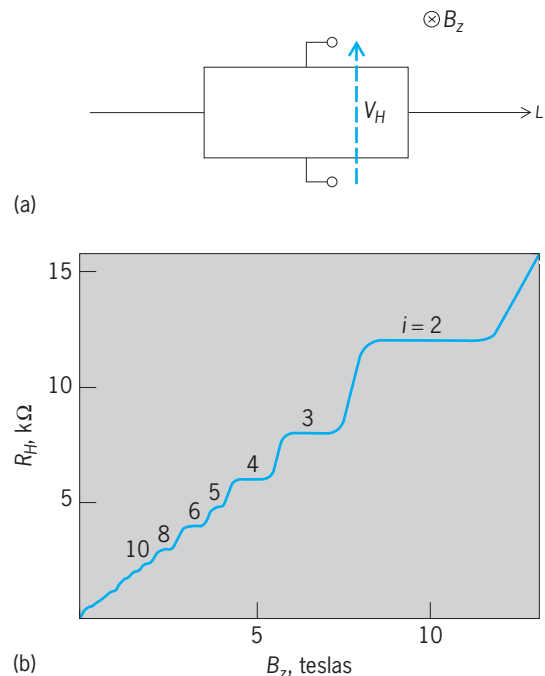


Fig. 4. Quantum Hall effect. (a) Geometry of semiconductor. (b) Hall resistance characteristic. $i =$ quantum number.

the ohm. This is not so, because measurements of e and h refer to the electrical standards and neither is known independently with sufficient accuracy. It is rather the case that these relations provide a new route for the measurement of e and h .

The values of k_J and R_H have therefore to be determined in terms of the absolute values of the electrical standards. Much effort has been devoted to these measurements and to obtaining international agreement on the values to be used. See FUNDAMENTAL CONSTANTS.

International Agreement

As the basis for a generally accepted system of measurement, it is essential that two requirements are satisfied: stability of the primary standards and international agreement on the values to be assigned to them. The first of these has been met for electrical standards by use of the Josephson and quantum Hall effects. The CCEM made an attempt to meet the second requirement in 1972, when it recommended that the value 483 594.0 gigahertz per volt should be used for the Josephson constant. However, this proved to be premature: several major national laboratories continued to use their own, different values, and subsequent measurements have shown this value to be in error by about 8 parts in 10^6 .

Recognizing this, the CCEM encouraged further research and in 1986 set up working groups that made critical assessments of the experimental results bearing on the values of the Josephson and von Klitzing constants. In September 1988 the CCEM recommended that from January 1, 1990, all laboratories should adopt the values of 483 597.9 GHz/V for the Josephson constant (to be designated K_{J-90}) and 25 812.807 Ω for the von Klitzing constant (to be designated R_{K-90}). It recommended that all laboratories should adjust the values assigned to their reference standards accordingly. To ensure consistency, it was recommended that these changes should take place on January 1, 1990, and not before. The CIPM subsequently endorsed the CCEM recommendations.

The CCEM estimated the uncertainties in the recommended values of the constants and expressed them by indicating that the Josephson effect together with the value of K_{J-90} can be used to establish a reference standard of emf having an uncertainty (1 standard deviation) relative to the SI volt of 4 parts in 10^7 and reproducibility significantly better; similarly, the quantum Hall effect, together with the value of R_{K-90} , can be used to establish a reference standard of resistance having an uncertainty (1 standard deviation) relative to the SI ohm of 2 parts in 10^7 and again a reproducibility significantly better.

Implementation of the CCEM recommendations required much work and expense on the part of national standards laboratories and national calibration services. The values of national standards of the volt and the ohm had to be adjusted by up to 10 parts per million for the volt and smaller amounts for the

ohm. These changes, of course, also affected the secondary standards of other derived electrical quantities and caused consequential changes in the calibration values assigned to accurate instruments.

Much work is continuing at the BIPM and the national laboratories to improve further the accuracy and availability of electrical standards. Josephson arrays operating at 10 V dc are available, and the BIPM has conducted international comparisons of the volt and the ohm using traveling Josephson array standards and quantum Hall resistance standards. These have led to results differing by no more than 2 and 4 parts in 10^9 . The BIPM also continues to carry out a variety of international comparisons of derived quantities. See ELECTRICAL MEASUREMENTS.

A. Earle Bailey

Bibliography. R. J. Bell (ed.), *SI: The International System of Units*, 6th ed., 1993; Comité Consultatif d'Electricité et Magnétisme, *Rapport de la 18^e Session, 1988, 19^e Session, 1992, 20^e Session, 1995, 21^e Session, 1997, 22^e Session, 2000, 23^e Session, 2002, 24^e Session, 2005*, BIPM, Sèvres, France; B. P. Kibble, Present state of electrical units, *Proc. IEE*, 138A:187-197, 1991; B. W. Petley, Electrical units: The last thirty years, *Metrologia*, 32:495-502, 1995; B. N. Taylor and T. J. Witt, New international electrical reference standards based on the Josephson and quantum Hall effects, *Metrologia*, 26:47-62, 1989.

Electricity

Those physical phenomena involving electric charges, their motions, and their effects. The motion of a charge is affected by its interaction with the electric field and, for a moving charge, the magnetic field. The electric field acting on a charge arises from the presence of other charges and from a time-varying magnetic field. The magnetic field acting on a moving charge arises from the motion of other charges and from a time-varying electric field. Thus electricity and magnetism are ultimately inextricably linked. In many cases, however, one aspect may dominate, and the separation is meaningful. See ELECTRIC CHARGE; ELECTRIC FIELD; MAGNETISM.

Historical development. The earliest observations of electrical effects were made on naturally occurring substances. Magnetism was observed in the attraction of metallic iron by the iron ore magnetite. The natural resin amber was found to become electrified when rubbed (triboelectrification) and to attract lightweight objects. Both of these phenomena were known to Thales of Miletus (640-546 B.C.). Jerome Cardan in 1551 first clearly distinguished the difference between the attractive properties of amber and magnetite, thus presaging the division of electrical and magnetic effects. He also envisioned electricity as a type of fluid, a viewpoint that was developed more extensively in the late eighteenth and early nineteenth centuries. In 1600 W. Gilbert observed variations in the amounts of electrification

of various substances. He divided substances into two classes, according to whether they did or did not electrify by rubbing. The division actually is into poor and good conductors, respectively. A two-fluid theory was first proposed by C. F. duFay in 1733. A one-fluid theory of electricity was propounded in 1747 by Benjamin Franklin, who called an excess of the fluid positive electrification, and a deficiency of fluid negative electrification. This theory fell into disrepute, but the choice of positive and negative remains. Although fluid theories of electricity were superseded at the end of the nineteenth century, the concept of electricity as a substance persists.

The quantitative development of electricity began late in the eighteenth century. J. B. Priestley in 1767 and C. A. Coulomb in 1785 discovered independently the inverse-square law for stationary charges. This law serves as a foundation for electrostatics. *See* COULOMB'S LAW; ELECTROSTATICS.

In 1800 A. Volta constructed and experimented with the voltaic pile, the predecessor of modern batteries. It provided the first continuous source of electricity. In 1820 H. C. Oersted demonstrated magnetic effects arising from electric currents. The production of induced electric currents by changing magnetic fields was demonstrated by M. Faraday in 1831. In 1851 he also proposed giving physical reality to the concept of lines of force. This was the first step in the direction of shifting the emphasis away from the charges and onto the associated fields. *See* ELECTROMAGNETIC INDUCTION; ELECTROMAGNETISM; LINES OF FORCE.

In 1865 J. C. Maxwell presented his mathematical theory of the electromagnetic field. This theory proposed a continuous electric fluid. It remains valid today in the large realm of electromagnetic phenomena where atomic effects can be neglected. Its most radical prediction, the propagation of electromagnetic radiation, was convincingly demonstrated by H. Hertz in 1887. Thus Maxwell's theory not only synthesized a unified theory of electricity and magnetism, but also showed optics to be a branch of electromagnetism. *See* ELECTROMAGNETIC RADIATION; MAXWELL'S EQUATIONS.

The developments of theories about electricity subsequent to Maxwell have all been concerned with the microscopic realm. Faraday's experiments on electrolysis in 1833 had indicated a natural unit of electric charge, thus pointing toward a discrete rather than continuous charge. Thus, the groundwork for exceptions to Maxwell's theory of electromagnetism was laid even before the theory was developed. H. A. Lorentz began the attempt to reconcile these viewpoints with his electron theory in 1895. He postulated discrete charges, called electrons. The interactions between the electrons were to be determined by the fields as given by Maxwell's equations. The existence of electrons, negatively charged particles, was demonstrated by J. J. Thomson in 1897 using a Crookes tube. The existence of positively charged particles (protons) was shown shortly afterward (1898) by W. Wien, who ob-

served the deflection of canal rays. Since that time, many particles have been found having charges numerically equal to that of the electron. The question of the fundamental nature of these particles remains unsolved, but the concept of a single elementary charge unit is apparently still valid. Of these many particles only two, the electron and the proton, exist in a stable condition on Earth. *See* BARYON; ELECTROLYSIS; ELECTRON; ELEMENTARY PARTICLE; HYPERON; MESON; PROTON; QUARKS.

A second departure from classical Maxwell theory was brought on by M. Planck's studies of the electromagnetic radiation emitted by "black" bodies. These studies led Planck to postulate that electromagnetic radiation was emitted in discrete amounts, called quanta. This quantum hypothesis ultimately led to the formulation of modern quantum mechanics. The most satisfactory fusion of electromagnetic theory and quantum mechanics was achieved in 1948 with the work of J. Schwinger and R. Feynman in quantum electrodynamics, which suppressed the particle aspect and emphasized the field. *See* HEAT RADIATION; QUANTUM ELECTRODYNAMICS; QUANTUM MECHANICS.

Sources. The sources of electricity in modern technology depend strongly on the application for which they are intended.

The principal use of static electricity today is in the production of high electric fields. Such fields are used in industry for testing the ability of components such as insulators and capacitors to withstand high voltages, and as accelerating fields for charged-particle accelerators. The principal source of such fields today is the Van de Graaff generator. *See* PARTICLE ACCELERATOR.

The major use of electricity today arises in devices using direct current and low-frequency alternating current. The use of alternating current, introduced by S. Z. de Ferranti in 1885-1890, allows power transmission over long distances at very high voltages with a resulting low-percentage power loss followed by highly efficient conversion to lower voltages for the consumer through the use of transformers. *See* ALTERNATING CURRENT; ELECTRIC CURRENT.

The development of materials that are superconducting at liquid-nitrogen temperatures may alter this concept. The cost of liquid nitrogen is trivial compared with the cost of the liquid helium previously required for maintaining superconductivity in materials. If the new materials can be formed into malleable wire, this economic factor could lead to the construction of lossless power lines and other applications in computers and in magnetically levitated transportation. *See* MAGNETIC LEVITATION; SUPERCONDUCTING DEVICES; SUPERCONDUCTIVITY.

Large amounts of direct current are used in the electrodeposition of metals, both in plating and in metal production, for example, in the reduction of aluminum ore. To avoid power transmission difficulties, such facilities are frequently located near sources of abundant power. *See* DIRECT CURRENT; ELECTROCHEMISTRY; ELECTROMETALLURGY; ELECTROPLATING OF METALS.

The principal sources of low-frequency electricity are generators based on the motion of a conducting medium through a magnetic field. The moving charges interact with the magnetic field to give a charge motion that is normal to both the direction of motion and the magnetic field. In the most common form, conducting wire coils rotate in an applied magnetic field. *See* ALTERNATING-CURRENT GENERATOR; DIRECT-CURRENT GENERATOR; GENERATOR.

The rotational power is derived from a water-driven turbine in the case of hydroelectric generation, or from a gas-driven turbine or reciprocating engine in other cases. In the latter situation the high-pressure gas that is required is obtained by heating a fluid or gas, using some thermal source, usually the burning of fossil fuels. Economic considerations, particularly the cost of natural gas and oil and the need to conserve them for petrochemical purposes, have pointed out the need for increased reliance on coal and nuclear reactors as heat sources. *See* ELECTRIC POWER GENERATION; GAS TURBINE; HYDROELECTRIC GENERATOR; NUCLEAR REACTOR.

An alternate method is magnetohydrodynamic generation. The requisite motion is that of a high-temperature gas flowing through the magnetic field. The gas consists of the combustion products of the fossil-fuel heat source with a small amount of potassium added to increase the conductivity. Extensive studies have been conducted, especially in the United States and Russia, to determine the feasibility of this method. Because the gas temperatures are very high, the exhaust gas can be used as a heat source for a conventional generator. It is estimated that a combined facility would be 50–55% efficient. This compares with an efficiency of 35–40% for the conventional facility. An added benefit appears to be a ready solution of the sulfur dioxide (SO₂) emission problem for coal-fueled systems. The potassium would efficiently interact with the sulfur dioxide and yield an easily removable solid residue. *See* ACID RAIN; ENERGY SOURCES; MAGNETOHYDRODYNAMIC POWER GENERATOR; MAGNETOHYDRODYNAMICS.

A more direct method of using fission or fusion reactors is the direct conversion of the energy released in the nuclear process into electricity. This has been achieved on a laboratory scale in the case of fission reactors.

Many high-frequency devices, such as communications equipment, television, and radar, involve the consumption of only moderate amounts of power, generally derived from low-frequency sources. If the power requirements are moderate and portability is needed, the use of ordinary chemical batteries is possible. Ion-permeable membrane batteries are a later development in this line. Fuel cells, particularly hydrogen-oxygen systems, are being developed. They have already found extensive application in earth satellite and other space systems. The successful use of thermoelectric generators based on the Seebeck effect in semiconductors has been reported in Russia and in the United States. In a partic-

ularly compact low-power device constructed in the United States, the heat needed for the operation of such a generator has been supplied by the energy release in the radioactive decay of suitably encapsulated isotopes produced in fission reactors. *See* BATTERY; FUEL CELL; ION-SELECTIVE MEMBRANES AND ELECTRODES; NUCLEAR BATTERY; THERMOELECTRIC POWER GENERATOR; THERMOELECTRICITY.

The solar battery, also a semiconductor device, has been used to provide charging current for storage batteries in telephone service and in communications equipment in artificial satellites. *See* SOLAR CELL.

Direct conversion of mechanical energy into electrical energy is possible by utilizing the phenomena of piezoelectricity and magnetostriction. These have some application in acoustics and stress measurements. Pyroelectricity is a thermodynamic corollary of piezoelectricity. *See* MAGNETOstriction; PIEZOELECTRICITY; PYROELECTRICITY.

Some other sources of electricity are those in which charged particles are released with some energy and collected in some manner. Charged particles are suitably released in radioactive decay, in the photoelectric effect, and in thermionic emission, among other ways. The photovoltaic effect may also be in this group. *See* ELECTRON EMISSION; PHOTOVOLTAIC EFFECT; RADIOACTIVITY; THERMIONIC EMISSION.

The differences of work functions of various materials can be used for energy conversion. The contact potential difference may be used to convert heat directly to electricity or to provide improved collection for currents arising from some other source such as radioactivity. *See* WORK FUNCTION (ELECTRONICS).

Other possible sources of electricity arise from the existence of electrokinetic potentials in flowing fluids and of phase-transition potentials such as occur in the Workman-Reynolds effect. The possibilities of combining several effects also exist as exemplified in thermogalvanic potentials. It also appears that organic materials (as distinguished from the inorganic materials for which most of the work already described was done) merit investigation. A primitive type of organic solar battery has been developed. *See* ATMOSPHERIC ELECTRICITY; CIRCUIT (ELECTRICITY); CONDUCTION (ELECTRICITY); ELECTRIC POWER MEASUREMENT; ELECTRICAL ENGINEERING; ELECTRICAL UNITS AND STANDARDS; ELECTROKINETIC PHENOMENA; ELECTRONICS; GEOELECTRICITY.

Walter Aron

Bibliography. P. H. Abelson (ed.), *Energy: Use, Conservation and Supply*, American Association for the Advancement of Science, 2 vols., 1974, 1978; B. I. Bleaney and B. Bleaney, *Electricity and Magnetism*, 2 vols., 3d ed., 1976, paper, 1989; R. P. Feynman et al., *Feynman Lectures on Physics: The Definitive and Extended Edition*, 3 vols., 2005; R. J. Fowler, *Electricity: Principles and Applications*, 6th ed., 2004; E. M. Purcell, *Electricity and Magnetism*, 2d ed., 1985; P. B. Visscher, *Fields and Electrodynamics*, 1988.

Electrochemical equivalent

The weight of a substance, according to Faraday's law, produced or consumed by electrolysis with 100% current efficiency during the flow of a quantity of electricity equal to 1 faraday or 96,485.309 coulombs per mole (1 coulomb corresponds to a current of 1 ampere during 1 second). Electrochemical equivalents are essential in the calculation of the current efficiency of an electrode process.

The electrochemical equivalent of a substance is equal to the gram-atomic or gram-molecular weight of this substance divided by the number of electrons involved in the electrode reaction. For example, the electrochemical equivalent of zinc, for which two electrons are required in order to deposit one atom, is $Zn/2$ or $65.37/2$ g. Thus, the faraday is equal to the product of the charge of the electron times the number of electrons (the Avogadro number) required to react with 1 atom- or molecule-equivalent of substance. The value of the faraday computed in this manner agrees with values obtained from electrochemical determinations. The relative error on the value 96,485.309 coulombs is smaller than $\pm 0.01\%$. See COULOMETER; ELECTROLYSIS. Paul Delahay

Electrochemical process

The principles of electrochemistry may be adapted for use in the preparation of commercially important quantities of certain substances, both inorganic and organic in nature.

Inorganic Processes

Inorganic chemical processes can be classified as electrolytic, electrothermic, and miscellaneous processes including electric discharge through gases and separation by electrical means. In electrolytic processes, chemical and electrical energy are interchanged. Current passed through an electrolytic cell causes chemical reactions at the electrodes. Voltaic cells convert chemicals into electricity. Electrothermic processes use electricity to attain the necessary temperature for reaction. For related information see ELECTROCHEMISTRY; ELECTROLYSIS; ELECTROLYTIC CONDUCTANCE; ELECTROMOTIVE FORCE (CELLS).

For a discussion of equipment used to convert alternating current to direct current for electrolytic plants. See MOTOR-GENERATOR SET; SEMICONDUCTOR RECTIFIER; SYNCHRONOUS CONVERTER.

Voltaic cells are used for the intermittent production of small amounts of electricity. When the chemicals involved are exhausted and must be replaced, the unit is called a primary cell. A special case of the primary cell is the fuel cell in which the fuel and oxidizer are fed continuously to the cell, converted to electricity, and the products removed. If exhausted components can be revived by passing electricity backward through the unit, it is called a secondary cell, storage battery, or accumulator. Cells

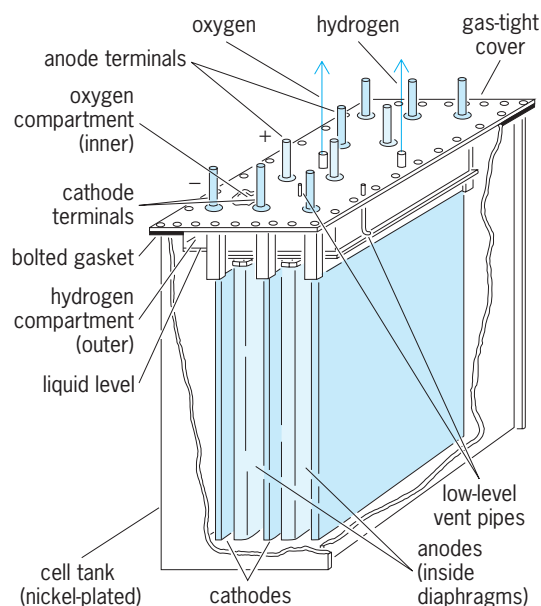


Fig. 1. Diagram of Stuart hydrogen-oxygen cell. (Electrolyser Corp. Ltd., Toronto, Canada)

may be connected in parallel or in series to form a battery.

For a discussion of the theory and description of commercial primary, secondary, and fuel cells see BATTERY; FUEL CELL.

Electrolysis in aqueous solutions. The electrolysis of water to form hydrogen and oxygen, according to the reaction $2H_2O \rightarrow 2H_2 + O_2$, may be considered as the simplest process for aqueous electrolytes. It does not compete with hydrogen from propane or from natural gas and with oxygen from liquid air, except in small installations. While simplicity, high hydrogen purity requirement, and lower capital cost (in small plants) have justified electrolytic plants, severely rising energy costs limited such applications. The electrolyte is 18–30% NaOH or KOH, the latter having a lower resistance but higher cost. The cathode is steel and the anode is nickel-plated steel separated by a diaphragm, usually of asbestos (Fig. 1). A cell voltage of 2.0–2.5 V is the summation of the decomposition voltage of water (1.23 V at room temperature) and the oxygen and hydrogen over voltages, plus the IR drop through the electrolyte, electrode contacts, and bus bars. The raw material is distilled or demineralized water. There are monopolar cells operating up to 20,000 A, bipolar or filter-press cells using 2000–5000 A, and 150 cells in series at 300–400 V overall. Cells are available which operate at 600 psi (4.1 megapascals). A. T. Kuh described 11 cells using 25–30% KOH and indicated designs operating at pressures up to 200 atm (20 MPa). See HYDROGEN; OXYGEN.

Heavy water, or deuterium oxide, used in moderating nuclear reactors is also a by-product of the electrolysis of water. Protium (H^1) is preferentially discharged, so that the electrolyte becomes richer in deuterium. Electrolysis must be combined with catalytic exchange or distillation processes when used

in primary or earlier stages of concentration; it is also used in final stages to produce 99–100% concentration. See DEUTERIUM; HEAVY WATER.

Metallurgical applications. Protective or decorative coatings on a base metal such as steel are obtained by electroplating. Plating may also be used to replace worn metal or to provide a wear-resistant surface. The final surface may require several layers of different metals or even layers of the same metal deposited under varying conditions. The metals plated are copper, cadmium, chromium, cobalt, gold, iron, lead, nickel, the platinum metals, silver, tin, and zinc, and many alloys. Electrogalvanizing is preferred over hot dipping for applying zinc to steel. Tin plate for containers is electrolytic. Factors affecting the resulting plate include pretreatment and cleaning of the metal surface, current density, concentration of metal ions, agitation, temperature, conductance of solution, pH, and addition agents. See ELECTROPLATING OF METALS; METAL COATINGS.

Electroforming is a method of forming or reproducing articles by electrodeposition. In contrast to electroplating, the product is removed from the base surface or mold. A nonconducting surface can be made conductive by metallizing or with graphite powder. A low-melting-point metal mold can be used, or a mold can be plated with a metal from which the final metal can be removed. The electrodeposits may be up to 0.5 in. (13 mm) thick, and after completion are removed from the mold. Phonograph records, electrotype, textile replicas, and brass instrument bells are common products of electroforming, which is also widely used in the automotive and electronic industries.

Electrodeposition of metal powders is used to produce particles in the 1–1000-micrometer range for use in powder metallurgy and metallic pigments. Powdered iron made electrolytically makes stronger parts than other types, and electrolytic powdered copper makes parts that are easier to machine and have improved wear.

In the anodizing of aluminum articles, a coating approximately 0.001 in. (25 μm) thick is applied, which can be dyed or made impervious. The article to be anodized is cleaned and made anodic in sulfuric acid solution.

Electrolytic polishing of metals is accomplished by making the article anodic in an electrolyte of mixed acids, such as phosphoric-chromic acids or phosphoric-sulfuric acids. High points on the surface apparently dissolve, to give a unique polish not achievable by mechanical polishing.

Electrolytic machining of metals is accomplished by making the metal part anodic in a suitable electrolyte. Metal dissolves at the anode and hydrogen discharges at the cathode. The cathode is contoured as a negative of the shape to be developed in the work piece. By circulating the electrolyte under very high pressure and using high current densities (30–2000 A/in.² or 5–300 A/cm²), practical machining rates are achieved. Very hard and very thin metal surfaces can be machined without changing the heat treatment.

Electrorefining is a process for purifying metals and recovering their impurities, which at times are more valuable than the original metal. Copper from its ore or scrap is purified of volatilizable impurities, cast into an anode, dissolved in an electrolyte in a cell, and deposited at a cathode as a very pure, highly conductive metal for electrical engineering purposes (Fig. 2). Gold, silver platinum, selenium, and tellurium are recovered as by-products. Nickel is freed of copper by using a diaphragm between anode and cathode, with cobalt and the platinum

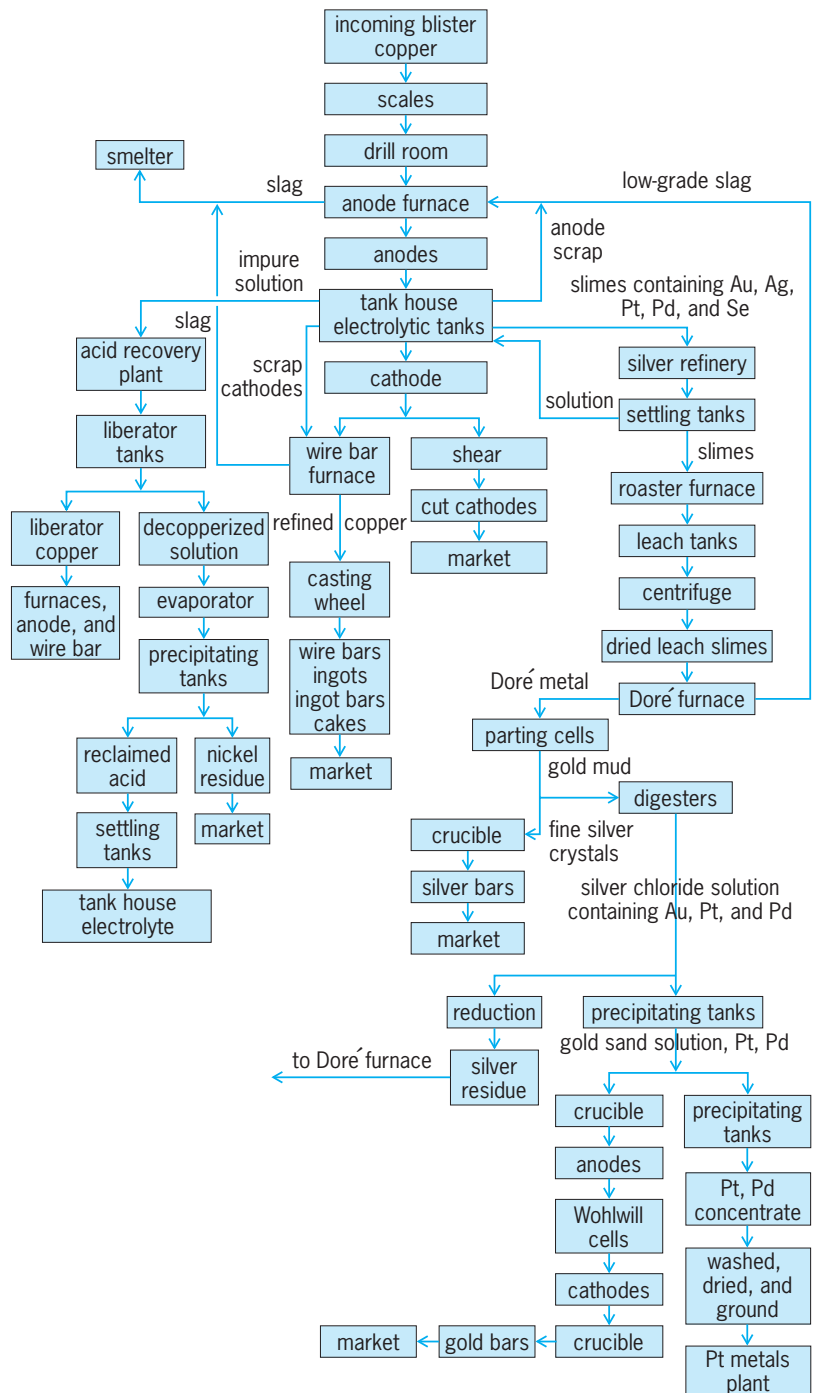


Fig. 2. Flowsheet for a copper plant and the refinery by-products resulting from a product of 99.97% conductor copper. (Ontario Refining Co., International Nickel Co. of Canada, Ltd.)

TABLE 1. Electrolytics for various individual metals

Metal	Source	Method	Electrolyte	Application
Antimony	Crude metal	Refining	Sulfates	Alloys, batteries, chemicals
Antimony	Antimony ores	Winning	Sulfides	
Bismuth	Lead refining slimes	Refining	Chlorides	Alloying agent, medicine
Cadmium	Zinc residues and slimes	Winning	Sulfates	Electroplating, alloys
Chromium	Chromium ores	Winning	Sulfates	High-temperature, alloys
Cobalt	Complex cobalt ores	Winning	Sulfates	Alloying agent, electronics
Cobalt	Complex nickel ores	Refining	Sulfates and chlorides	Alloying agent, electronics
Copper	Copper ores, complex ores	Winning	Sulfates	Electrical conductors, wire, brass alloys
Copper	Crude metal, secondary materials, waste	Refining	Sulfates	Electrical conductors, wire, brass alloys
Copper powder	Refined metal	Plating	Sulfates	Powders for powder metallurgy, oilless bearings
Copper sheet	Purified electrolyte	Plating	Sulfates	Sheet for printed circuits, electronics
Gallium	Sodium aluminate liquors	Winning	Caustic	Low-melting-point metals and alloys
Gold	Copper refining slimes	Refining	Chlorides	Jewelry, dentistry, plating
Indium	Residues, wastes, ores	Refining	Chlorides	Silver alloys, jewelry, television
Iron	Low-carbon steel	Refining	Sulfates	Powder metallurgy
Lead	Crude lead	Refining	Fluosilicates	Separation of bismuth etc., alloys, fittings
Manganese	Manganese ores	Winning	Sulfates	Stainless steels, carbon-free metal, alloys
Nickel	Crude metal, nickel matte	Refining	Sulfate-chloride	Plating, alloys, stainless steel
Silver	Copper refinery slimes, lead residues, crude metal, ore	Refining	Nitrates	Jewelry, electrical applications, alloys
Silver	Silver alloys, photographic wastes	Winning	Nitrates	Jewelry, electrical applications, alloys
Solder	Waste and crudes	Refining	Fluosilicate	Metal joining, electronic components
Tin	Crude metal	Refining	Cresol-sulfonates	Tin plate, solder, alloys
Zinc	Zinc ores, complex ores	Winning	Sulfates	Die-casting alloys, battery cups, brass, galvanizing
Zinc	Ores and residues	Winning	Caustic	Chemicals, paint

metals as by-products. Lead is freed from bismuth; tin from lead, antimony, and bismuth with silver as a by-product; and zinc from copper and lead, with cadmium as a by-product, to make commercial a 99.99% pure metal for diecasting. The last operation is usually referred to as electrowinning. This has rapidly replaced pyrometallurgy or fire processing because it can eliminate sulfur dioxide atmospheric pollution and contamination by particulates.

Electrowinning, sometimes termed aqueous electrometallurgy, involves processing of metallic ores, usually of very low metal content but large in volume, by leaching solutions, usually sulfates, to ob-

tain metal-containing electrolytes which can be processed with insoluble anodes and metal cathodes. Examples are found in cadmium, copper, cobalt, manganese, and zinc.

Electrorefining, electrowinning, and electroforming are summarized as to the individual metals in **Table 1**, while **Table 2** gives anodes, cathodes, and diaphragms for commercially successful operation. **Table 3** lists energy consumption of aqueous electrochemical operations. See ELECTROMETALLURGY.

Electrolytic corrosion of metals. This occurs because some parts of the surface of metals act as

TABLE 2. Anodes, cathodes, and diaphragms used on commercial processes

Metal	Method	Anode	Anolyte-catholyte	Diaphragm	Cathode	Voltage
Antimony	Winning	Insoluble	Same	No	Steel	2.5–3
Cadmium	Winning	Insoluble	Same	No	Aluminum	4.0
Chromium	Winning	Insoluble	More acid–Acid	Yes	Hastelloy	4.2
Cobalt	Refining	Soluble	Same	No	Cobalt-stainless steel	2.5
Copper	Refining	Soluble	Same	No	Copper	0.2–0.3
Copper	Winning	Insoluble	Same	No	Copper	2–2.1
Gold	Refining	Soluble	Same	No	Gold	0.5–2.8
Lead	Refining	Soluble	Same	No	Lead	0.35–0.45
Manganese	Winning	Insoluble	Acid–Alkaline	Yes	Stainless steel	
Nickel	Refining	Soluble	Same–Pure	Yes	Nickel	2.4
Silver	Refining	Soluble	Same	Yes	Stainless steel or carbon	1.3–5.4
Tin	Refining	Soluble	Same	No	Tin	0.3
Zinc	Winning	Insoluble	Same	No	Aluminum	3.25–3.7

anodes and corrode, whereas other parts act as cathodes and do not corrode.

Cathodic protection is provided if the whole surface is made cathodic to a separate anode and sufficient voltage is available between the two electrodes. This type of protection is used to inhibit corrosion of boilers, condensers, underground pipelines, ships, and water tanks. Sacrificial anodes of zinc, magnesium, or aluminum alloy may provide the potential, or inert anodes such as graphite, stainless steel, or platinum-plated titanium may be used with power supplied from a rectifier.

Anodic protection can be used to create a passive

layer on the surface of some metals, such as steel and stainless steel in some environments. It is a practical method of controlling corrosion of tanks in the chemical industry, but is not feasible for copper or brass vessels. The tank is made anodic. An inert cathode, such as platinum-clad metal, is installed in the liquid in the tank. Current is applied so as to maintain a predetermined voltage between the anodic surface and a reference electrode, such as silver-silver chloride, in the liquid. Equipment to maintain precise potential control at high current output has made anodic protection practical. For example, a current of 0.0015 A/ft² (0.016 A/m²) at +0.900 V to a

TABLE 3. Energy consumption of electrochemical products*

Industry	kWh/lb	lb/kWh	Voltage/tank, cell, or furnace	Range, A/unit	Line voltage	Range, kVA/cell or furnace	kVA/line
Electrolytic refining							
Copper: multiple system, series system	0.09–0.2	5–13.3	0.2–0.3	6000–15,000	80–200	1.2–4.5	480–3000
Gold (troy lb)	0.074	13.5	120–200				
Lead	0.15	6.6	0.5–2.8	150–500	3–10	0.1–1.4	0.5–5
Nickel	0.07–0.09	11.1–14.3	0.35–0.5	5000–6000	100–185	1.75–3.0	500–1200
Silver, Moebius (troy lb)	1.1	0.9	2.4–2.6	5000–6000	220–230	12–16	1300–1500
Silver, Thum (troy lb)	0.27–0.3	3.3–3.7	2.3–2.8	400–500	45–250	1–1.4	18–125
Solder	0.33–0.6	1.67–3.0	1.3–5.4	150–200	200–220	0.3–1.2	30–50
Tin	0.08	12.5	0.34				
Tin	0.085	11.8	0.3–0.35	5000	100–200	1.5–2.0	500–1000
Electrowinning							
Antimony			2.5–3	1500	300	4	450
Cadmium	0.65–0.97	1.03–1.54	2.5–3.1				
Chromium	5.0–8.4	0.12–0.2	4.2–4.3	10,000	250–300	40–45	2500–3000
Cobalt	1.2–1.56	0.64–0.83	3.7–4.5				
Copper	0.89–1.34	0.77–1.12	2.0–2.12	10,000–25,000	150–200	20–30	1500–5000
Manganese	4–4.5 dc	0.22–0.25	5–5.4	6000	600	30–40	3600
Silver	5–5.3 dc	0.19–0.2					
Zinc	0.58–1.0	1–1.73	1–1.5	300		0.3–0.5	0.3–0.5
Zinc	1.4–1.6	0.61–0.71	3–3.7	5000–10,000	600–800	15–40	3000–8000
Metal melting							
Copper	0.12–0.15	6.67–8.33	85–225	12,000–30,000		6000	
Copper alloys	0.15–0.3	3.33–6.67					
Steel: cold charge, hot charge	0.237–0.35	2.86–4.22	80–250	30,000–100,000		25,000–33,000	
Zinc	0.05–0.2	5–20	80–250	30,000–100,000		25,000–33,000	
Zinc	0.045	22.2					
Metal powders							
Copper	0.3–0.5	2–3.3	0.3–0.5	10,000–15,000			
Iron	1.20	0.83	2.5				
Nickel	1.13	0.83	1.4–1.5				
Zinc	1.37	0.73	3.4				

*1 kWh/lb = 7.937 MJ/kg; 1 lb/kWh = 0.1260 kg/MJ.

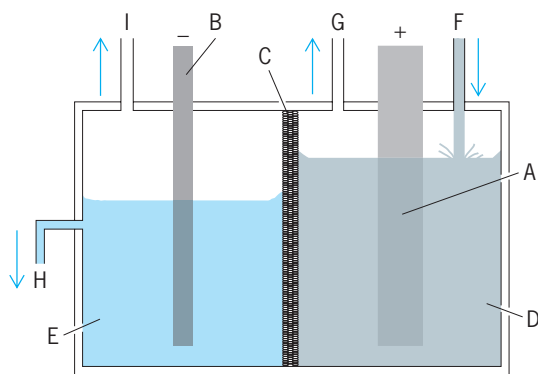
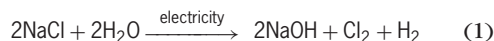


Fig. 3. Diagram of diaphragm cell for chlorine and caustic soda. A = graphite anode, B = iron screen cathode, C = asbestos diaphragm, D = anode compartment for brine and chlorine, E = cathode compartment for NaOH-NaCl cell liquor, F = brine inlet, G = chlorine outlet, H = cell liquor (caustic soda) outlet, I = hydrogen outlet.

silver-silver chloride electrode will maintain passivity of a carbon-steel tank holding 93% sulfuric acid at 27°C (80°F). See CORROSION.

Alkali-chlorine processes. Electrolysis of alkali-halides is the basis of the alkali-chlorine and chlorate industries. Chlorine, Cl_2 , and caustic soda, NaOH (or caustic potash, KOH), are made by electrolysis of brine, a solution of sodium chloride, NaCl, in water. This is represented by reaction (1).



Two processes are used to prevent the products from the diaphragm cell and the mercury cell. In the diaphragm cell process (Fig. 3) an asbestos diaphragm is interposed between a graphite anode and an iron screen cathode. Saturated purified brine fed around the anode passes through the diaphragm to the cathode. Chlorine is formed at the anode. Hydrogen is released at the cathode, leaving NaOH as a 10-15% solution and 10-15% residual NaCl

crystallizes out and is recycled. The decomposition voltage of brine to form chlorine and hydrogen is 2.3 V. At 0.75 A/in.² (11.6 A/dm²) the average voltage components of a diaphragm cell are as follows:

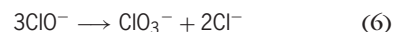
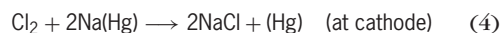
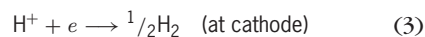
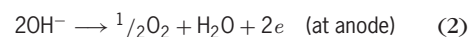
Anode potential	1.50 V
Cathode potential	1.25
Anolyte	0.47
Diaphragm	0.30
Conductors	<u>0.18</u>
Total	3.70 V

The current efficiency is 95.5-95.5%, due to some oxygen discharge at the anode and some chlorine being carried through the diaphragm. An installation is pictured in Fig. 4.

In the mercury cell process brine is electrolyzed between graphite anodes and a flowing mercury cathode, forming a dilute (0.2-0.4%) sodium amalgam which is decomposed in another compartment by water in contact with graphite surfaces to form H_2 and NaOH (Figs. 5-7). The products of the mercury cell are purer than those of the diaphragm cell. To offset the cost of mercury, a much higher current density is used in mercury cells. Typical components of voltage in a mercury cell at 5.12 A/in.² (80 A/dm²) are as follows:

Anode potential, reversible	1.34 V
Cathode potential, reversible	<u>1.76</u>
Decomposition voltage	3.10 V
Anode polarization	0.35 V
Cathode polarization	0.06
Electrolyte	0.60
Conductors and contacts	<u>0.29</u>
Total cell voltage	4.40 V

Economic factors dictate the use of higher current densities, equal to or exceeding 6.5 amperes/in.² (100 A/dm²). Cell voltage at these higher current densities can be calculated for good cell designs on the market from the equation: $V = 3.20 + 0.015 C$, where C is the cathode current density in amperes per square decimeter. The current efficiency is approximately 95%. Inefficiency reactions are demonstrated in reactions (2) through (6). Adverse con-



ditions can increase these inefficiencies. See CHLORINE.

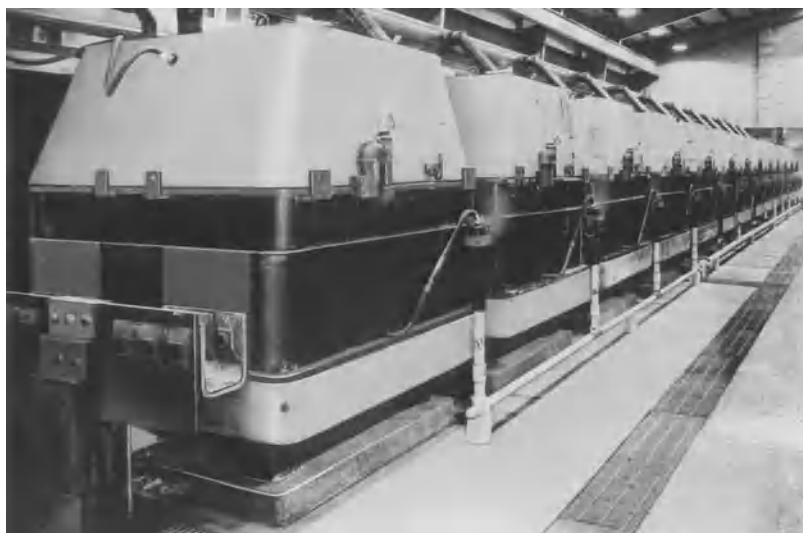


Fig. 4. Photograph of an installation of diaphragm alkali-chlorine cells. (Hooker Chemical Corp.)

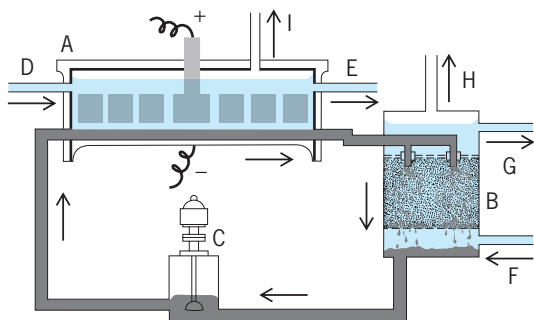


Fig. 5. Diagram of a mercury alkali-chlorine cell. A = electrolyzer, B = decomposer with graphite packing, C = mercury pump, D = feed brine, E = spent brine, F = water, G = 50% caustic soda, H = hydrogen, I = chlorine.

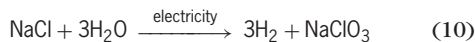
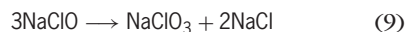
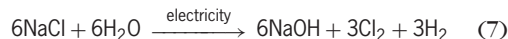
Two factors have been important in the development of chlorine technology: (1) The Nafion diaphragm of the synthetic resin type replacing the deposited asbestos diaphragm. (2) The dimensionally stable anode replacing the graphite anode. The dimensionally stable anode is a titanium substrate with a platinum-group coating of metals and oxides. Its use eliminates the continuously necessary voltage increase or anode-cathode spacing adjustment needed because of graphite-anode wear.

In addition, there has been a swing away from mercury cells because of widespread publicity of so-called mercury poisoning by mercury discharges. These discharges have been reduced by better hydrogen cooling and recycle of metal dross, cutting mercury losses as well as permitting better "house-keeping."

Sodium hypochlorite is formed when the products of the electrolysis of brine are mixed. Electrolytic cells have been built for this purpose, but have limited or special use, such as for sterilization of swimming pools and algae control in power plant condensers. Sodium hypochlorite is usually made chemically.

Sodium chlorate is made in cells with graphite or lead peroxide anodes and steel cathodes. When mixing is encouraged, changes take place according

to reactions (7)–(9). The overall reaction is (10).



The temperature is kept below 40°C (104°F) in cells using graphite anodes to prevent excessive attack. The optimum efficiency is at pH 6.8; hydrochloric acid is added as required. Sodium dichromate prevents reduction of chlorate at the steel cathode. The conversion of hypochlorite to chlorate is a somewhat slow chemical reaction, occurring partly in the cells and partly in a rundown tank. Salt is added and electrolysis is continued until the sodium chlorate is down to about 100 g/liter and the chlorate has reached the desired concentration. It is then recovered by crystallization. Cells operate at 3–3.5 V, 1–3 A/dm², and 80–85% current efficiency. Energy consumption is about 2.5 kWh/lb (20 megajoules/kg).

Hydrochloric acid electrolysis is of interest for recovery of chlorine from HCl resulting as a by-product from organic chlorinations. A filter-press electrolyzer is used which has 30–50 unit cells with polyvinyl-cloth diaphragms and graphite electrodes. Hydrochloric acid of 30–33% concentration is fed to the anode compartment. Weak acid is withdrawn at about 20% and reconcentrated by absorption of HCl gas. The graphite anode is not attacked as long as the concentration is kept at 20% HCl or higher. The current efficiency is 92–96%, the loss due to electrical leakage. Energy consumption is 1800 kWh/2000 lb (7.1 MJ/kg) chlorine (direct current). The voltage balance of a unit cell is as follows:

Anode potential	1.02 V
Cathode potential	0.28
Anode polarization	0.2
Cathode polarization	0.5
Electrolyte, diaphragm	0.3
Total	2.30 V

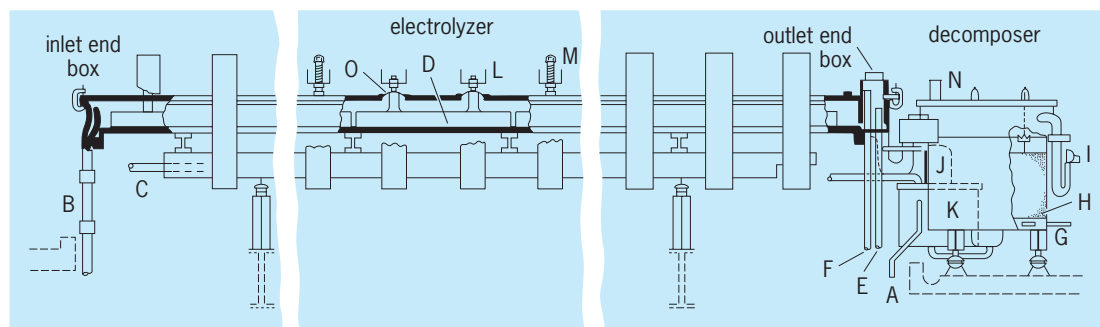


Fig. 6. Longitudinal section of Olin Mathieson E-11 mercury cell. A = dilute caustic outlet, B = brine inlet, C = mercury return, D = anode, E = brine-chlorine outlet, F = outlet end box vent, G = water inlet, H = graphite packing, I = caustic outlet, J = mercury pump, K = mercury pump sump, L = anode support bus, M = lifting screws, N = hydrogen outlet, O = anode seal. (Olin Mathieson Chemical Corp.)

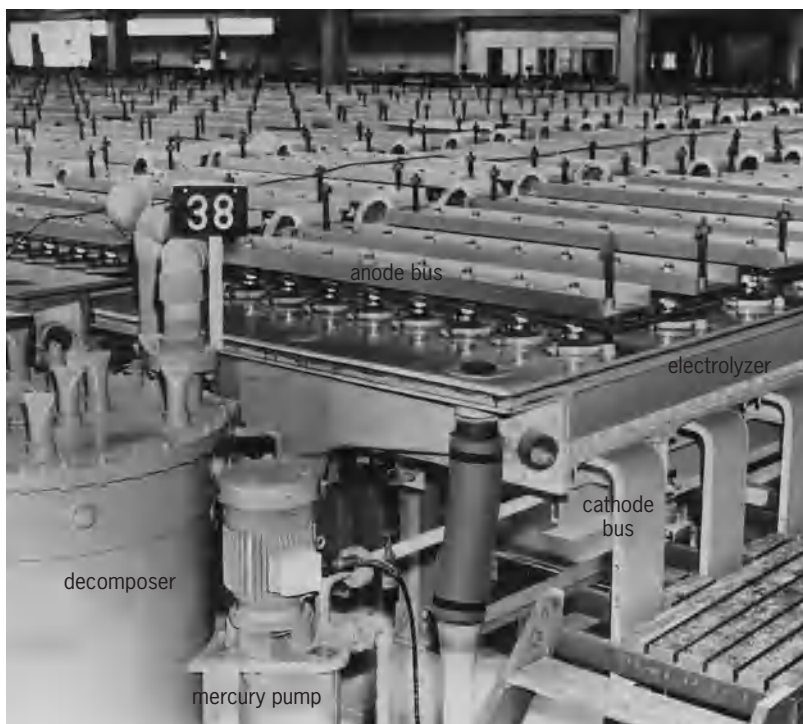
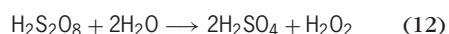
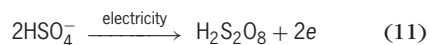


Fig. 7. Chlor-alkali mercury cell; 300,000-A capacity. (Olin Mathieson)

Oxidations and reductions. These reactions occur in all cells, but in a narrower sense oxidation reactions are those in which oxygen or chlorine at the anode oxidizes some material to form a new compound; reduction reactions are those in which hydrogen, liberated at the cathode, reduces a material to a new product. There are no commercial applications of inorganic electrochemical reductions by this narrow definition.

Sodium perchlorate is made by oxidation of a solution containing NaClO_3 at pH 6.1–6.4 by use of a platinum anode and an iron cathode, with chromate in the electrolyte. A lead dioxide anode may be used with a stainless-steel or nickel cathode, with no chromate in the electrolyte. Energy consumption is 1.4–1.6 kWh/lb (11–13 MJ/kg) NaClO_4 (direct current). Other perchlorates are made by metathesis with NaClO_4 .

Persulfuric acid, $\text{H}_2\text{S}_2\text{O}_8$, is made by oxidizing sulfuric acid as an intermediate in the production of hydrogen peroxide. The reactions for the process are shown in reactions (11) and (12). Alkali persulfates



can be made in the same way. The cell has smooth platinum anodes, a porous stoneware diaphragm, and a lead cathode cooled to 30°C (86°F). The reversible potential is 2.18 V and the operating voltage 5.0–5.5 V. The energy released as heat in the cell must be removed by cooling with stoneware or glass coils in the cell. Hydrogen peroxide is recovered by distillation. See PEROXIDE.

Lead peroxide anodes are used in sodium chlorate, sodium perchlorate, sodium bromate, and periodic or periodic acid regeneration cells. The material is dense and made from a lead nitrate solution. For some uses it is produced on a steel base, from which it is removed mechanically and chemically. It is also applied to tantalum, platinum-clad tantalum, or graphite.

Periodic acid is used in producing dialdehyde starch, and spent solution from oxidation can be regenerated in a cell by use of lead dioxide anode, a porous ceramic diaphragm, and an iron cathode.

Electrolytic manganese dioxide for batteries is made by electrolyzing hot MnSO_4 solutions at pH 6.5–7.5 by use of graphite electrodes. The MnO_2 deposited on the anode is pulverized with the graphite and separated mechanically. Energy consumption is 1 kWh/lb (8 MJ/kg) MnO_2 . The quality of the MnO_2 for battery use depends on cell temperature and anode current density.

Ion-permeable membrane cells. These utilize diaphragms made of ion-exchange resins. Cation-permeable membranes permit cations to pass through but not anions, whereas the reverse holds for anion-permeable membranes. **Figure 8** shows the movement of ion and water in an electric membrane stack. Purification of seawater is the most important application. Salt has been recovered from seawater that has been concentrated in this way. See ION-SELECTIVE MEMBRANES AND ELECTRODES; WATER DESALINATION.

Fused-salt electrolysis. Aluminum, barium, beryllium, cerium and misch metal, fluorine, lithium, magnesium, sodium, molybdenum, thorium, titanium, uranium, and zirconium are obtained by electrolysis of fused salts, because water interferes with the desired reaction. Raw materials must all be purified before addition to fused-salt electrolytes as aqueous electrolytes. Metallizing is a process of depositing a metal as an alloy on a substrate from a fused complex metal salt.

Aluminum is produced in steel pots, lined with carbon, graphite, or silicon carbide, containing an electrolyte of alumina dissolved in fused cryolite, $\text{AlF}_3 \cdot 3\text{NaF}$ at $950\text{--}1000^\circ\text{C}$ ($1740\text{--}1830^\circ\text{F}$). The pool of aluminum in the bottom of the pot is cathode. Contact is made by iron bars buried in the carbon lining or by titanium diboride contacts. Anodes are prebaked carbon blocks or the Söderberg type made from carbon paste baked in place. The arrangement of the apparatus for aluminum production is shown in **Fig. 9**. Aluminum is siphoned out periodically. Oxygen released at the carbon anode forms carbon monoxide. See ALUMINUM.

Barium, used in the electronics industry and in alloys, is obtained by aluminum reduction at low pressures and high temperatures.

Beryllium is made by batch electrolysis of fused salt, starting with 25% BeCl_2 and 75% NaCl in a chrome-iron pot, which acts as cathode to a graphite anode. Beryllium is deposited on the wall of the pot and is cleaned out and broken up when cold. Salt is

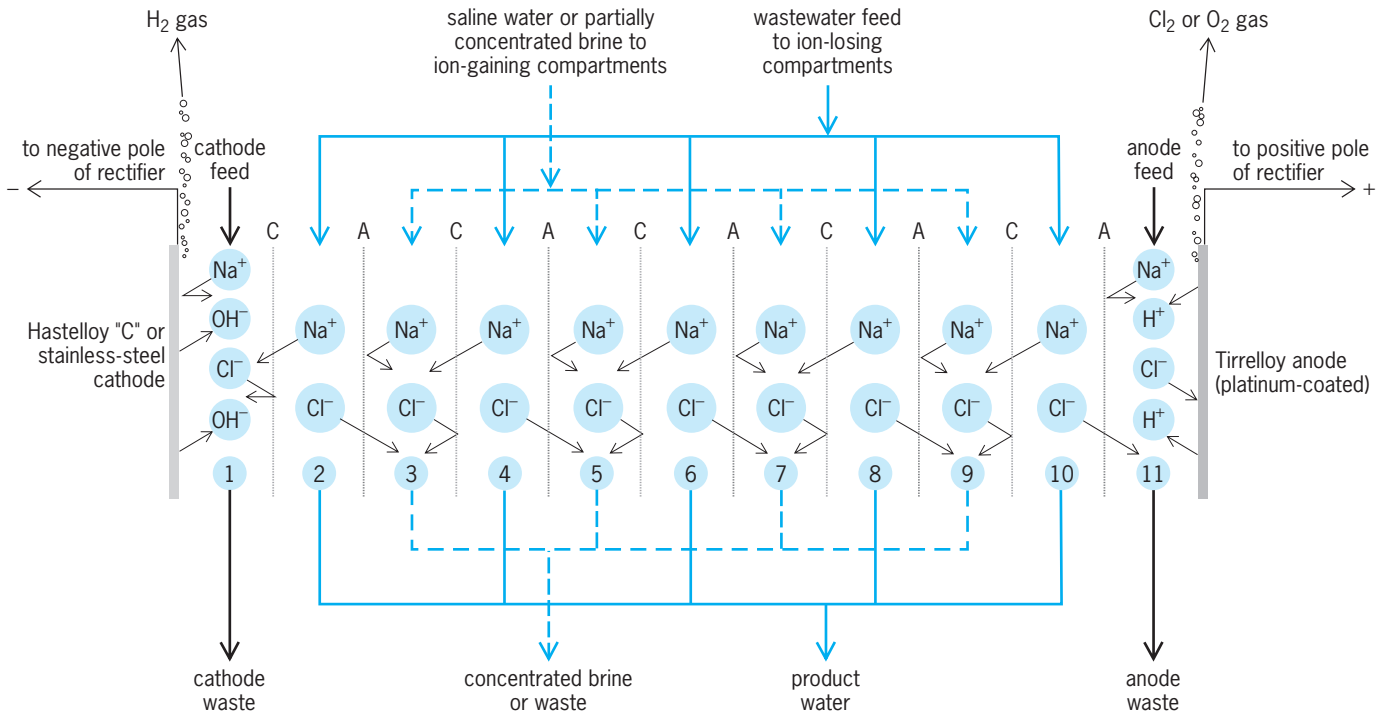


Fig. 8. Diagram of the basic ion and water flow in electric membrane stack. C = cation membrane; A = anion membrane; Na⁺ = any cation, such as sodium; Cl⁻ = any anion, such as chloride. Numbers = compartments.

washed out with water. The metal is in the form of bright crystalline flakes. A beryllium-copper eutectic can be made by using a copper cathode. Beryllium alloys containing copper and nickel are made from BeO in arc furnaces.

Calcium has been made by electrolysis of pure fused calcium chloride about 800°C (1470°F). In this case the cathode is solid calcium, and is mechanically withdrawn from the cell as a “carrot.” Calcium is now made by aluminum reduction at low pressures and high temperatures.

Cerium and misch metal are made from CeCl₃ or mixtures of chlorides of cerium, lanthanum,

and neodymium in fused-salt electrolysis with NaCl. Misch metal is used for lighter flints.

Fluorine for separation of uranium isotopes is produced by electrolysis of 40% HF in KF between carbon anodes and steel cathodes at 88–100°C (190–212°F). A diaphragm of Monel screen keeps the products H₂ and F₂ separated. Dry HF gas bubbled continuously into the electrolyte. At a current density of 1 A/in.² (15 A/dm²) the cell operates at 9–12 V and 96% current efficiency. Energy consumption is 3.0 kWh/lb (24 MJ/kg) fluorine. The theoretical decomposition potential is 2.85 V.

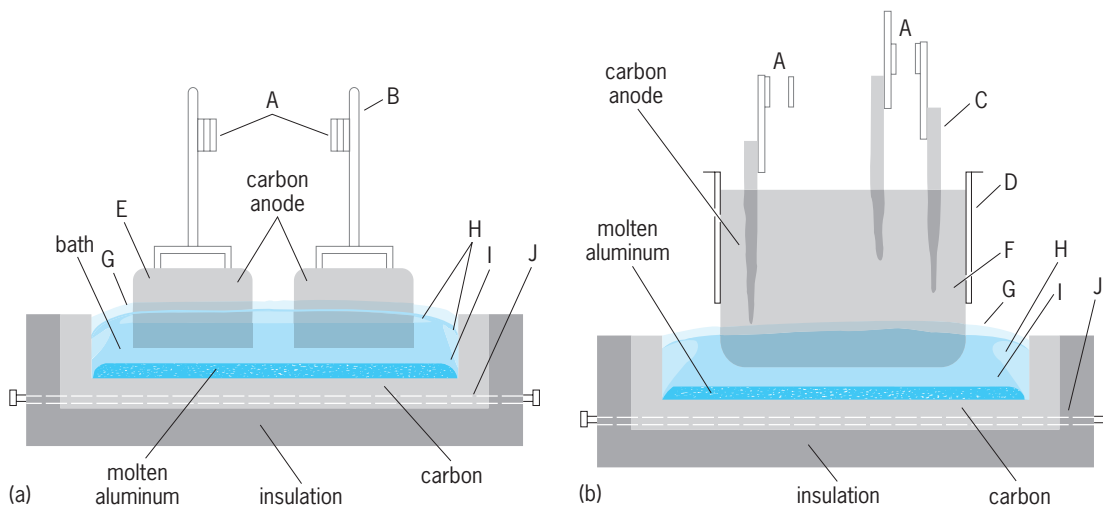


Fig. 9. Two types of aluminum cells, (a) utilizing prebaked carbon anode and (b) utilizing Söderberg carbon anode. A = anode bus, B = anode rod, C = anode stub, D = anode casing, E = prebaked carbon anode, F = Söderberg carbon anode baked in place, G = crust of frozen bath and alumina, H = frozen bath, I = bath (electrolyte), J = steel pin cathode collector.

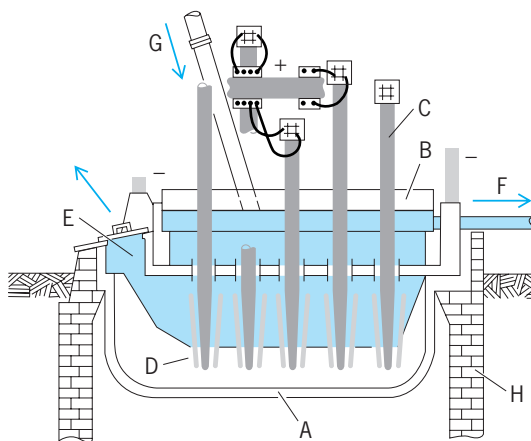


Fig. 10. Cross section of the Dow magnesium cell. A = steel container, B = ceramic cover, C = graphite anodes, D = steel cathodes, E = magnesium collecting well, F = chlorine outlet, G = magnesium chloride feed, H = furnace setting.

Lithium is made by electrolysis of fused 60% LiCl and 40% KCl at 450–500°C (840–930°F) in a cell similar to the Downs sodium cell.

Magnesium is produced by electrolysis of fused 25% MgCl₂ and 75% NaCl at around 700°C (1290°F). The Dow process for making magnesium from seawater uses material approximating MgCl₂·2H₂O, which is fed around the graphite anodes, where dehydration occurs (Fig. 10). Gas from the anode compartment is wet chlorine, air, and hydrogen chloride; the latter is used to make fresh magnesium chloride from magnesium hydroxide. Magnesium metal is deposited on steel cathodes, which direct the metal to a collecting zone. The cell is a cast-steel pot in a furnace setting. Other cells use molten anhydrous magnesium chloride feed. They have brick-lined steel bodies with graphite anodes (Fig. 11). Magnesium chloride, with other chlorides, is separated in vacuum crystallizers from brines, dehydrated, and melted in electric resistance cells from which molten MgCl₂ is tapped periodically to feed the cells. Molten magnesium is ladled from the cells

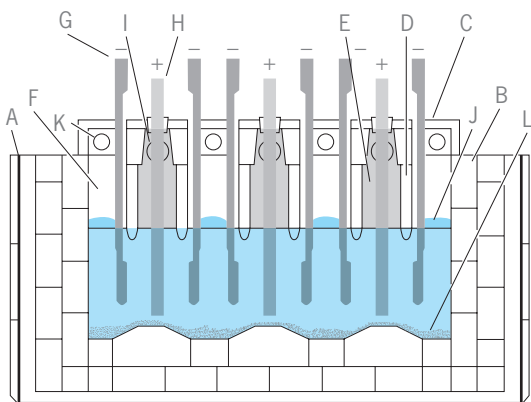
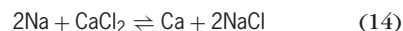
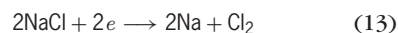


Fig. 11. Cross section of anhydrous magnesium chloride cell. A = steel box, B = ceramic lining, C = ceramic cover, D = ceramic separators, E = anode (chlorine) compartment, F = cathode (metal) compartment, G = iron cathode, H = graphite anode, I = chlorine outlet, J = magnesium, K = cathode compartment vent, L = mud.

and cast into molds. The cell operates at 6–7 V and 80–88% current efficiency and uses 8–8.5 kWh/lb (63–67 MJ/kg) metal.

Sodium was once made by electrolysis of fused NaOH, but since 1929 it has been made by electrolysis of NaCl in the Downs cell. The electrolyte is 40% NaCl and 60% CaCl₂ at 590°C (1090°F). The cell consists of brick-lined steel vessels. Four graphite anodes project upward from the bottom. The cathode is made of steel cylinders concentric with the anodes and supported from iron arms extending through the sides of the cell, which also conduct current. A diaphragm of 26-mesh-per-inch (10 mesh-per-centimeter) iron screen directs the sodium into an inverted trough leading to a riser pipe, which cools the metal and conducts it to a collecting tank beside the cell. Chlorine is collected in a nickel cone inverted over the anode. Pure dry salt is fed to the cell. The reactions in the cell are Eqs. (13) and (14). The



metal at cell temperature is 5% Ca, but as it cools, the second equation is reversed so that the cool metal is about 1% Ca. It is filtered just above the melting point of the sodium, and the final product is under 0.04% Ca. Cells of 38,000 amp operate at 7 V and 83% current efficiency. Energy consumption is about 4 kWh/lb (32 MJ/kg) metal.

Molybdenum, thorium, titanium, uranium, and zirconium can all be made by electrolysis of their complex halides, K₃MoCl₆, ThF₄·KF, K₂TiF₆, KUF₅, and K₂ZrF₆, in molten NaCl. Inert atmospheres are required, and the cell is usually a graphite crucible acting as anode with a graphite or molybdenum cathode, on which the metal deposits as crystals or powder. The batch of metal on the cathode is cooled in an inert atmosphere, then broken off, pulverized, and leached with water. The metal powder is used as such or melted in a vacuum arc furnace.

In the case of tantalum, pure K₂TaF₇ is heated to 900°C (1650°F) in a graphite pot acting as an anode with a removable metal cathode. When the cathode is loaded with deposited metal, it is removed and quickly replaced. The bath is replenished with K₂TaF₇. The cathode deposit is pulverized and washed with acid. The metal powder can be compacted by sintering.

Electrothermics. The manufacture of many products requires temperatures higher than can be obtained by combustion methods. Electric heat can usually be developed at, or close to, the point where it is required, so that it is relatively quick. It permits easy control of the atmosphere for oxidizing, reducing, or neutral conditions.

Products of the electric furnace include iron and steel; ferrous; nonferrous metals and alloys; the exotic metals titanium, zirconium, hafnium, thorium, and uranium; and nonmetallic products such as calcium carbide, calcium cyanamide, sodium cyanide, silicon carbide, boron carbide, graphite,

fused alumina, magnesium oxide, quartz, silica, thorium, zirconia, lime, spinel, kyanite, sodium aluminate, dolomite, boric acid and borides, carbides of zirconium and titanium and related metals, graphite, phosphorus and phosphoric acid, and chlorides of magnesium, boron, zirconium, and titanium. An electric smelting process converts ilmenite into iron and a titanium slag for a pigment manufacture. Zinc metallurgy uses an electric furnace. Steam is generated in electric boilers where economically feasible.

Large amounts of United States steels, stainless steels, and special alloys are made in electric furnaces. **Table 4** gives values for electrode consumption, furnace voltage, secondary amperage (from transformer to electrodes), and current density.

Methods of electric heating utilize resistance, arcs, or induction. Resistance furnaces (**Fig. 12**) may use the substance being heated as a resistor, or auxiliary resistors may be used. Arc furnaces (**Fig. 13**) may have a direct arc between an electrode and the material being heated, for example, steel scrap or the arc may be between two or more electrodes. The arc may be indirect or it may be submerged in the material being melted. Induction furnaces (**Fig. 14**) use the crucible or its charge as the closed secondary circuit of a transformer with low- or high-frequency alternating current. The temperature in the carbon arc is approximately 3100°C (5600°F). Graphite or amorphous carbon are used for electrodes, depending on economics. Metals such as titanium, zirconium, and hafnium may be compacted into electrodes which are consumed by resistance or arc melting or combinations of both in a nonreactive-atmosphere or vacuum furnace. Some types of furnaces are adaptable to several products, but usually the design is developed for an individual product. See ELECTRIC FURNACE; ELECTRIC HEATING; PYROMETALLURGY; STEEL MANUFACTURE.

Zone refining of metals for the electronics industry, such as silicon for diodes and transistors, is accomplished by induction melting of the metal in a narrow zone and slow movement of the molten in the metal ingot from one end to the other in an evacuated or inert gas-filled enclosure. Impurities move

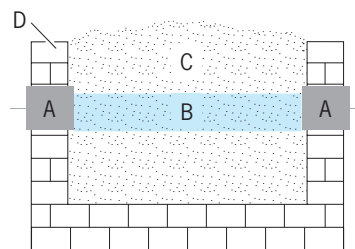


Fig. 12. Resistance furnace. A = current conductors, B = conducting core of granulated carbon, C = granular charge, D = furnace wall of brick and refractory.

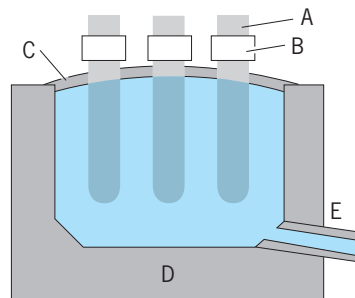


Fig. 13. Three-phase arc furnace. A = carbon or graphite electrodes, B = current clamps, C = refractory roof, D = refractory shell, E = tap hole.

toward the end of the ingot. The operation is repeated until the desired purity is obtained. See ZONE REFINING.

High-melting point metals, carbides, oxides, and nitrides are melted with a plasma-arc torch by coating objects with molten droplets carried in a jet of inert gas passed through the torch. Raw material, such as a powder or wire, is fed through a dc arc between a tungsten cathode and water-cooled copper anode. A strong bond with the base material is obtained. See TORCH.

Processes in gases. Electrical discharge through gases has industrial application in ozone production and nitrogen fixation. Ozonizers consist of two metal electrodes with an air gap and a dielectric, such as a gas, between them. Very dry air passed through an air gap will then contain 10–12 mg ozone per liter.

TABLE 4. Representative values for electrode consumption, voltage, current, and current density for electric furnace products

Product	Electrode consumption*		Voltage range, phase to phase	Current range, secondary amperes	Current density, A/in. ²
	Approx. lb/ton of product	Approx. lb/1000 kWh			
Ferrosilicon, 50%	30–40	7	130–195	35,000–55,000	35–55
Ferrosilicon, 65%	60–90	10	125–170	35,000–50,000	35–60
Ferrosilicon, 75%	90–130	12	120–170	35,000–50,000	35–60
High-carbon ferrochrome	30–50	10	110–200	25,000–40,000	25–40
Silicomanganese	80–100	23	110–140	35,000–70,000	35–70
Ferromanganese	30–50	16	85–100	30,000–60,000	35–65
Calcium carbide	40–60	17	120–200	30,000–80,000	40–70
Phosphorus	35–70	3	250–400	13,000–30,000	20–30
Refined copper	4–6	19	80–190	12,000–30,000	80–110
Nonferrous castings	4–5	15	100–110	800–5000	110–300

*1 lb/ton = 0.5 kg/metric ton. 1 lb/k Wh = 0.1260 kg/MJ.

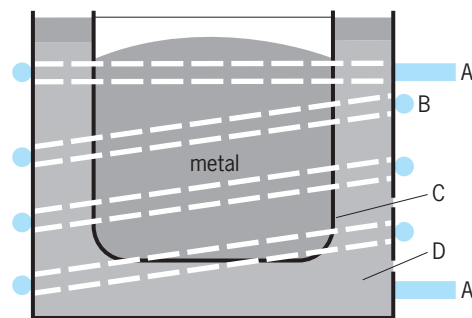


Fig. 14. Induction furnace. A = terminals of high frequency, B = copper coil, C = crucible, D = refractory.

Fixation of nitrogen by passing air through an arc furnace, thus forming oxides of nitrogen, was once practiced when power was cheap in Norway, France, and Italy, but has been replaced by conventional processes. See ELECTRICAL CONDUCTION IN GASES; NITROGEN FIXATION; OZONE.

Electromagnetic separation. Magnetic separation removes tramp iron from mixtures of granular solids, suspensions, and solutions. Magnetic separation is also used to separate solids of various magnetic susceptibilities, such as mineral fractions. See MECHANICAL SEPARATION TECHNIQUES.

Electrodialysis. This is the separation of low-molecular-weight electrolytes from aqueous solutions by migration of the electrolyte through semipermeable membranes in an electric field. It is used on an industrial scale for deashing starch hydrolyzates and whey, and in many municipalities for producing potable water from saline water. Its uses also include the concentration of liquid foods such as dairy products and citrus juices, the recovery of sulfite pulp waste and pickling acid, and the isolation of proteins. See COLLOID; DIALYSIS.

Electrophoretic deposition. This is the deposition of a nonconductive material in a finely divided state from a suspension in an inert medium. Electrophoresis is the migration of colloidal particles, which acquire positive or negative charges in an electric field. The process is useful in electropainting; for instance, electropainting of automobile bodies and other objects has now been adopted on a large scale. Rubber latex is an example of negatively charged colloid which can be plated on an anode. Electronic components can be coated with inorganic salts, oxides, and ceramics suspended in organic media. Bitumen can be electrodeposited out of an aqueous dispersion onto the anodic surface of steel pipe by use of an axially placed cathode. See ELECTROPHORESIS.

Electroendosmosis. This is the movement of a liquid with respect to an immobilized colloid in an electric field. The process is used in the dehydration of peat, dye pastes, and clay. Dies in a clay extrusion press can be "lubricated" by making the die cathodic which attracts a film of water to it from the wet clay. It is also used commercially for dewatering soils in mining, road building construction, and other civil engineering works.

Electrostatic technique. The deposition of charged particles from suspension in gases has many useful applications. The Cottrell electrostatic precipitator removes dusts and mists from gases. In one form, a fine wire is axial in a pipe and insulated from it. The pipe is grounded, and the wire is negative in a high-voltage dc circuit. Particles in the gas that is passing through the pipe become electrically charged and move to the pipe. Liquid particles drain off, while solids are periodically vibrated off. See ELECTROSTATIC PRECIPITATOR.

Spray painting with a high voltage between the gun and the work is particularly effective in providing an even coating with an economical use of paint or irregular and open surfaces, such as a screen. Complete coverage is achieved because an uncovered spot or pore has a higher field strength and attracts the charged droplets to it.

In xerography a sheet of plain paper is electrically sensitized in those areas corresponding to an original so that colored resin particles carrying an opposite charge are attracted and retained only on the sensitized areas, thus producing a visible image corresponding to the original.

Abrasive paper and cloth are coated with an adhesive and abrasive powers attracted to the base material in an electrostatic field. Pile fabrics can be produced in a similar manner, with the short fibers oriented by the electric field. Powdered adhesive can be applied to paper.

Plant design for electro process. The greatest changes have taken place in the electrical distribution system portion of electrolytic plants. Copper bus-bars have taken for granted for years. They show good corrosion resistance, low contact voltages, and contacts can be kept clean readily. Aluminum is not as corrosion-resistant, and contacts are poor unless specifically designed. Aluminum and copper are competitive conductors when the ratio of the cost of a given amount of aluminum to that of copper is about 1.7. Aluminum is below the equality price, copper is above. Therefore, electrolytic plants cannot afford the luxury of copper bus-bars, but still want the advantages of copper contacts—namely, the copper-surface aluminum bus-bar and the copper-shoe aluminum crossbar, achieved by plating, machined sections, explosive forming, and duplex preparation.

Electrolytic engineering encompasses cell design; electrical isolation; use of sumps instead of sewers; economic selection of bus-bars; solution circuit isolation; under-floor distribution of lines, supports, conductors, crystallizing liquors, and plastic pipe and conduits; lined and unlined cells; contacts and their inspection; and avoidance of pollution.

Early in plant development, the problems of impurities in solutions and the effects of traces of the electrolytic deposit were evident, appreciated, and solved, largely by Edisonian methods. Maximum tolerance were served up, and purification procedures were developed to a high level in copper and zinc. Crystallizing solutions were a problem for manganese. In this area, control of deposition potential,

voltage-amperage and amperage-time recording instruments gave the concept of impurity concentrations affecting the deposition potential and the current density needed to reach this potential. There was an explanation for receding or resolution of deposits. In chromium the problem was solved, in effect, by converting the ore during purification into a chemical compound purified by crystallization. Purification of electro-winning solutions of magnesium, sodium, potassium, barium, calcium, iron, aluminum, antimony, arsenic, tin, cadmium, silica, lead nickel, and cobalt by chemical methods is a common occurrence.

The reconditioning of cathodes for metal deposition and ease of product removal in the electro-winning of manganese and chromium is a repetitive, demanding, and skilled-labor operation. Lines have been designed and put in operation resembling the continuous, conditioning cleaning, pickling, and preparation lines of the electroplating phases of the automotive industry, a long-delayed cross-fertilization. When electrolytic copper sheet needed separate plants for printed electronic circuits, auxiliary electrodes were developed to control edges and ensure more uniform thickness. Perhaps manufacture of starting sheets can be mechanized by adopting continuous electroplating lines.

The concept of electrolytic cell as an insulated mechanism, supported on structures and insulated from them, with separate support for conductors, supply lines, circulating systems, all insulated and electrically broken so that they are not shunt circuits, has been well established for decades. Tanks were wood, but are now unreinforced concrete, lined or unlined, and in few cases glass fiber or cloth-reinforced plastics, such as poly(vinyl chloride) or its competitors. In general the copper industries studies have shown that the investment cost for a lead-lined concrete cell is a little more than half that of a ribbed, equivalent-strength plastic. Because of local availability, tin refining cells are poly(vinyl chloride)-lined concrete units, but chromium metal is made in plastic cells. A few foreign plants employ "painted" or thin-lined concrete.

Selenium rectifiers are used, to a greater extent than copper oxide, for small power supplies, control and signaling circuits, and for applications in which their relative immunity to voltage surges is desired.

There have been no mercury-arc rectifiers built in the United States since about 1960, and very few in Europe for electrolytic or industrial applications, except for high-voltage direct current, welding, and transmission. Monocrystalline diode rectifiers, first germanium and then silicon, have completely replaced mercury-arc rectifiers in industrial applications. Manufacturers have converted to silicon in many multianode mercury-arc rectifiers and excitron and ignition units. All new installations and expansions utilize silicon rectifiers.

The only type of power rectifiers for electrolytic loads, such as aluminum reduction, metal refining, electro-winning, chlorine and chlorate cells, and so on, are silicon rectifiers. Silicon rectifiers range in

sizes up to 90,000 A and up to nearly 1000 V, with a maximum rating of 30,000 kVA.

Organic Processes

Organic electrochemistry was once regarded as a tantalizing area with many important laboratory achievements but few successes in commercial practice. This situation has changed, however, in that electroorganic processes are commercially advantageous if they can fulfill either of two conditions: (1) performance under conditions of voltage corresponding thermodynamically to the conversion of an organic group to a reduced or oxidized group, with the cell products relatively easy to isolate and purify; (2) performance of a highly selective, specific technique to make an addition at a double bond, or to split a particular bond (for example, between carbon atoms 17 and 18 of a complex molecule having 25 carbon atoms).

Selectivity and specificity are highly important in electroorganic processes for the manufacture of complicated molecules of vitamins and hormones—as well as for the medicinal products whose action on pathogenic organisms is a function of their spatial arrangement, steric forms, and resonance.

The electrolytic approach can also be competitive for some low-cost, tonnage products. Here continuous processing is important, and only a single phase should be present, that is, a solution rather than an emulsion, dispersion, or mechanical mixture. Only for fairly valuable products is it practical to find a conducting solvent and then to engineer around it.

The electrolytic oxidation and reduction of organic compounds differ from the corresponding and more familiar inorganic reactions only in that organic reactions tend to be complex and have low yields. The electrochemical principles are precisely those of inorganic reactions, while the procedures for handling the chemicals are precisely those of organic chemistry.

Most organic molecules are insoluble in the aqueous solutions that are the best electrical conductors. Unless solubility can be increased, the only other approach is to use organic solvents. These make relatively poor conductors; hence are encountered power loss, heat build-up, chemical inversion, and often, stepwise, complex reactions.

Oxidations. Commercial success in inorganic electrochemistry has come about by well-engineered combinations of organic and inorganic techniques in areas where strictly chemical methods are either impossible or inefficient, for example, in catalytic hydrogenation or oxidation.

The conventional oxidation reagents of the organic chemist are expensive. There is no market for the oxidant once it is reduced, and chemical regeneration is prohibitive in cost. Accordingly, these reagents are avoided. Electrolytic regeneration, however, can be relatively inexpensive if linked to a carefully controlled organic operation. Typical of this approach is manganese dioxide oxidation of anthraquinone with electrolytic regeneration of the oxidant. In the electrolytic oxidation of anthracene

there is a cost-efficient process that utilizes a 20% sulfuric acid suspension with a small quantity of ceric sulfate as a catalyst. Other examples include chromic acid oxidation of oleic acid to perlargonic and azelaic acids, where the oxidant is regenerated electrolytically, and the electrolytic regeneration of periodic acid (a costly oxidizing reagent) in the dialdehyde starch process. These involve savings not only in the purchase price of the oxidants but also in the disposal cost of products that cannot be marketed.

In electrolytic reactions the cell surface is the only source of reductants or oxidants, and these must be produced at highest efficiencies. The mass action effects, concentrations, temperature of reaction, reaction velocities, diffusion, and equilibria between the initial and final products of the reaction all apply to electrochemical reduction and oxidation in the same manner as do the corresponding reactions carried on outside the electrolytic cell.

Materials that are reduced absorb hydrogen at the cathode, and may be considered cathodic depolarizers; those that are oxidized absorb oxygen at the anode, and are anodic depolarizers. Oxidation reactions may involve substances other than oxygen, such as chlorine.

Anodes are selected with a high oxygen or halogen overvoltage, and cathodes are selected with a high hydrogen overvoltage. (Because the accumulation of electrolysis products at anode and cathode causes polarization, the overvoltages are needed to move the products away and keep the process going.)

Reductions. Substances that are easy to reduce may be acted on, at the interface of cathodes, with low hydrogen overvoltage. Hard-to-reduce materials may require much higher overvoltages, which are reached through either the cathode composition or the current density.

The aromatic nitrogen-containing organic compounds have been extensively studied by many researchers. As early as 1900, F. Haber and K. Elbs showed that the nitro compounds of the type RNO_2 , where R is an organic radical, could either be directly reduced to the amine RNH_2 (that is, aniline), or successively reduced to the nitroso product RNO , the beta aryl hydroxylamine RNHOH , the azoxy product $\text{RN}=\text{O}=\text{NR}$, and the material containing an azo group $\text{RN}=\text{NR}$, which in turn is reduced to the hydrazoform $\text{RNH}=\text{NHR}$ or the idene type as $\text{H}_2\text{NR}=\text{RNH}_2$.

Although commercial processes based on these reactions were successful, they were eventually replaced by improved nonelectrolytic processes.

Oxidation of starch. The original process for the oxidation of starch to dialdehyde starch was by periodic acid, continuously regenerated at the anode of an electrochemical cell. Dialdehyde starch is an important polymeric aldehyde used in paper manufacture and leather tanning. In this process a lead anode, coated with lead dioxide, is immersed in a suspension of starch; the cathodes are steel, enclosed in porous aluminum oxide diaphragms. One difficulty with this process is that free iodine, formed

by the migration of iodine ions to the anode, tends to react with starch, despite the presence of the diaphragm.

In 1960 an improved two-stage process was devised, where just enough periodic acid is added to oxidize the starch in a straightforward chemical reaction. The iodic acid produced during the reaction is then separated from the dialdehyde starch and sent to an electrolytic cell for reconversion to periodic acid. In the electrolytic cell, this reconversion is achieved with current efficiencies of 80–90%. Materials of construction are structural plastics that are unaffected by the reagents or reactants; only the electrodes are metal. Continuous operation is achieved, with the instrumentation permitting good control of temperature, concentration, conversions per pass, gas venting or collection, external recirculation rates, and other rate factors. The wide range over which such control can be exercised allows the cell to function like a reactor or catalytic converter. In fact, the cell has many of the characteristics of the catalytic reactors, the main difference being that in the cell the variables are related to the type and surface of the electrodes, while in the catalytic reactor they are related to the surface of the catalyst and its support.

All piping is manifolded exterior to the cell (unlike the hydrogen-oxygen cell), and the feeds are valve controlled externally. The electrode spacing is prearranged, and all reaction chambers are in parallel both relative to flow and electrically. Because valve closure can shut off a fraction of the parallel streams through the cell, shut downs because of mechanical or other failures are minimized. The electrolytes are simple solutions, not emulsions or suspensions, and do not contain stabilizers or salting-in components. Gas evolution, separation, and collection are separately handled in anolyte and catholyte streams. These two streams may be blended in any proportion desired for maintenance of pH and control of concentration.

Adiponitrile production. In 1965, an electroorganic process for adiponitrile production was developed and commercialized. Power was used as a reagent in the electrolytic reductive coupling (or hydrodimerization) of acrylonitrile to adiponitrile, a component of nylon-6,6. The adiponitrile molecule is a six-carbon dimer, essentially made up of two molecules of three-carbon acrylonitrile.

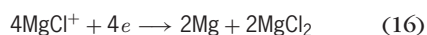
In the commercial plant, each cell consists of a lead cathode and an alloy anode separated by a sulfonated polystyrene resin membrane. Spacers support the membrane and prevent it from touching the electrodes. The cells are made of polypropylene with neoprene gaskets and are grouped in banks of 24. The reactants, the catholyte and anolyte, circulate separately in an upflow direction through a bank of cells in parallel; the direct current passes through the bank in series. The catholyte contains a quaternary ammonium salt, that is, a McKee salt. This provides ions for electrical conductivity and increases the acrylonitrile's solubility above 10%; lesser concentrations would lead to the formation of propionitrile.

The catholyte is circulated to remove the film of product at the cathode surface.

Lead alkyl production. In one electroorganic process for the electrolytic production of lead alkyls, electrochemical synthesis of tetramethyllead or tetraethyllead is achieved by employing a Grignard reagent and lead shot. The electrolyte consists of alkyl magnesium halides, which ionize in ether solutions. The solution is fed to 8000-gallon (30,280-liter) electrolytic cells, with lead pellet anodes, the walls of the cell as cathode, and a nonconducting low-permeability diaphragm. The reaction is shown in reaction (15), where R represents the methyl or ethyl groups.



The $MgCl^+$ ions migrate to the cathode cell walls. Magnesium chloride and metallic magnesium are formed according to reaction (16).



In this process a variable-voltage approach is used. As the rate of reaction drops and the resistance of the cell increases, more and more voltage is applied to the cell, thus maintaining a constant rate of reaction per unit of time.

Mixed alkyl lead compounds may be made by electrolyzing ethylmagnesium chloride and adding methyl chloride in the cells; or, conversely, by starting with methylmagnesium chloride and then adding ethyl chloride.

The electrolysis cell bank is made up of 10 cells, each with an associated recirculation surge drum. In the cell the copper bus handles the low-voltage, high-amperage current loads. The lead pellet storage hopper and tubular conveyor equipment are located directly above the cell dome. The upper third of the recirculating drum is an insulated refrigerant-storage vessel that supplies coolant to carry off heat developed during the course of electrolysis. Refrigerant vapors are returned to centrifugal compressors.

Charles L. Mantell

Bibliography. A. Bard and L. R. Faulkner, *Electrochemical Methods, Fundamentals and Applications*, 2d ed., 2000; D. R. Crow, *Principles and Applications of Electrochemistry*, 4th ed., 1994; A. Kuhn (ed.), *Industrial Electrochemical Processes*, 1971; H. Lund and M. M. Baizier (eds.), *Organic Electrochemistry: An Introduction and a Guide*, 3d ed., 1991; D. Pletcher and F. C. Walsh, *Industrial Electrochemistry*, 2d ed., 1993; K. Scott, *Electrochemical Reaction Engineering*, 1991.

Electrochemical series

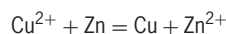
A series in which the metals are listed in the order of their chemical reactivity, the most active at the top and the less reactive or more "noble" metals at the bottom. In a broader sense such an activity series need not be limited to the metals but may be

Electrochemical series of the elements*			
Element	Symbol	Element	Symbol
Lithium	Li	Zinc	Zn
Potassium	K	Chromium	Cr
Rubidium	Rb	Gallium	Ga
Cesium	Cs	Iron	Fe
Radium	Ra	Cadmium	Cd
Barium	Ba	Indium	In
Strontium	Sr	Thallium	Tl
Calcium	Ca	Cobalt	Co
Sodium	Na	Nickel	Ni
Lanthanum	La	Molybdenum	Mo
Cerium	Ce	Tin	Sn
Magnesium	Mg	Lead	Pb
Scandium	Sc	Germanium	Ge
Plutonium	Pu	Tungsten	W
Thorium	Th	Hydrogen	H
Beryllium	Be		
Uranium	U	Copper	Cu
Hafnium	Hf	Mercury	Hg
Aluminum	Al	Silver	Ag
Titanium	Ti	Gold	Au
Zirconium	Zr	Rhodium	Rh
Manganese	Mn	Platinum	Pt
Vanadium	V	Palladium	Pd
Niobium	Nb	Bromine	Br
Boron	B	Chlorine	Cl
Silicon	Si	Oxygen	O
Tantalum	Ta	Fluorine	F

*According to standard oxidation potentials E° at 25°C or 77°F.

carried on through the electronegative (nonmetallic) elements as well. See the **table** for a list of common elements.

The electrochemical series as it applies to metals was first established by laboratory experiments in which the purpose was to determine which metals would displace others from solutions of their salts. Thus a clean strip of zinc immersed in a solution of copper sulfate is soon found to be covered by a deposit of copper, while zinc in turn goes into solution from the strip as zinc ions. By definition, then, zinc is a more reactive metal than copper, since it will displace copper from a solution of Cu^{2+} ions. The reaction is readily seen to be an oxidation-reduction transfer of electrons, which can be summarized by the equation below. Similarly, copper will displace



silver from a solution containing Ag^+ ions, depositing crystals of metallic silver and coloring the solution with Cu^{2+} ions. From these observations an activity series may be set up in the order Zn, Cu, Ag. By exhaustive experiments with other metals it becomes possible to draw up a complete list in the order of chemical activity, in which the metals at the top of the list are those which are found to give up their electrons most readily (that is, are the most electropositive elements). Such a list is shown in the table, where lithium exhibits the most reactivity as a metal.

The ease with which an isolated atom of an element gives up an electron, known as the first ionization potential, is a precise physical quantity which can be measured by electrical experiments

on gases or vapors at low pressure. The replacement experiments which determine the order of the electrochemical series take place in a very different environment, since they involve solid phases and also aqueous solutions with their consequent hydration effects. Moreover, it might well be expected that displacement reactions in solution would depend upon the concentrations of the reagents used, and also upon the presence or absence of other dissolved substances. *See* IONIZATION POTENTIAL.

To obtain a more accurate and reproducible activity series, it is best to turn to the more exact quantity called electrode potential, or oxidation-reduction potential, which is defined as the voltage developed by a sample of pure metal immersed in a solution of one of its salts (at unit activity and at 25°C or 77°F) versus a hydrogen electrode immersed in hydrochloric or sulfuric acid of equivalent concentration. For details about this measurement *see* ELECTRODE POTENTIAL.

It is evident that by confining the experimental conditions to a standard concentration and temperature the hydration and concentration effects can be kept quite constant, making possible a more exact listing of metals according to their activity. Hence any present-day electrochemical series must rely on the measurements of oxidation potential and should be in agreement with the accepted values determined from such electrochemical cells. Such reliance has the further advantage that the series need not then be confined to metals but may be extended to the nonmetals, or electronegative elements. As before, those metals which will liberate hydrogen from dilute acids (such as hydrochloric or sulfuric) will stand above hydrogen in the series, while those metals and nonmetallic elements which will not liberate hydrogen from such dilute acids will stand below hydrogen in the list. Since the oxidation potentials are also related to the equilibrium constants for reversible reactions, it becomes possible to calculate oxidation potentials from other information when direct experiments are inconvenient, as in the case of the alkali metals versus aqueous solutions of their salts. *See* ELECTROCHEMISTRY; ELECTRONEGATIVITY; OXIDATION-REDUCTION. Eugene G. Rochow

Bibliography. F. A. Cotton, G. Wilkinson, and P. L. Gaus, *Basic Inorganic Chemistry*, 3d ed., 1994; C. A. Hampel (ed.), *Encyclopedia of Electrochemistry*, 1964, reprint 1972.

Electrochemical techniques

Experimental methods developed to study the physical and chemical phenomena associated with electron transfer at the interface of an electrode and solution. The objective is to obtain either analytical or fundamental information regarding electroactive species in solution. Fundamental electrode characteristics may be investigated also.

The physical and chemical phenomena important in electrode processes generally occur very close to

the electrode surface (usually within a few micrometers). Mass transfer of species involved in an electrode process to and from the bulk of solution is one important aspect. Inclusion of a large excess of inert electrolyte in most electrochemical systems eliminates electrical migration (current flow as a result of movement of ions of the substrate toward or away from the electrode) as an important means of mass transfer for electroactive species, and only convection and diffusion are considered.

Important chemical aspects of electrode processes include the oxidation or reduction occurring as a result of electron transfer, and coupled chemical reactions. Coupled reactions are initiated by production or depletion of the primary products or reactants at the electrode surface. Identification of the nature and mechanism of such coupled reactions is of particular importance in studies of electrode reactions of organic compounds, where multireaction and multielectron cascades are often found to be initiated by electron transfer to or between the electrode and the electroactive substance.

The primary experimental variables involved in electrochemical techniques are the potential E , the current I , and the time t . The potential or current at the working electrode is controlled and the other is observed as a function of time. The many ways in which either may be controlled give rise to the wide variety of controlled-potential or controlled-current techniques. In all such techniques it is necessary to specify whether mass transport of the electroactive species to the electrode is by convection or diffusion, since mathematical treatments of these two processes are quite different. Mass transport by convection is more efficient than diffusion by several orders of magnitude but is much more difficult to model mathematically.

The general scheme for electron transfer at an electrode in solution is shown in reaction (1), where O



and R are the oxidized and reduced forms of the electroactive species and n is the number of electrons transferred. When k_f and k_b , the rate constants for the forward and back reaction, are very fast and O and R are not involved in preceding or following chemical reactions, the system is called reversible and the Nernst equation (2) holds. In this relation,

$$E = E^{\circ'} + \frac{0.059}{n} \log \frac{C_O}{C_R} \quad (2)$$

E is the electrode potential, $E^{\circ'}$ is the formal standard potential for the redox couple, and C_O and C_R are concentrations at the electrode surface. In the following discussions, only reversible reduction processes will be considered, although oxidations are equally applicable. *See* ELECTROCHEMICAL PROCESS; ELECTRODE; ELECTRODE POTENTIAL; ELECTROLYSIS; OXIDATION-REDUCTION.

Controlled potential. A variety of methods of this type have been developed, depending on whether the electrode potential is held constant or varied

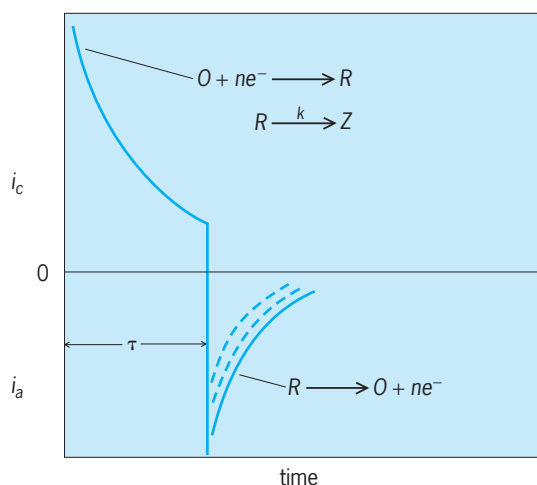


Fig. 1. Current-time behavior for double-potential-step experiment at a stationary working electrode. Broken lines in the reverse current-time curve represent cases where the rate of chemical reaction is fast compared to the time scale of the experiment. $i_a/i_c = f(k)$.

during the experiment and on whether mass transport is by convection (stirring) or diffusion.

Constant potential with convection. In this technique, known as controlled-potential electrolysis, the electrode potential is held constant as the solution is stirred or the electrode is rotated at a constant rate. The current is controlled by the concentration of electroactive substance and the stirring rate. Complete conversion of the electroactive substance takes place at a first-order rate whose constant is determined by cell geometry and stirring efficiency. With good cell geometry and efficient stirring, electrolysis is complete in a manner of minutes.

Constant potential with diffusion (chronoamperometry). When a reducing potential is imposed instantaneously on a stationary working electrode in quiescent solution, current will rise sharply and then decay as the electroactive species near the electrode is depleted by electrolysis. The magnitude of the current is proportional to the bulk concentration of electroactive species, and if the imposed potential is sufficiently negative of $E^{\circ'}$, the Nernst equation demands complete conversion to the reduced form. Under these conditions, the current for this chronoamperometric, or potential-step, experiment is diffusion-controlled and decays with $1/t^{1/2}$. This is shown by Eq. (3), where F is the faraday, A is the electrode

$$I = nFAC^* \left(\frac{D}{\pi t} \right)^{-1/2} \quad (3)$$

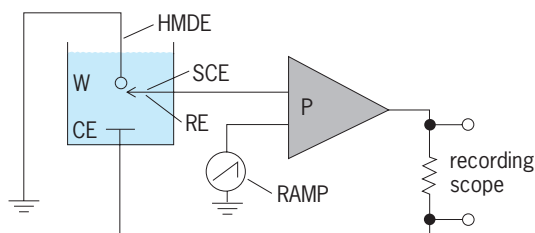
area, D is the diffusion coefficient for the electroactive species, and C^* is the bulk concentration of electroactive species.

An important variant on this experiment is double-potential-step chronoamperometry, in which a second potential step is applied. During the initial step, electrolysis occurs, depleting the oxidized form O , but producing the reduced form R in the immediate vicinity of the electrode. If the potential is instantaneously

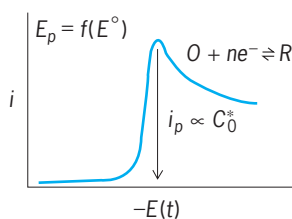
switched back to the initial value after a time τ , species R will be reoxidized, and an inverted $I-t$ curve will be obtained (Fig. 1). This figure also illustrates the common convention that reduction and oxidation currents are plotted in the positive (upward) and negative (downward) directions, respectively. The value of this technique for studying coupled chemical reactions is illustrated by the broken lines in Fig. 1, which demonstrate that while coupled processes have no effect upon the first step, the current after the second step is diminished by a following chemical reaction (decomposition of R to Z). The ratio of forward and back currents, measured at $t = \tau$ and $t = 2\tau$, respectively, provides a convenient experimental measure of this additional chemical complication. The degree to which the current is diminished by coupled reactions can be used to obtain both rate constants and the molecularity of such reactions. A wide dynamic range of rate constants can be measured by varying τ . See ELECTROCHEMISTRY; ELECTROLYSIS.

Variable potential. Several procedures of this type have been developed.

Linear sweep voltammetry (LSV). In this technique, the potential of a stationary electrode in quiescent solution is varied by applying a linear voltage ramp to the electrode. The resulting signal is recorded as a plot of current versus potential. A peak-shaped voltammogram is observed (Fig. 2). The decay of current after the peak arises, as in the experiment at fixed potential, because of control of the current by diffusion; indeed, the decay after the peak has a $1/t^{1/2}$ dependence. The height of the peak is dependent upon the concentration of the electroactive species, and the peak potential E_p is related to $E^{\circ'}$. Analytically, the experiment is applicable in the 10^{-3} M



(a)



(b)

Fig. 2. Linear sweep voltammetry. (a) Generalized instrumentation. (b) Current-voltage curve obtained on the scope. HMDE = hanging mercury drop electrode; SCE = saturated calomel reference electrode; CE = auxiliary counter electrode; P = potentiostat; W = working electrode; RE = reference electrode; and RAMP = linearly increasing voltage signal.

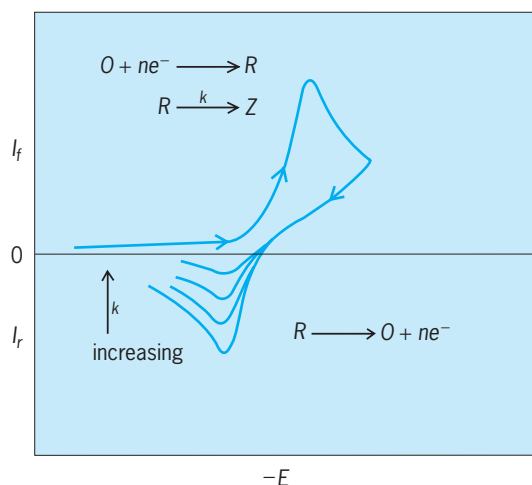


Fig. 3. Current-voltage curves in cyclic voltammetry.

range, but this can be extended to the 10^{-5} M range by plotting dI/dE , the derivative of the current with respect to potential, as a function of potential (so-called derivative LSV); to a good approximation, the background current rises constantly with potential in the region of E_p and is therefore eliminated by differentiation.

Cyclic voltammetry. The most widely used electrochemical analytical and mechanistic technique is cyclic voltammetry which bears the same relation to linear sweep voltammetry as double-potential-step chronoamperometry does to the single-step experiment. That is, at the end of the first linear ramp, the potential is swept linearly back to the initial potential. Typical cyclic voltammograms are shown in Fig. 3 for the case of reduction of O to R and subsequent decomposition of R to Z . The figure shows a series of cyclic voltammograms corresponding to different values of the rate constant k . As in the double-step experiment, the appearance of the forward sweep is unaffected by the coupled chemical process. However, as k increases, the height of the peak on the reverse sweep due to reoxidation of R decreases. Sweep rates may be varied over a range of roughly 10^{-2} to 10^2 V/s in conventional cyclic voltammetry. This may be extended considerably, as high as 10^5 V/s or greater by using microelectrodes (electrodes with diameters in the micrometer range) and appropriate cell shielding to avoid stray currents.

Voltammetry in stirred solution. The potential in a cell configured for controlled potential electrolysis is varied linearly with time (about 5 mV/s), and the current is measured as a function of potential. The current increases sharply when the potential passes through the region of $E^{o'}$ for an electroactive species in the solution. Between reduction steps, current plateaus are established, the heights of which are proportional to the concentrations of each electroactive species. A typical voltammogram is shown in Fig. 4a, where I_L is the limiting (plateau) current which is proportional to the bulk concentration of electroactive species C^* . The point designated as $E_{1/2}$ is called

the half-wave potential and is approximately equal to $E^{o'}$. A small-amplitude (1–15 mV, 10–100 Hz) alternating potential (Fig. 4b) has been superimposed on the linear voltage ramp applied to the electrode of Fig. 4a. The resulting alternating-current response is peak-shaped, and the heights of the peaks are proportional to C^* . This experiment, alternating-current voltammetry, can be used to determine rates of heterogeneous electron transfer to or from the electrode to the electroactive species; in this case, it is generally done in quiescent solution and at a solid or dropping mercury electrode.

Stripping analysis. These techniques can sometimes be combined in useful ways. Analytical sensitivity is enhanced by the method known as stripping analysis. A constant reducing potential is applied to an electrode in a stirred solution containing one or more metal ions for a fixed time (about 1–60 min) sufficient to deposit substantial amounts of the metal on the electrode. After this controlled-potential electrolysis, a linear potential sweep to positive (oxidizing) potentials is applied to the electrode. Because of the preconcentration of the metal onto the electrode, the current is not limited by mass transport, and thus the stripping current is considerably larger than would have been observed in a voltammogram in stirred or quiescent solution. Sensitivity as high as 10^{-9} M is achieved. The technique is useful for measuring trace-metal analyses in seawater.

Square-wave polarography. This technique resembles conventional polarography at a dropping mercury electrode, except that a small-amplitude (about 1–25 mV) square wave (frequency about 225 Hz) is imposed on the controlling dc potential. The current measured is an alternating component obtained by measuring the current at times near the end of each half-cycle. This avoids the large capacitive background current associated with instantaneous electrode potential changes and allows greater

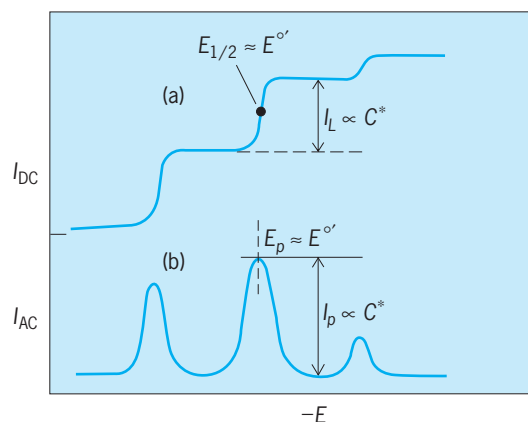


Fig. 4. Plots of the variation of current with constant applied potential in reduction processes. (a) Direct current versus voltage at a rotated working electrode (or stirred solution at a stationary electrode), where three successive reductions of electroactive species occur. (b) Alternating-current versus direct-current voltage applied at a working electrode in an alternating-current voltammetric experiment under the conditions of Fig. 4a.

sensitivity. Analyses over the range of 10^{-3} to 10^{-7} M are possible.

Pulse polarography. This is a variation of square-wave polarography, where a single small potential step is imposed on the controlling dc potential during the life of each mercury drop. The change in current with each step is measured and plotted as a function of the dc potential. It is convenient to make current measurements near the end of the mercury drop life, where the drop is growing slowest and the capacitive current is minimal. For this reason, the drop time is usually controlled mechanically and electronically synchronized with the current measurement. The sensitivity is similar to that of square-wave polarography. Like alternating current and square-wave polarography, a peak-shaped curve resembling Fig. 4b is obtained. All of these techniques may use stationary electrodes of carbon or a variety of metals in place of the dropping mercury electrode; in this case, they are referred to as voltammetry rather than polarography. See POLAROGRAPHIC ANALYSIS.

Controlled current. Controlled-current techniques, although easier to implement than those in which the potential is controlled, are not as selective and therefore have fallen out of favor in recent years. In chronopotentiometry, a constant current is imposed at the working electrode, and its potential is monitored with time. The electrode must assume a potential which will cause electrolysis sufficient to maintain the imposed current. Thus, the electrode initially adopts the potential of the most easily reduced species. As this species is depleted near the electrode, the potential shifts to that of the next most easily reduced species. This continues until reduction of solvent or electrolyte occurs. Characteristic curves are plotted in Fig. 5. The transition time τ is related to the bulk concentration of the electroactive species. However, it is difficult to establish an accurate background correction, thus rendering void the use of chronopotentiometry for analytical purposes.

Coulostatic analysis involves application of a very large, short pulse of current to the electrode, after which the cell circuit is opened and the return of the working electrode potential to its initial value

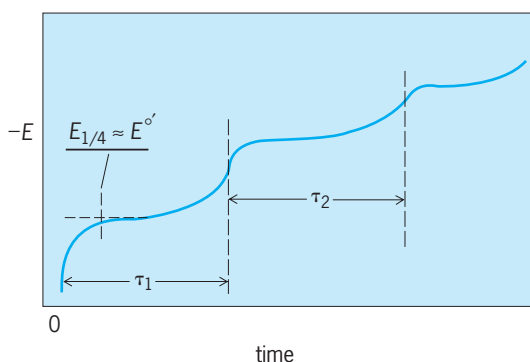


Fig. 5. Potential-time curves for chronopotentiometric experiments. Curves represent successive reduction of three different electroactive species in solution.

is monitored. The open-circuit condition requires that the current necessary to discharge the electrode interface comes from electrolysis of electroactive species in solution. The current pulse charges up the electrode-solution interface to a new potential, and the electrode discharges, and returns to its original potential, by reducing the electroactive species in solution. The time dependence of the change in electrode potential is proportional to C^* , but use of this dependence for analytical purposes requires the independent and inconvenient measurement of the electrode-solution capacitance. The change in electrode potential ΔE versus $t^{1/2}$ results in a straight-line plot, the slope of which is proportional to concentration, as in Eq. (4), where C is the electrical

$$\pm \Delta E = \frac{2nFD^{1/2}C^*}{\pi^{1/2}C} t^{1/2} \quad (4)$$

capacitance associated with the electrode-solution interface.

Digital simulation. Computer simulation of electrode reactions involving one or more electron transfer steps together with several coupled chemical reactions is a challenging task because concentration gradients of the various chemical species involved vary as a function of both time and distance from the electrode during the experiment. Major advances in the necessary mathematical and computational techniques, as well as increases in computing power even at the desktop level, have led to corresponding major advances in the ease with which electrode reactions can be digitally simulated. Many modern electrochemical investigations, especially in the areas of organic and organometallic chemistry, rely on digital simulation to confirm mechanistic postulates and to extract kinetic and thermodynamic parameters.

Thin-layer electrochemistry. Electrochemical techniques can be applied in cells where a solution only 10–100 μm thick is electrolyzed. The principal advantage is that, for experiments lasting more than about 1 s, mass transfer limitations can be ignored. Thus, correlations between observed parameters and system characteristics are more straightforward. For example, for a chronopotentiometric experiment the transition time represents complete depletion of the electroactive species. Thus, C^* is proportional to τ for each of a sequence of reduction steps. Also, if a linear potential sweep is applied, a symmetrical current peak is obtained, the height or area of which is related to the concentration of electroactive species.

Albert J. Fry

Bibliography. A. J. Bard and L. R. Faulkner, *Electrochemical Methods: Fundamentals and Applications*, 2d ed., 2000; A. J. Fry, *Synthetic Organic Electrochemistry*, 2d ed., 1989; P. T. Kissinger and W. R. Heineman (eds.), *Laboratory Techniques in Electroanalytical Chemistry*, 2d ed., 1996; K. B. Oldham and J. C. Myland, *Fundamentals of Electrochemical Science*, 1994; D. T. Sawyer, A. Sobkowiak, and J. L. Roberts, Jr., *Electrochemistry for Chemists*, 2d ed., 1995.

Electrochemistry

The science dealing with the chemical changes accompanying the passage of an electric current, or the reverse process in which a chemical reaction is used as the source of energy to produce an electric current, as in a battery or fuel cell. Electric conduction occurs through the motion of charged particles. The charged particles may be electrons (as in metals or semiconductors) or ions, which are electrically charged atoms, molecules, or molecular aggregates. Ionic conduction in electrolytes (liquid solutions, molten salts, and certain ionically conductive solids) is a phase of electrochemistry, as is corrosion. Conduction in metals, semiconductors, and gases is generally considered a portion of physics. *See* ELECTROLYTIC CONDUCTANCE.

Galvanic cells. These are better known as electric batteries. Many chemical reactions can be arranged to produce electrical energy by physically separating the reaction into half-reactions: one supplying electrons to an electrode forming the negative terminal of the cell, the other removing the electrons from the positive terminal. In the lead storage battery, electrons are supplied to the negative terminal by a half-reaction in which metallic lead is oxidized to form lead sulfate. At the positive terminal, lead dioxide is reduced to lead sulfate.

The electrons flowing in the external circuit from the negative to the positive terminal constitute the desired electric current. Charging of the lead storage battery by forcing a current to flow in the reverse direction results in the reversal of both half-reactions, and the storage of electric energy in the form of lead and lead dioxide. Such a cell is called a secondary cell, in contrast to a primary cell, such as the Leclanché cell or dry cell, which is not designed to be recharged. In the Leclanché cell, electric energy is produced by the oxidation of zinc and the reduction of manganese dioxide at a carbon electrode. The electrolyte is a moist mixture of zinc chloride, ammonium chloride, and powdered carbon. The fuel cell is designed for the continuous production of electric current through the consumption of oxidant and reductant at separate electrodes. The most common fuel cell is the hydrogen-oxygen (or air) cell with alkaline electrolyte. Many other systems, including hydrocarbon-chlorine and sodium-sulfur, have been proposed. *See* BATTERY; FUEL CELL.

Electrodeposition. The most important type of chemical reaction brought about by the passage of electric current is the deposition of a metal at a cathode from a solution of its ions. Electroplating of many metals, such as silver, cadmium, nickel, and chromium, is used for protective and decorative coatings. Electroforming is a variety of electrodeposition in which the electrically nonconductive prototype of an article to be produced is rendered conductive by spraying with a thin metallic coating, next electroplated with a metallic deposit that is then stripped from its substrate and filled with backing to reproduce the original article. Electrowinning is another application of electrodeposition used for the

commercial production from molten salts of active metals such as aluminum, magnesium, and sodium, and from aqueous solution of other metals such as copper, manganese, and antimony. Electrorefining is commonly used to purify metals such as silver, lead, and copper. The impure metal is used as the anode, and purified metal is deposited at the cathode. *See* ELECTROMETALLURGY; ELECTROPLATING OF METALS.

Electrolytic processes. Many electrode reactions other than metal deposition are of commercial or scientific use. Electrolysis of brine to yield chlorine at the anode, hydrogen at the cathode, and sodium hydroxide in the electrolyte is a very important industrial process, consuming a significant fraction of the national electrical production. Many organic compounds can be prepared electrolytically. *See* ELECTROLYSIS.

Electrothermics. While not strictly electrochemical, electrothermics is generally recognized as a part of the field. It includes high-temperature processes involving electric arc or resistance furnaces. *See* ELECTRIC FURNACE.

Electroanalytical chemistry. Many electrochemical measurements are useful for analytical purposes. Electrodes are commonly used for analytical purposes through measurement of their potentials and include the glass electrode for pH measurements, and ion-selective electrodes for certain ions, such as sodium or potassium ion (special glass compositions), calcium ion (liquid membrane), and fluoride ion (doped lanthanum fluoride single crystals). Solid electrodes, especially those composed of platinum and carbon, are generally used to examine current-voltage behavior of substances in solution. Polarography involves the use of a dropping mercury electrode as one electrode of an electrolytic cell. Its use is declining except for analytical applications and those involving operation at highly negative potentials. Qualitative analysis is carried out by measurement of characteristic potentials (half-wave potentials) for electrode processes, and quantitative analysis by measurement of diffusion-controlled currents. Coulometry involves the application of Faraday's law for analytical purposes. *See* COULOMETER; ELECTRODE; ION-SELECTIVE MEMBRANES AND ELECTRODES; POLAROGRAPHIC ANALYSIS.

Several methods involving electrolysis during short periods of electrolysis permit the application of diffusion theory (in the absence of convection) to calculate mass-transport rates and to identify reactive intermediates and reaction mechanisms in multistep electrode processes. These techniques include linear sweep and cyclic voltammetry (measurement of currents with linear voltage scan) and chronopotentiometry (measurement of potential-time transients under constant current conditions). Several titration methods involve electrochemical measurements, for example, conductometric, potentiometric, and amperometric titrations. *See* ELECTROCHEMICAL TECHNIQUES; TITRATION.

Electrode kinetics. Studies of the kinetics of electrode processes are valuable not only for

understanding the mechanisms of electrode reactions but also for studying homogeneous reactions that occur in solutions preceding or following the charge-transfer step. Such studies are made by pulse or transient techniques, or by steady-state methods involving dynamic systems, such as rotating-disc electrodes or flowing solutions.

Miscellaneous phenomena. Electrochemical transport of ions through synthetic or natural membranes is important for processes such as desalination of water and electro dialysis. In biological systems, the transmittal of nerve impulses and the generation of electrical signals such as brain waves are basically of electrochemical origin. A set of related phenomena can be grouped together under electrokinetic behavior, including the motion of colloidal particles in an electric field (electrophoresis), the motion of the liquid phase relative to the stationary solid under the influence of a potential gradient (electroosmosis), and the inverse generation of a potential gradient caused by a flowing liquid (streaming potential). Alternating-current phenomena, such as dielectric behavior, double-layer charging, and faradaic rectification, may also be included in a general definition of electrochemistry. Corrosion and passivation of metals are electrochemical in nature. *See* COLLOID; CORROSION; ELECTROENCEPHALOGRAPHY; ELECTROKINETIC PHENOMENA; ELECTROPHORESIS; STREAMING POTENTIAL.

Albert J. Fry; Herbert A. Laitinen

Organic electrochemistry. Organic electrochemistry involves the study of the chemical reactions that take place when an electric current is passed through a solution containing one or more organic compounds. It is a highly interdisciplinary science because understanding and fully developing a given organic electrochemical reaction may involve not only techniques of synthesis, purification, and identification of organic products but also theory and practice of a variety of sophisticated electroanalytical techniques, surface science, cell design, electronics, engineering scaleup, and materials modification. In recent years, organic electrochemistry has been of increasing interest for industrial applications because the costs of electrical current have been rising more slowly than the costs of conventional chemical reagents and because electrochemical procedures can be environmentally less intrusive than conventional chemical processes. A number of fine chemicals are made electrochemically on scales ranging from several kilograms for expensive specialty compounds to several tons per day. A few chemicals are made on a much larger scale; the best known is adiponitrile [$\text{NC}(\text{CH}_2)_4\text{CN}$], a key intermediate in the synthesis of nylon.

Types of processes. An extremely wide range of organic chemical conversions can be carried out electrochemically. A single electrochemical cell and power supply can generate either anode potentials powerful enough to oxidize saturated alkanes or cathode potentials sufficiently negative to reduce benzene to cyclohexene. If necessary, such processes could be done simultaneously in the same

electrochemical cell. An inorganic analogy illustrates this, whereby it is possible to oxidize fluoride ion to elemental fluorine at the anode of an electrochemical cell while reducing sodium ion to metallic sodium at the cathode. Of course, it is necessary to separate the two compartments so that the two products do not immediately react with each other. Almost any functional group can be modified electrochemically. *See* OXIDATION-REDUCTION.

The wide range of redox potentials that can be attained electrochemically and the precision with which this can be done represents only one aspect of organic chemical reactivity. A variety of other processes are possible in electrochemical cells. For example, it is possible to generate a chemical reagent at an electrode under carefully controlled conditions for reaction with another constituent of the medium. Reagents that have been generated in this manner include bromine, chlorine, iodine, iodine ion, hydrogen ion, hydroxide ion, solvated electrons, the powerful oxidant ruthenium oxide (RuO_4), and a variety of low- and high-valent metal ions. For many substrates, electrolysis is done at constant current, in which a fixed current is passed through the solution until the reaction is complete. In this technique, a simple and inexpensive power supply can be used. Selective electrolysis of a complex substance may require controlled-potential electrolysis, which uses a more sophisticated power source (a potentiostat) to maintain the electrode potential at a constant value. *See* ELECTROCHEMICAL PROCESS.

Methodology. The way that an electrochemical reaction is done depends upon the complexity of the process and the amount of material desired. In the laboratory, small amounts of material, from a few milligrams to a hundred grams, can be prepared in relatively simple cells or even in open beakers. Larger quantities are usually prepared in flow cells, often with many of these arranged in series or parallel. Industrial-scale electrochemical reactions are usually done in flow cells to reduce power costs and simplify operating conditions. An emerging technology is the use of microreactors for rapid (<1 s) reaction of electrochemically generated intermediates.

Electrode materials. A wide variety of substances have been used as electrodes, including the metals platinum, mercury, and iron, as well as carbon, which is proving to be the most versatile and least expensive electrode material. For many applications, the material of which the electrode is composed has little effect on the nature of the reaction taking place. The principal requirement is simply that the material be a conductor of electricity. Some applications, particularly those involving adsorbed intermediates and reagents, do take place more efficiently at electrodes made of certain materials. Electrodes made of an intimate mixture of carbon and sulfur, selenium, or tellurium have been used to introduce the latter elements into an organic substrate. Electrolysis at a platinum electrode in acidic media can effect catalytic hydrogenation of unsaturated substances. Other applications depend upon a special coating on the electrode surface to change

the nature of the electrode reactions occurring there. This coating can be a conductive polymer, a chemical redox agent, or an enzyme, where a high degree of selectivity is required. *See* ORGANIC CONDUCTOR.

Reaction mechanisms. The electrochemical behavior of thousands of organic substances has been studied. Consequently, it is possible to make generalizations about the kinds of processes likely to take place with a given substrate. Cathodic organic reactions fall into several categories: cleavage of single bonds, reduction of functional groups, and reduction of large conjugated systems such as activated alkenes and aromatic compounds. Anodic reactions are equally diverse. The oldest and best-known anodic reaction is the Kolbe reaction, in which electrochemical oxidation of carboxylate ions yields dimeric products with the evolution of carbon dioxide. A wide range of functional groups can be accommodated in this versatile reaction. Anodic oxidation of other substrates generally involves initial formation of cation radicals, which may undergo loss of a proton or alkyl group, reaction with solvent or a second component of the medium, or dimerization.

The powerful yet precisely controllable oxidizing and reducing conditions that can be achieved electrochemically are useful for generating novel organic intermediates, including carbocations, carbanions, radicals, radical ions, carbenes, nitrenes, and arynes. These reactive intermediates can frequently be trapped by other organic substances added to the medium to extend the synthetic utility of an electrode reaction. *See* ORGANIC REACTION MECHANISM; REACTIVE INTERMEDIATES.

Voltammetric methods are commonly used to provide detailed mechanistic information. In these, the electrode potential is varied in a controlled fashion in the vicinity of the redox potential of the substrate, and the current response is measured as a function of experimental variables such as scan rate, concentration, and added reagents. Other methods include identification of the products from a preparative electrolysis (often at controlled potential in mechanistic experiments); coulometry, that is, measurement of the actual amount of current passed in the electrolysis; comparison of experimental voltammetric or other responses with computer simulations of the theoretical behavior for various mechanisms; studies of changes in the structure of the electrode surface during electrolysis; and spectroscopic identification of intermediates during electrolysis. Microelectrodes, which have micrometer diameters, extend the speeds with which electrochemical measurements can be made, and thus permit measurement of even very fast rates of chemical reactions associated with electron transfers at electrodes.

Albert J. Fry

Bibliography. A. J. Fry, *Synthetic Organic Electrochemistry*, 2d ed., 1989; H. Lund and O. Hammerich (eds.), *Organic Electrochemistry*, 4th ed., 2000; T. Shono, *Electroorganic Synthesis*, 1991; S. Torii (ed.), *Novel Trends in Electroorganic Synthesis*, 1998.

Electrochromic devices

Self-contained, hermetically sealed, two-electrode electrolytic cells that change their ability to transmit (or reflect) light in response to a small bias (typically 1–2 V) applied across the two electrodes. The operation of electrochromic devices relies upon their electrochromic material content. These materials are organic or inorganic substances that are able to interconvert between two or more color states upon oxidation or reduction, that is, upon electrolytic loss or gain of electrons. The electrochromic materials that are appropriate for most practical applications are strong light absorbers in one redox state but colorless in another. Historically, the phenomenon of electrochromism based on physical phenomena, such as the Franz-Keldish and Stark effects, dates back to 1932. Since the mid-1970s, however, development has been focused on color changes carried out electrochemically, so that in modern usage the terms electrochromic devices and materials refer to devices and substances defined above. *See* STARK EFFECT.

Assembly. A typical electrochromic device is a sandwichlike structure with two glass plates and an electrolyte (**Fig. 1**). Each glass plate is coated on the inside with a transparent electrically conducting layer of indium-tin oxide, which operates as an electrode. Electrochromic mirrors include an additional reflective coating (for example, aluminum) on the outside of one of the glass plates. The electrolyte carries the ionic current inside the cell between the two electrodes, and it can be as simple as a salt (for example, sodium chloride, NaCl) dissolved in a dissociating solvent such as water. However, development has focused on gel and solid electrolytes, because they offer several advantages: they are easier to confine in the space between the electrodes; they function as

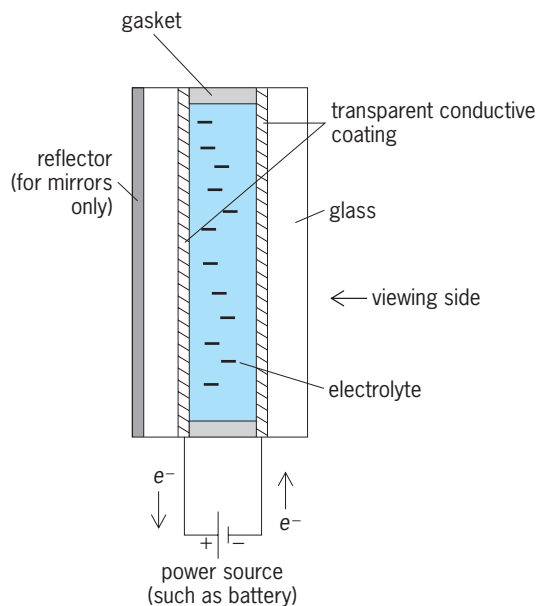


Fig. 1. Electrochromic device. The electrochromic materials can be either dissolved in the electrolyte or coated on the transparent electrodes.

laminators holding the two glass plates together; and their use minimizes the hydrostatic pressure that can cause substrate deformation and leakage problems, particularly in large-area devices such as smart windows.

State-of-the-technology electrochromic devices utilize two electrochromic materials with complementary properties: the first electrochromic material is normally reduced ($\text{ECM}_1^{\text{red}}$) and undergoes a colorless-to-colored transition upon oxidation (loss of electrons), while the second electrochromic material is normally oxidized (ECM_2^{ox}) and undergoes a similar transition upon reduction (gain of electrons). The electrochromic materials $\text{ECM}_1^{\text{red}}$ and ECM_2^{ox} are selected so that they do not react with each other. The oxidation of $\text{ECM}_1^{\text{red}}$ and the reduction of ECM_2^{ox} then are forced by the external power source (Fig. 1), which operates as an electron pump consuming energy in order to transfer electrons from one electrode to the other. Oxidation of $\text{ECM}_1^{\text{red}}$ occurs at the positive electrode (anode) and is a source of electrons, while reduction of ECM_2^{ox} occurs at the negative electrode (cathode) and is a sink of electrons. This approach, known as complementary counterelectrode technology, has two distinct advantages. First, the long-term operating stability of the electrochromic cell is greatly enhanced, because providing both a source and a complementary sink of electrons within the same cell prevents any electrolytic decomposition of the electrolyte. Second, the reinforcing effect of two electrochromic materials changing color simultaneously enhances the contrast difference between the color states per unit charge consumed. Depending on the location of the two electrochromic materials within the electrochromic devices, three main types of such devices exist: solution, precipitation, and thin-film (Fig. 2). See ELECTRODE; ELECTROLYTE; OXIDATION-REDUCTION.

Solution type. In these devices both electrochromic materials, $\text{ECM}_1^{\text{red}}$ and ECM_2^{ox} , are dissolved in the electrolyte, and they move to the electrodes by diffusion driven by the concentration gradients that are generated by the electrolytic depletion of $\text{ECM}_1^{\text{red}}$ and ECM_2^{ox} in the vicinity of the electrodes (Fig. 2a). The main advantage of the solution-type electrochromic devices is the variety of materials that can be used; any material that can be oxidized or reduced and is electrochromic is a potential candidate. Also, unlike the other types of electrochromic devices, the maximum color intensity can be controlled easily by varying the thickness of the electrolyte layer, typically between 0.1 and 5 mm (0.004 and 0.2 in.). More than 90% of the electrochromic antiglare rearview mirrors for automobiles sold worldwide are based on solution-type devices. A typical device employs N,N,N',N' -tetramethyl-*p*-phenylenediamine (TMPD) as $\text{ECM}_1^{\text{red}}$, and N,N' -diheptyl-4,4'-bipyridinium dication (HV^{2+}) as ECM_2^{ox} . Electric current forced through the electrode-solution interface is responsible for the electrolysis of the electrochromic materials to their colored forms. See ELECTROLYSIS.

In organic solvent-based electrolytes the blue radicals TMPD^+ and HV^+ are also soluble and diffuse away from the electrodes back into the bulk electrolyte, where they meet, react, and annihilate each other back to their respective colorless states. This property is a desirable feature that renders electrochromic rearview mirrors self-erasing, thereby providing a fail-safe system that automatically recovers the colorless, fully reflective state under power-failure conditions. However, in order for these devices to remain colored, power must be applied continuously. This property increases the energy consumption and makes these devices rather unsuitable for application in architectural glaze. Finally, another potential drawback of this approach is the speed of coloration, which is relatively slow (typically on the order of a few seconds) because it is controlled by diffusion of the electrochromic materials in the bulk electrolyte, a slow process. Implementation of gel or solid electrolytes with these devices slows the diffusion processes even further. See DIFFUSION.

Precipitation type. In these devices at least one of the electrochromic materials is originally dissolved in the electrolyte, but upon oxidation or reduction the colored product is electrodeposited onto the corresponding electrode (Fig. 2b). One of the first examples (1962) of electrolytic cells in light modulation is a precipitation device in that it employs the reversible electroplating of silver. Later (1973) it was discovered that in aqueous electrolytes the one-electron reduction product, HV^+ , of the N,N' -diheptyl-4,4'-bipyridinium dication, HV^{2+} , forms a blue precipitate on the electrode. That discovery significantly shifted attention from physical to electrochemical electrochromism.

The switching speed of precipitation-type devices is comparable to that of solution types, because both are controlled by diffusion. But precipitation devices offer inherently higher resolution. For example, with patterned-segment electrodes they can be used potentially for high-resolution displays. However, such devices have not been commercialized because precipitation-type electrochromic materials usually suffer poor cycling lifetimes; certain recrystallization effects within the colored precipitates progressively render them electrically inaccessible by the electrode.

Thin-film type. In these devices, both electrochromic materials are immobilized and confined as thin-film coatings (0.1–5 micrometers thick) on the surfaces of the two electrodes (Fig. 2c). The main advantage of this arrangement is that physical separation of the two electrochromic materials prevents annihilation of the colored forms, and provides an open-circuit memory that significantly decreases the consumption of power. A typical device consumes less than 10 millicoulombs/cm² for full coloration. In a sense, a thin-film-type electrochromic device is a rechargeable battery, in which the color of the electrodes depends upon the state of charge. Discharging (decolorizing) can be carried out either by reversing the bias or by simply short-circuiting the two

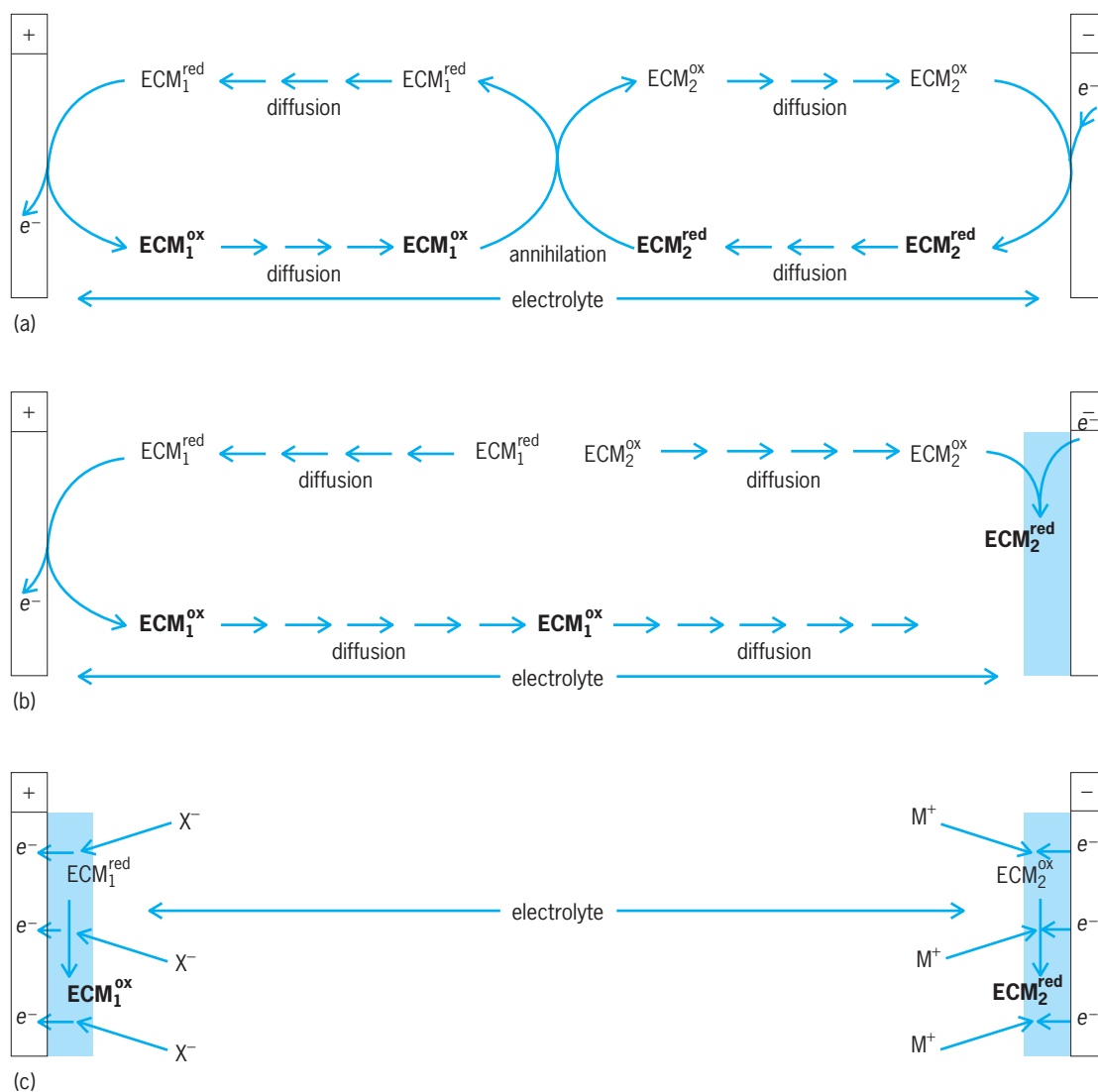


Fig. 2. Operating principles of electrochromic devices. (a) Solution type. (b) Precipitation type. (c) Thin-film type. The color-bearing forms of the electrochromic materials are shown in bold type.

electrodes. Furthermore, the switching speed of this type of device can be greatly improved because coloration does not depend on electrolysis of electrochromic materials dissolved in the bulk electrolyte. In practice, however, coloration depends on diffusion of ions in thin solid films, which is a much slower process than diffusion in the bulk electrolyte. Therefore, depending on the nature of the electrochromic materials, the switching speed of thin-film devices can vary from a fraction of a second to a few tens of seconds.

Several organic and inorganic materials are available for surface confinement on electrodes as thin films. A film of tungsten trioxide (WO_3), for instance, is an inorganic electrochromic material which, upon reduction, uptakes hydrogen ions (H^+) reversibly from an aqueous electrolyte [or lithium ions (Li^+) and sodium ions (Na^+) from a nonaqueous electrolyte] to form a blue tungsten bronze (M_nWO_3 , where M is H, Li, or Na). Certain electrochromic rearview mirrors used in trucks employ WO_3 as ECM_2^{ox} , and a nickel hydroxide [$\text{Ni}(\text{OH})_2$] film as $\text{ECM}_1^{\text{red}}$. These are man-

ual, switch-operated devices which are preset at a level with which the driver feels comfortable. Another inorganic electrochromic material that is complementary to WO_3 and suitable for thin-film electrochromic devices is Everitt's salt, $\text{K}_4\text{Fe}_4[\text{Fe}(\text{CN})_6]_3$, which is the colorless reduced form of Prussian blue, $\text{Fe}_4[\text{Fe}(\text{CN})_6]_3$.

Among organic materials, conductive polymers such as poly(pyrrrole), poly(aniline), and poly(thiophene) have been evaluated as thin-film electrochromic materials for oxidative coloration, and polymeric films derived from diquaternized 4,4'-bipyridine for reductive coloration. Generally, organic electrochromic polymer films are able to switch color in less than a second, and they seem to be inherently more durable than their inorganic counterparts because they better accommodate the structural changes induced upon oxidation and reduction. See ORGANIC CONDUCTOR.

Properties. Electrochromic devices are analogous to liquid-crystal devices in that they do not generate their own light but modulate the ambient light.

Liquid-crystal devices require use of polarizers; consequently, their viewing angle is limited, and lateral size limitations are imposed because the spacing between the electrodes (thickness) must be controlled within a few micrometers over the entire device area. Electrochromic devices do not require polarizers, thereby allowing a viewing angle approaching 180°, and contrast ratios similar to black ink on white paper (20:1 or better); moreover, control of the thickness is not important. Other desirable features of electrochromic devices include inherent color, continuous gray scale, and low average power consumption for the thin-film-type devices. Furthermore, it has been shown that electrochromic thin films can be patterned with a 2–5- μm resolution to form a large number of display elements that can be matrix-addressed. Nevertheless, even though there is no apparent intrinsic limitation, the best cycling lifetimes claimed for electrochromic materials are of the order of 10–20 million cycles, while the lifetime of liquid-crystal devices is of the order of several hundred million cycles. This long lifetime has made liquid-crystal devices a very successful technology in matrix-addressed, flat-panel displays.

The larger tolerance in thickness variation for electrochromic devices renders them better suited than liquid-crystal devices for large-area light modulation applications, such as smart windows, space dividers, and smart mirrors. Another possible application is in large-area displays that do not need frequent refreshing, such as signs and announcement boards. Reconfigurable optical recording devices (for example, disks) have been proposed as a high-resolution application that is within the presently available lifetimes of electrochromic materials. *See* ELECTROCHEMICAL PROCESS; ELECTROCHEMISTRY; ELECTRONIC DISPLAY; ELECTROOPTICS; LIQUID CRYSTALS. Nicholas Leventis

Bibliography. A. R. Kmetz and E. K. von Willisen (eds.), *Nonemissive Electrooptic Displays*, 1976; C. M. Lampert and C. G. Granqvist (eds.), *Large-Area Chromogenics: Materials and Devices for Transmittance Control*, vol. IS 4, SPIE Institute Series, 1990; N. Leventis and Y. C. Chung, New complementary electrochromic system based on polypyrrole-prussian blue composite, a benzylviologen polymer, and poly(vinylpyrrolidone)/potassium sulfate aqueous electrolyte, *Chem. Mater.*, 4:1415–1422, 1992; S. Matsumoto (ed.), *Electronic Display Devices*, 1990.

Electroconvulsive therapy

The controlled induction of a grand mal seizure under anesthesia for treating individuals suffering from certain psychiatric illnesses. Electroconvulsive therapy (ECT) was introduced as a treatment for psychiatric disorders by the Italian neurologist Ugo Cerletti in 1938. Cerletti developed a method of electrically inducing seizures in laboratory animals as a means of studying epilepsy. Aware of the observation that individuals with epilepsy and depression

experienced an improvement in mood following a seizure, he postulated that inducing a seizure in a depressed individual might improve his or her condition. In 1938 he reported the first case of a psychiatric patient treated with ECT with dramatic, positive results.

Since 1938 ECT has developed considerably. The major changes in the treatment since its inception include the use of general anesthesia to modify the motor component of the seizure; the use of a pulse wave ECT apparatus, which allows careful dosing of electrical energy; the delineation of specific pathological conditions that respond to ECT; and a clearer understanding of how electrode placement and energy dosing affects both therapeutic results and side effects.

Procedure. ECT is performed under general anesthesia. While often performed on hospitalized patients, it is increasingly used in an outpatient setting. The patient is evaluated medically to determine the risk of anesthesia. The required medical workup consists of a physical examination, an electrocardiogram, and laboratory tests, usually including a blood count and determination of electrolytes concentrations. The patient is not allowed to eat or drink after midnight prior to the procedure. Both cardiovascular monitoring and electroencephalographic (EEG) monitoring are done during the procedure. An intravenous line is started, and the patient is anesthetized using a short-acting barbiturate. After the patient is fully asleep, he or she is paralyzed using a short-acting muscle relaxant. When the patient is fully asleep and paralyzed, a grand mal seizure is induced using a pulse wave ECT apparatus. The seizure typically lasts 30 to 60 s. Patients regain consciousness in 5 to 15 min. They are monitored until they are fully recovered and able to eat and drink, which usually takes an additional 30 min.

Treatment dosing. Although all ECT treatment induces a grand mal seizure, it is now known that not all seizures are alike. Both the efficacy of the treatment and the side effects are influenced by stimulus electrode placement and the quantity of energy. Bilateral and right unilateral electrode placements are the most widely used treatment placements and the most researched. Bilateral placement with a dosing slightly above seizure threshold is marginally more effective than right unilateral placement with an energy dosing 4–6 times above seizure threshold. Bilateral placement yields a slightly faster response but also slightly more cognitive side effects. Recently, bifrontal electrode placement with an energy dosing slightly above seizure threshold has been introduced with results comparable to bilateral placement. However, this placement requires more study before it is widely introduced into clinical practice.

ECT-responsive disorders. The use of ECT to treat serious psychiatric disorders has been refined over the past half century. It has shown efficacy in major depression with and without psychotic features; in both the manic and depressed phases of bipolar disorder; in the treatment of acute schizophrenia, particularly when confusion is a prominent symptom;

and in the treatment of catatonia (a psychiatric condition characterized by unresponsiveness to external stimuli) regardless of the psychiatric etiology. Due to the expansion of medication options over the last 20 years, pharmacotherapy is almost always tried before ECT is initiated. The most notable exceptions to this are in circumstances in which a fast response is needed, such as in the treatment of individuals who are dangerously suicidal or severely nutritionally impaired. ECT can show results in days rather than the weeks it typically takes with medication treatment.

The efficacy of ECT in these disorders varies. It appears to be most effective in the treatment of depressed individuals with psychotic symptoms, yielding about a 92% response rate. This compares with an approximately 62% response rate in individuals with major depression without psychotic symptoms. ECT is also extremely effective in the treatment of catatonia, with response rates of about 90%. The efficacy of medication in the treatment of mania and acute schizophrenia is so high that ECT is rarely used in these disorders currently. ECT has been used to treat obsessive-compulsive disorder (OCD) with disappointing results. However, there is a high rate of comorbid (coexisting) depression with OCD, and in this circumstance ECT can be very effective at treating the depression.

ECT has also shown efficacy in the treatment of neuroleptic malignant syndrome (NMS), a confusional state believed to result from an idiosyncratic reaction to antipsychotic medication. ECT has been shown to limit the symptoms and shorten the course of NMS, particularly when it is initiated early.

ECT also has efficacy in treating the depression of Parkinson's disease. In addition, patients with Parkinson's without depressive symptoms have shown improvement in their motor symptoms after receiving a course of ECT. *See* AFFECTIVE DISORDERS; OBSESSIVE-COMPULSIVE DISORDER; PARKINSON'S DISEASE; SCHIZOPHRENIA.

Side effects. The side effects of ECT are divided into two categories: those resulting from the induction of a grand mal seizure and those resulting from the administration of anesthesia.

Grand mal seizure. The induction of a grand mal seizure is accompanied by numerous physiologic changes that can cause side effects. When the electrical stimulus is administered, there is an immediate centrally mediated vagal nerve response, which results in sinus bradycardia (a slow heart rate). This is followed by sinus tachycardia (a fast heart rate) and increase in blood pressure during the seizure. When the seizure ends, there is a peripherally mediated vagal nerve response, which again results in bradycardia. As a consequence, the heart is subjected to rapid changes in rate and blood pressure. This can result in both ventricular and atrial arrhythmias. Typically these arrhythmias consist of several premature atrial contractions or premature ventricular contractions that are self-limited. The potential does exist for the development of significant and potentially life-threatening arrhythmias. For this reason alone, ECT

is always done in the presence of an anesthesiologist who is prepared to treat any arrhythmia that might occur. *See* CARDIAC ELECTROPHYSIOLOGY.

The occurrence of a grand mal seizure is associated with short-term memory loss and confusion. The seizure itself causes both anterograde and retrograde amnesia, that is, loss of memory for events occurring both after and before the seizure, respectively. Typically this effect is cumulative so that individuals report memory problems worsening over the course of treatment. These effects are most intense during the treatment course and the several weeks surrounding it. On rare occasion individuals have reported autobiographical memory loss as far back as 6 months prior to the treatment. Evaluating the deleterious cognitive effects from ECT is confounded by the fact that depression causes significant memory problems, as do many of the medications used in its treatment. When assessing the risks and benefits of ECT in regard to memory problems, it is important to weigh the fact that tests show memory functioning—as measured by attention, concentration, and the capacity to acquire and use new information—actually improves after ECT.

The grand mal seizure also increases intracranial pressure. A minor consequence of this is headache. A potentially major consequence is that any central nervous system condition that is at risk for worsening from this increase in pressure needs to be carefully evaluated prior to the individual's undergoing treatment. For example, aneurysms can potentially rupture, or tumors could result in herniations. *See* AMNESIA; MEMORY; SEIZURE DISORDERS.

Anesthesia. Anesthesia-related side effects are similar to those associated with anesthesia used in any minor surgical procedure. The most serious risk is an idiosyncratic reaction to the anesthetic agent, which is associated with a mortality rate of approximately 1 in 100,000. More common side effects are minor headaches, muscle pain, dry mouth, and, rarely, urinary retention. *See* ANESTHESIA.

Mechanism of action. The therapeutic effects of ECT are related to the energy level, the grand mal seizure, and the electrode placement. The treatment produces several effects on neurotransmitters that can be responsible for its efficacy. There is an increase in both neurotransmitter production and release and neurotransmitter receptor sensitivity. There is also an increase in seizure threshold over the course of the treatments that is associated with a clinical response. Energy dosing appears to be important with respect to the electrode placement. Energy dosing at seizure threshold is effective with bilateral electrode placement, but the dosing needs to be four to six times above threshold when right unilateral placements are used. The reason for this is not precisely understood. There is a relative sparing of the speech and language center when right unilateral electrode placement is used, which may explain the reduction in cognitive side effects with this technique.

ECT has been shown to increase long-term potentiation (LPT) pathways in laboratory animals. LPT is

involved in the formation of neural nets related to memory formation. When this effect is blocked in laboratory animals using chemical agents, the memory side effects of ECT are reduced. Work is now under way with human subjects to see if blocking this effect can result in decreased memory side effects.

Robert Ostroff

Bibliography. American Psychiatric Association Committee on Electroconvulsive Therapy, *The Practice of Electroconvulsive Therapy*, 2d ed., American Psychiatric Press, Washington, DC, 2001; S. H. Lisanby et al., The effects of electroconvulsive therapy on memory of autobiographical and public events, *Arch. Gen. Psychiat.*, 57:581-590, 2000; K. G. Rasmussen, S. M. Sampson, and T. A. Rummans, Electroconvulsive therapy and newer modalities for the treatment of medication-refractory mental illness, *Mayo Clinic Proc.*, 77:552-556, 2002; H. A. Sackeim et al., A prospective, randomized, double-blind comparison of bilateral and right unilateral electroconvulsive therapy at different stimulus intensities, *Arch. Gen. Psychiat.*, 57:425-435, 2000.

Electrode

An electrical conductor through which an electric current enters or leaves a conducting medium, whether it be an electrolytic solution, solid, molten mass, gas, or vacuum. For electrolytic solutions, many solids, and molten masses, an electrode is an electric conductor at the surface of which a change occurs from conduction by electrons to conduction by ions. For gases and vacuum, the electrodes merely serve to conduct electricity to and from the medium. See ELECTRODE POTENTIAL; ELECTROLYSIS; ELECTROMOTIVE FORCE (CELLS).

Walter J. Hamer

Electrode potential

The equilibrium potential difference between two conducting phases in contact, most often an electronic conductor such as a metal or semiconductor on the one hand, and an ionic conductor such as an electrolyte solution (a solution containing ions) on the other. Electrode potentials are not experimentally accessible, but the differences in potential between two electronic conductors making contact with the same ionic conductor (that is, the difference between two electrode potentials) can be measured. A useful scale of electrode potentials can therefore be obtained when a particular electrode potential is set equal to zero by definition. There are several conventions, based on different definitions of the zero point on the scale of electrode potential, but all tables use the so-called standard hydrogen convention. See ELECTRODE; REFERENCE ELECTRODE.

Mechanism. The interfacial potential difference is usually the consequence of the transfer of some

charge carriers from one conducting phase to the other. For example, when a piece of silver, which contains silver ions and free, so-called conduction electrons, is in contact with an aqueous solution of silver nitrate, the only species common to the two phases are the silver ions. Their concentration (volume density) is constant in the metal but variable (from zero to the solubility limit of the silver salt used) in the solution. When more silver ions transfer from the solution to the metal than in the opposite direction, an excess of negatively charged nitrate ions remains in the solution, which therefore acquires a negative charge. However, the metal gains more silver ions than it loses, and therefore acquires a positive charge. Such a charge separation leads to a potential difference across the boundary between the two phases. The continued buildup of such charges makes the potential of the metal more and more positive with respect to that of the solution. This effect in turn leads to electrostatic repulsion of the silver ions in the solution phase immediately adjacent to the metal; these are the very metal ions that are candidates for transfer across the boundary. Consequently, the electrostatic repulsion decreases the tendency of silver ions to move from the solution to the metal, and eventually the process reaches equilibrium, at which point the tendency of ions to transfer is precisely counterbalanced by the repulsion of the candidate ions by the existing potential difference. At that potential, there is no further net transfer of charges between the contacting phases, although individual charges can still exchange across the phase boundary, a process which gives rise to the exchange current.

With a small number of metals, especially silver and mercury, interfacial equilibrium is often established rapidly, well within milliseconds. However, with most metals, the equilibrium state is reached much more slowly, especially when the electrode process involves bond breaking or other complicated reaction mechanisms. In some cases, equilibrium is never attained, and the equilibrium potentials can be known only through calculation.

Concentration dependence. For a metal in contact with its metal ions of valence z , the potential difference E can be expressed in terms of a standard potential E° (describing the affinity of the metal for its ions) and the concentration c of these ions in solution through the Nernst equation (1), where R is the gas

$$E = E^\circ + (RT/zF) \ln c \quad (1)$$

constant, T is the absolute temperature, and F is the Faraday.

When the metal ions in solution form a sparingly soluble salt, the solubility equilibrium can be used to convert a metal electrode responding to its own metal cations (positive ions) into an electrode responding to the concentration of anions (negative ions). Typical examples are the silver/silver chloride electrode, based on the low solubility of silver chloride (AgCl), and the calomel electrode, based on the poor solubility of calomel (Hg₂Cl₂).

An equilibrium potential difference between a metal and an electrolyte solution can also be established when the latter contains a redox couple, that is, a pair of chemical components that can be converted into each other by the addition or withdrawal of electrons, by reduction or oxidation respectively. In that case the metal often merely acts as the supplier or acceptor of electrons. When metal electrons are donated to the solution, the oxidized form of the redox couple is reduced; when the metal withdraws electrons from the redox couple, its reduced component is oxidized. Again, the buildup of a charge separation generates a potential difference, which counteracts the electrochemical charge transfer and eventually brings the process to equilibrium, a state in which the rate of oxidation is exactly equal to the rate of reduction. The dependence of the equilibrium electrode potential on the concentrations of c_{ox} and c_{red} of the oxidized and reduced forms respectively is described by a Nernst equation of the form (2), where n denotes the number of electrons

$$E = E^\circ + (RT/nF) \ln(c_{\text{ox}}/c_{\text{red}}) \quad (2)$$

(e^-) transferred between the oxidized species (Ox) and the reduced species (Red) in the reactions $\text{Ox} + ne^- \rightleftharpoons \text{Red}$. A typical example is the reduction of hydrogen ions H^+ to dissolved hydrogen molecules H_2 , and vice versa, in which case the reactions are $2\text{H}^+ + 2e^- \rightleftharpoons \text{H}_2$, and for which the Nernst equation is of the form (3).

$$E = E^\circ + (RT/2F) \ln(c_{\text{H}^+}^2/c_{\text{H}_2}) \quad (3)$$

Redox potentials involving a gas are often established slowly, if at all. For determinations of such a redox potential, platinum is often used as the metal, because it is chemically and electrochemically stable. See OXIDATION-REDUCTION.

Electrode potentials can also be established at double phase boundaries, such as that between two aqueous solutions separated by a glass membrane. This glass electrode is commonly used for measurements of the pH, a measure of the acidity or basicity of solutions. The mechanism by which the glass electrode operates involves ion exchange of hydrogen ions at the two glass-solution interfaces. See ION EXCHANGE.

In all the above examples, the two contacting phases can have only one type of charge carrier in common. Usually, no equilibrium potential difference is established when more than one type of charge carrier can cross the interface, but often (depending on the nature of the metal and of the chemical components of the solution, and sometimes also depending on the geometry of the contact region) an apparently stable potential can still be obtained, which corresponds to zero net charge transfer. This can be a so-called mixed potential, important in metal corrosion, or a junction potential, which figures in most measurements of electrochemical potentials and usually limits the accuracy and precision of such measurements, including that of the pH. See CORROSION; ELECTROCHEMICAL SERIES.

Complications. In determining electrode potentials, there are several complications. In the first place, it follows from thermodynamics that the Nernst equation should be written in terms of activities rather than concentrations. The difference between these two parameters is often small but seldom completely negligible. The problem is that the activities of ions (in contrast to those of neutral molecules) are not experimentally accessible, although there are theoretical models that allow their estimation in sufficiently dilute solutions. See ACTIVITY (THERMODYNAMICS).

Second, measurements of potential differences always involve the potential difference between two metals rather than that between a metal and a solution. Therefore, electrode potentials as defined above cannot be measured either. The latter complication is circumvented by setting the electrochemical potential of one particular metal-solution interface equal to zero by definition. W. Nernst suggested that a normal hydrogen electrode be used as the common zero point; he defined this as a metal in contact with a solution saturated with hydrogen gas at 1 atm partial pressure and containing 1 mole of a monoprotic acid per liter. However, it was soon found that the nature of the acid used, and the presence of other components in the solution, would affect the potential of such an electrode—a reflection of the fact that the electrode responds to the activity rather than the concentration of the hydrogen ions. The problem here is that the activity of ions (including that of hydrogen ions) is not known, except at or near infinite dilution. Thus the standard hydrogen electrode, defined in terms of hydrogen ion activity rather than concentration, is not experimentally realizable. An alternative definition, in terms of an acid activity, has an anion sensitivity similar to the original definition in terms of concentrations. Moreover, use of the hydrogen electrode can lead to a high liquid junction potential, whereas use of the standard calomel electrode tends to reduce it. In practice, therefore, other electrodes are often used, such as silver-silver chloride electrode or a calomel electrode. In many areas of electrochemistry, the de facto reference electrode is the saturated calomel electrode, which consists of mercury plus calomel in a saturated aqueous solution of potassium chloride (KCl). The standard hydrogen electrode is used mostly for physico-chemical calculations.

Because these measurements involve a potential difference, there has been considerable confusion about the definition of that difference; that is, whether it is defined as the potential of the metal minus that of the solution, or the other way around. This is a matter of a sign convention. The problem is usually framed in terms of oxidation potentials versus reduction potentials. In 1953 the international community settled on the latter and simply called them electrode potentials.

Another zero point is often used to anchor the scale of electrode potentials of semiconductor electrodes. Here the reference point is the work function, that is, the electrical work necessary to remove an

electron from the interior of the phase to a far-away position in vacuum. This is, theoretically, a more meaningful scale. Unfortunately, the work function is not known with sufficient accuracy and precision to make this scale practical. *See* WORK FUNCTION (THERMODYNAMICS).

There are other factors besides charge transfer that can cause interfacial potential differences. Oriented dipolar molecules at interfaces can give rise to interfacial potential differences, even though they represent a charge separation within a layer of only molecular dimensions. Similarly, adsorbed molecules and ions can affect the interfacial potential difference when their adsorption involves a partial charge transfer. None of these phenomena affect the thermodynamic potentials, but they have effectively blocked attempts to define an absolute potential that is both practical and theoretically meaningful.

Applications. There are four main applications for measurements of electrode potentials: (1) in the establishment of the oxidative and reductive power of redox systems, the so-called electromotive series; (2) as concentration probes, such as in pH measurements; (3) as sources of chemical equilibrium data; and (4) as the primary (or independent) variable in studies of electrode reactions.

The electromotive series is a listing of redox couples in the order of their potentials. The most positive potentials correspond to the strongest oxidants, while most negative potentials identify the strongest reductants. This series therefore serves to organize redox systems according to their oxidative or reductive power. *See* ELECTROMOTIVE FORCE (CELLS).

A combination of two electrodes can be employed to determine the concentration of a particular species in a sample solution. One of the two electrodes is the indicator electrode, sensitive to a particular ionic species in the sample solution. The other, reference electrode typically comprises a separate compartment containing a metal in contact with its sparingly soluble salt and an excess of the anion of that salt (as in the calomel and silver/silver chloride electrodes), or with a redox couple such as the combination of iodide and triiodide ions. The solution in the reference electrode compartment then makes contact with the sample solution through a very constrictive contact, such that solution mixing is minimized. Moreover, the nature of the solution in the junction is chosen so that the potential difference across that liquid junction will be as small as possible. This can be achieved by using a high concentration of a salt with, as much as possible, equitransferent ions (ions that have near-equal mobilities, such as in KCl or NH_4NO_3). Under those conditions the liquid junction potential can often be neglected, and the Nernst equation can be written in terms of measurable potentials and concentrations. This is the basis of the measurement of pH using the combination of a glass electrode and a calomel, silver/silver chloride or iodide/triiodide reference electrode. Many other indicator electrodes are available, responding in a similar fashion to different cations and anions. There are also electrodes that respond to neutral compounds.

Such electrodes are often based on the presence of an immobilized enzyme. For instance, when an enzyme reacts with its substrate, it may release or consume hydrogen ions. The electrode then responds to the resulting change in the concentration of hydrogen ions and, therefore, indirectly to the enzyme substrate. *See* ION-SELECTIVE MEMBRANES AND ELECTRODES.

Because equilibrium potentials can be correlated with chemical equilibrium constants, measurements of the potentials of electrochemical cells are a major source of numerical information about equilibrium constants. Measurements of electrode potentials not involving liquid junctions can yield data accurate to about 10 microvolts, from which salt activities (but not the corresponding ionic activities) can be calculated. Measurements involving liquid junctions (such as pH measurements using reference electrodes with liquid junctions) are much less reliable, because of the uncertainties in the liquid junction potentials, of the order of one or a few millivolts.

The rates of electrode reactions depend strongly (typically, exponentially) on the electrode potential, and most measurements of the rates of electrode reactions are therefore performed under tight control of this potential, often using a specialized instrument called a potentiostat or, in biophysical applications, a voltage clamp. *See* BIOPOTENTIALS AND IONIC CURRENTS; ELECTROCHEMISTRY.

Robert de Levie

Bibliography. M. S. Antelman, *The Encyclopedia of Chemical Electrode Potentials*, 1982; A. J. Bard and H. Lund (eds.), *Encyclopedia of the Electrochemistry of the Elements and Compounds*, vols. 1-15, 1973-1984; A. J. Bard, R. Parsons, and J. Jordan (eds.), *Standard Potentials in Aqueous Solutions*, 1985.

Electrodermal response

A transient change in certain electrical properties of the skin, associated with the sweat gland activity and elicited by any stimulus that evokes an arousal or orienting response. Originally termed the psychogalvanic reflex, this phenomenon became known as the galvanic skin response. Electrodermal response (EDR) has replaced galvanic skin response as the collective term.

Mechanism. The skin of a relaxed person has a low electrical conductance (high resistance), and the skin surface is some 40 mV negative with respect to interior tissues. Sweat gland activity changes these electrical properties by increasing skin conductance and by changing the balance of positive and negative ions in the secreted fluid. Human skin is perforated by the ducts of 2-3 million eccrine sweat glands that are distributed most densely on the chest and forehead and, especially, on the palms and soles, where the orifices can be seen (under 10 \times magnification) as a row of tiny craters lining the dermatoglyphic ridges. On most of the body surface, sweat glands subserve the function of thermoregulation, but on the palms and soles, sweating varies with

psychological arousal. In many persons this “mental sweating” can spread to regions of the head or chest, especially if the person is highly aroused or if the ambient temperature is high.

Significance. Tonic skin conductance varies with psychological arousal, rising sharply when the subject awakens and rising further with activity, mental effort, or especially stress. Phasic skin conductance responses are wavelike increases in skin conductance that begin 1–2 s after stimulus onset and peak within about 5 s. The amplitude of the skin conductance response varies with the subjective impact of the eliciting stimulus, which in turn varies with the intensity of the stimulus, its novelty or unexpectedness for the subject, and its meaning or signal value. Aroused subjects display spontaneous skin conductance responses, generated apparently by mental events or other internal stimuli; their frequency, like the tonic skin conductance level, increases with the level of arousal. Skin potential phenomena are more complicated, and it is doubtful that there is any psychophysiological information present in skin potential level or skin potential responses that cannot be more easily obtained from measurements of skin conductance. Volar sweating, the parent phenomenon, increases the tactile sensitivity, the resistance to abrasion, and the adhesive qualities of these grasping surfaces; thus its function seems to be preparatory to manipulation.

Measurement. A low-impedance voltage source (0.5 V, 500 ohms) is applied to the skin through nonpolarizing electrodes attached typically to the ventral surface of the distal phalange, where sweat glands are most numerous. Skin contact is effected by means of an electrolyte paste containing 0.5% KCl or NaCl. A small (500-ohm) resistor in series with the subject generates a voltage that varies directly with changes in skin conductance. This signal is amplified and recorded on the moving paper chart of a polygraph. *See* LIE DETECTOR.

Applications. Electrodermal responses are measured in studies of emotion and stress, conditioning, habituation, and cognitive processing: that is, when it is desired to assess the differential or changing impact of a series of stimuli. *See* ELECTROENCEPHALOGRAPHY; SYMPATHETIC NERVOUS SYSTEM.

David T. Lykken

Electrodiagnosis

Diagnosis employing equipment that can be used clinically to measure the intrinsic electrical activity of the heart (electrocardiography), the brain (electroencephalography), peripheral nerves (nerve conduction), and skeletal muscles (electromyography). Electrodiagnostic procedures depend upon the difference in electrical potential between the interior and external surfaces of living cells. This potential difference is created by use of cellular metabolic energy to establish and maintain the concentration gradient of ions across the cell membrane. A discharge of the electrical potential (depolarization) by stimulation

or injury occurs in a wave proceeding from cell to cell. The ensuing electrical currents are conducted through the tissues to the surface of the body where electrodiagnostic equipment amplifies and records the changes in potential.

In peripheral nerve conduction studies, signals elicited by electrical stimulation of a nerve are measured either over a muscle supplied by the nerve or along the nerve itself. In another electrodiagnostic test, termed electromyography, electrodes are inserted into a muscle to record its electrical activity. A set of diagnostic tests called evoked potentials records the electrical responses of the brain that follow visual or auditory stimuli.

Electrocardiography. The electrocardiogram (ECG) employs electrodes on the extremities and at various points on the chest wall to record voltage changes due to depolarization and repolarization of heart muscle (myocardium). The ECG is most useful for the diagnosis and treatment of disturbances in heart rhythm (arrhythmias) or in the conduction of electrical impulses within the heart. Both inadequate delivery of oxygen to the myocardium (myocardial ischemia) and localized death of myocardium (infarction) produce characteristic changes in the ECG pattern. The ECG also detects hypertrophy, or overgrowth (of the myocardium), and certain abnormalities of blood electrolytes. ECG recordings during a graded exercise test or over a 24-h period of normal activities (Holter monitoring) may detect ischemic changes or arrhythmias that are not present at rest. *See* CARDIAC ELECTROPHYSIOLOGY; INFARCTION.

Electroencephalography. The electroencephalogram (EEG) records waves of electrical activity from the surface of the brain through the placement of electrodes on the scalp in a standardized pattern. Electroencephalography is the most useful procedure for evaluating individuals with epileptic seizures. Because of their greater sensitivity, computerized tomography (CT) scanning and magnetic resonance imaging (MRI) have largely replaced the EEG for the detection and localization of structural brain abnormalities such as tumors or infarcts. The most important use of EEG is for evaluation of brain abnormalities due to metabolic disorders, drug intoxication, and brain infections that can produce changes in the EEG pattern. The total absence of electrical activity (noted by a flat line on the EEG) for 30 min is an important criterion for diagnosing brain death. *See* DEATH; ELECTROENCEPHALOGRAPHY.

Peripheral nerve conduction studies. Nerve conduction velocity is determined by supramaximally stimulating the surface over a nerve and measuring the time it takes for the signal to reach a recording electrode over the innervated muscle (for a motor nerve) or at another spot along the path of a sensory nerve. Nerve damage results in slowed conduction or decreased amplitude of electrical signals along a sensory nerve. Slowed conduction is a common feature of demyelinating familial neuropathies; reduced amplitudes of action potentials along sensory nerves

are observed in diabetic neuropathy. The delivery of rapid, repetitive shocks to a motor nerve results in a swift reduction in the amplitude of the electrical response of the muscle in individuals with myasthenia gravis. *See* MYASTHENIA GRAVIS.

Electromyography. Electromyography (EMG) assesses the function of muscles with measurements of their electrical activity both at rest and during contraction. The electrical activity of muscles is examined by inserting an electrode into them. The amplified signals are displayed on an oscilloscope screen or heard through a loudspeaker. Abnormal EMG patterns can result from damage to the muscles or their nerve supply. Normal skeletal muscles are electrically silent at rest, but damage to or interruption of the nerve supply to a muscle can lead to scattered electrical spikes caused by spontaneous depolarization of single nerve fibers. During their contraction, damaged muscles produce electrical signals that are diminished in amplitude and duration. Electromyography is useful in individuals with amyotrophic lateral sclerosis (Lou Gehrig's disease) and muscular dystrophy. *See* ELECTROMYOGRAPHY.

Evoked potentials. Visual, auditory, or other sensory stimulation produces small transient potentials in the brain. These evoked potentials require hundreds or thousands of rapidly repeated stimuli to generate the evoked potential or response that requires a computer to summate and average the signals. Frequently used for the diagnosis of optic neuritis or multiple sclerosis is a recording over the back of the skull of the visual evoked potential generated by the display of a light signal to one eye at a time. The auditory evoked potential recorded over the scalp in response to high-frequency clicks delivered by an earphone is also useful for the diagnosis of multiple sclerosis. *See* BIOPOTENTIALS AND IONIC CURRENTS.

Simeon Margolis

Bibliography. M. J. Aminoff, *Electrodiagnosis in Clinical Neurology*, 4th ed., Churchill Livingstone, New York, 1999; K. H. Chiappa (ed.), *Evoked Potentials in Clinical Medicine*, 3d ed., Lippincott-Raven, Philadelphia, 1997; A. L. Goldberger, *Clinical Electrocardiography: A Simplified Approach*, 6th ed., Mosby, St. Louis, 1999; A. H. Ropper and R. H. Brown, Electrophysiologic and laboratory aids in the diagnosis of neuromuscular disease, Chap. 45 in *Adams and Victor's Principle of Neurology*, 8th ed., McGraw-Hill, New York, 2005; A. H. Ropper and R. H. Brown, Special techniques for neurologic diagnosis, Chap. 2 in *Adams and Victor's Principle of Neurology*, 8th ed., McGraw-Hill, New York, 2005.

Electrodynamics

The study of the relations between electrical, magnetic, and mechanical phenomena. This includes considerations of the magnetic fields produced by currents, the electromotive forces induced by changing magnetic fields, the forces on currents in magnetic fields, the propagation of electromag-

netic waves, and the behavior of charged particles in electric and magnetic fields. Classical electrodynamics deals with fields and charged particles in the manner first systematically described by J. C. Maxwell, whereas quantum electrodynamics applies the principles of quantum mechanics to electrical and magnetic phenomena. Relativistic electrodynamics is concerned with the behavior of charged particles and fields when the velocities of the particles approach that of light. Cosmic electrodynamics is concerned with electromagnetic phenomena occurring on celestial bodies and in space. *See* ELECTROMAGNETISM; ELECTRON MOTION IN VACUUM; MAXWELL'S EQUATIONS; QUANTUM ELECTRODYNAMICS; RADIO-WAVE PROPAGATION; RELATIVISTIC ELECTRODYNAMICS.

John W. Stewart

Electroencephalography

The biomedical technology and science of recording the minute electric currents produced by the brains of human beings and other animals. Electroencephalography (EEG) was discovered by Hans Berger in the late 1920s, and was found to have important clinical significance for the diagnosis of brain disease. The interpretation of EEG records has become a clinical specialty for neurological diagnosis.

Technology and physiology. The recording machine, the electroencephalograph, usually produces a 16-channel ink-written record of brain waves, the electroencephalogram. It is interpreted by an electroencephalographer. The placement of about 20 equally spaced electrodes pasted to the surface of the scalp is in accordance with the standard positions adopted by the International Federation of EEG, and is called the 10/20 system. Electrode positions are carefully measured so that subsequent EEGs from the same person can be compared. About 10 patterns or montages of combinations of electrode pairs are selected for transforming the spatial location from the scalp to the channels which are traced on the EEG pen writer.

The aggregate of synchronized neuronal activity from hundreds of thousands or millions of neurons acting together form the electrical patterns on the surface of the brain. The cellular basis of the EEG depends on the spontaneous fluctuations of post-synaptic membrane potentials between the inside and the outside of the dendritic processes of post-synaptic cells. The brain waves do not derive from the synchronized axonal spike discharges of cerebral cortical nerve cells.

Electrical voltage is transduced from the scalp by differential input amplifiers and amplified about a million times in order to drive the pens for the paper record. The recording usually takes 30–60 min during a relaxed waking state, and also during sleep when possible. Often, activating procedures are used, such as a flickering light stimulator and hyperventilation or overbreathing for about 3 min.

Normal EEG. EEG waves are defined by form and frequency. Various frequencies are given Greek

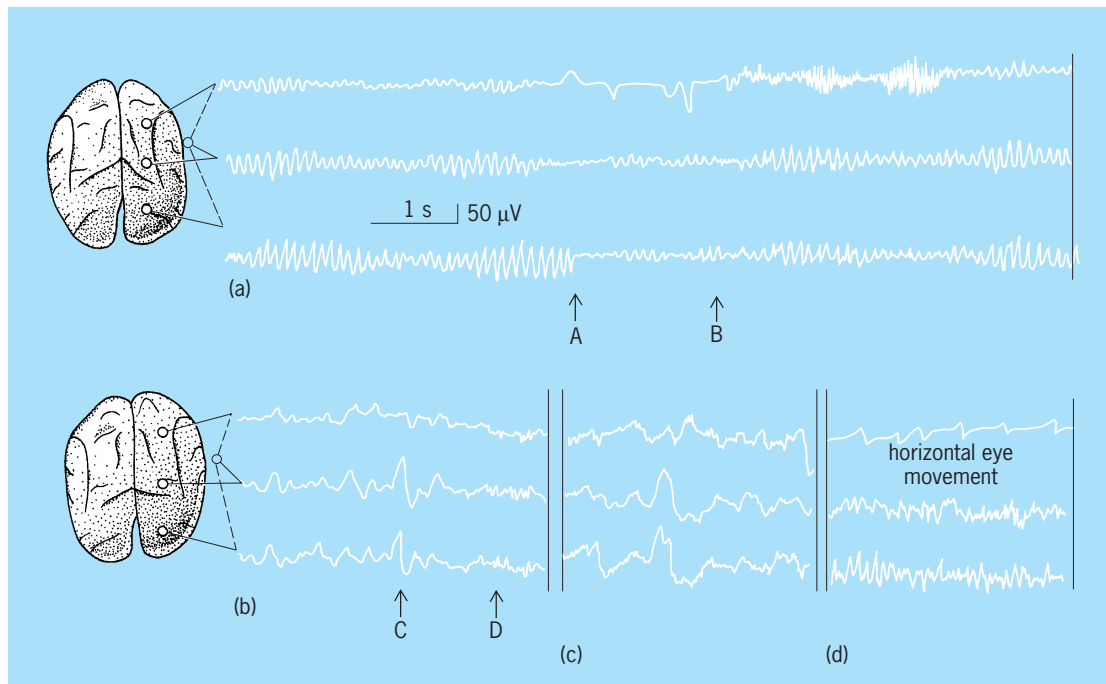


Fig. 1. Normal EEG waves. (a) The alpha rhythm in the waking record blocks at eye opening; eyes opened at A and closed at B. (b) A light-sleep EEG has theta waves and a sleep hump at C and spindles at D. (c) Deep sleep reveals constant, irregular slow delta waves. (d) In the dream state, alpha and theta activities are seen, almost like in the waking record.

letter designations. Alpha rhythm is defined as 8–12-Hz sinusoidal rhythmical waves (**Fig. 1a**). Alpha waves are normally present during the waking and relaxed state and enhanced by closing the eyes. They are suppressed or desynchronized when the eyes are open, or when the individual is emotionally aroused or doing mental work. They may be synchronized by bright light flashes and driven over a wide range of frequencies by repetitive visual stimulation (alpha driving). They are of highest amplitude in the posterior regions of the brain. The alpha rhythm develops with age, reaching maturity by about 12 years, stabilizes, and then declines in frequency and amplitude in old age (over 65).

Beta rhythms are faster, low-voltage sinusoidal waves, usually about 14–30 Hz. They are more prominent in the frontal areas. They are often synchronized and prevalent during sedation with phenobarbital or with the use of tranquilizers and some sedative drugs.

Slower rhythms are theta and delta waves. Theta waves of 4–7 Hz usually replace the alpha rhythm during drowsiness and light sleep (**Fig. 1b**). Delta waves of 0.5–4 Hz are present during deep sleep in normal people of all ages (**Fig. 1c**), and they are the primary waves present in the records of normal infants. Delta waves are almost always pathological in the waking records of adults.

Sleep stages are objectively defined by the use of EEG records. Theta waves are present in drowsiness and light sleep (stage I). Light sleep (stage II) is characterized by negative sharp waves at the vertex (top of head), called sleep humps, and spindle bursts of 13–15-Hz waves in the central regions of the head (over the Rolandic fissures). These are inter-

mixed with theta waves, and seen only during light sleep, sometimes called the hump-and-spindle stage of sleep.

Low-voltage, desynchronized delta waves are present during deep sleep (stage III), and higher-amplitude, slower, and more synchronized delta waves are seen in the deeper sleep stage (stage IV). Stage V sleep is called the dream stage or rapid-eye-movement (REM) stage. That is when dreams take place, accompanied by rapid horizontal eye movements, as if the sleeper were scanning a real scene. The EEG pattern is more like waking, with alpha rhythm present (**Fig. 1d**). The REM state is the one in which it is hardest to wake up the subject, so that it is sometimes called paradoxical sleep. *See SLEEP AND DREAMING.*

Clinical uses of EEG. The EEG reveals functional abnormalities of the brain, whether caused by localized structural lesions, essential paroxysmal states such as epilepsy, or toxic and abnormal metabolic conditions. The three major classes of abnormalities are asymmetries between the hemispheres, slow rhythms, and very sharp waves or spikes. Slow waves represent a depression of cerebral cortical activity or injury in the projection pathways beneath the recording electrodes. Sharp waves or spikes often indicate a hyperexcitable or irritable state of the cortex. During a full epileptic seizure attack, spikes become repetitive and synchronized over the whole surface of the brain. *See SEIZURE DISORDERS.*

The EEG is frequently used for the evaluation of comatose states. The record is slowed in all areas in coma, with delta waves predominating. If the EEG becomes isoelectric or flat for several hours, brain function is not recoverable and the coma may be

considered terminal. "Brain death" is indicated by a flat EEG, recorded at the highest gain with widely spaced electrode positions and the absence of cerebral reflexes and spontaneous respiration. *See DEATH.*

Cerebral evoked potentials. Computer advances in the analysis of EEG signals that are emitted by the brain during sensory stimulation and motor responses have led to the discovery and measurement of electrical waves known as event-related potentials or evoked potentials. These responses are averaged by a computer to enhance the small signals and increase the signal-to-noise ratio, so that they may be graphed and seen.

During the first 10 ms after a click is presented to the ear, electrical responses can be detected from the cochlea, the brainstem, and midbrain structures. Responses to 1000 clicks must be averaged to detect those small, rapid, earliest evoked responses.

This example of a far-field evoked potential is termed the brainstem auditory evoked response (BAER), usually derived from a vertex lead at the top of the scalp compared to a lead at one or the other earlobe as a reference electrode. Each ear is stimulated by clicks at a loudness level set at a certain number of decibels above the auditory threshold. The vertex electrode is far from the origin of the initial electrical activity coming from the cochlear nerve and the brainstem and the midbrain thalamic auditory nuclei, which are far-field or several inches from the active vertex electrode, which is recording field potentials generated by volume conduction rather than by cortical electrical activity from synaptic neural transmission pathways. Neuronal activity from the site of the recording electrode is termed near-field potentials. BAERs are clinically useful in diagnosing peripheral hearing loss due to damage to the structures of the ear, as distinct from hearing disorders involving the neural transmission pathways to the cortex by way of the brainstem and nuclei in the thalamus. Also, evoked potential audiometry has become an active area of clinical research for the evaluation of hearing ability in subjects who cannot respond voluntarily, such as very young children and language-impaired adults.

Somatosensory evoked responses (SERs) are far-field, fast low-voltage brainstem evoked responses derived from stimulating peripheral nerves, usually in the arms or legs, and recorded also by vertex electrodes referred to the earlobes as neutral sites. Like BAERs, the SERs are recorded to trace the somatosensory levels of neural functioning from the spinal cord to the brainstem, thalamus, and cortical areas. Abnormally long latencies or severely reduced amplitudes of the normal waves are diagnostic of neural damage at various levels of conduction of the neural impulses at various synaptic levels.

Fewer signals can be averaged to detect and measure the amplitudes and delays of the larger, slower signals that arrive at the sensory and association cortex 50–150 ms after sensory stimulation. Slower responses have been recorded from association cortex of the brain when psychological factors such as at-

tention, decision, and surprise are indicated. In general, fast potentials relate to the normal physiological functions of direct anatomical pathways, and slower potentials relate to the elaboration of significance and meaning of the responses mediated by nonspecific pathways and association cortex.

Cerebral evoked response patterns have been determined for the major senses under a variety of stimulating conditions. Repetitive electrical stimuli to the skin activate the sensory nerves. Strong stimuli to the eyes such as repeated bright flashes or reversing light and dark squares in a checkerboard pattern, which looks like a checkerboard moving back and forth from left to right and back again, evoke visual responses from the primary and association areas of the visual cortex (VERs). Sinusoidal variations in brightness several times per second are also utilized to evoke visual responses, synchronizing occipital-area brain rhythms, the alpha rhythm, to the stimulating light frequency. The clinical applications of VERs include estimations of visual acuity in nonresponsive subjects, including young children and adults who are not capable of giving verbal or voluntary responses. Visual and other evoked responses are valuable for the diagnosis of early multiple sclerosis, often before any other definite symptoms appear. Optic nerve atrophy or lesions in the visual pathway are detected by these methods.

Evoked responses have also been discovered in the motor systems of the brain preceding voluntary movement by similar averaging methods, using muscle movement outputs as the synchronizing signals to the computer, and tracing voltage changes prior to movements or even speech.

Psychological personality types or psychiatric illnesses may be indicated by modifications of these slow components of the evoked potentials. A negative slow shift, called the contingent negative variation (CNV), is seen when the subject associates a relationship of expectancy between two sensory stimuli when a response is required. If a warning signal is given 1 s before an imperative signal to which the subject must press a button, the contingent negative variation is seen between the two stimuli. The contingent negative variation has been found to be absent or disrupted in some behavioral disorders.

Both normal EEG waves and the fast and slow evoked responses are depressed on the side of a brain injury in one hemisphere compared to the normal hemisphere. Other computer techniques involve the power spectrum analysis of the EEG, such as the Fourier transform, and also cross-correlation between equivalent channels of the two hemispheres.

Evoked potentials are quantified by computer averaging, which enhances the signal-to-noise ratio by the square root of the number of trials. Various evoked potential components are described by the electrical sign, positive or negative, and by their latency in milliseconds to peak amplitude. Factor analytic methods, such as the principal components analysis, are also used by some investigators.

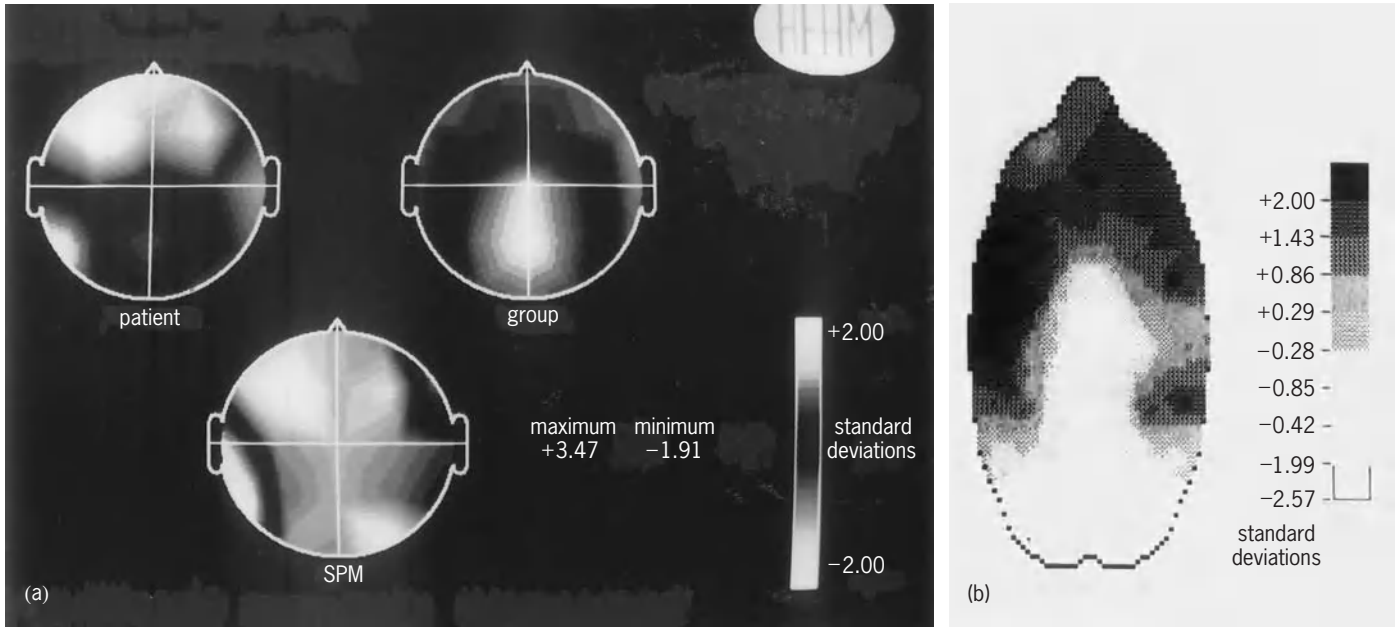


Fig. 2. Two examples of display conventions for topographic mapping. (a) Gray-scale technique, depicting amplitude of electrical discharge at particular points in time across the cortex. The maps display comparison of an individual's electroencephalogram with those of a control group, each normalized by subtracting the mean scores and dividing by the standard deviations based on values recorded at each electrode on 3 min of eyes-closed recording. Patient's map shows response to an auditory stimulus at time 156–192 ms after stimulus onset. Group map depicts the electrical response of a comparable control group to the same stimulus at the same point in time. Significance probability map depicts regional differences between patient and control group based upon standard deviation from mean amplitude (courtesy of F. H. Duffy). (b) Modification of the standard 10–20-electrode placement system, developed by M. Buchsbaum, using 32 electrodes and a four-point linear interpolation. Map shown is a top-down view, with a gray scale depicting brain activity within a specified frequency range.

Topographic mapping. The complexity of evoked potential and EEG analysis makes interpretation difficult in relation to where various components originate and their pattern of spread through time along the neural transmission pathways. Both of these considerations are especially important in determining the way that the brain processes cognitive information as well as sensory signals, that is, how the brain performs at work rather than at rest. In the early 1950s, scientists devised a method of displaying the distribution of EEG wave frequencies over the surface of the head. About 10 oscilloscope displays represented the various head positions and presented the dominant frequency distributions at those spatial locations. It was not until the 1980s, with the development of minicomputers and color graphics screens, that the presentation of topographic information could be analyzed in sophisticated statistical ways for research and clinical purposes by electroencephalographers and neurophysiologists. This method is best known as brain electrical activity mapping (BEAM) and is used in many research investigations of brain activity patterns in learning and language dysfunctions, psychiatric disorders, aging changes and dementia, and studies of normal and impaired child development. Difficult neurological diagnostic problems that do not show anatomical deformities by brain scan methods may often be clarified by these new electrographic procedures.

To create topographic images of electroencephalography, numeric values are assigned to fre-

quency bands, and these are translated into a gray scale or a rainbow scale. Data are collected from a number of electrodes, usually 20, and an interpolation algorithm is used to develop intermediate values between electrodes (Fig. 2).

Topographic images of evoked potentials can be made at millisecond intervals over a period of time and then viewed, like a movie or cartoon, as the computer displays these images on screen at the terminal. The spectral, temporal, and spatial analyses required in electroencephalography and evoked potential interpretation are performed by the computer. When reading the traditional electroencephalogram, the neurophysiologist may take advantage of these imaging techniques to support an interpretation. **Figure 3** provides a simple block diagram of the elements used in generating topographic images of brain electrical activity.

When viewing a topographic map, the clinician is still confronted with the critical question of whether the image (or electroencephalographic recordings associated with the image) represents the neurological activity of a normal, healthy individual or that of one with a pathological condition. A statistical technique called significance probability mapping (SPM) has been developed to aid in this discrimination. Significance probability mapping compares an individual's topographic image to an image based on neurological data of a comparable reference group, and generates a new image that depicts the deviation of the individual from the normal group.

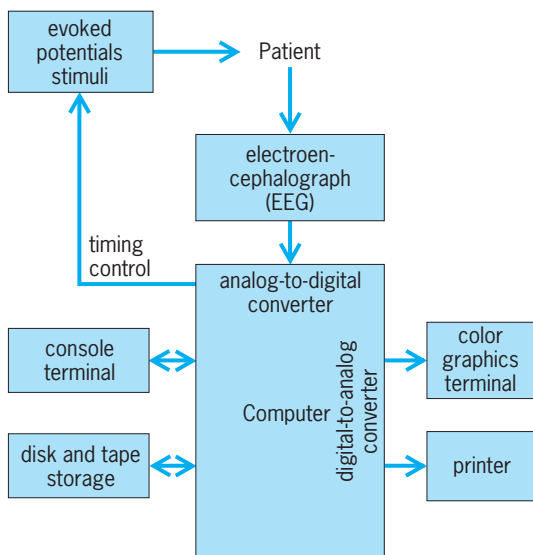


Fig. 3. Primary elements of equipment required for topographic mapping of brain electrical activity.

Such comparisons provide the final complicated step in clinical evaluation of electroencephalographic evoked potential data, and have proven extremely valuable in diagnosing neurologic abnormalities. Additionally, significance probability mapping techniques are valuable in research applications comparing two or more groups of subjects under a number of resting or activated (evoked) conditions. See BIOPOTENTIALS AND IONIC CURRENTS; BRAIN; NEUROBIOLOGY.

Jerome Cohen

Bibliography. F. H. Duffy, *Topographic Mapping of Brain Electrical Activity*, 1986; D. Giannitrapani and L. Murri (eds.), *The EEG of Mental Activities*, 1988; R. Johnson, Jr., et al. (eds.), *Current Trends in Event-Related Potential Research*, 1988; E. Niedermeyer, F. H. Lopes da Silva, and L. Fernando (eds.), *Electroencephalography: Basic Principles, Clinical Applications, and Related Fields*, 3d ed., 1993.

Electrokinetic phenomena

Phenomena associated with the movement of charged particles through a continuous medium or with the movement of a continuous medium over a charged surface. The four main electrokinetic phenomena are electrophoresis, electroosmosis, streaming potential, and sedimentation potential, or Dorn effect. These phenomena are related to one another through the zeta potential ζ of the electrical double layer that exists in the neighborhood of the charged surface.

Electrically charged layers. The distribution of electrolyte ions in the neighborhood of a negatively charged surface and the variation of potential ψ with distance from the surface are shown in Fig. 1. According to O. Stern, two different layers of ions are associated with the charged surface. The layer of ions immediately adjacent to the surface is called

the Stern layer. The ions of this layer are held to the charged surface by a combination of electrostatic attraction and specific adsorption forces, such as short-range van der Waals interactions and chemical bonds. The thickness δ of this layer is assumed to be equal to the ionic radius of the adsorbed ion species. The second layer of ions is the Gouy layer. The boundary between the two layers is the limiting Gouy plane. The ions in the Gouy layer are acted upon only by electrostatic forces and thermal motions of the liquid environment (Brownian motion), and they form a diffuse atmosphere of opposite charge (positive charge in Fig. 1) to the net charge at the limiting Gouy plane. The net charge density of the diffuse ion atmosphere of the Gouy layer decreases exponentially with distance from the limiting Gouy plane. The Gouy layer forms the positive half of an electrical double layer, and the charged surface plus the Stern layer form the negative half. The effective distance of separation $1/\kappa$ between the two halves of the double layer is determined by the concentration of electrolyte (ionic strength). For an electrolyte of univalent ions in water at 25°C (77°F), the relationship for $1/\kappa$ from the Debye-Hückel theory is Eq. (1), where c is the concentration of electrolyte (moles/liter).

$$\frac{1}{\kappa} = \frac{3 \times 10^{-8}}{\sqrt{c}} \quad (1)$$

Variation of potential ψ with the distance x from the charged surface is shown by a solid curve in Fig. 1. Here ψ_0 represents the thermodynamic reversible electrode potential which is independent of the properties of the electrical double layer and dependent only on the activity of the ion which is in reversible electrochemical equilibrium with the substance of the charged surface. The potential ψ

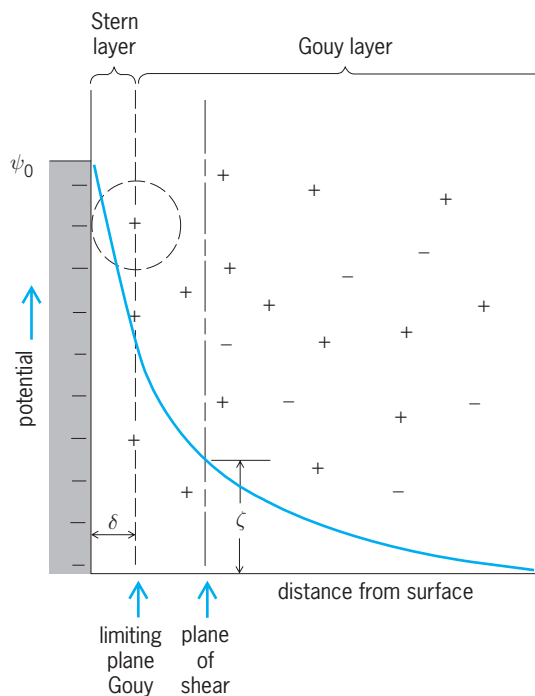


Fig. 1. Electrical double layer.

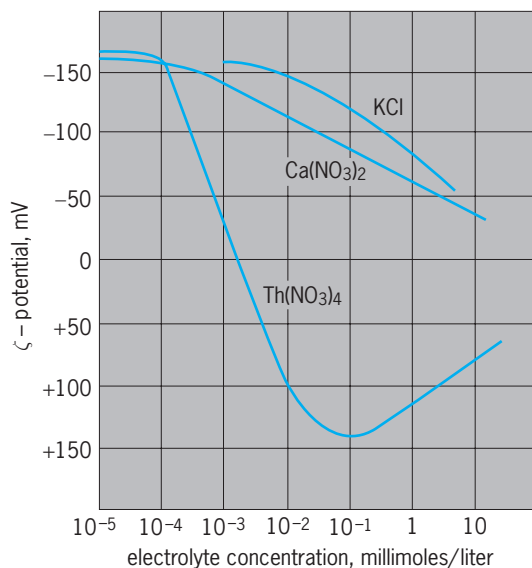


Fig. 2. Effect of electrolyte concentration on the ζ -potential of glass-water interfaces.

decreases linearly with increasing distance x in the region of the Stern layer. In the region of the Gouy layer, ψ decreases exponentially with increasing distance x .

Displacement of charged layers. In the four listed electrokinetic phenomena, a displacement occurs at some plane (plane of shear) between the charged surface and its atmosphere of ions. The position of the slipping plane in Fig. 1 is shown to be located in the Gouy layer. The potential of the plane of shear is the ζ -potential. From the theories of Gouy and Chapman, for spherical particles Eq. (2) holds. Here

$$\zeta = \frac{q}{Da} \frac{1}{1 + \kappa a} \quad (2)$$

$1/\kappa$ is the effective thickness of the double layer, q the net charge of the particle inside the plane of shear, D the dielectric constant of the liquid, and a the particle's radius at the plane of shear. For flat surfaces

$$\zeta = \frac{4\pi e}{D\kappa} \quad (3)$$

Eq. (3) holds where e is the charge per unit area of surface. Equations (2) and (3) show that ζ -potential is determined by the net charge at the plane of shear and $1/\kappa$, the effective thickness of the ion atmosphere. In turn, ζ -potential controls the rate of transport between the charged surface and the adjacent liquid. The relationship between rate of transport v_E and ζ -potential which is valid for all four electrokinetic phenomena is Eq. (4), where v_E is the velocity

$$v_E = \frac{D\zeta E}{4\pi\eta} \quad (4)$$

of the liquid at a large distance from the charged surface, E is the field strength (V/cm), and η is the viscosity of the liquid. The conditions for validity of Eq. (4) are that the double layer thickness ($1/\kappa$) must be small compared to the radius of curvature of the

surface; the substance of the surface must be nonconducting; and the surface conductance of the interface must be negligible. The equations which relate ζ -potential to electroosmotic flow rate and streaming potential may be obtained from Eq. (4) by use of Poiseuille's law for laminar flow through a capillary. For electrophoresis and sedimentation potential (Dorn effect), v_E is the velocity of the particles. E is the applied field strength for electrophoresis, whereas it is the gradient of potential developed by the sedimentation of charged particles in the Dorn effect.

Electrophoresis, electroosmosis, and streaming potential experiments have been shown to yield identical ζ -potentials for several different interfaces, particularly glass-water and protein-water systems. The sedimentation potential has not been significantly studied.

The effect of electrolytes on the ζ -potential of glass-water interfaces is shown in Figs. 2 and 3. As shown in Fig. 2, an increase in electrolyte concentration produces a decrease in ζ -potential, and ions of high charge of opposite sign to that of the surface can completely reverse the sign of the ζ -potential. The explanations for these two effects are given in Fig. 3, where the variation in ψ with distance from the surface is shown for low concentration of electrolyte in curve 1; moderate concentration of electrolyte in curve 2; and charge reversal by adsorption of ions (Th^{4+} on glass) in curve 3. Curves 1 and 2 show that an increase in electrolyte concentration reduces ζ -potential by reducing $1/\kappa$, as indicated by Eqs. (1), (2), and (3). Curve 3 shows that reversal of charge by ion adsorption occurs in the Stern layer and that

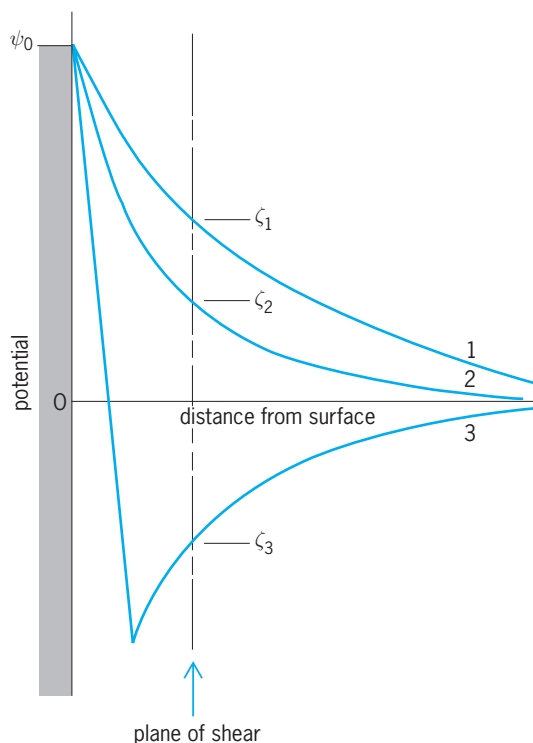


Fig. 3. Reversal of ζ -potential by ion adsorption.

this gives rise to a ζ -potential of opposite sign to the original value. See COLLOID; ELECTROPHORESIS; STREAMING POTENTIAL. Quentin Van Winkle

Electroless plating

A chemical reduction process which, once initiated, is autocatalytic. The process is similar to electroplating except that no outside current is needed. The metal ions are reduced by chemical agents in the plating solutions, and deposit on the substrate. Electroless plating is used for coating non-metallic parts. Decorative electroless plates are usually further coated with electrodeposited nickel and chromium. There are also applications for electroless deposits on metallic substrates, especially when irregularly shaped objects require a uniform coating. Electroless copper is used extensively for printed circuits, which are produced either by coating the non-metallic substrate with a very thin layer of electroless copper and electroplating to the desired thickness or by using the electroless process only. Electroless iron and cobalt have limited uses. Electroless gold is used for microcircuits and connections to solid-state components. Deeply recessed areas which are difficult to plate can be coated by the electroless process.

Nonmetallic surfaces and some metallic surfaces must be activated before electroless deposition can be initiated. Activation on nonmetals consists of the application of stannous and palladium chloride solutions. Once electroless plating is begun, it will continue to a desired thickness; that is, it is autocatalytic. The process thus differs from a displacement reaction, in which a more noble metal is deposited while a less noble one goes into solution; this ceases when the more noble deposit, if pore-free, covers the less noble substrate.

Electroless copper and gold deposits consist of very small crystals. Electroless nickel deposits are really highly supersaturated alloys containing phosphorus or boron, depending on the reducing agent used. These deposits, which are so fine-grained that they are almost amorphous, can be precipitation-hardened. Electroless nickel deposits are generally harder and more brittle than the electroplated variety. Corrosion resistance of electroless nickel depends, among other factors, on the uniformity of distribution of the phosphorus or boron. Adhesion to nonmetallic substrates is achieved primarily by mechanical means—by plating into pores which are created by selectively etching the substrate. An advantage of electroless plating over plating with current is the more uniform thickness of the surface coating. See ELECTROPLATING OF METALS. Rolf Weil

Electroluminescence

A general term for the luminescence excited by the application of an electric field to a system, usually in the solid state. Solid-state electroluminescent sys-

tems can be made quite thin, leading to applications in thin-panel area light sources and flat screens to replace cathode-ray tubes for electronic display and image formation. See LUMINESCENCE.

Destriau effect. Modern interest in electroluminescence dates from the discovery by G. Destriau in France in 1936 that when a zinc sulfide (ZnS) phosphor powder is suspended in an insulator (oil, plastic, or glass ceramic) and an intense alternating electric field is applied with capacitorlike electrodes, visible light is emitted. The phosphor, prepared from zinc sulfide by addition of a small amount of copper impurity, was later shown to contain particles of a copper sulfide (Cu_2S) phase in addition to copper in its normal role as a luminescence activator in the zinc sulfide lattice. The intensification of the applied electric field by the sharp conductive or semiconductive copper sulfide inhomogeneities is believed to underlie the mechanism of Destriau-type electroluminescence. Minority carriers are ejected from these high-field spots into the low- or moderate-field regions of the phosphor, where they recombine to excite the activator centers. This interpretation is supported by microscopic observations showing that the light is emitted from a few small spots in each phosphor grain. The structure of a Destriau-type electroluminescent cell is shown in Fig. 1; the light is observed through the transparent indium-tin oxide electrode.

The brightness L of lamps made with Destriau-type phosphors increases rapidly with the applied voltage V according to the equation below, where A and B

$$L = A \exp \frac{-B}{V^{1/2}}$$

are constants. The brightness increases slightly less than linearly with increase in frequency, until a saturation is reached at about 100,000 Hz. Maximum efficiency is not obtained under the same operating conditions as maximum brightness. Although flat-panel lamps were developed for general lighting application at 120 V with an active layer as thin as 1 mil (25 micrometers), the efficiency of these devices is low compared to conventional light sources; for example, electroluminescence, 3–5 lumens/watt; incandescent lamp, 16 lm/W; and fluorescent lamp,

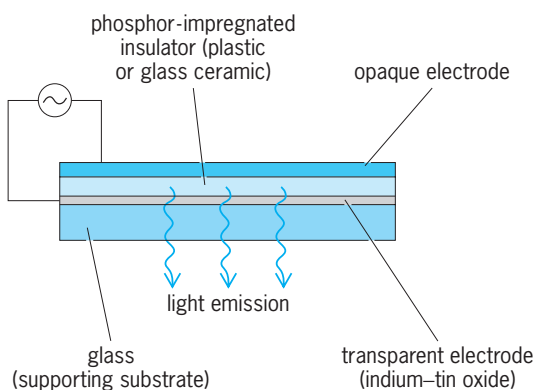


Fig. 1. Structure of powdered-phosphor (Destriau) electroluminescent cell, edge view (not to scale).

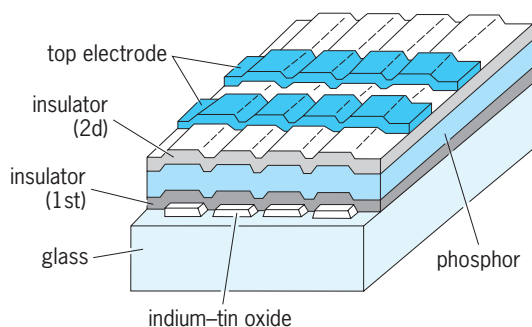


Fig. 2. Structure of a thin-film electroluminescent matrix display. (After J. Kuwata et al., *Feasibility of a $40 \times 40 \text{ cm}^2$ TFEL display with 2000×2000 lines operating at 50 W*, *Society for Information Display Digest*, pp. 297–300, 1988)

80 lm/W. Electroluminescent lamps are therefore not competitive for general illumination, but they have many specialized uses. See LIGHT PANEL.

Thin-film electroluminescence. The application of electroluminescence to display and image formation received great impetus from work in the late 1960s and mid-1970s on thin-film electroluminescence (TFEL), giving rise to devices that are different in structure and mechanism from the Desriau conditions. The phosphor in these devices is not a powder but a thin (about 500 nanometers) continuous film prepared by sputtering or vacuum evaporation. The luminescence activators are manganese or rare-earth ions, atomic species with internal electronic transitions that lead to characteristic luminescence. The phosphor film does not contain copper sulfide or any other separate phase, and is sandwiched between two thin (about 200 nm) transparent insulating films also prepared by evaporative means. Conducting electrodes are applied to the outside of each insulating film; one of the electrodes is again a transparent coating of indium-tin oxide on glass, which serves as supporting substrate. If an imaging matrix is desired, both electrodes consist of grids of parallel lines, with the direction of the grid on one insulator (row) orthogonal to the other grid (column). By approximate circuitry the entire matrix can be scanned, applying voltage where desired to a phosphor element that is located between the intersection of a row and column electrode (Fig. 2).

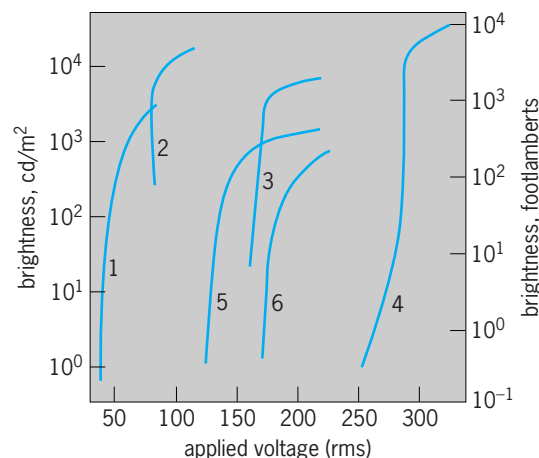
A thin-film electroluminescent device is thus a symmetrical structure (conductor-insulator-phosphor-insulator-conductor). It acts like a pure capacitor at low applied voltage; no light is emitted until the voltage reaches a threshold value determined by the dielectric properties of the insulator and phosphor films (Fig. 3). Above this threshold the device is no longer purely capacitive, a dissipative current flows, and light emission occurs. The brightness increases very steeply with the applied voltage but is finally saturated. The light output, or average brightness, is roughly proportional to the frequency up to at least 5 kHz, and also depends on the waveform of the applied voltage.

Thin-film electroluminescent devices are believed

to operate as follows. As voltage is applied, electrons trapped at the insulator-phosphor interface are released when the electric field reaches a critical value of about 5 MV/in. (2 MV/cm). The electrons are accelerated through the phosphor film, causing impact excitation of the activator centers, which then emit their characteristic luminescence. After the electrons complete their passage through the phosphor, they are retrapped at the second insulator-phosphor interface; on reversal of polarity the action is repeated and the electrons return to the original interface.

The best thin-film electroluminescent phosphor is manganese-activated zinc sulfide, which emits yellow light peaking at 585 nm. A brightness of approximately 10^4 footlamberts (3.5×10^4 candelas/m²) has been achieved with this phosphor at saturation under 5 kHz excitation, and it exhibits a brightness of 25 footlamberts (85 cd/m²) and an efficiency of 2 lm/W in a conventional thin-film electroluminescent matrix application at a frequency of 60 Hz. Activation of zinc sulfide and certain alkaline earth sulfides with different rare earths has yielded many other promising electroluminescent phosphors emitting blue, green, red, and white, and making full-color matrix-addressed thin-film electroluminescent displays possible. The light output of thin-film electroluminescent displays has been very reliable, with typically only 10% loss after tens of thousands of hours of operation.

Hysteresis in brightness versus voltage curves offers the possibility of using thin-film electroluminescence in memory devices. In conjunction with



Curve	Emission color	Activator	Insulator	Frequency (sine wave)
1	Yellow	Manganese	Lead titanate	5 kHz
2	Yellow	Manganese	Aluminum oxide	5 kHz
3	Yellow	Manganese	Aluminum oxide	10 kHz
4	Yellow	Manganese	Yttrium oxide	5 kHz
5	Green	Terbium	Yttrium oxide	5 kHz
6	Red	Samarium	Yttrium oxide	5 kHz

Fig. 3. Brightness-voltage characteristics of some zinc sulfide thin-film electroluminescent devices as a function of activator, insulating film, and operating frequency. (After R. Mach and G. O. Müller, *Physical concepts of high-field thin-film electroluminescent devices*, *Phys. Stat. Sol. (a)*, 69:11–66, 1982)

photoconductors, electroluminescence can also be employed in light amplifiers and logic circuits for computers, although semiconductor logic circuits are faster. See ELECTRONIC DISPLAY; LIGHT AMPLIFIER.

Injection electroluminescence. Injection electroluminescence results when a semiconductor *pn* junction or a point contact is biased in the forward direction. This type of emission was first observed from silicon carbide (SiC) in England by H. J. Round in 1907. It is the result of radiative recombination of injected minority carriers with majority carriers in material. Such emission has been observed in a large number of semiconductors, including silicon (Si), germanium (Ge), diamond, cadmium sulfide (CdS), zinc sulfide (ZnS), zinc selenide (ZnSe), zinc telluride (ZnTe), zinc oxide (ZnO), and many of the so-called III-V compounds, such as gallium phosphide (GaP), gallium arsenide (GaAs), indium phosphide (InSb), boron phosphide (BP), boron nitride (BN), and aluminum nitride (AlN). The wavelength of the emission corresponds to an energy equal, at most, to the forbidden band gap of the material, and hence in most of these materials the wavelength is in the infrared region of the spectrum. In some cases the efficiency is high, approaching one emitted photon. In suitable structures the excitation intensity can become so high that stimulated rather than spontaneous emission predominates, and laser action results, with spectral narrowing and coherent emission. Direct electrical excitation is more convenient than optical excitation, but the beam divergence of injection lasers is normally much greater than that of gas lasers or the optically pumped solid type. See JUNCTION DIODE; LASER; SEMICONDUCTOR; SEMICONDUCTOR DIODE.

Other effects. If a *pn* junction is biased in the reverse direction, so as to produce high internal electric fields, other types of emission can occur, but with very low efficiency. The presence of very energetic ("hot") carriers can result in emission at energies greater than the band gap of the material (avalanche emission). The emission in this case may be correlated with small active regions called microplasmas. Light emission may also occur when electrodes of certain metals, such as Al or Ta, are immersed in suitable electrolytes and current is passed between them. In many cases this galvanoluminescence is electroluminescence generated in a thin oxide layer formed on the electrode by electrolytic action. In addition to electroluminescence proper, other interesting effects (usually termed electrophotoluminescence) occur when electric fields are applied to a phosphor which is concurrently, or has been previously, excited by other means. These effects include a decrease or increase in steady-state photoluminescence brightness when the field is applied, or a burst of afterglow emission if the field is applied after the primary photoexcitation is removed. See PHOTOLUMINESCENCE.

James H. Schulman; Clifford C. Klick

Bibliography. H. K. Henisch, *Electroluminescence*, 1962; H. F. Ivey, *Electroluminescence and Related*

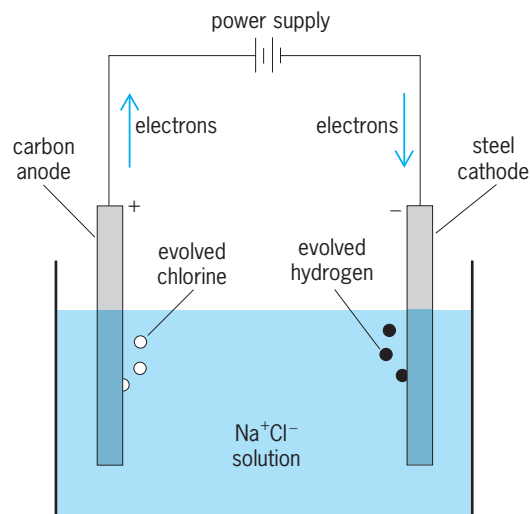
Effects, 1963; A. H. Kitai (ed.), *Solid State Luminescence: Theory, Materials, and Devices*, 1993; J. I. Pankove (ed.), *Electroluminescence*, 1977; S. Shionoya (ed.), *Electroluminescence*, 1989.

Electrolysis

A means of producing chemical changes through reactions at electrodes in contact with an electrolyte by the passage of an electric current. Electrolysis cells, also known as electrochemical cells, generally consist of two electrodes connected to an external source of electricity (a power supply or battery) and immersed in a liquid that can conduct electricity through the movement of ions. Reactions occur at both electrode-solution interfaces because of the flow of electrons. Reduction reactions, where substances add electrons, occur at the electrode called the cathode; oxidation reactions, where species lose electrons, occur at the other electrode, the anode. In the cell shown in the **illustration**, water is reduced at the cathode to produce hydrogen gas and hydroxide ion; chloride ion is oxidized at the anode to generate chlorine gas. Electrodes are typically constructed of metals (such as platinum or steel) or carbon. Electrolytes usually consist of salts dissolved in either water or a nonaqueous solvent, or they are molten salts. See ELECTROCHEMISTRY; ELECTRODE; ELECTROLYTE; OXIDATION-REDUCTION.

Applications include industrial synthesis of chemicals, electroplating of metals, metallurgical extraction and refining of metals, metal finishing and electromachining, and the production of electricity in batteries. Corrosion of metals often occurs through electrolytic processes. Electrolytic cells are used in analytical chemistry and in laboratory studies of reaction mechanisms. See BATTERY; CORROSION.

Principles. The flow of current, measured in amperes (where 1 A equals the passage of 1 coulomb of charge per second), represents the rate of flow of



Schematic diagram of an electrolysis cell in which the electrolyte is a solution of sodium chloride.

electrical charge through the electrolysis cell. The quantity of a substance produced or consumed in an electrode reaction is proportional to the quantity of electricity (coulombs) passed during the electrolysis: a constant current of 1 A passed for 1 h (equivalent to 3600 coulombs) will produce 0.018656 mole or 1.3228 g of chlorine gas. The quantity of electricity is the integral of the current over the duration of the electrolysis and can be determined with a coulometer. *See* COULOMETER; ELECTROCHEMICAL EQUIVALENT.

In bulk electrolysis, for example, in the production of chemicals, the flow of current produces appreciable changes in the concentration of species in the electrolyte. Cells for bulk electrolysis usually employ relatively large electrodes with the liquid kept in constant motion. For analytical applications or the characterization of chemical systems, the electrolysis occurs only near the surface of the electrodes. Here, unstirred solutions and small electrodes are typically used.

When only one reaction occurs at an electrode, it is said to occur with 100% current efficiency. When two or more reactions occur at the same electrode, for example, when both oxygen and chlorine are evolved at the anode of a cell, the current efficiency of each reaction is given by the fraction of the quantity of electricity consumed for each process. For large-scale electrolytic processes, the power consumption (or energy) efficiency is of interest because it is related to the cost of electrical energy needed to produce a given amount of product. The energy efficiency is the ratio of the ideal theoretical energy for the cell reaction to the actual energy (directly proportional to the cell voltage) consumed in the electrolysis. The cell voltage depends upon a number of factors, including the thermodynamic energy requirements for the electrode reactions, the energy needed to drive the reactions at a given rate, and the resistance of the electrodes and electrolyte. *See* OVERVOLTAGE.

In analytical applications the electrode reaction at only one of the electrodes is of interest, and a third electrode, called a reference electrode, is usually introduced into the cell so that the potential of that electrode can be determined with respect to a known reference potential. The measured potential is also less perturbed by solution resistance effects in a three-electrode cell. *See* REFERENCE ELECTRODE.

Applications. There are many industrial applications for the production of important inorganic chemicals. Chlorine and alkali are produced by the large-scale electrolysis of brine (the chloralkali process) in cells carrying out the same reactions as those shown in the illustration. Other chemicals produced include hydrogen and oxygen (via water electrolysis), chlorates, peroxydisulfate, and permanganate. Other processes are carried out with molten salts as solvents, because the electrode reactions of interest would be masked in aqueous solutions by the electrolysis of water to produce hydrogen and oxygen. Electrowinning of aluminum, magnesium, and

sodium metals is carried out by molten salt electrolysis. Fluorine is produced by electrolysis of a 2:1 mixture of anhydrous hydrofluoric acid and potassium fluoride.

The major electrolytic processes involving organic compounds are the hydrodimerization of acrylonitrile to produce adiponitrile and the production of tetraethyllead. Many other organic compounds have been studied on the laboratory scale.

Electroplating involves the electrochemical deposition of a thin layer of metal on a conductive substrate, for example, to produce a more attractive or corrosion-resistant surface. Chromium, nickel, tin, copper, zinc, cadmium, lead, silver, gold, and platinum are the most frequently electroplated metals. Metal surfaces can also be electrolytically oxidized (anodized) to form protective oxide layers. This surface-finishing technique is most widely used for aluminum but is also used for titanium, copper, and steel. Electrolytic capacitor production involves anodization of aluminum, tantalum, and niobium. Anodic dissolution of a metal is used in electrochemical machining to produce a desired structure. It is mainly applied in the machining of very hard alloys or in producing complex structures that would be difficult to make by conventional machining methods. *See* ELECTROPLATING OF METALS.

Metals can be purified by electrorefining. Here, the impure metal is used as the anode, which dissolves during the electrolysis. The metal is plated, in purer form, on the cathode. Copper, nickel, cobalt, lead, and tin are all purified by this technique. *See* ELECTROMETALLURGY.

Electroanalysis involves the use of electrolytic processes to identify and quantitate a species. Coulometric methods are based on measuring the quantity of electricity used for a desired process. Voltammetric methods allow characterization of species through an analysis of the effect of potential and electrolysis conditions on the observed currents. Interest in photo-electrolysis is growing; this involves the utilization of radiant (for example, solar) energy to produce electricity or to drive chemical reactions in electrolytic cells. In such cells, which have not yet found practical application, the irradiated electrodes are usually made of semiconductor materials, for example, titanium dioxide, gallium arsenide, or cadmium selenide. Electrolytic cells can also be used to purify waste streams, for example, by the cathodic deposition and recovery of metals and the oxidation of organic pollutants. *See* POLAROGRAPHIC ANALYSIS.

Allen J. Bard

Bibliography. A. J. Bard and H. Lund (eds.), *The Encyclopedia of the Electrochemistry of the Elements*, 1973–1986; A. J. Bard and L. R. Faulkner, *Electrochemical Methods: Fundamentals and Applications*, 2d ed., 2000; D. R. Crow, *Principles and Applications of Electrochemistry*, 4th ed., 1994; H. Lund and O. Hammerich (eds.), *Organic Electrochemistry*, 4th ed., 2000; D. Pletcher and F. C. Walsh, *Industrial Electrochemistry*, 2d ed., 1990; M. Schlesinger and M. Paunovic (eds.), *Modern Electroplating*, 4th ed., 2000.

Electrolyte

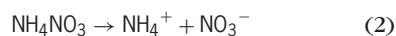
A material that conducts an electric current when it is fused or dissolved in a solvent, usually water. Electrolytes are composed of positively charged species, called cations, and negatively charged species, called anions. For example, sodium chloride (NaCl) is an electrolyte composed of sodium cations (Na^+) and chlorine anions (Cl^-). The ratio of cations to anions is always such that the substance is electrically neutral. If two wires connected to a light bulb and to a power source are placed in a beaker of water, the light bulb will not glow. If an electrolyte, such as sodium chloride, is dissolved in the water, the light bulb will glow because the solution can now conduct electricity. The amount of electric current that can be carried by an electrolyte solution is proportional to the number of ions dissolved. Thus, the bulb will glow more brightly if the amount of sodium chloride in the solution is increased. See ELECTRIC CURRENT; ION.

Hydration. Water is a special solvent because its structure has two different sides. On one side is the oxygen atom, and on the other are two hydrogen atoms. Covalent molecules, such as water, are held together by covalent bonds, which are formed when two atoms share a single pair of electrons. However, when two different atoms form a covalent bond, the sharing of electrons is not always equal. An electron in a covalent bond between two different atoms might spend more time near one atom or the other. In water, the electrons from the O-H bonds spend more time near the oxygen atom than near the hydrogen atoms. As a result, the oxygen has somewhat more negative charge than the hydrogen atoms. This phenomenon is not the same as ionization, where the electron is completely transferred from one atom to the other; so to indicate the difference, the oxygen is said to have a partial negative charge (δ^-) and the hydrogen to have a partial positive charge (δ^+). See CHEMICAL BONDING.

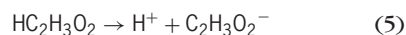
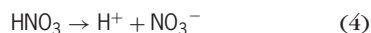
Since water has two different sides with opposite partial charges, it is said that water is polar. When ions are dissolved in polar solvents such as water, the solvent molecules are organized around the ion to create the maximum amount of electrostatic attraction. Thus, if a cation is dissolved in water, the water molecules will organize themselves around the ion with the oxygen atoms directed toward the cation, because the partial negative charges on the oxygen atoms in the water molecules are electrostatically attracted to the positive charge on the cation. Likewise, water molecules will organize about anions with the partially positive hydrogen atoms directed toward the negatively charged anion. When water molecules surround dissolved ions, the ions are said to be hydrated. The electrostatic attraction associated with hydration provides the energetic driving force for dissolving ions. Since water is polar, it will also dissolve polar molecules that are not charged, such as ethanol ($\text{CH}_3\text{CH}_2\text{OH}$), which contains a polar O-H bond. Thus, water can dissolve both ionic sub-

stances and polar molecules that are not charged. See SOLVENT; WATER.

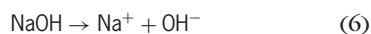
Strong and weak electrolytes. Any substance that produces ions when dissolved is an electrolyte. These substances include ionic materials composed of simple monatomic ions, such as sodium chloride, or substances composed of polyatomic ions, such as ammonium nitrate (NH_4NO_3). When these substances are dissolved, hydrated ions are generated, as in reactions (1) and (2).



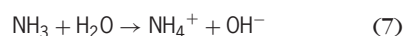
A special type of electrolyte is an acid, in which the cation is H^+ . When acids are dissolved in water, H^+ ions are produced along with an anion, which can be either a monatomic ion, as in hydrochloric acid (HCl), or a polyatomic ion, as in nitric acid (HNO_3) or acetic acid ($\text{HC}_2\text{H}_3\text{O}_2$), as in reactions (3)–(5).



Electrolytes such as sodium hydroxide (NaOH) that yield the hydroxyl (OH^-) anion when dissolved in water are called bases [reaction (6)]. Some



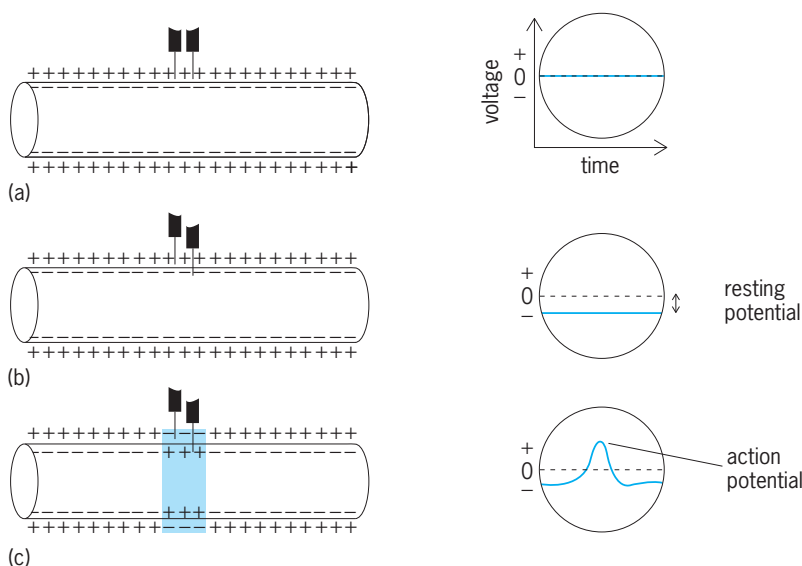
molecules that do not contain ions, such as ammonia (NH_3), generate OH^- ions when dissolved in water as in reaction (7); these bases are also electrolytes.



Polar covalent molecules, such as ethanol, dissolve in water but do not generate any ions and are called nonelectrolytes.

Many electrolytic substances completely dissociate into ions when they are dissolved in water. In these cases, the reactions shown above would proceed completely to the right, leaving only dissolved ions and no associated, electrically neutral molecules. For example, when sodium chloride is dissolved in water, all of the dissolved material is present as Na^+ and Cl^- ions, with no dissolved NaCl molecules. Substances that completely dissociate are called strong electrolytes, because every molecule dissolved generates ions that contribute to the electrical conductivity. Ionic substances, such as NaCl and NH_4NO_3 , are strong electrolytes. Acids and bases that completely dissociate are called strong acids and strong bases, and these substances are also strong electrolytes. Hydrochloric acid (HCl) and HNO_3 are strong acids, so every molecule of HCl or HNO_3 that dissolves generates a H^+ cation and a Cl^- or NO_3^- anion, respectively.

Some substances dissolve in water but do not dissociate completely. For example, when acetic acid is dissolved in water, some molecules dissociate to form H^+ and $\text{C}_2\text{H}_3\text{O}_2^-$ ions, while others remain



Measurement of the electric potential across a neuronal membrane by using an oscilloscope. (a) With both electrodes outside the membrane, no potential is measured. **(b)** With one electrode inside the membrane and the neuron in the resting state, the resting potential is measured. **(c)** When a nerve impulse travels along the neuron, the charge is reversed, and an action potential that is the opposite of the resting potential is measured for a brief period. The shaded region represents the portion where the polarity is reversed. (After H. Curtis and N. S. Barnes, *Biology*, 5th ed., Worth, 1989)

associated as $\text{HC}_2\text{H}_3\text{O}_2$ units, which contain polar O-H bonds and are therefore readily soluble in water. Acids such as acetic acid that do not dissociate completely are called weak acids. Similarly, the reaction of ammonia with water in reaction (7) proceeds to only a small extent. Thus, some molecules of ammonia react with water to generate OH^- and NH_4^+ ions, but many simply remain as NH_3 . Since each dissolved molecule of ammonia does not generate an OH^- anion, ammonia is said to be a weak base. The electrical conductivity of an electrolyte solution depends on the number of ions present, so a solution of acetic acid or ammonia will not conduct as well as a solution of an equal amount of a strong electrolyte. Thus, when a light bulb is connected to a solution of acetic acid, it will not glow as brightly as a light bulb connected to a solution of sodium chloride. In fact, when 100 acetic acid molecules are dissolved in water, 99 exist as associated $\text{HC}_2\text{H}_3\text{O}_2$ molecules, while only one dissociates to form H^+ and $\text{C}_2\text{H}_3\text{O}_2^-$ ions. Therefore, in order to get an equal amount of electrical conductivity as a solution of a strong electrolyte, 100 times more acetic acid must be dissolved. For this reason, substances that dissociate somewhat but not completely are called weak electrolytes. Weak acids and weak bases are the most common types of weak electrolytes. See IONIC EQUILIBRIUM.

Physiological roles. In biological systems, electrolytes play important roles in regulating kidney function and the retention of water. Electrolytes are also vital for providing the electric current needed for nerve impulses in neurons. When the ratio of cations to anions is different on two sides of a cell membrane, the electric charge on either side of the membrane is different, leading to the establishment

of an electric potential. The cell membranes of neurons are impermeable to anions, but they contain special proteins that maintain a higher concentration of cations outside the neuron by either actively transporting sodium (Na^+) and potassium (K^+) ions into and out of the cell or by changing the permeability of the membrane to these cations. Since anions cannot follow the cations out of the neuron, a negative charge builds up inside the neuron and a positive charge builds up outside the neuron, generating an electric potential. When no impulse is being transmitted by the neuron, it is said to be in its resting state and the resulting electric potential is called the resting potential. If two electrodes are placed outside the neuron in its resting state, no potential is measured (illus. a); however, if one electrode is inside the neuron and the other electrode is outside, the resting potential can be measured (illus. b). See ELECTRODE; OSMOREGULATORY MECHANISMS.

When a nerve impulse is generated, a localized section of the neuron suddenly becomes permeable to Na^+ cations. Since the concentration of cations is higher outside the neuron, Na^+ ions rapidly flow into the neuron because of attraction to the excess negative charge inside and the concentration gradient. An excess of positive charge is built up inside the neuron, and the electrical charge is reversed. As a result, a new action potential that is opposite of the resting potential is generated (illus. c). When the Na^+ ions enter the neuron, the membrane becomes permeable to K^+ cations, which rapidly flow out of the neuron, restoring the resting potential. The action potential therefore persists for only a half millisecond. This process is repeated at adjacent sites in a single direction along the neuron, propagating the nerve impulse. See BIOPOTENTIALS AND IONIC CURRENTS; SYNAPTIC TRANSMISSION. H. Holden Thorp

Bibliography. A. J. Bard and L. R. Faulkner, *Electrochemical Methods: Fundamentals and Applications*, 2000; R. Chang, *Chemistry*, 1991; H. Curtis and N. S. Barnes, *Biology*, 1989; L. Stryer, *Biochemistry*, 4th ed., 1995; G. Wulfsberg, *Principles of Descriptive Inorganic Chemistry*, 1992.

Electrolytic conductance

The transport of electric charges, under electric potential differences, by particles of atomic or larger size. This phenomenon is distinguished from metallic conductance, which is due to the movement of electrons. The charged particles that carry the electricity are called ions.

Positively charged ions are termed cations; the sodium ion, Na^+ , is an example. The negatively charged chloride ion, Cl^- , is typical of anions. The negative charges are identical with those of electrons or integral multiples thereof. The unit positive charges have the same magnitude as those of electrons but are of opposite sign. Colloidal particles, which may have relatively large weights, may be ions, and may carry many positive or negative

charges. Electrolytic conductors may be solids, liquids, or gases. Semiconductors have properties that are intermediate between the metallic conductors and insulators.

Measurement. Conductances are usually reported as specific conductances κ , which are the reciprocals of the resistances of cubes of the materials, 1 cm in each dimension, placed between electrodes 1 cm square, on opposite sides. These units are sometimes called mhos, that is, ohms spelled backward. Conductances of solutions are usually measured by Friedrich Kohlrausch's method, in which a Wheatstone bridge is employed. Such a bridge is shown diagrammatically in Fig. 1. The resistances R_3 and R_4 (usually of the same value) form two arms of the bridge. Resistance R_2 is adjustable, and the remaining arm is the cell holding the electrolytic conductor, or as is usually stated, solution of electrolyte. Direct current and the usual galvanometers cannot be used because of an apparent failure of Ohm's law. Passage of direct current produces chemical reactions and a back electromotive force (emf) is generated by the galvanic action of the products. By using an alternating current, the electrochemical reactions occurring when the current is briefly passed in one direction may be reversed when the direction of the current is changed. When a small alternating-current input signal is used, practically all the electric charge passed during each half cycle is stored in the electric double layer, which acts as a capacitor. The electrodes are usually made of platinum and are platinized, that is, coated with finely divided platinum. The surface area, and hence the electrode capacitance, is thereby greatly increased. By making measurements at several frequencies and extrapolating to infinite frequency, the effect of electrode reactions can be eliminated. For less exact measurements, a fixed frequency of 60-1000 Hz is commonly used.

To determine the conductance C , that is, the reciprocal resistance $1/R$ of the cell of Fig. 1, the resistance R_2 of the bridge is adjusted until a minimum of sound is heard in the telephone.

Greater sensitivity may be obtained by electronic amplification of the off-balance signal and by using an oscilloscope to detect the point of balance. When the bridge is in balance, the conductance is given by the relation $C = R_4/R_2R_3$. From this the specific conductance κ may be obtained from the equation $\kappa = KC$, in which K is the cell constant. Occasionally,

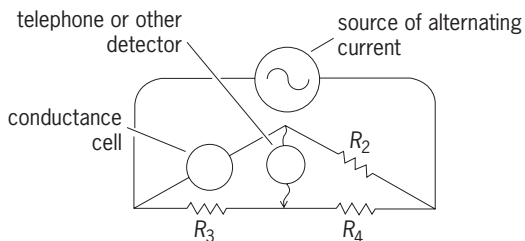


Fig. 1. Wheatstone bridge circuit for the measurement of electrolytic conductance. R_3 and R_4 are fixed resistances; R_2 is a variable resistance.

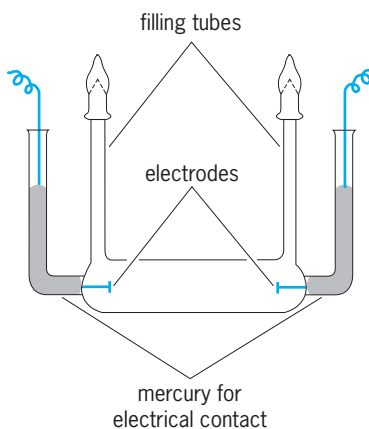


Fig. 2. Diagram of conductance cell.

this constant can be computed from the dimensions of the cell. Usually, however, it is determined by using a solution whose κ value is accurately known from measurements in such a cell, or by comparison with the specific conductance of mercury.

For precision work, care must be taken to avoid errors due to electrical reactances. This has been done in specially designed bridges. A typical, properly designed conductance cell is shown in Fig. 2. The cell is filled with solution through the center tubes. Electrical contact is made with the electrodes by platinum wires sealed through the glass wall. These connect the mercury in the outside tubes, which are widely spread to avoid errors due to electrical capacity.

Equivalent conductance. Although many substances and mixtures show electrolytic conductance, the greater part of the research on the subject has been on aqueous solutions of salts, acids, and bases. There have been considerable data accumulated for solutions of such electrolytes in nonaqueous solvents, such as alcohols. The data are usually given in terms of equivalent conductance Λ , which is defined by Eq. (1), in which κ is the

$$\Lambda = \frac{1000\kappa}{c} \quad (1)$$

specific conductance and c is the concentration in equivalents per liter. Values of Λ change with the concentration and, in general, increase as the solutions measured are made more dilute, that is, as c is decreased. A plot of values of the equivalent conductance λ against \sqrt{c} for some typical electrolytes is shown in Fig. 3. Svante Arrhenius, who was the first to assume that electrolytic conductance is due to freely moving charged ions, explained the decrease of Λ with increasing c by assuming that the number of ionic carriers gets smaller as the concentration increases, and he computed a degree of dissociation α by formula (2). The term Λ_0 is obtained by

$$\alpha = \frac{\Lambda}{\Lambda_0} \quad (2)$$

determining Λ at a series of low concentrations and extrapolating to a limiting value, termed the equivalent conductance at infinite dilution. Though

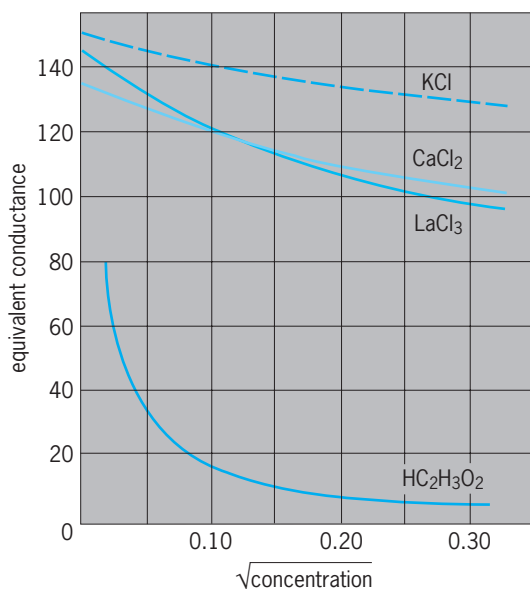


Fig. 3. Equivalent conductance at infinite dilution for some typical electrolytes.

Eq. (2) has been shown by later work to give nearly the right values of α for certain poorly conducting solutions, it is now considered to be much in error for the so-called strong electrolytes. These include most salts, such as potassium chloride, KCl, and sodium sulfate, Na_2SO_4 , and inorganic acids and bases, such as hydrochloric acid, HCl, and sodium hydroxide, NaOH. For an electrolyte which yields two types of ion, it can be shown that Eq. (3) holds,

$$\Lambda = F\alpha (U^+ + U^-) \quad (3)$$

in which F is the faraday and U^+ and U^- are the mobilities, or speeds, under unit potential difference of the positive and negative ions, respectively. For Eq. (2) to hold, these mobilities must be constant from the equivalent concentration c at which Λ is measured to infinite dilution.

Since the advent of the Debye-Hückel theory of interionic attractions, strong electrolytes have been considered to be completely dissociated, that is, the term α of Eq. (3) is equal to unity for these substances. The decreases observed in the values of the equivalent conductances Λ , with increases in concentration, are assumed to be due to reductions in the values of the ionic mobilities U^+ and U^- . According to the theory of P. Debye and E. Hückel, the ion possesses an ionic atmosphere distributed with radial symmetry around the ion as center. This is due to the fact that interionic attractions and repulsions, together with thermal vibrations, tend to produce a slight preponderance of negative ions around a positive ion, and vice versa. The presence of this atmosphere leads to the lowering of ionic mobilities with increasing ion concentrations. The adaptation of the Debye-Hückel theory for conducting solutions is due to Lars Onsager. His equation for very dilute uni-univalent electrolytes, such as sodium chloride, is shown in Eq. (4), in which θ and σ are given for uni-

univalent electrolytes by Eqs. (5) and (6), in which D

$$\Lambda = \Lambda_0 - (\theta \Lambda_0 + \sigma) \sqrt{c} \quad (4)$$

$$\theta = \frac{8.16 \times 10^5}{(DT)^{3/2}} \quad (5)$$

$$\sigma = \frac{8.28}{\eta(DT)^{1/2}} \quad (6)$$

is the dielectric constant at the absolute temperature T and η is the viscosity.

Equation (4) yields accurate values of the data for strong electrolytes up to concentrations of about 0.001 M , above which there are small deviations. Modifications of Eq. (4) for solutions of salts of higher valence types, such as calcium chloride, CaCl_2 , and lanthanum chloride, LaCl_3 , are available and have also been found to agree with the data for dilute solutions.

Onsager, in his derivation of Eq. (4), treated the ions as point charges. Later Raymond Fuoss and Onsager extended the theory to include the radii of the ions and also the effects of higher concentrations. The ion sizes obtained from conductance data agree closely with those calculated from activity measurements. See ACTIVITY (THERMODYNAMICS).

Equation (4), or its empirical and theoretical extensions, can be used to obtain values of Λ_0 from the data on equivalent conductances of dilute solutions. Some typical figures for limiting equivalent conductances Λ_0 of typical strong electrolytes in aqueous solution at 25°C (77°F) are listed below.

HCl	426.16	NH_4Cl	149.86	MgCl_2	129.40
NaOH	247.8	KNO_3	144.96	NaI	126.94
KBr	151.9	CaCl_2	135.84	NaCl	126.45
KI	150.38	Na_2SO_4	129.9	LiCl	115.03

Ionic conductances. The fraction of the current carried by a given type of ion is known as the transference number, t , which depends on the relative mobilities of the ions. This relation is shown in Eqs. (7). The ionic mobilities, and therefore the trans-

$$t^+ = \frac{U^+}{U^+ + U^-} \quad t^- = \frac{U^-}{U^+ + U^-} \quad (7)$$

ference numbers, vary somewhat with concentration, but the limiting equivalent conductances Λ_0 for various electrolytes can be accurately represented as the sum of the limiting ionic conductances. Thus, for potassium chloride, $\Lambda_{0,\text{KCl}} = \lambda_{0,\text{K}^+} + \lambda_{0,\text{Cl}^-}$, and the value of λ_{0,Cl^-} is the same whether it is derived from measurements on HCl, NaCl, or KCl solutions. This additive relation is known as Kohlrausch's law of the independent mobility of ions. However, it is necessary to obtain the value of λ_0 for at least one ion constituent independently in order to establish the ion conductances of the other ions. The relation used is shown in Eqs. (8), in which λ_0 is the

$$t_0^+ \Lambda_0 = \lambda_0^+ \quad \text{or} \quad t_0^- \Lambda_0 = \lambda_0^- \quad (8)$$

limiting equivalent conductance of an electrolyte and t_0^+ and t_0^- are the limiting transference numbers of the positive and negative ion constituents, respectively.

The same value of λ_{0,Cl^-} , within 0.02%, is obtained from precision conductance and transference measurements on solutions of hydrogen, lithium, sodium, and potassium chlorides. Values of the limiting ionic conductance at 25°C (77°F) are given below for some ions.

H ⁺	349.82	1/2Ba ²⁺	63.64
OH ⁻	198	Ag ⁺	61.92
SO ₄ ²⁻	79.8	1/2Ca ²⁺	59.50
Br ⁻	78.4	1/2Mg ²⁺	53.06
I ⁻	76.8	Na ⁺	50.11
Cl ⁻	76.34	HCO ₃ ⁻	44.48
K ⁺	73.52	CH ₃ CO ₂ ⁻	40.9
NH ₄ ⁺	73.4	CH ₂ ClCO ₂ ⁻	38.7
NO ₃ ⁻	71.44	Li ⁺	38.69

Nonaqueous systems. In addition to the study of water solutions of electrolytes, considerable study has been given to electrolytes in nonaqueous and mixed solvents. In general, the same principles as those outlined above apply to the interpretation of the results. However, fewer of the electrolytes are completely dissociated, and the degrees of dissociation of the weaker acids and bases are lower. This is due to the fact that, in general, the dielectric constants of nonaqueous solvents are smaller than those of water, so that the attractions between positive and negative ions are greater.

It will be observed that in this article discussion is confined to quite dilute solutions of electrolytes. For concentrated solutions few generalizations of any value can be given.

Molten salts exhibit a wide range of conductivities, depending upon their structures. Salts of alkali and alkaline earth metals usually are largely ionic in character and are highly conductive in the molten state, whereas heavy metal salts may be essentially covalent and exhibit little or no conductivity. Thus the conductivities of liquid arsenic chloride (AsCl₃) and bismuth chloride (BiCl₃) near their melting points are approximately 10⁻⁶ and 0.44 ohm⁻¹ cm⁻¹, respectively, reflecting the more ionic structure of BiCl₃.

If, instead of using quite low potentials in the measurement of electrolytic conductances, voltages of the order of 100,000 are employed, the conductances observed are no longer constant but tend to increase with the potential used. Under these conditions Ohm's law evidently is not valid. This increase of conductance with high potentials is called the Wien effect. This effect is in accord with the interionic attraction theory. When the velocity of the ions becomes sufficiently great, the ion atmospheres do not have time to form to their full extent, so that both the electrophoretic and time of relaxation effects exert less influence on the conductance. How-

ever, a large Wien effect is also found for weak acids and bases. It would appear that the high potentials produce, temporarily, additional ionization of these substances. This explanation has been proposed and discussed theoretically by Onsager. If very high frequencies are used in the measurements, an increase in the conductance, termed the Debye-Falkenhagen effect, is observed. This can also be explained by the interionic attraction theory.

Solids. Electrical conduction in solids can range from purely electronic (in metals and semiconductors) to purely ionic (in solids that are electronic insulators but have ions with appreciable mobilities). Examples of purely ionic solid conductors are silver iodide and silver sulfide, in which all of the current is carried by silver ions. In other cases, such as metallic oxides, both ions may serve as charge carriers. Ionic conduction must lead to transport of matter and to chemical reactions at the electrodes, whereas electronic conduction leads only to a transfer of charge. See ELECTROCHEMISTRY; ELECTROMOTIVE FORCE (CELLS).
Herbert A. Laitinen

Bibliography. J. L. Copeland, *Transport Properties of Ionic Liquids*, 1974; D. Inman and D. G. Lovering (eds.), *Ionic Liquids*, 1981; R. F. Snipes, *Statistical Mechanical Theory of the Electrolytic Transport of Non-Electrolytes*, 1974.

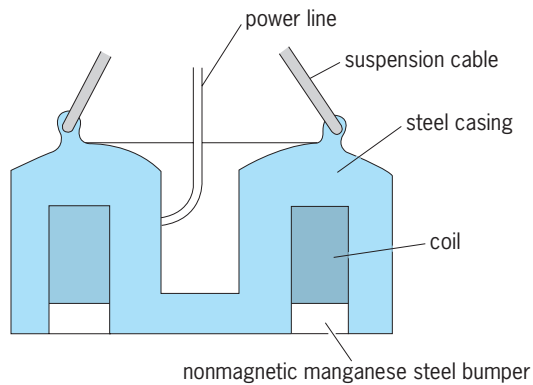
Electromagnet

A soft-iron core that is magnetized by passing a current through a coil of wire wound on the core. Electromagnets are used to lift heavy masses of magnetic material and to attract movable magnetic parts of electric devices, such as solenoids, relays, and clutches.

The difference between cores of an electromagnet and a permanent magnet is in the retentivity of the material used. Permanent magnets, initially magnetized by placing them in a coil through which current is passed, are made of retentive (magnetically "hard") materials which maintain the magnetic properties for a long period of time after being removed from the coil. Electromagnets are meant to be devices in which the magnetism in the cores can be turned on or off. Therefore, the core material is nonretentive (magnetically "soft") material which maintains the magnetic properties only while current flows in the coil. All magnetic materials have some retentivity, called residual magnetism; the difference is one of degree. See MAGNETIZATION.

A magnet, when brought near other susceptible material, induces magnetic poles in the susceptible material and so attracts it. A force is developed in the susceptible material that tends to move it in a direction to minimize the reluctance of the flux path of the magnet. The reluctance force may be expressed quantitatively in terms of the rate of change of reluctance with respect to distance. See MAGNETISM.

In an engineering sense the word electromagnet does not refer to the electromagnetic forces incidentally set up in all devices in which an electric current



Cross section of circular lifting electromagnet.

exists, but only to those devices in which the current is primarily designed to produce this force, as in solenoids, relay coils, electromagnetic brakes and clutches, and in tractive and lifting or holding magnets and magnetic chucks.

Electromagnets may be divided into two classes: traction magnets, in which the pull is to be exerted over a distance and work is done by reducing the air gap; and lifting or holding magnets, in which the material is initially placed in contact with the magnet. For examples of the first type *see* BRAKE; CLUTCH; RELAY; SOLENOID (ELECTRICITY).

Examples of the latter type are magnetic chucks and circular lifting magnets. The illustration shows a typical circular lifting magnet. The outer rim makes up one pole and the inner area is the opposing pole. The coil is wound cylindrically around the center pole. Manganese steel, used as a protective cover plate for the coil, is nonmagnetic and thus forces the flux through the magnetic member being lifted.

The mechanical force between two parallel surfaces is given by Maxwell's equation, shown below,

$$F = B^2 A / 2\mu_0 \text{ (newtons)}$$

where B is the flux density (in webers/m²), A is the cross-sectional area (in m²) through which the flux passes, and μ_0 is the permeability of free space ($4\pi \times 10^{-7}$ henry/m). When two poles are active, the force produced by each is calculated to find the total force. An interesting result of this relation is that the force is not simply the result of the total flux (BA) but also of the flux density. Thus if the same flux can be forced through one-half the cross-sectional area, the net pull will be doubled. In practice it is difficult, if not impossible, to calculate the actual lifting capacity of the magnet by using Maxwell's equation since the capacity varies with the shape and kind of material lifted, how the material is stacked, and other factors. Therefore, lifting magnets are usually rated on their all-day average lifting capacity.

Since currents are large and the circuit is highly inductive, control of a lifting magnet is a problem. If the line switch were simply opened, a destructive arc would result due to the coil's inductance not allowing the current to change instantaneously to a

zero value. Therefore, the controller employed with a lifting magnet usually does the following things automatically: (1) reduces magnet current after an initially high value to a lower holding value to reduce heating of the magnet, (2) introduces a shunt discharge resistor across the magnet coil before allowing the line to be opened when the operator turns the magnet off, and (3) causes a reduced current of reverse polarity to flow in the magnet coil for a short time after the operator turns the switch off. Thus the residual magnetism is canceled and scraps and small chunks that might have continued clinging to the magnet are released.

Jerome Meisel

Bibliography. D. G. Fink and H. W. Beaty (eds.), *Standard Handbook for Electrical Engineers*, 14th ed., 2000; G. R. Jones, M. A. Laughton, and M. G. Say (eds.), *The Electrical Engineers' Reference Book*, 15th ed., 1993; E. M. Purcell, *Electricity and Magnetism*, 2d ed., 1984.

Electromagnetic compatibility

The situation in which electrical and electronic devices and systems work as intended, both within themselves and in their electromagnetic environment.

Electromagnetic interference (EMI) is said to exist when unwanted voltages or currents are present so that they adversely affect the performance of a device or system. Such voltages or currents may reach the victim circuit or device by conduction or by non-ionizing radiation. In all cases, electromagnetic interference arises because of a combination of three factors: a source, a transmission path, and a response, at least one of which is unplanned. Electromagnetic interference control refers to the process of making design changes or adjustments of signal or noise levels in order to achieve electromagnetic compatibility (EMC). Examples of electromagnetic interference often experienced by the public are herringbone patterns on the television screen and buzzing and clicking sounds in the radio.

Sources of interference. Designers usually do not intend that their devices be sources of interference. However, what is a desired signal in one path may be an undesired signal, that is, noise, in a path into which it is inadvertently coupled. Interference may be unintentional, or nonfunctional, such as an arc discharge, radiation from a lightning stroke, a corona discharge from high-tension power lines, or a noise caused by a sudden change in current flow in a conductor. Functional interference, while not usually produced for interference purposes, often includes sine waves, computer clock pulses, speech or video waves, or pulses forming data trains. An example of functional interference is signal leakage from cable television systems. Fluorescent lamps, commutators, automotive ignition systems, and industrial, scientific, and medical equipment such as diathermy machines all constitute sources of interference, as does the electromagnetic pulse (EMP) which accompanies a nuclear detonation. Short-wave radio

listeners often experience interference from high-frequency radar from Russia. Interference also exists across national borders because of incomplete international agreement on the use of the radio-frequency spectrum. *See* ELECTRICAL INTERFERENCE; ELECTROMAGNETIC PULSE (EMP); RADIO SPECTRUM ALLOCATION.

Interference coupling. Electromagnetic interference always starts with a flow of current through a conductor (which may include a conducting gas), and always appears at the victim in the form of a current or a voltage. The coupling path, however, may be a conduction or a radiation path. Thus the actual paths include common wiring, capacitance between devices, mutual inductance between adjacent wiring, nonionizing radiation, or wires in an electromagnetic field. Such coupling is aided by the fact that all conductors, whether wires, printed circuit board traces, or the conductors within an integrated circuit, exhibit both resistance and inductance. Likewise, between wires, mutual inductance and capacitance exist, resulting in a variety of possible coupling paths. *See* CAPACITANCE; COUPLED CIRCUITS; INDUCTANCE.

Grounds and bonds. Grounding is the establishment of an electrically conductive path between an electrical or electronic element of a system and a reference point or plane known as the ground. Grounding thus is a circuit concept, whereas bonding refers to the physical implementation of that concept. Grounding may also refer to an electrical connection made to the earth.

Grounds for currents whose wavelength is long compared with circuit dimensions are best made directly to a single point in a system. For currents whose wavelength is comparable to, or shorter than, circuit dimensions, however, the multipoint ground must be used to avoid the situation in which the ground lead might be a significant fraction of a quarter-wavelength, and thus not serve as a low-impedance conductor. *See* ELECTRONIC EQUIPMENT GROUNDING; GROUNDING.

Shields. The purpose of shielding is to confine radiated energy to a specific region, or to prevent radiated energy from entering a specific region. The most effective shield is a solid metallic enclosure, made of a permeable metal (for example, iron or steel) if frequencies below 100 kHz are to be shielded, or of any metal if higher frequencies are to be shielded. However, the solid shield does not permit light, air, water, or other substances to be passed through it, so shields with holes, including screens, braids, and honeycomb arrangements, as well as conductive glass may be needed. The widespread use of plastic enclosures has made thin film shields vital in achieving the needed shielding effectiveness in the use of such enclosures. *See* ELECTRICAL SHIELDING.

Filters. An electrical filter offers relatively little opposition to the passage of certain frequencies or direct current (dc) while blocking the passage of other frequencies. Accordingly, filters play a significant role in reducing conducted interference to the extent that such interference has a

spectral content different from that of the desired signals.

A filter may be either reflective or lossy. Reflective filters present an impedance mismatch to unwanted frequencies, thereby returning them to the input, whereas lossy filters absorb unwanted frequencies. A filter may be designed on a time-domain basis as well as on a frequency-domain basis. *See* ELECTRIC FILTER; IMPEDANCE MATCHING.

Mathematical models. The complexity of electronic circuitry no longer allows the use of simple trial-and-error methods in achieving an electromagnetically compatible system design. System behavior can be expressed in terms of equations that describe source outputs, transmission or coupling characteristics, and susceptor responses. Such models have been incorporated into system analysis programs which not only predict system compatibility but also can be used to evaluate the effect of parameter changes on overall system performance, as may be needed in analyzing a proposed subsystem specification waiver, or in generating interface specifications for use within an overall system design. *See* SIMULATION.

Standards. Numerous standards have been developed for military as well as for industrial and consumer products. For military equipment, MIL-STD-461B, *EMI Characteristics Requirements for Equipment*, defines the emission and susceptibility limits applicable to units under various test conditions; MIL-STD-462, *Measurement of EMI Characteristics*, describes how the measurements are to be made.

Such organizations as the Institute of Electrical and Electronics Engineers (IEEE), the Society of Automotive engineers (SAE), the Electronic Industries Association (EIA), and the American National Standards Institute (ANSI) are active in promulgating standards for electronic equipment, including electromagnetic compatibility aspects.

The Federal Communications Commission's (FCC) *Rules and Regulations*, Part 15, as amended by FCC Docket 20780, establishes emission limits on "computing devices," defined as any electronic devices or systems that use digital techniques. A class A computing device is destined for commercial, industrial, or business use. It must meet certain radiation limits (for frequencies of 30 to 100 MHz) at a 96-ft (30-m) distance, and conduction limits (for frequencies of 0.45 to 30 MHz) with respect to its power leads. A class B computing device is destined for residential use. This device must meet certain radiation limits at a 10-ft (3-m) distance. Both its radiation limits and its conduction limits are more stringent than those of a class A device. *See* DIGITAL COMPUTER; MICROCOMPUTER.

Mutual effects. Digital systems, such as computers, tend to interfere with analog systems, such as voice and video communications, more readily than analog systems interfere with digital systems. Therefore, data streams to be transmitted over analog voice circuits are converted to a quasianalog tone form first. Computer clocks also may have to be shielded and

their output circuits may have to be filtered to prevent interference to communication equipment. In addition, personal computers must be connected to television receivers so that the video output of the personal computer does not reach the television receiving antenna, which then would radiate such signals. See ELECTRICAL NOISE. Bernhard E. Keiser

Bibliography. J. Goedbloed, *Electromagnetic Compatibility*, 1993; B. E. Keiser, *Principles of Electromagnetic Compatibility*, 3d ed., 1987; C. R. Paul, *Introduction to Electromagnetic Compatibility*, 1992; D. A. Weston, *Electromagnetic Compatibility Principles and Applications*, 1991.

Electromagnetic field

A changing magnetic field always produces an electric field, and conversely, a changing electric field always produces a magnetic field. This interaction of electric and magnetic forces gives rise to a condition in space known as an electromagnetic field. The characteristics of an electromagnetic field are expressed mathematically by Maxwell's equations. See ELECTRIC FIELD; ELECTROMAGNETIC RADIATION; ELECTROMAGNETIC WAVE; MAXWELL'S EQUATIONS. Jesse W. Beams

Electromagnetic induction

The production of an electromotive force either by motion of a conductor through a magnetic field in such a manner as to cut across the magnetic flux or by a change in the magnetic flux that threads a conductor.

Motional electromotive force. A charge moving perpendicular to a magnetic field experiences a force that is perpendicular to both the direction of the field and the direction of motion of the charge. In any metallic conductor, there are free electrons, electrons that have been temporarily detached from their parent atoms.

If a conducting bar (Fig. 1) moves through a magnetic field, each free electron experiences a force due to its motion through the field. If the direction of the motion is such that a component of the force

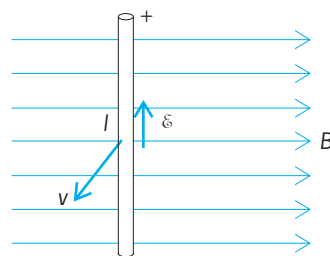


Fig. 1. Flux density B , motion v , and induced emf \mathcal{E} when a conductor of length l moves in a uniform field. (After R. L. Weber, M. W. White, and K. V. Manning, *Physics for Science and Engineering*, McGraw-Hill, 1957)

on the electrons is parallel to the conductor, the electrons will move along the conductor. The electrons will move until the forces due to the motion of the conductor through the magnetic field are balanced by electrostatic forces that arise because electrons collect at one end of the conductor, leaving a deficit of electrons at the other. There is thus an electric field along the rod, and hence a potential difference between the ends of the rod while the motion continues. As soon as the motion stops, the electrostatic forces will cause the electrons to return to their normal distribution.

From the definition of magnetic induction (flux density) B , the force on a charge q due to the motion of the charge through a magnetic field is given by Eq. (1), where the force F is at right angles to a

$$F = Bqv \sin \theta \tag{1}$$

plane determined by the direction of the field, and the component $v \sin \theta$ of the velocity is perpendicular to the field. When B is in webers/m², q is in coulombs, and v is in m/s, the force is in newtons.

The electric field intensity E due to this force is given in magnitude and direction by the force per unit positive charge. The electric field intensity is equal to the negative of the potential gradient along the rod. In motional electromotive force (emf), the charge being considered is negative. Thus, Eqs. (2)

$$E = \frac{F}{-q} = -Bv \sin \theta = -\frac{\mathcal{E}}{l} \tag{2}$$

$$\mathcal{E} = Blv \sin \theta$$

hold. Here l is the length of the conductor in a direction perpendicular to the field, and $v \sin \theta$ is the component of the velocity that is perpendicular to the field. If B is in webers/m², l is in meters, and v is in m/s, the emf \mathcal{E} is in volts.

This emf exists in the conductor as it moves through the field whether or not there is a closed circuit. A current would not be set up unless there were a closed circuit, and then only if the rest of the circuit does not move through the field in exactly the same manner as the rod. For example, if the rod slides along stationary tracks that are connected together, there will be a current in the closed circuit. However, if the two ends of the rod were connected by a wire that moved through the field with the rod, there would be an emf induced in the wire that would be equal to that in the rod and opposite in sense in the circuit. Therefore, the net emf in the circuit would be zero, and there would be no current.

EMF due to change of flux. When a coil is in a magnetic field, there will be a flux Φ threading the coil the magnitude of which will depend upon the area of the coil and its orientation in the field. The flux is given by $\Phi = BA \cos \theta$, where A is the area of the coil and θ is the angle between the normal to the plane of the coil and the magnetic field. Whenever there is a change in the flux threading the coil, there will be an induced emf in the coil while the change is taking place. The change in flux may be caused by

a change in the magnetic induction of the field or by a motion of the coil. The magnitude of the induced emf, Eq. (3), depends upon the number of turns of

$$\mathcal{E} = -N \frac{d\Phi}{dt} \quad (3)$$

the coil N and upon the rate of change of flux. The negative sign in Eq. (3) refers to the direction of the emf in the coil; that is, it is always in such a direction as to oppose the change that causes it, as required by Lenz's law. If the change is an increase in flux, the emf would be in a direction to oppose the increase by causing a flux in a direction opposite to that of the increasing flux; if the flux is decreasing, the emf is in such a direction as to oppose the decrease, that is, to produce a flux that is in the same direction as the decreasing flux. See FARADAY'S LAW OF INDUCTION; LENZ'S LAW.

Consider the case of a flat coil of area A rotating with uniform angular velocity ω about an axis perpendicular to a uniform magnetic field of flux density B . For any position of the coil, the flux threading the coil is $\Phi = BA \cos \theta = BA \cos \omega t$, where the zero of time is taken when θ is zero and the normal to the plane of the coil is parallel to the field. Then the emf induced as the coil rotates is given by Eq. (4).

$$\mathcal{E} = -N \frac{d\Phi}{dt} = -NBA \frac{d(\cos \theta)}{dt} = NBA\omega \sin \omega t \quad (4)$$

The induced emf is sinusoidal, varying from zero when the plane of the coil is perpendicular to the field to a maximum value when the plane of the coil is parallel to the field.

Self-induction. If the flux threading a coil is produced by a current in the coil, any change in that current will cause a change in flux, and thus there will be an induced emf while the current is changing. This process is called self-induction. The emf of self-induction is proportional to the rate of change of current. The ratio of the emf of induction to the rate of change of current in the coil is called the self-inductance of the coil.

Mutual induction. The process by which an emf is induced in one circuit by a change of current in a neighboring circuit is called mutual induction. Flux produced by a current in a circuit A (Fig. 2) threads or links circuit B . When there is a change of current in circuit A , there is a change in the flux linking coil B , and an emf is induced in circuit B while the change is taking place. Transformers operate on the principle of mutual induction. See TRANSFORMER.

The mutual inductance of two circuits is defined as the ratio of the emf induced in one circuit B to the rate of change of current in the other circuit A . For a detailed discussion of self- and mutual inductance See INDUCTANCE.

Coupling coefficient. This refers to the fraction of the flux of one circuit that threads the second circuit. If two coils A and B having turns N_A and N_B , respectively, are so related that all the flux of either threads both coils, the respective self-inductances are given by Eqs. (5) and the mutual inductance of the pair is

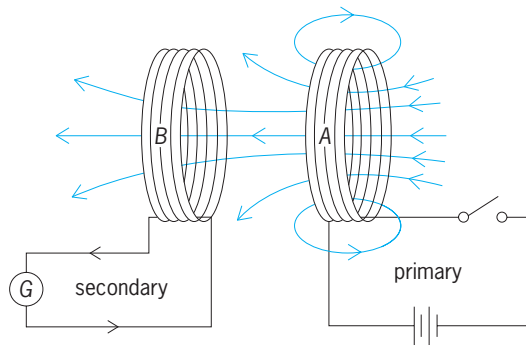


Fig. 2. Mutual induction. An emf is induced in the secondary when the current changes in the primary. (After R. L. Weber, M. W. White, and K. V. Manning, *Physics for Science and Engineering*, McGraw-Hill, 1957)

given by Eq. (6). Then Eqs. (7) hold.

$$L_A = \frac{N_A \Phi_A}{I_A} \quad L_B = \frac{N_B \Phi_B}{I_B} \quad (5)$$

$$M = \frac{N_A \Phi_B}{I_B} = \frac{N_B \Phi_A}{I_A} \quad (6)$$

$$M^2 = \frac{N_A N_B \Phi_A \Phi_B}{I_A I_B} = \frac{N_A \Phi_A}{I_A} \frac{N_B \Phi_B}{I_B} = L_A L_B \quad (7)$$

$$M = \sqrt{L_A L_B}$$

In general, not all the flux from one circuit threads the second. The fraction of the flux from circuit A that threads circuit B depends upon the distance between the two circuits, their orientation with respect to each other, and the presence of a ferromagnetic material in the neighborhood, either as a core or as a shield. It follows that for the general case that Eq. (8) holds.

$$M \leq \sqrt{L_A L_B} \quad (8)$$

The ratio of the mutual inductance of the pair to the square root of the product of the individual self-inductances is called the coefficient of coupling K , given by Eq. (9). The coupling coefficient has a maximum

$$K = \frac{M}{\sqrt{L_A L_B}} \quad (9)$$

value of unity if all the flux threads both circuits, zero if none of the flux from one circuit threads the other. For all conditions, K has a value between 0 and 1.

Applications. The phenomenon of electromagnetic induction has a great many important applications in modern technology. For example, see COUPLED CIRCUITS; GENERATOR; INDUCTION HEATING; MICROPHONE; MOTOR; SERVOMECHANISM.

Kenneth V. Manning

Bibliography. W. J. Duffin, *Electricity and Magnetism*, 4th ed., 1990; D. Halliday and R. Resnick, *Physics*, pt. 2, 4th ed., 1992; A. Mottershead, *Introduction to Electricity and Electronics*, 2000; E. M. Purcell, *Electricity and Magnetism*, 2d ed., 1984.

Electromagnetic pulse (EMP)

A transient electromagnetic signal produced by a nuclear explosion in or above the Earth's atmosphere. Though not considered dangerous to people, the electromagnetic pulse (EMP) is a potential threat to many electronic systems.

Discovery. The existence of a nuclear-generated EMP has been known for many years. Originally predicted by scientists involved with the early development of nuclear weapons, it was not considered to be a serious threat to people or equipment. Then in the early 1960s some of the high-altitude nuclear tests conducted in the Pacific led to some strange occurrences many miles from ground zero. In Hawaii, for example, some 800 mi (1300 km) from the Johnston Island test, EMP was credited with setting off burglar alarms and turning off street lights. In later tests that were conducted in Nevada, significant EMP-induced signals were coupled to cables. These experimental results gave credibility to the potential military use of EMP and encouraged investigators to more accurately describe its origin, electromagnetic characteristics, and coupling to systems, and to develop an affordable method for system protection.

Initial nuclear radiation. In a typical nuclear detonation, parts of the shell casing and other materials are rapidly reduced to a very hot, compressed gas, which upon expansion gives rise to enormous amounts of mechanical and thermal energy. At the same time the nuclear reactions release tremendous amounts of energy as initial nuclear radiation (INR). This INR is in the form of rapidly moving neutrons and high-energy electromagnetic radiation, called x-rays and prompt gamma rays. About a minute after detonation, the radioactive decay of the fission products gives rise to additional gamma rays and electrons (or beta particles), known as residual nuclear radiation (RNR). The distribution of the total explosive energy of a hypothetical fission detonation in the atmosphere below an altitude of 6 mi (10 km) is 50% blast, 35% thermal, 10% RNR, 5% INR. At higher altitudes where the air is less dense, the thermal energy increases and the blast energy decreases proportionately. See BETA PARTICLES; GAMMA RAYS; NUCLEAR FISSION; NUCLEAR REACTION; RADIOACTIVITY.

EMP is associated with the INR output, which is a small percentage of the total explosive energy. Nevertheless, EMP is still capable of transferring something of the order of 0.1–0.9 joule/m² (0.007–0.06 ft-lbf/ft²) onto a collector, more than enough to cause upset or damage to normal semiconductor devices.

Early research first developed the physics of high-altitude EMP. Then in 1978 a consistent explanation was set forth of how EMP is related to the total gamma source strength regardless of the detonation height above ground. **Figure 1** identifies six contributions to the gamma source strength for a hypothetical 1-megaton surface burst. As the detonation height is elevated, gamma source contributions from neutrons interacting with the ground and air decrease; for a high-altitude detonation (higher than

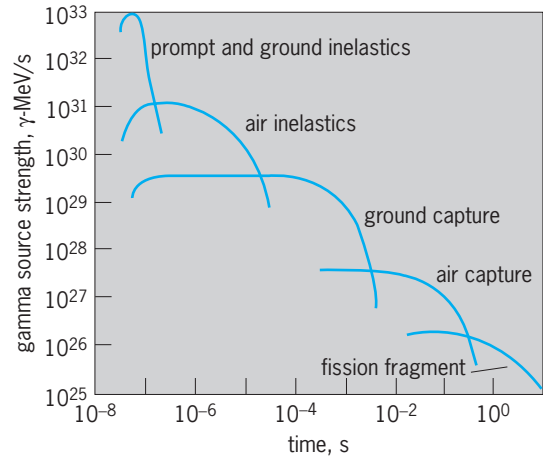


Fig. 1. Total gamma source strength versus time for nominal 1-megaton surface burst. Both horizontal and vertical scales are logarithmic. (After C. L. Longmire, *On the electromagnetic pulse produced by nuclear explosions*, *IEEE Trans. Antennas Propag.*, AP-26(1):3–13, 1978)

37 mi or 60 km) the gamma source strength becomes essentially the prompt gammas from the nuclear burst. The gamma source strength wave shape then approaches a smooth curve, approximating a double exponential, and since EMP is generated from the gamma source strength as described below, it too approaches a double exponential wave shape.

EMP generation in a high-altitude burst. As the prompt gammas move away from a high-altitude nuclear detonation (**Fig. 2**), those gamma rays moving toward the Earth penetrate a more dense region of the atmosphere called the source or deposition region. In this 6-mi thick (10-km) region, approximately 15–21 mi (25–35 km) above the Earth, the highly energetic gamma rays interact with the air molecules to form Compton electrons (with energies starting at 1 MeV) and less energetic gamma rays, which then proceed in the same general direction as the original gamma rays. The fast Compton electrons eventually slow down by stripping other electrons from air molecules to form secondary electron-ion pairs. (Though these secondary electrons and ions

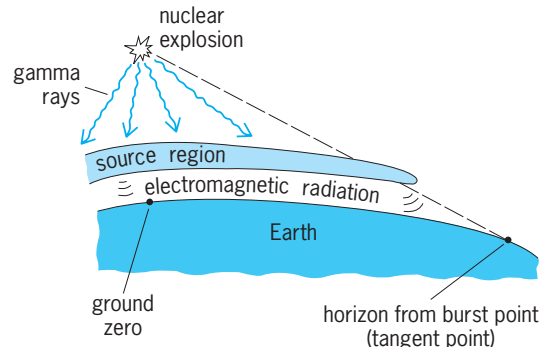


Fig. 2. Schematic representation of the EMP in a high-altitude burst. The extent of the source region varies with the altitude and the yield of the explosion. (After S. Glasstone and P. J. Dolan, eds., *The Effects of Nuclear Weapons*, U.S. Department of Defense and the Energy Research and Development Administration, 3d ed., 1977)

do not contribute to the generation of the EMP, they do cause the region to become highly conductive, and therefore play an important role in determining the EMP wave shape and amplitude.) While slowing down the very intense, short-duration flux of Compton electrons is also deflected by the Earth's geomagnetic field, according to Eq. (1), where the

$$\vec{F} = -q(\vec{v} \times \vec{B}) \quad (1)$$

deflection force \vec{F} is perpendicular to the geomagnetic field \vec{B} and the velocity \vec{v} of a Compton electron of charge q . The Compton electrons then spiral about the geomagnetic lines, radiating electromagnetic energy in the form of EMP until they eventually recombine with local, positively charged ions. See COMPTON EFFECT; SYNCHROTRON RADIATION.

Characteristics for a high-altitude burst. For a high-altitude nuclear detonation, the radiated EMP observed at large distances from the source region can be represented in time t as an electric field $E(t)$ and a magnetic field $H(t)$, given by Eqs. (2) and (3), with

$$E(t) = E_0(e^{-\alpha t} - e^{-\beta t}) \quad t \geq 0 \quad (2)$$

$$H(t) = H_0(e^{-\alpha t} - e^{-\beta t}) \quad t \geq 0 \quad (3)$$

$E_0 = 5.2 \times 10^4$ V/m (1.6×10^4 V/ft), $H_0 = 1.4 \times 10^2$ A/m (4.2×10 A/ft), $\alpha = 4.0 \times 10^6$ s $^{-1}$, and $\beta = 4.76 \times 10^8$ s $^{-1}$. (In air, E and H are related by the impedance of free space, 377 ohms.) If the observer is directly below the detonation (ground zero), the polarization of both fields is predominantly horizontal; if the observer is at the horizon, the fields can have both horizontal and vertical components. The magnitude of these field components varies according to the latitude and longitude of the observer and the direction of the burst. See POLARIZATION OF WAVES.

Surface-burst EMP. Should the nuclear detonation occur closer to the Earth, the EMP generation process becomes far more complex and the electric and magnetic fields become very complicated. The most dramatic change occurs with a surface burst (Fig. 3).

When the observer is somewhere within the source region of a near-surface burst, where the air conductivity varies between 10^{-4} and 10^{-2} siemen/m (3×10^{-5} and 3×10^{-3} siemen/ft), the resultant electric field is predominantly vertical, and the resultant magnetic field is polarized perpendicular to the plane of the figure. The electric field polarization is due to the Compton-electron and ion-charge-separation fields which tend to become perpendicular to the conducting earth, leaving a resultant vertical electric field near the ground. As these Compton electrons move radially away from the detonation, they curve earthward and return to the detonation point through the conducting earth. This Compton current loop then gives rise to a resultant magnetic field. When the observer is far from the detonation (outside the source region), the electric and magnetic fields begin to approximate fields from a vertical dipole, decaying with distance r as $1/r$.

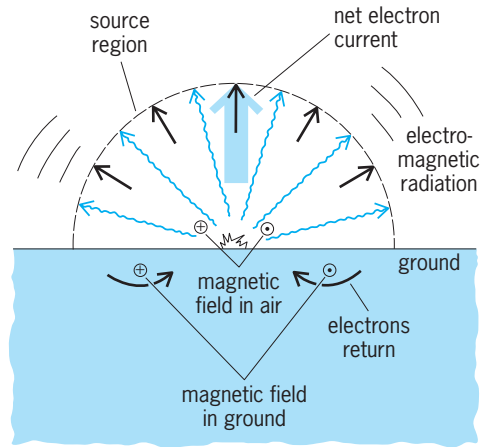


Fig. 3. Schematic representation of the EMP in a surface burst. (After S. Glasstone and P. J. Dolan, eds., *The Effects of Nuclear Weapons*, U.S. Department of Defense and the Energy Research and Development Administration, 3d ed., 1977)

Internal EMP. It is also possible for INR (both x-rays and gamma rays) to directly interact with systems, causing EMP signals internal to structures. This phenomenon has been called internal or system-generated EMP and is potentially a serious problem for satellites in orbit and for electronics in metallic enclosures on or near the ground. These forms of EMP are generated by x-rays interacting with satellites and gamma rays impinging on ground-based enclosures, producing currents of Compton electrons internally that then produce internal electromagnetic waves. They are very dependent upon the nuclear detonation, the system topology, and the relative position of one to the other.

Coupling. The coupling of these rather wide-band (10^4 to 10^8 Hz) signals to systems of different topologies can be significant. For example, it is not unusual to predict voltage and current levels of hundreds of thousands of volts and a few thousand amperes coupled by high-altitude EMP onto extended systems. The exact level of coupling depends upon the size of the system and its orientation with respect to the incident field, and upon whether or not it is near an earth ground. For ground-based systems the incident field components can actually be enhanced or degraded, depending upon the polarization. Solving Maxwell's equations at the interface between an ideal conductor and a dielectric shows that the net horizontal electric field is zero at the interface and the net vertical electric field doubles. (For the magnetic field components the opposite is true.) Since the Earth ground is not a perfect conductor, the field components do not completely cancel or quite double in amplitude. See MAXWELL'S EQUATIONS.

An estimate of about 1 joule (0.7 ft-lbf) of EMP-coupled energy is considered reasonable for many systems. Even if the coupling onto circuits is inefficient, as little as 10^{-13} J can upset some semiconductor devices and 10^{-6} J can cause damage. The potential for such upset and damage in critical electronic circuits has led to the incorporation of EMP protection in many system designs. This protection

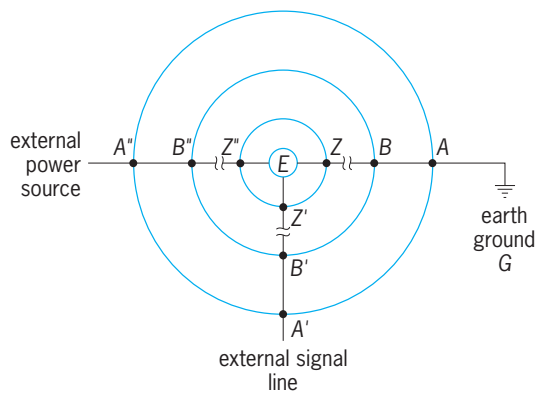


Fig. 4. Typical system protection scheme.

is most prevalent in communications systems whose disruption by EMP is considered an important civil and military vulnerability. See ELECTROMAGNETIC RADIATION; NUCLEAR EXPLOSION.

Protection. The most common form of protection incorporated in system designs is a combination of shielding and penetration control. Figure 4 illustrates a protection scheme in which a system's electronics *E* is isolated from the external environment by one or more nested, shielded enclosures (often called Faraday cages). Penetration control is then maintained by minimizing the number of shield penetrations (in this case, a power line, a signal line, and a ground wire connecting *E* to earth ground *G*) and by applying terminal protection devices, such as spark gaps, high-power Zener diodes, or metal oxide varistors, at selected shield penetration points (*A*, *A'*, *A''*, *B*, *B'*, and *B''*; or *Z*, *Z'*, and *Z''*; or both). In this way, system protection can be designed not only for EMP but also for other electromagnetic transients (such as near-strike lightning and electromagnetic interference). Furthermore, cost-effective, field maintainable protection can be achieved by properly selecting off-the-shelf shielding techniques and terminal protection devices and applying them to systems in accordance with such military standards as MIL-STD 461. See ELECTRIC PROTECTIVE DEVICES; ELECTRICAL INTERFERENCE; ELECTRICAL SHIELDING; ELECTROMAGNETIC RADIATION; LIGHTNING AND SURGE PROTECTION; NUCLEAR EXPLOSION. Robert A. Pfeffer

Bibliography. S. Glasstone and P. J. Dolan (eds.), *The Effects of Nuclear Weapons*, U.S. Department of Defense and the Energy Research and Development Administration, 3d ed., 1977, reprint 1983; K. S. Lee (ed.), *EMP Interaction: Principles, Techniques and Reference Data*, rev. ed, 1986; J. R. Pierce et al., *Evaluation of Methodologies for Estimating Vulnerability to Electromagnetic Pulse Effects*, National Research Council Rep., 1984.

Electromagnetic pump

A pump that operates on the principle that a force is exerted on a current-carrying conductor in a magnetic field. The high electrical conductivity of the

liquid metals pumped (liquid metals are used as the heat-transfer media in some nuclear reactors and magnetohydrodynamic systems) allows a pumping force to be developed within the metals when they are confined in a duct or channel and subjected to a magnetic field and to an electric current. These pumps are designed principally for use in liquid-metal-cooled reactor plants where liquid lithium, sodium, potassium, or sodium-potassium alloys are pumped. Other metallic and nonmetallic liquids of sufficiently high electrical conductivity, such as mercury or molten aluminum, lead, and bismuth, may also be pumped in nonnuclear applications. The absence of moving parts within the pumped liquid eliminates the need for seals and bearings that are found in conventional mechanical pumps, thus minimizing leaks, maintenance, and repairs, and improving reliability. In liquid-metal-cooled nuclear reactor plants, electromagnetic pumps with a capacity of up to several thousand gallons per minute have operated without maintenance for decades. See MAGNETOHYDRODYNAMICS; NUCLEAR POWER; NUCLEAR REACTOR.

The major classifications of electromagnetic pumps, either conduction or induction, is based on the method employed to cause the electric current to flow in the pumped liquid.

Direct-current conduction pumps. In this simplest form of electromagnetic pump the liquid-containing pump duct has two heavy electrodes to provide a path for direct current to be conducted at a right angle to the direction of pumping force in the duct. A magnetic field generated by either a permanent magnet or electromagnet is impressed across the duct at the location of the electrodes in a direction at right angles to both the current flow and the pumping direction. This orthogonal arrangement of current, magnetic field, and pump duct produces a force in the liquid which is dependent on the magnitude of the current, magnetic field, and the length of the current path in the liquid. It is also possible to arrange the duct so that the current flowing in the liquid also produces the required magnetic field (that is, the self-induced pump). The physical arrangement of the conduction pump allows only a single passage of current between electrodes, resulting in a relatively short length of current path in the liquid. To compensate, the current must be very high (thousands of amperes), thus requiring special dc power supplies such as homopolar generators. However, the very low voltage allows special electrical insulation to be used, which is suitable for high-temperature (1500°F or 800°C) operation.

In addition to the externally energized dc conduction pump described above, there are two other unique pump design concepts of the dc conduction pump family in use or under development. In the thermoelectric electromagnetic pump, the current in the duct is generated by thermoelectric elements attached to the duct electrodes which are driven by the temperature difference between the duct and a heat sink. The major application of thermoelectric electromagnetic pumps is in the

liquid-metal loops of space reactors. See THERMOELECTRICITY.

The second design concept is the flow coupler. In this design a dc conduction pump duct in a secondary (driven) loop is directly coupled to a magnetohydrodynamic generator located in an adjacent parallel duct of a primary (driving) loop to which it is connected by appropriate buswork. Both ducts lie in a magnetic field produced by a permanent magnet or electromagnet. The initial use of this concept was to drive small auxiliary loops in experimental liquid-metal facilities. However, with the advent of liquid-metal pool reactors the concept is being studied as a means of circulating the primary coolant in the pool.

Alternating-current conduction pump. This pump uses a pump duct, magnetic field, and electrode arrangement similar to the dc conduction pump, except that the magnetic field is produced by an ac (single-phase) electromagnet and the current is produced by the transformer action of the same or separate single-phase electromagnet with a secondary (step-down) winding of usually one turn. The heavy secondary winding connects to the electrodes. With the sinusoidal magnetic field and duct current in phase with each other, a pumping force will be developed. This design can be operated directly from normal single-phase ac power equipment since it has an integral transformer to develop the necessary high current. It too can have a high-temperature application. See ELECTROMAGNET; TRANSFORMER.

Helical induction pump. This pump closely resembles a polyphase squirrel-cage ac induction motor in principle and operating characteristics. The fluid channel is formed by helical vanes encased in a concentric cylinder and is analogous to the rotor of a conventional motor. Arranged around the periphery of the outer cylinder is a polyphase multipole winding which is similar to the stator of an induction motor. The polyphase winding creates a rotating magnetic field which induces a current in the liquid metal, and this current in turn reacts with the magnetic field to produce a force on the fluid tending to rotate it in the duct. However, the fluid is channeled by the helical vanes, forcing it to produce an axial pressure and resultant axial flow in the duct. The helical induction pump is used in applications requiring higher pressures and lower flows than the linear induction pumps described below. See INDUCTION MOTOR.

Linear induction pumps. These pumps use a stationary polyphase ac winding to produce a traveling magnetic field that induces a current in the liquid. The reaction between the induced current and the traveling magnetic field causes the liquid to follow the traveling magnetic field. The force generated in the liquid is confined in the duct, producing a linear flow.

Two principal designs are the flat linear and the annular linear induction pumps. The flat linear induction pump has a rectangular duct connected to inlet and outlet nozzles. Several liquid-metal-cooled reactors use such pumps in their main heat transport systems. The annular induction pump design has a cylindrical duct formed by two concentric tubes.

These pumps are designed for loop application. They are temperature-limited with conventional electrical insulation and in most cases require a separate cooling system for the electrical windings. However, for application to pool- or tank-type liquid-metal-cooled reactors, there is considerable interest in the design and development of large submersible self-cooled (windings cooled by the liquid being pumped) linear induction pumps. Leslie R. Dahl

Bibliography. R. S. Baker and M. J. Tessier, *Handbook of Electromagnetic Pump Technology*, 1987; J. Collett et al., *Thermoelectric Electromagnetic Pump Design for Sp-100, Symposium on Space Nuclear Power Systems*, 1988; *Fast Breeder Reactors—Experience and Trends: Proceedings of a Symposium*, Lyons, 1986; *International Conference on Design and Safety of Advanced Nuclear Power Plants*, Tokyo, 1992; *International Conference on Fast Reactors and Related Fuel Cycles*, Kyoto, 1991; A. E. Walter and A. B. Reynolds, *Fast Breeder Reactors*, 1981.

Electromagnetic radiation

Energy transmitted through space or through a material medium in the form of electromagnetic waves. The term can also refer to the emission and propagation of such energy. Whenever an electric charge oscillates or is accelerated, a disturbance characterized by the existence of electric and magnetic fields propagates outward from it. This disturbance is called an electromagnetic wave. The frequency range of such waves is tremendous, as is shown by the electromagnetic spectrum in the **table**. The sources given are typical, but not mutually exclusive, as is shown by the fact that the atomic interstellar hydrogen radiation whose wavelength is 0.210614 m falls in the radar region. The other monochromatic radiation listed is that from positron-electron annihilation whose wavelength is 2.42626×10^{-12} m.

Detection of radiation. In theory, any electromagnetic radiation can be detected by its heating effect. This method has actually been used over the range from x-rays to radio. Ionization effects measured by cloud chambers, photographic emulsions, ionization chambers, and Geiger counters have been used in the gamma- and x-ray regions. Direct photography can be used from the gamma-ray to the infrared region. Fluorescence is effective in the x-ray and ultraviolet ranges. Bolometers, thermocouples, and other heat-measuring devices are used chiefly in the infrared and microwave regions. Crystal detectors, vacuum tubes, and transistors cover the microwave and radio frequency ranges.

Free-space waves. A charge in simple harmonic (linear sinusoidal) motion in a vacuum generates a simple wave which becomes spherical at distances from the source much larger than the amplitude of the motion and so great that many oscillations have occurred before the disturbance arrives. The wave is plane when the dimensions of the area observed are very small compared with the radius of spherical

Electromagnetic spectrum			
Frequency, Hz	Wavelength, m	Nomenclature	Typical source
10^{23}	3×10^{-15}	Cosmic photons	Astronomical
10^{22}	3×10^{-14}	γ -rays	Radioactive nuclei
10^{21}	3×10^{-13}	γ -rays, x-rays	
10^{20}	3×10^{-12}	X-rays	Atomic inner shell, positron-electron annihilation
10^{19}	3×10^{-11}	Soft x-rays	Electron impact on a solid
10^{18}	3×10^{-10}	Ultraviolet, x-rays	Atoms in sparks
10^{17}	3×10^{-9}	Ultraviolet	Atoms in sparks and arcs
10^{16}	3×10^{-8}	Ultraviolet	Atoms in sparks and arcs
10^{15}	3×10^{-7}	Visible spectrum	Atoms, hot bodies, molecules
10^{14}	3×10^{-6}	Infrared	Hot bodies, molecules
10^{13}	3×10^{-5}	Infrared	Hot bodies, molecules
10^{12}	3×10^{-4}	Far-infrared	Hot bodies, molecules
10^{11}	3×10^{-3}	Microwaves	Electronic devices
10^{10}	3×10^{-2}	Microwaves, radar	Electronic devices
10^9	3×10^{-1}	Radar	Electronic devices, interstellar hydrogen
10^8	3	Television, FM radio	Electronic devices
10^7	30	Short-wave radio	Electronic devices
10^6	300	AM radio	Electronic devices
10^5	3000	Long-wave radio	Electronic devices
10^4	3×10^4	Induction heating	Electronic devices
10^3	3×10^5		Electronic devices
100	3×10^6	Power	Rotating machinery
10	3×10^7	Power	Rotating machinery
1	3×10^8		Commutated direct current
0	Infinity	Direct current	Batteries

curvature. In this case the choice of the rectangular coordinates x and z as the directions of the oscillation and of the observation or field point, respectively, permits the electric intensity \mathbf{E} and the magnetic flux density \mathbf{B} to be written as Eq. (1). The field

$$E_x = vB_y = E_0 \cos [\omega(t - v^{-1}z)] \quad (1)$$

amplitude E_0 is constant over the specified area and not dependent on z if the z range is small compared with the source distance, as in stellar radiation. The angular frequency of the source is ω radians per second, which is the frequency ν in hertz multiplied by 2π . The velocity of the wave is v , the direction of propagation z , and the time t . The wavelength λ is $2\pi v/\omega$. If t is in seconds and z is in meters, then v is in meters per second and λ is in meters. It is found that in a lossless, isotropic, homogeneous medium Eq. (2) holds; here μ is the permeability, and ϵ the

$$v = (\mu\epsilon)^{-1/2} \quad (2)$$

capacitance, or dielectric constant. This wave is transverse because \mathbf{E} and \mathbf{B} are normal to z . It is plane-polarized because E_x and B_y are parallel to fixed axes. The plane of polarization is taken as that defined by the electric vector and the direction of propagation.

Plane waves. An electromagnetic disturbance is a plane wave when the instantaneous values of any field element such as \mathbf{E} and \mathbf{B} are constant in phase over any plane parallel to a fixed plane. These planes are called wavefronts. In empty unbounded space, \mathbf{E} and \mathbf{B} lie in the wavefront normal to each other; if the wave is unpolarized, their direction fluctuates in this plane in random fashion. If the plane waves

are bounded, as on transmission lines and in waveguides, the amplitudes may vary over the wavefront, and in the case of waveguides and crystals some of the elements will not in general lie in the wavefront. The equation for an undamped plane wave whose front is normal to z is Eq. (3), where F is one of the

$$F = \Phi_1(x, y)f_1(z - vt) + \Phi_2(x, y)f_2(z + vt) \quad (3)$$

field elements such as \mathbf{E} or \mathbf{B} . Note that if an observer sees a certain value of $\Phi_1(x, y)$ at z and then jumps instantaneously in the z direction to a point $z + \Delta z$, the observer will, after waiting a time $\Delta z/v$, see the same value $\Phi_1(x, y)$ because Eq. (4) is valid.

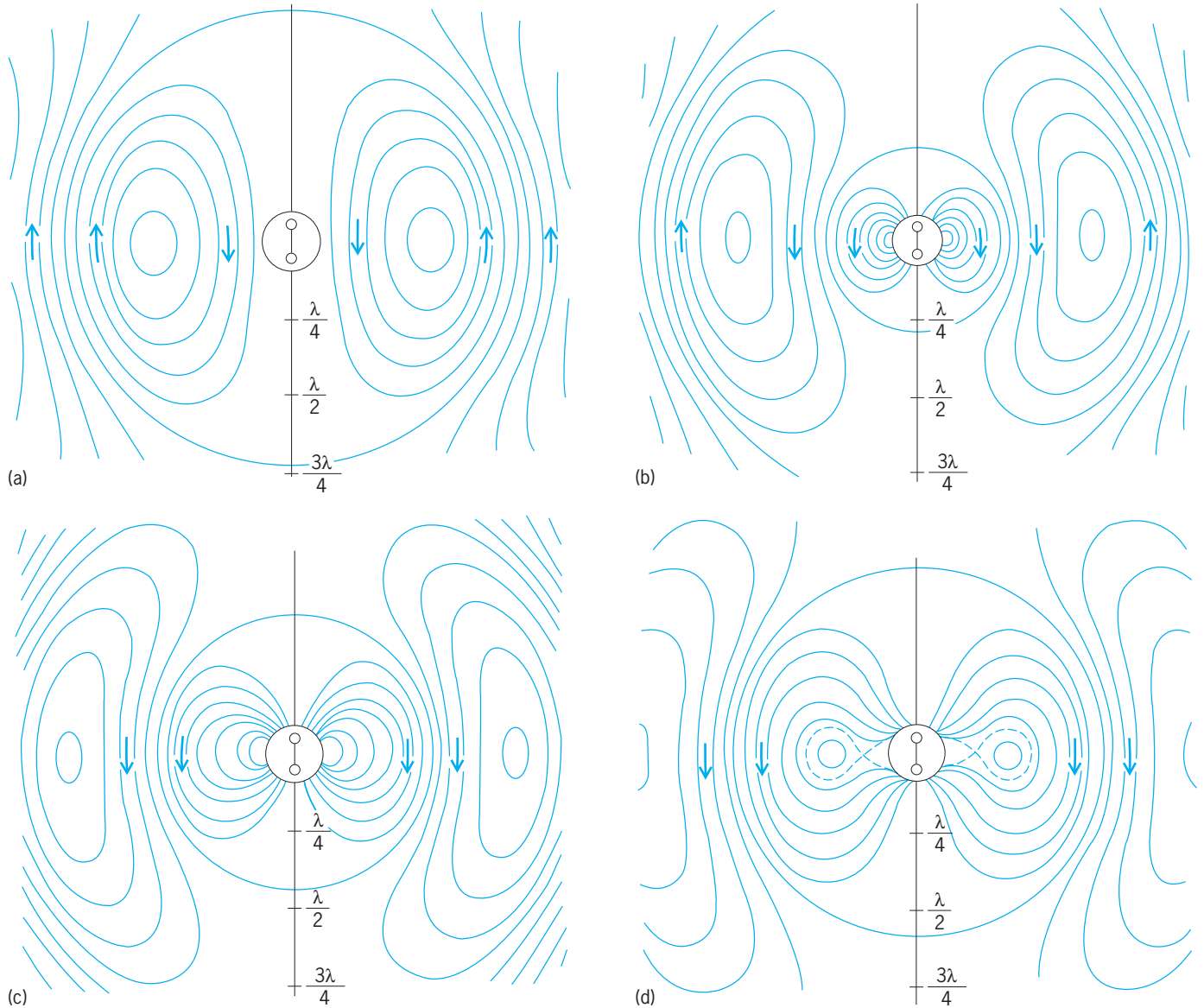
$$f(z - vt) = f[z + \Delta z - v(t + \Delta z/v)] \quad (4)$$

Thus, the first term represents a wave moving in the z direction with a velocity v . The second term represents a wave in the negative z direction. The form of $\Phi_1(x, y)$ and $\Phi_2(x, y)$ depends on the boundary conditions. See WAVE EQUATION.

Spherical waves. A wave is spherical when the instantaneous value of any field element such as \mathbf{E} or \mathbf{B} is constant in phase over a sphere. The radiation from any source of finite dimensions becomes spherical at great distances in an unbounded, isotropic, homogeneous medium. The equation for an undamped spherical wave is Eq. (5). The first term represents a

$$F = r^{-1}\Phi_1(\theta, \varphi)f(r - vt) + r^{-1}\Phi_2(\theta, \varphi)f(r + vt) \quad (5)$$

diverging and the second a converging wave. Again, the form of $\Phi_1(\theta, \varphi)$ and $\Phi_2(\theta, \varphi)$ depends on the nature of the source and other boundary conditions.



Diagrams of electric dipole. The outward moving electric field lines generated by the hertzian oscillator are shown at successive eighth-period intervals. (a) $t = 0$; (b) $t = T/8$; (c) $t = T/4$; (d) $t = 3T/8$.

Damped waves. If there are energy losses which are proportional to the square of the amplitude, as in the case of a medium of conductivity γ which obeys Ohm's law, then the wave is exponentially damped, and Eq. (1) becomes Eq. (6). The symbol α is called

$$E_x = E_0 e^{-\alpha z} \cos(\omega t - \beta z) \quad (6)$$

the attenuation constant, and β the wave number or phase constant which equals ω/v' , where v' is the damped-wave velocity. The electric wave amplitude at the origin has been taken as E_0 . The ratio of E_0 to B_0 , as well as that of α to β , depends on the permeability μ , the capacitivity ϵ , and the conductivity γ of the medium. In terms of the phasor \check{E}_x , Eq. (6) may be written as the real part of Eq. (7).

$$E_x = \check{E}_x e^{j\omega t} = E_0 e^{-(\alpha + j\beta)z} e^{j\omega t} \quad (7)$$

This is exactly the form for the current on a trans-

mission line. (Phasors are complex numbers of form such that, when multiplied by $e^{j\omega t}$, the real part of the product gives the amplitude, phase, and time dependence.)

Wave impedance. Those trained in transmission line theory find it useful to apply the same techniques to wave theory. Consider an isolated tubular section of the wave in Eq. (1) bounded by $x = 0$, $x = 1$, and $y = 0$, $y = 1$ as a transmission line. The potential across the line between $x = 0$ and $x = 1$ is E . The line integral of \mathbf{B} around the $x = 0$ boundary from $y = 0$ to $y = 1$ is μI by Ampère's law and equals \mathbf{B} because \mathbf{B} is zero on the negative side. Thus, the impedance of the line is, making use of Eqs. (1) and (2), given by Eq. (8). This depends only on the properties of

$$\check{Z}_k = \frac{V}{I} = \frac{\mu E}{B} = \frac{E}{H} = \left(\frac{\mu}{\epsilon}\right)^{1/2} = \eta \quad (8)$$

the medium and is known as the wave impedance.

In transmission line theory the ratio μ/ϵ would be replaced by the ratio of the series impedance $Z_L = j\omega L$ to the shunt admittance $Y = j\omega C$, where L is the inductance per unit length, and C the capacitance per unit length across the line. If there is a resistance R per unit length across the line, then $1/R$ must be added to Y . This resistance is $1/\gamma$ for the tubular section. Thus, for a conducting medium, Eq. (8) becomes Eq. (9). The last term is a common

$$\check{Z}_k = \left(\frac{j\omega\mu}{\gamma + j\omega\epsilon} \right)^{1/2} = \frac{j\omega\mu}{\alpha + j\beta} \quad (9)$$

transmission line form. The reflection and refraction of plane waves at plane boundaries separating different mediums may be calculated by transmission line formulas with the aid of Eqs. (8) and (9).

Electric dipole. A charge undergoing simple harmonic motion in free space is a dipole source when the amplitude of the motion is small compared with the wavelength. The term is loosely applied to the hertzian oscillator, usually pictured as a dumbbell-shaped conductor in which the electrons oscillate from one end to the other, leaving the opposite end periodically positive. An electric dipole of moment M is defined as the product qa when two large, equal and opposite charges, $+q$ and $-q$, are placed a small distance a apart. A dipole is oscillating when M is periodic in time and is the simplest source of spherical waves. Much can be learned by a study of H. Hertz's picture of the outward moving electric field lines at successive time intervals of one-eighth period in a plane which passes through the hertzian oscillator axis, shown in the **illustration**. The most striking feature of the pictures is that, after breaking loose from the dipole, all electric field lines are closed, which means that the divergence of \mathbf{E} is zero. This is true of all unbounded waves. It is also noteworthy that the waves become truly spherical with a fixed wavelength λ only in a direction perpendicular to the dipole and at a distance which greatly exceeds the dipole dimensions. This distance is beyond the edges of the picture. Lengths $\lambda/4$, $\lambda/2$, and $3\lambda/4$ are marked off on the axis for comparison. The magnetic field lines are circles coaxial with the oscillator, so they intersect the plane of the diagram normally. They are most dense where the electric lines are closely spaced. The radiant energy emitted by atoms and molecules is essentially radiation of the dipole type. See ABSORPTION OF ELECTROMAGNETIC RADIATION; ANTENNA (ELECTROMAGNETISM); DIFFRACTION; ELECTROMAGNETIC WAVE TRANSMISSION; GAMMA RAYS; HEAT RADIATION; INFRARED RADIATION; INTERFERENCE OF WAVES; LIGHT; MAXWELL'S EQUATIONS; MICROWAVE; POLARIZATION OF WAVES; RADIATION; RADIO-WAVE PROPAGATION; REFLECTION OF ELECTROMAGNETIC RADIATION; REFRACTION OF WAVES; SCATTERING OF ELECTROMAGNETIC RADIATION; TRANSMISSION LINES; ULTRAVIOLET RADIATION; WAVEGUIDE; WAVE MOTION; X-RAYS.

William R. Smythe

Bibliography. M. Born and E. Wolf, *Principles of Optics*, 7th ed., 1999; N. A. Dyson, *X-Rays in*

Atomic and Nuclear Physics, 2d ed., 1990; B. D. Guenther, *Modern Optics*, 1990; A. Ishimaru, *Electromagnetic Wave Propagation, Radiation, and Scattering*, 1996; M. F. Iskander, *Electromagnetic Fields and Waves*, 2000; J. A. Kong, *Electromagnetic Wave Theory*, 3d ed., 2000; S. Ramo, J. R. Whinnery, and T. Van Duzer, *Fields and Waves in Communications Electronics*, 3d ed., 1993; V. V. Sarwate, *Electromagnetic Fields and Waves*, 1993; W. R. Smythe, *Static and Dynamic Electricity*, 3d ed., 1968, reprint 1989.

Electromagnetic wave

A disturbance, produced by the acceleration or oscillation of an electric charge, which has the characteristic time and spatial relations associated with progressive wave motion. A system of electric and magnetic fields moves outward from a region where electric charges are accelerated, such as an oscillating circuit or the target of an x-ray tube. The wide wavelength range over which such waves are observed is shown by the electromagnetic spectrum. The term electric wave, or hertzian wave, is often applied to electromagnetic waves in the radar and radio range. Electromagnetic waves may be confined in tubes, such as wave guides, or guided by transmission lines. They were predicted by J. C. Maxwell in 1864 and verified experimentally by H. Hertz in 1887. See ELECTROMAGNETIC RADIATION; MAXWELL'S EQUATIONS.

William R. Smythe

Electromagnetic wave transmission

The transmission of electrical energy by wires, the broadcasting of radio signals, and the phenomenon of visible light are all examples of the propagation of electromagnetic energy. Electromagnetic energy travels in the form of a wave. Its speed of travel is approximately 186,000 mi/s (3×10^8 m/s) in a vacuum and is somewhat slower than this in liquid and solid insulators. An electromagnetic wave does not penetrate far into an electrical conductor, and a wave that is incident on the surface of a good conductor is largely reflected.

Electromagnetic waves originate from accelerated electric charges. For example, a radio wave originates from the oscillatory acceleration of electrons in the transmitting antenna. The light that is produced within a laser originates when electrons fall from a higher energy level to a lower one. See LASER.

The waves emitted from a source are oscillatory in character and are described in terms of their frequency of oscillation. Local telephone lines (not using carrier systems) carry electromagnetic waves with frequencies of about 200–4000 Hz. Medium-wave radio uses frequencies of the order of 10^6 Hz, radar uses frequencies of the order of 10^{10} Hz, and a ruby laser emits light with a frequency of

4.32×10^{14} Hz. The method of generating an electromagnetic wave depends on the frequency used, as do the techniques of transmitting the energy to another location and of utilizing it when it has been received. See RADAR.

The communication of information to a distant point is generally accomplished through the use of electromagnetic energy as a carrier. A familiar example is the telephone, in which sound waves in the range of frequencies from a few hundred to a few thousand hertz are converted into corresponding electromagnetic waves, which are then guided to their destination by a pair of wires. Another familiar example is radio, in which the signals are caused to modify an identifiable characteristic, such as the amplitude or frequency, of an electromagnetic carrier wave. The electromagnetic wave, thus modified, or modulated, is radiated from an antenna and can be received over a considerable region. See MODULATION; RADIO BROADCASTING.

Features of electromagnetic waves. Figure 1 illustrates schematically some of the essential features of an electromagnetic wave. Shown in Fig. 1 are vectors that represent the electric field intensity E and the magnetic field intensity H at various points along a straight line taken in the direction of propagation of the wave. The electric field is in a vertical plane and the wave is said to be vertically polarized. The magnitude of the field, at a given instant, varies as a sinusoidal function of distance along the direction of propagation. The magnetic field intensity H lies in a plane normal to that of E and, at each point, is proportional in magnitude to E , as shown in Eq. (1),

$$\frac{E}{H} = \sqrt{\frac{\mu}{\epsilon}} \quad (1)$$

where H is the magnetic field intensity in amperes/m, E is the electric field intensity in volts/m, ϵ is the permittivity, or absolute dielectric constant, of the medium, and μ is the absolute permeability of the medium. For a vacuum, $\epsilon = 8.854 \times 10^{-12}$ farad/m and $\mu = 4\pi \times 10^{-7}$ henry/m; therefore for a vacuum the ratio E/H is approximately 377 ohms. This ratio is termed the wave impedance of the medium. See MAXWELL'S EQUATIONS; WAVE EQUATION.

The E and H waves travel along a straight line, as suggested in Fig. 1. Of the two possible directions along this line, the actual direction of travel can be determined by imagining a screw with a right-hand thread placed along the axis and turned from E toward H ; then the longitudinal direction of travel of the screw is the direction of propagation of the energy.

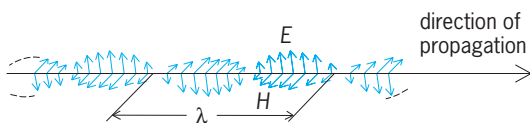


Fig. 1. Representation of an electromagnetic wave at a particular instant of time.

The velocity of travel of the wave is shown in Eq. (2). In a vacuum this is approximately $3 \times$

$$v = 1/\sqrt{\mu\epsilon} \quad (2)$$

10^8 m/s. The velocity in air is only slightly smaller.

The wavelength is the distance between two successive similar points on the wave, measured along the direction of propagation. The wavelength is denoted by λ in Fig. 1.

As the wave travels past a stationary point, the values of E and H at the point vary sinusoidally with time. The time required for one cycle of this variation is termed the period, T seconds. The number of hertz is the frequency f ; and $f = 1/T$. In one cycle the wave, traveling at the velocity v , moves one wavelength along the axis of propagation. Therefore, $\lambda = vT$, or may be calculated by Eq. (3). Assuming a velocity

$$\lambda = \frac{v}{f} \quad (3)$$

of 3×10^8 m/s, an electromagnetic wave having a frequency of 60 Hz has a wavelength of 5×10^6 m, or approximately 3100 mi. At a frequency of 3 MHz (3×10^6 Hz), λ is 100 m, and at 3000 MHz, λ is 10 cm. Visible light has a frequency of the order of 5×10^{14} Hz and a wavelength of approximately 6×10^{-5} cm.

The density of energy in an electric field is $\epsilon E^2/2$ joules/m³, and that in a magnetic field is $\mu H^2/2$ J/m³. With the aid of Eq. (1), the relationships become those shown in Eq. (4). Therefore the

$$\frac{\mu H^2}{2} = \frac{\mu(\sqrt{\epsilon/\mu}E)^2}{2} = \frac{\epsilon E^2}{2} \quad (4)$$

electric and magnetic fields carry equal energies in the electromagnetic wave. The total energy density at any point is equal to ϵE^2 J/m³. Since this is transported with a velocity equal to $1/\sqrt{\mu\epsilon}$, the rate of flow of energy per square meter normal to the direction of propagation is $\epsilon E^2 v$ or $E^2 \sqrt{\epsilon/\mu}$ watts/m². In radio broadcasting a field of 50 millivolts/m is considered to be strong. An electromagnetic wave with this intensity has an average energy density of 2.2×10^{-14} J/m³, and the average rate of energy flow is 6.6×10^{-6} W/m².

Radiation from an antenna. Figure 2 illustrates the configuration of the electric and magnetic fields about a short vertical antenna in which flows a sinusoidal current of the form $i = I_{\max} \sin 2\pi ft$ amperes. The picture applies either to an antenna in free space (in which case the illustration shows only the upper half of the fields), or to an antenna projecting above the surface of a highly conducting plane surface. In the latter case the conducting plane represents to a first approximation the surface of the Earth. The fields have symmetry about the axis through the antenna. For pictorial simplicity only selected portions of the fields are shown in Fig. 2. The magnetic field is circular about the antenna, is perpendicular at every point to the direction of the electric field, and is proportional in intensity to the magnitude of the electric field, as indicated by Eq. (1). All parts of the wave

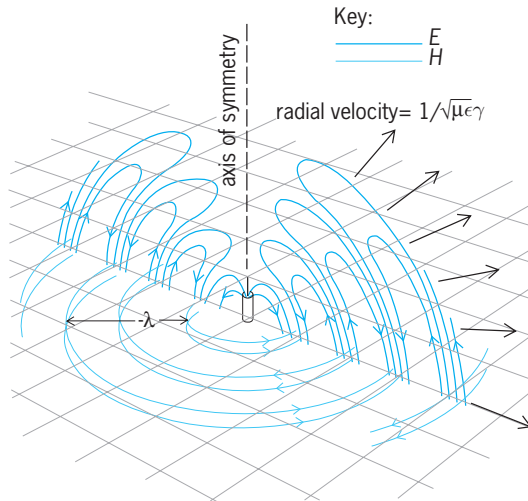


Fig. 2. Configuration of electric and magnetic fields about a short vertical antenna.

travel radially outward from the antenna with the velocity given by Eq. (2); the wave is described as spherical, with the antenna located at the center of the wave. The wavelength of the radiation is given by Eq. (3).

If a short antenna that projects above a highly conducting plane surface carries a current of $i = I_{\max} \sin 2\pi ft$ that is uniform throughout the length of the antenna, the intensity of the electric field in the radiated wave is that shown in Eq. (5), where l is the

$$E_{\max} = \sqrt{\frac{\mu}{\epsilon}} \frac{I_{\max}}{r} \frac{l}{\lambda} \cos \theta \quad \text{V/m} \quad (5)$$

length of the antenna, r is the radial distance from the antenna, and θ is 0, and varies inversely with the distance from the antenna.

If the rate of flow of energy per unit area $E^2\sqrt{\epsilon/\mu}$ is integrated over an imaginary spherical surface about the antenna, the average power radiated is that shown by Eq. (6). The factor $(l/\lambda)^2$ is of par-

$$P_{\text{av}} = \frac{2\pi}{3} \sqrt{\frac{\mu}{\epsilon}} I_{\max}^2 \left(\frac{l}{\lambda}\right)^2 \quad \text{W} \quad (6)$$

ticular importance, for it indicates that a longer antenna is required at the longer wavelengths (lower frequencies). The radiation of appreciable energy at a very low frequency requires an impractically long antenna.

The foregoing relations assume a uniform current throughout the length of the antenna. An approximation to this can be achieved in practice by connecting a long horizontal conductor to the top of the antenna. Where some such construction is not utilized, the current will not be uniform in the antenna and will, in fact, be zero at the tip. The results given by Eqs. (4) and (5) must be modified, but the qualitative features of the radiation remain as shown in Fig. 2. See ANTENNA (ELECTROMAGNETISM).

Often it is desired to concentrate the radiated energy into a narrow beam. This can be done either by

the addition of more antenna elements or by placing a large reflector, generally parabolic in shape, behind the antenna. The production of a narrow beam requires an antenna array, or alternatively a reflector, that is large in width and height compared with a wavelength. The very narrow and concentrated beam that can be achieved by a laser is made possible by the extremely short wavelength of the radiation as compared with the cross-sectional dimensions of the radiating system.

Propagation over the Earth. The foregoing discussion shows some of the important features of the radiation of electromagnetic energy from an antenna, but is oversimplified insofar as communication to and from positions on or near the Earth is concerned. The ground is a reasonably good, but not perfect, conductor; hence, the actual propagation over the surface of the Earth will show a more rapid decrease of field strength than that indicated by the factor of $1/r$ in Eq. (5). Irregularities and obstructions may interfere. In long-range transmission the spherical shape of the Earth is important. Inhomogeneities in the atmosphere refract the wave somewhat. For long-range transmission, the ionized region high in the atmosphere known as the Kennelly-Heaviside layer, or ionosphere, can act as a reflector. The electric field of the wave produces oscillation of the charged particles of the region, and this causes the refractive index of the layer to be smaller than that of the atmosphere below. The result is that, if the angle of incidence is not too near the normal and if the frequency of the wave is not too high, the wave may be refracted back toward the Earth. Successive reflections between ionosphere and Earth can provide communication for long distances around the periphery of the Earth. See RADIO-WAVE PROPAGATION.

Hollow waveguides. When an electromagnetic wave is introduced into the interior of a hollow metallic pipe of suitably large cross-sectional dimensions, the energy is guided along the interior of the pipe with comparatively little loss. The most common cross-sectional shapes are the rectangle and the circle. The cross-sectional dimensions of the tube must be greater than a certain fraction of the wavelength; otherwise the wave will not propagate in the tube. For this reason hollow waveguides are commonly used only at wavelengths of 10 cm or less (frequencies of 3000 MHz or higher).

A single wave of the type in Fig. 1 cannot propagate longitudinally inside a tubular conductor since, at some portions of the inner surface of the conducting tube, the E vector of the wave necessarily would have a component tangential to the surface. This is impossible because an electric field cannot be established along a good conductor, such as the wall of the tube. An electromagnetic wave can propagate along the interior of the tube only by reflecting back and forth between the walls of the tube. This reflection is a comparatively simple one between the plane surfaces of a rectangular tube, but is a complex reflection in tubes of other cross-sectional shapes.

A dielectric rod can also be used as a waveguide. Such a rod, if of insufficient cross-sectional

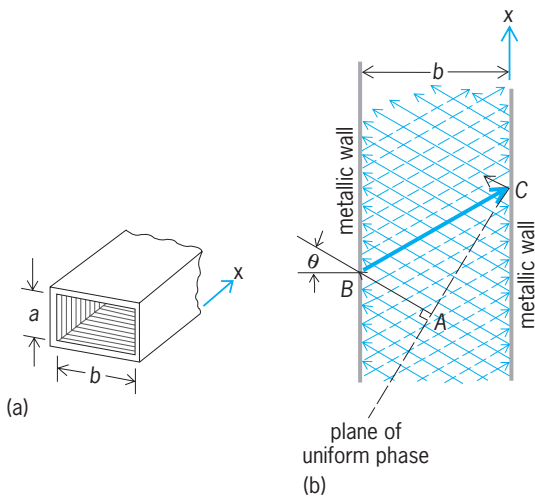


Fig. 3. Hollow metallic waveguide of rectangular cross section. (a) Guide. (b) Paths of electromagnetic energy in the simplest mode of propagation.

dimensions, can contain the electromagnetic wave by the phenomenon of total reflection at the surface.

A hollow metallic waveguide of rectangular cross section is shown in Fig. 3a. The simplest mode of propagation is indicated in Fig. 3b. The entire space is filled with a plane electromagnetic wave which moves obliquely to the left in the direction shown by the solid arrows. This wave has its E vector normal to the paper and its H vector in the plane of the paper. Any plane normal to the direction of propagation is a plane of uniform phase (thus the name plane wave), and one such plane is indicated in the illustration by a broken line. The wave strikes the wall at an angle θ from the normal and is reflected at an equal angle. As the wave is reflected, the direction of its E vector reverses so as to make the tangential component of the electric field equal to zero at the conducting wall. The wave incident on the left wall thus is reflected to the right, where it is again reflected and moves to the left. By successive reflections the energy propagates longitudinally along the interior of the guide. As the wave incident upon the wall reflects and reverses the direction of its E vector, electric currents are caused to flow in the conducting wall. Since the wall is not a perfect conductor, some of the energy of the wave is transformed into heat. Consequently the amplitude of the wave diminishes exponentially as it passes down the guide; this phenomenon is termed attenuation. For an electromagnetic wave with a frequency of 3000 MHz (wavelength of 10 cm) propagating down the interior of a rectangular copper waveguide with cross-sectional dimensions of 1.4 by 3.2 in (or 4 by 8 cm) half the power is lost in a distance of approximately 450 ft (150 m). Hollow waveguides are used chiefly for short-distance transmission, as from a transmitter to an antenna. See WAVEGUIDE.

The requirements on the reflection of the wave, as outlined above, restrict the wavelength that can be propagated in a hollow guide. Consider the ray ABC in Fig. 3. The wave propagates from A to B ,

where it is reflected with reversal of the E vector; thereupon it propagates from B to C , where it is again reflected with another reversal of the E vector. But AC is a line of equal phase, and so the wave emerging from C must have the same phase as that at A . Thus the distance ABC must be an integral multiple of a wavelength, or $n\lambda$, where n is a positive integer. The distance ABC is $2b \cos \theta$ where b is the breadth of the guide; hence, $n\lambda = 2b \cos \theta$. The condition for propagation down the axis of the guide is that $\theta > 0$; hence $\cos \theta < 1$, and the restriction on wavelength is $\lambda < 2b/n$. The greatest ratio of wavelength to breadth of guide is obtained when $n = 1$, whence $\lambda/b < 2$. Therefore, the breadth of the waveguide must be somewhat greater than $\lambda/2$.

In the simple mode of propagation described above, the fields are independent of distance in the direction of the dimension a , and this dimension has no influence on the propagation. The net electric vector caused by the sum of the two waves is everywhere transverse to the longitudinal axis of the guide, and so the mode is described as transverse electric (TE).

If the wavelength of the radiation is small enough in comparison with the cross-sectional dimensions of the guide, more complex modes of propagation are possible, in which the wave reflects obliquely against a side wall, proceeds to the top of the guide, and reflects from there to the other side wall, then to the bottom wall, and so on. With this type of reflection it is possible to have both transverse-electric and transverse-magnetic (TM) modes. In the latter the net H vector is everywhere transverse to the axis of the guide.

When the dimensions of the guide are such that complex modes are possible, so also are the simple ones. The transmission of energy by a combination of modes introduces complications in abstracting the energy from the guide at the receiving end. Propagation in only the simplest mode is ensured by selecting the dimension b to be greater than $\lambda/2$ but not as large as λ , and also by restricting the dimension a so as to render complex modes impossible.

Waveguides of circular cross section are sometimes used. Analysis of these shows that the first TE mode is propagated if the diameter of the guide is greater than 0.586λ , and that the first TM mode is propagated if the radius is greater than 0.766λ .

Two-conductor transmission lines. Electromagnetic energy can be propagated in a simple mode along two parallel conductors. Such a waveguiding system is termed a transmission line. Three common forms are shown in Fig. 4. If the spacing between conductors is a small fraction of the wavelength of the transmitted energy, only one mode of propagation is possible. This corresponds to the wave of Fig. 1, with the direction of propagation taken longitudinally along the line. The E and H vectors are in the plane of the cross section, and the mode is termed transverse electromagnetic (TEM). The E vector must be at right angles to a highly conducting surface, and the oscillating H vector must be parallel with such a surface. With two separated conductors, there is for each

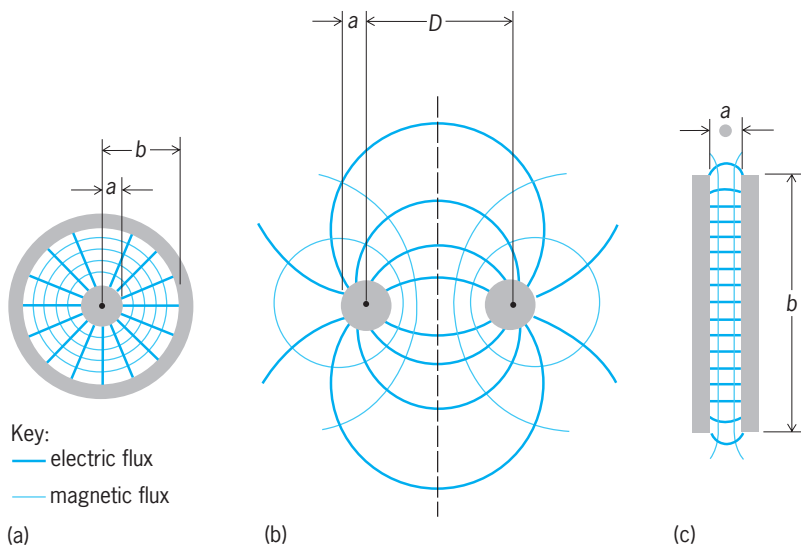


Fig. 4. Cross sections of common two-conductor transmission lines. (a) Coaxial cable. (b) Two-wire line. (c) Parallel-strip line.

geometrical arrangement of conductors one and only one cross-sectional field configuration which will satisfy the boundary conditions at the metal surface. The field configurations for coaxial and two-wire lines are shown in Fig. 4. At each point the ratio of E to H is as given by Eq. (1), and the velocity of propagation of the wave is as given by Eq. (2). Half of the propagated energy is contained in the electric field and half in the magnetic field. This mode of propagation is in contrast with the more complex modes required in a hollow metal pipe, where the conditions required at the boundaries can be satisfied only by means of reflections at the metal walls. As a result, the two-conductor transmission line does not have the upper limit on wavelength that was imposed on the hollow waveguide by the requirement of reflections; in fact, the two-conductor line operates completely normally at zero frequency (direct current).

At wavelengths that are small enough to be comparable with the cross-sectional dimensions of the line, more complex modes, involving reflections from the surfaces of the conductors, become possible. High-frequency energy can thus be propagated in several modes simultaneously. In a coaxial cable a rough criterion for the elimination of higher modes is that the wavelength should be greater than the average of the circumferences of the inner and outer conductors.

As the wave propagates along the line, it is accompanied by currents which flow longitudinally in the conductors. These currents can be regarded as satisfying the boundary condition for the tangential H field at the surface of the conductor. The conductors have a finite conductivity, and so these currents cause a transformation of electrical energy into heat. The energy lost comes from the stored energy of the wave, and so the wave, as it progresses, diminishes in amplitude. The conductors are necessarily supported by insulators which are imperfect and cause additional attenuation of the wave. In a typical open-wire telephone line operating at voice frequencies,

half the energy is lost in a distance of perhaps 60 mi (95 km). The losses increase with frequency, and for a typical air-insulated coaxial line operating at 5 MHz, half the energy is lost in a distance of less than 1 mi (1.6 km). At a frequency of 3000 MHz ($\lambda = 10$ cm), typical distances in which half the energy is lost are, for air-insulated coaxial cable, 82 ft (25 m); for flexible coaxial cable insulated with polyethylene, 33 ft (10 m).

Noise. In a transmission line intended for the transmission of large amounts of power, such as the cross-country lines joining electrical generating stations to centers of population, the loss of an appreciable proportion of the power enroute is a serious matter. In a communication system, however, the average rate of flow of energy is rather small, and the intrinsic value of the energy itself is not of prime importance. The important characteristic of such a system is the accurate transmission of information, and the limiting factor is noise. Noise is always present in a transmission channel. Two common causes are thermal agitation and nearby electrical discharges. In a transmission system conveying information by an electromagnetic wave, the loss of energy in transmission becomes a serious matter if the wave is attenuated to the point where it is not large enough to override the noise. Amplifiers must be inserted in the transmission system at sufficiently close intervals, so that the signal never falls into the noise level, from which it could not be recovered and interpreted accurately.

Circuit analysis of transmission lines. Because the conductors of a transmission line are almost always spaced much closer together than a quarter wavelength of the electromagnetic energy that they are guiding, it is possible to analyze their performance quantitatively by circuit theory. It is then possible to deal with the voltages between the conductors and the currents flowing along the conductors, instead of with the electric and magnetic fields that exist in the insulating medium.

The waveguiding properties of the transmission line can be examined most conveniently if losses of energy are ignored. If L is defined as the inductance of the pair of conductors per unit length and C the capacitance between the conductors per unit length, field theory shows that L is μF_g and C is ϵ/F_g , where F_g is a geometrical factor that depends on the cross-sectional configuration of the conductors. For a coaxial line (Fig. 4a), $F_g = (1/2\pi) \log_e(b/a)$. For a two-wire line (Fig. 4b), $F_g = (1/\pi) \log_e(D/a)$. For a parallel-strip line (Fig. 4c), neglecting edge effects, $F_g = a/b$.

In the circuit analysis of a transmission line, the line can be visualized as being composed of a cascaded set of sections, each of short length Δx , as shown in Fig. 5b. The partial differential equations which describe the voltage e and the current i are shown by Eqs. (7) and (8).

$$\frac{\partial e}{\partial x} = -L \frac{\partial i}{\partial t} \tag{7}$$

$$\frac{\partial i}{\partial x} = -C \frac{\partial e}{\partial t} \tag{8}$$

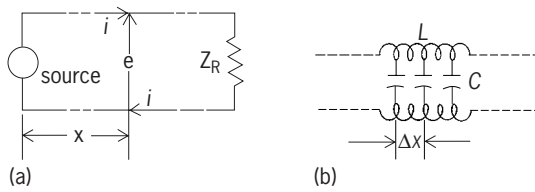


Fig. 5. Schematic representation of a transmission line. (a) Circuit diagram. (b) Visualization of L and C .

The solution of these equations is shown by Eqs. (9) and (10), where f_1 and f_2 are any finite, single-

$$e = f_1 \frac{x-t}{\sqrt{LC}} + f_2 \frac{x+t}{\sqrt{LC}} \quad (9)$$

$$i = \frac{1}{\sqrt{L/C}} \left(f_1 \frac{x-t}{\sqrt{LC}} - f_2 \frac{x+t}{\sqrt{LC}} \right) \quad (10)$$

valued functions of the arguments $x - t/\sqrt{LC}$ and $x + t/\sqrt{LC}$, respectively. These are interpreted physically as traveling waves, the first travelling in the positive x direction with the speed $1/\sqrt{LC}$ and the second travelling in the negative x direction at the same speed. Substitution of the values for L and C for any configuration of conductors yields the velocity $1/\sqrt{LC}$, which equals $1/\sqrt{\mu\epsilon}$.

The quantity $\sqrt{L/C}$ has the dimensions of ohms, termed the characteristic impedance Z_0 of the line, as shown by Eq. (11). Thus, Z_0 is a real quantity

$$Z_0 = \sqrt{L/C} = \sqrt{\mu/\epsilon} F_g \quad (11)$$

(a resistance) and is equal to the wave impedance of the insulating medium, $\sqrt{\mu/\epsilon}$, multiplied by the geometrical factor, F_g , characteristic of the particular configuration of conductors. For the traveling waves of voltage and current, Eqs. (9) and (10), the ratio of voltage to current of the forward-traveling wave is Z_0 ; that of the backward-traveling wave is $-Z_0$.

In Fig. 5a, a source of electrical energy is connected at one end of a transmission line and an electrical load is connected to the other. Electromagnetic energy is propagated from the sending end to the receiving end, and a portion of the energy is reflected back toward the sending end if the load impedance Z_R is different from the characteristic impedance Z_0 of the line. If Z_R equals Z_0 , there is no reflection of energy at the load, and in Eqs. (9) and (10), the function f_2 , representing leftward-traveling energy, is absent. This is the condition desired when the purpose of the line is to deliver energy from the source of the load. The sending-end impedance of the line is then equal to Z_0 .

In addition to impedance matching to reduce reflections (echoes) along a transmission line, it is also necessary to minimize signal distortion, which consists of amplitude and phase (delay) distortion. If the line attenuation is frequency-dependent, then a signal consisting of a group of different-frequency components will undergo amplitude distortion due to the unequal attenuation of each component of the signal. Similarly, if the velocity of propagation along the line is frequency-dependent, then a delay in phase

of each component will result in associated phase distortion of the signal.

Signal distortion can be minimized by the use of line loading, which is the addition of series impedances along the line and which is used to adjust the line parameters to obtain the so-called distortionless condition. Under distortionless operation the attenuation and velocity of propagation are independent of frequency. For a discussion of the distortionless line See TRANSMISSION LINES.

Instead of loading a line, one may employ equalizing circuits to compensate for the phase distortion along the line.

Short sections of transmission line are sometimes used to provide low-loss reactive impedances and resonant circuits at high frequencies. This is done by open-circuiting or short-circuiting the receiving end of the line to provide complete reflection of the incident energy. A short-circuited low-loss line provides the sending-end impedance shown by Eq. (12),

$$Z_s = jZ_0 \tan(2\pi fl/v) \quad (12)$$

measured in ohms. When l equals $v/4f$, the line is a quarter wavelength long, the argument $(2\pi fl/v)$ of the tangent function in Eq. (12) is $\pi/2$, and Z_s approaches an infinite value. In actual practice, losses keep Z_s to a finite value. However, at high frequencies the quarter wavelength is short and the losses are small, and such a short-circuited quarter-wave section can be used successfully as a low-loss insulator. Such a section is a resonant one and can be used as a substitute for a parallel-resonant LC circuit, for example, a tank circuit for a high-frequency oscillator. At low frequencies the required quarter wavelength is so large that losses impair the performance; also, the length becomes inconveniently great. At a frequency lower than $v/4l$, the sending-end impedance of the short-circuited line is inductive, and at frequencies between $v/4l$ and $v/2l$ the impedance is capacitive. This provides the possibility of using sections of short-circuited line as reactive elements in circuits.

Walter C. Johnson

Bibliography. F. A. Benson and T. M. Benson, *Fields, Waves, and Transmission Lines*, 1991; C. W. Davidson, *Transmission Lines for Communications*, 2d ed., 1989; C. T. A. Johnk, *Engineering Electromagnetic Fields and Waves*, 2d ed., 1987; J. A. Kong, *Electromagnetic Wave Theory*, 3d ed., 2000; P. C. Magnusson et al., *Transmission Lines and Wave Propagation*, 2000; S. Ramo, J. R. Whinnery, and T. Van Duzer, *Fields and Waves in Communication Electronics*, 3d ed., 1993; W. Sinnema, *Electronic Transmission Technology: Lines, Waves, and Antennas*, 2d ed., 1988.

Electromagnetism

The branch of science dealing with the observations and laws relating electricity to magnetism. Electromagnetism is based upon the fundamental observations that a moving electric charge produces a

magnetic field and that a charge moving in a magnetic field will experience a force.

The magnetic field produced by a current is related to the current, the shape of the conductor, and the magnetic properties of the medium around it by Ampère's law. *See* AMPÈRE'S LAW.

The magnetic field at any point is described in terms of the force it exerts on a moving charge at that point. The electrical and magnetic units are defined in terms of the ampere, which in turn is defined from the force of one current upon another.

The association of electricity and magnetism is also shown by electromagnetic induction, in which a changing magnetic field sets up an electric field within a conductor and causes the charges to move in the conductor. *See* EDDY CURRENT; ELECTRICITY; ELECTROMAGNET; ELECTROMAGNETIC INDUCTION; FARADAY'S LAW OF INDUCTION; HALL EFFECT; INDUCTANCE; LENZ'S LAW; MAGNETISM; MAXWELL'S EQUATIONS; RELUCTANCE.

Kenneth V. Manning

Electrometallurgy

The branch of process metallurgy dealing with the use of electricity for smelting or refining of metals. The electrochemical effect of an electric current brings about the reduction of metallic compounds, and thereby the extraction of metals from their ores (electrowinning) or the purification of the metals (electrorefining).

In other metallurgical processes, electrically produced heat is utilized in smelting, refining, or alloy manufacturing. For a discussion of electrothermics, that is, the theory and applications of electric heating to metallurgy, *see* ELECTRIC FURNACE; ELECTRIC HEATING; ELECTROCHEMICAL PROCESS; STEEL MANUFACTURE.

Electrowinning. This metallurgical process involves the recovery of a metal, usually from its ore, by dissolving a metallic compound in a suitable electrolyte and reducing it electrochemically through passage of a direct electric current.

Following beneficiation and, sometimes, chemical pretreatment, the metal-bearing constituent of the ore is dissolved in an aqueous solution or in a molten salt. The electrolysis of the purified electrolyte with direct current yields the reduced metal at the cathode (negative electrode of an electrolytic cell), and the nonmetal is the oxidation product at the anode (positive electrode). Since the process is very selective, the purity of the electrowon metal is high, and no further refining is needed.

Water is a suitable solvent for electrolysis of metals less active than zinc and manganese, and sulfuric acid is the preferred leaching agent. The high mobility of the solvated proton provides high-conductivity solutions, and the sulfate ion is electrochemically inert; acid is regenerated by anodic oxidation and recycled for leaching.

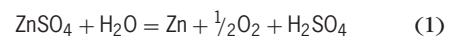
Zinc. The process of electrowinning of zinc was developed in 1915 for treating complex ores not amenable to thermal processing, and it is presently

the process of choice for most new plants. Changes in the chemistry have been minor, but steady progress has been made in the areas of process instrumentation and control, automation, and pollution-free operation.

The zinc sulfide ores are concentrated, roasted, and leached in several stages with sulfuric acid, and the electrolyte is purified. Lead and silver are insoluble through the acid leach; iron precipitates by oxidation in the neutral leach, and it coprecipitates impurities such as arsenic, antimony, and germanium. Other impurities more noble than zinc, such as copper, cadmium, cobalt, and nickel, are removed by galvanic precipitation on zinc dust.

The reduction potential of zinc is far more cathodic (-0.763 V) than that of the hydrogen ion, and its efficient electrodeposition depends upon maintaining a high hydrogen overvoltage. Metals more noble than zinc deposit on the cathode, and most promote hydrogen evolution and zinc corrosion. This reduces the yield of zinc, which is expressed by the current efficiency—the ratio of the amount of electricity (coulombs) theoretically required to yield a given quantity of metal to the amount actually consumed. Current efficiencies equal to 90% are achieved by bringing the concentration of the most harmful impurities such as germanium, arsenic, and antimony down to 0.01 mg/liter (1.3×10^{-6} oz/gal). *See* OVERVOLTAGE.

The zinc electrolyte, containing 120–160 g/liter (16–21 oz/gal) of zinc as sulfate, is circulated in lead-lined or plastic-lined [for example, poly(vinyl chloride)] concrete or wooden tanks in which the electrolysis proceeds between vertically suspended electrodes. The insoluble anodes are lead that contains small amounts of silver, and the cathodes are aluminum sheets from which the zinc deposits are stripped every 24–28 h. The submerged cathode area varies from 1.3 to 2.6 m². The current density varies from 300 to 600 A/m². From the Gibbs free-energy change for chemical reaction (1), a reversible voltage



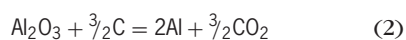
of 2 V can be calculated. The operating cell voltage is about 3.5 V; the difference is due to the anodic and cathodic overvoltages and the ohmic voltage drop in the electrolyte and at the electrode contacts. The power consumption averages 3.2 kWh/kg zinc for modern plants. This relatively high energy cost is acceptable since special high-grade zinc is produced (99.995%). *See* ZINC METALLURGY.

Copper. Copper electrowinning is similar to zinc electrowinning, but the pyrometallurgical reduction is preferred for sulfide ores. Most of the electrowon copper comes from the direct leaching of low-grade sulfide ores, or from oxide ores. The recent development of liquid ion-exchange reagents specially for copper makes possible an upgrading, through solvent extraction of dilute solutions, from dump or in-place leaching in order to obtain a concentrated solution suitable for electrolysis. The electrowinning of copper is carried out in sulfate solutions containing

from 20 to 40 g/liter (2.7–5.4 oz/gal) of copper and about 100 g/liter (13 oz/gal) of acid at current densities from 160 to 200 A/m² (15–19 A/ft²). The anodes are made of lead, alloyed with antimony, and copper is deposited on copper starting sheets.

Aluminum. Electrowinning has been the only process used for the commercial production of aluminum since it was developed independently by C. M. Hall and by P. L. T. Héroult in 1885. Alumina extracted from bauxite ores and purified by digestion in a caustic soda solution is dissolved in a mixture of cryolite (AlF₃ · 3NaF) and other fluorides (for example, CaF₂ and LiF). The molten-salt electrolyte floats on the molten aluminum, which is the cathode, and carbon anodes dip into the bath from above.

The carbon reacts with the oxygen produced by anodic oxidation. The overall cell reaction (2) corre-



sponds to a theoretical reversible voltage of 1.2 V, but the cell operates under a total voltage of about 4.2 V. The heat generated by irreversible electrochemical and ohmic losses maintains the cell temperature at about 960°C (1760°F).

Current densities are about 1 A/cm², with total currents of 100,000 A for a modern cell. Alumina is periodically added to the melt to maintain its concentration between 2 and 5%. The current efficiency averages 89%; the main losses are attributed to the recombination of aluminum by reduction of carbon dioxide to carbon monoxide. The energy consumption is about 14 kWh/kg aluminum. See ALUMINUM.

Magnesium. The electrowinning of magnesium, another well-established process, accounts for about 80% of the metal output. The electrolyte is a mixture of chlorides of potassium, sodium, calcium, and magnesium. Since the reduced metal has a lower specific gravity than the electrolyte has, it floats to the surface; the cells are designed so that it does not come in contact with air or with the chlorine gas produced at the anode.

Other metals. Cadmium, and in smaller quantities, thallium and indium, are recovered by electrolysis in sulfate solutions. Cobalt is also electrowon by a similar process, and some of the nickel production is achieved by electrolysis in sulfate-chloride solutions. Pure metallic chromium and manganese are produced by electrowinning; these metals are rather difficult to electrodeposit, and the acid produced at the anode must be kept away from the cathode by a diaphragm. Antimony is electrowon from a sodium sulfide electrolyte, and gallium is produced by electrolysis of a caustic solution.

Metals more active than manganese cannot be reduced cathodically from aqueous solutions, but most can be electrowon from mixtures of molten salts. Usually, a compound of the metal is dissolved in salts of still more active metals in order to improve the conductivity and lower the melting point of the melt. Solid metallic deposits are difficult to recover, and all the commercially successful processes yield the liquid metal. The solubility and reactivity of the reduced

metal in the melt is often important, and the cells must be designed to protect the reduced metal from contact with the anodic reaction product or with air. The cell electrolyte cannot be easily recycled, and the feed must be carefully purified and dehydrated.

All the sodium metal produced today comes from electrolysis of a molten chloride mixture between a steel cathode and a graphite anode separated by a perforated steel diaphragm. Lithium is produced by a similar process. The other alkali and alkaline-earth metals, as well as the rare earths, can be electrowon from molten chlorides. Processes have been developed for electrowinning titanium, tantalum, and niobium in molten salts, and uranium, thorium, and zirconium have been obtained in the laboratory. Beryllium and boron are also produced by electrolysis.

Electrorefining. This is a purification process in which an impure metal anode is dissolved electrochemically in a solution of a salt of the metal being refined; the pure metal is recovered by electrodeposition at the cathode.

Electrorefining is the most economical method for securing the high purity required for many uses of nonferrous metals. Since the same electrochemical reaction proceeds in opposite directions at the anode and the cathode, the overall chemical change is a small change in the activity of the metal at the two electrodes, and the reversible cell voltage is practically zero. The anodic and cathodic overvoltages, required for an economical reaction rate, are small, and the operating cell voltage consists mainly of the ohmic voltage drop through the electrolyte and at the electrode contacts. The power consumption is moderate, and electrorefining is less sensitive to the cost of electric power than other electrometallurgical processes are.

Electrochemical refining is a more efficient purification process than other chemical methods are because it is very selective. In particular, for metals such as copper, silver, gold, and lead, the operating electrode potential is close to the reversible potential, and a sharp separation is accomplished, both at the anode, where more noble metals do not dissolve, and at the cathode, where more active metals do not deposit.

The first commercial electrolytic copper refinery was established in 1871, and today most of the copper is electrorefined. Minor quantities of some impurities lower the electrical conductivity of copper markedly, and the refining process ensures that the metal will meet the specifications of the electrical industry. Furthermore, silver, gold, and other precious metals are removed by the electrolytic refining, and their recovery is thus an economic asset of the process.

The impure copper is cast in the shape of anodes approximately 3.0 × 3.3 ft (0.9 × 1.0 m) and 1.4–1.8 in. (3.5–4.5 cm) thick. The starting cathodes are pure copper sheets, prepared by electrodeposition. The electrodes are placed vertically in cells which are lead-lined or plastic-lined wooden or concrete tanks. The anodes are replaced at regular intervals of

20–28 days, and two successive cathodes are produced during the same period. The solution contains about 45 g/liter (6 oz/gal) of copper as copper sulfate and approximately 200 g/liter (26 oz/gal) of sulfuric acid. The temperature is maintained at 55–60°C (131–140°F) to lower the resistance of the electrolyte which is circulated through the cell. The current density is usually 200–250 A/m², and the current efficiency averages 95%. The cell voltage is 0.2–0.3 V, and the power consumption varies from 0.18 to 0.25 kWh/kg copper.

Metals less noble than copper, such as iron, nickel, cobalt, zinc, and manganese, dissolve from the anode. Lead is oxidized, and it precipitates as lead sulfate. Other impurities such as arsenic, antimony, and bismuth partly remain as insoluble compounds in the slimes and partly dissolve as complexes in the electrolyte. Precious metals do not dissolve at the potential of the anode, and they remain as metals in the slimes, along with insoluble selenium and tellurium compounds and some metallic copper particles that fall from the anode. Metals less noble than copper do not deposit at the cathode. In particular, nickel and arsenic accumulate in the electrolyte, and their concentration is controlled by partial withdrawal and purification of the solution.

Nickel is electrorefined in a sulfate-chloride electrolyte. Its anodic dissolution requires a significant overvoltage, and both copper and nickel dissolve at the anode potential. To prevent plating of the dissolved copper at the cathode, a diaphragm cell is used, and the anodic solution is circulated through a purification circuit before entering the cathode compartment. Cobalt is electrorefined by a similar process.

Very high-purity lead is produced by electrorefining in which a fluosilicate electrolyte is used, and a sulfate bath is used for purifying tin. Silver is electrorefined in a copper and silver nitrate electrolyte, and gold is refined by electrolysis in a chloride bath.

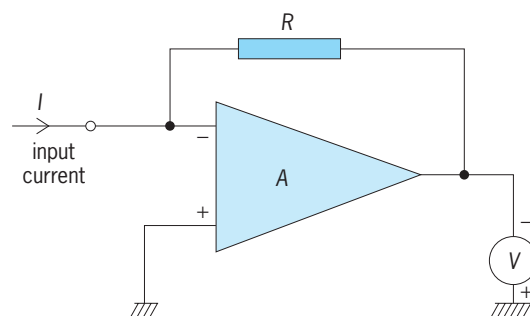
The electrorefining of many metals can be carried out with molten-salt electrolytes, but the processes are usually expensive, and the only industrial application is the electrorefining of aluminum by the “three-layer” process. The density of the electrolyte is adjusted so that a pure molten aluminum cathode will float on the molten salt, which in turn will float on an anode consisting of a molten mixture of impure aluminum and copper. See PYROMETALLURGY.

Paul Duby

Bibliography. A. T. Kuhn, *Industrial Electrochemical Processes*, 1971; A. T. Kuhn (ed.), *Techniques in Electrochemistry, Corrosion, and Metal Finishing: A Handbook*, 1987; D. Pletcher and F. C. Walsh, *Industrial Electrochemistry*, 2d ed., 1993.

Electrometer

A highly sensitive instrument which measures all or some of the following variables: current, charge, voltage, and resistance. There are two classes of electrometers, mechanical and electronic. The me-



A common arrangement for an electronic electrometer.

chanical instruments have been largely replaced by electronic types.

Mechanical. These instruments, with origins in the nineteenth century and earlier, rely for their operation on the mechanical forces associated with electrostatic fields. The negligible current drawn by many electrometers contributes to their high sensitivity, particularly in the context of charge and current measurement.

In common with other electrostatic instruments, electrometers are classified as either attracted-disk or symmetrical types. Attracted-disk instruments, in which the attractive force between two plates, with a potential difference between them, is measured in terms of the fundamental units of mass and length, are sometimes termed absolute electrometers. Attracted-disk instruments are widely used as electrostatic voltmeters for measuring potentials greater than 1 kilovolt. The quadrant electrometer was a particularly sensitive instrument based on a horizontal vane suspended by a torsion fiber so that it could rotate between fixed electrodes. See ELECTROSCOPE; ELECTROSTATICS; VOLTMETER.

Electronic. Electronic types utilize some form of electronic amplifier, typically an operational amplifier with a field-effect-transistor input stage to minimize the input current. In the illustration the amplifier *A* has a current measuring resistor *R* providing negative feedback. The unknown input current *I* can be found from the voltmeter reading *V*, as $I = V/R$. Charge *Q* can be measured if *R* is replaced by a capacitor *C*. For this case, $Q = VC$. See AMPLIFIER; TRANSISTOR.

In the most sensitive applications, problems arise due to drift in the amplifier characteristics and electrical noise present in the circuit components. To obviate these effects, electrometers employing a vibrating capacitor or varactor diodes are used. The signal to be measured is converted to an alternating-current (ac) signal and subsequently amplified by an ac amplifier which is less susceptible to drift and noise. The amplified signal is finally reconverted to direct current (dc). See ELECTRICAL NOISE; VARACTOR.

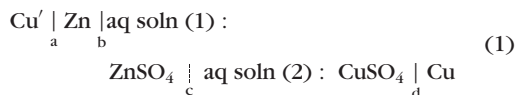
The most sensitive electrometers commercially available can resolve currents of 0.2 femtoampere under favorable circumstances. See CURRENT MEASUREMENT; ELECTRICAL MEASUREMENTS; VOLTAGE MEASUREMENT.

R. W. J. Barker

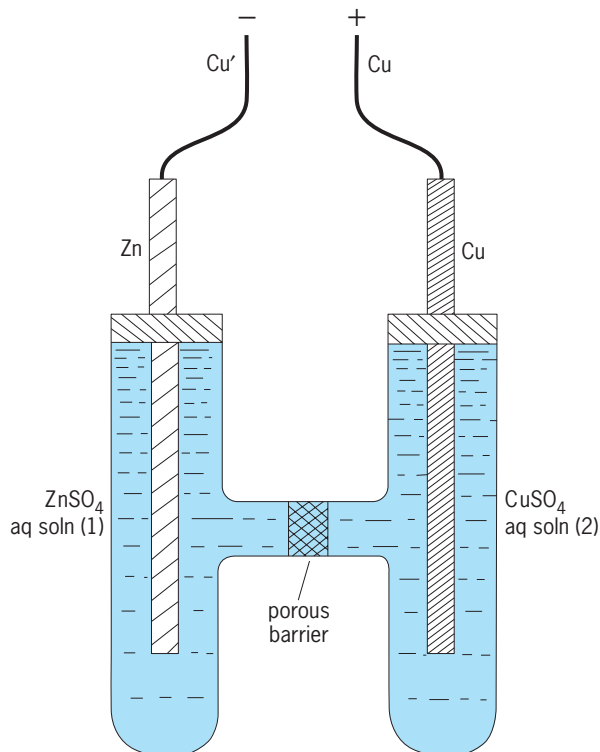
Electromotive force (cells)

The voltage or electric potential difference across the terminals of a cell when no current is drawn from it. The electromotive force (emf) is the sum of the electric potential differences produced by a separation of charges (electrons or ions) that can occur at each phase boundary (or interface) in the cell. The magnitude of each potential difference depends on the chemical nature of the two contacting phases. Thus, at the interface between two different metals, some electrons will have moved from the metal with a higher free energy of electrons to the metal with a lower free energy of electrons. The resultant charge separation will produce a potential difference, just as charge separation produces a voltage across a capacitor; at equilibrium this exactly opposes further electron flow. Similarly, potential differences can be produced when electrons partition across a metal|solution interface or metal|solid interface, and when ions partition across a solution|membrane|solution interface.

Origin. How a cell emf is composed of the sum of interfacial potential differences is shown by the Daniell cell (see **illus.**) and scheme (1);



aq soln denotes an aqueous solution, a solid line indicates a phase boundary, and the broken line a porous barrier permeable to all the ions in the adjacent solutions. The copper connector to the zinc electrode is



Cross-section diagram of Daniell cell. (After W. J. Moore, *Physical Chemistry, 5th ed., Longman, 1974*)

denoted Cu'. The barrier prevents physical mixing of the zinc sulfate (ZnSO₄) and copper sulfate (CuSO₄) solutions.

The cell emf is the open-circuit (that is, zero current) potential difference measured between the two Cu leads (any potential-measuring device must ultimately measure the potential difference between two chemically identical phases, in this example the Cu and Cu' phases). An expression can be written for the cell emf in terms of the electrochemical potentials of the species involved in the interfacial equilibria. The electrochemical potential $\bar{\mu}_j^\alpha$ of species j in phase α is defined by Eq. (2),

$$\bar{\mu}_j^\alpha = \mu_j^\alpha + z_j F \phi^\alpha \quad (2)$$

where μ_j^α is the chemical potential of species j in phase α , z_j is the number of electronic charges (with sign) on the species, F is the Faraday constant (96485.3 coulombs for Avogadro's number of electrons), and ϕ^α is the inner potential of phase α ; $z_j F \phi^\alpha$ is the electrical work required to move $z_j F$ coulombs of charge into phase α from infinite distance in a vacuum. For a species in solution, Eq. (2) can be written as Eq. (3),

$$\bar{\mu}_j^\alpha = \mu_j^{\circ\alpha} + RT \ln a_j^\alpha + z_j F \phi^\alpha \quad (3)$$

where $\mu_j^{\circ\alpha}$ is the standard chemical potential of species j in phase α , R is the gas constant (8.3145 joules/mole/degree Kelvin), T is the temperature in Kelvin, and a_j^α is the activity of species j in phase α . For a neutral species (that is, when $z_j = 0$), Eq. (4) applies.

$$\bar{\mu}_j^\alpha = \mu_j^\alpha = \mu_j^{\circ\alpha} + RT \ln a_j^\alpha \quad (4)$$

For a fully ionized solution species at low concentrations, a_j^α is well approximated by the concentration c_j^α ; the activity of a pure phase (for example, zinc and copper in the Daniell cell) is defined as unity.

The value of $\bar{\mu}_j$ must be identical for any species j equilibrated between two different phases. Thus for equilibrium of species j between phases α and β , Eq. (5) applies.

$$\bar{\mu}_j^\alpha = \bar{\mu}_j^\beta \quad (5)$$

More generally, for any interfacial equilibrium between phases α and β , Eq. (6) applies; v_j is the sto-

$$\sum_i v_j \bar{\mu}_j^{\text{phase}} = 0 \quad (6)$$

ichiometric number of species j (v_j is positive for products and negative for reactants), and $\bar{\mu}_j^{\text{phase}}$ is the electrochemical potential of species j in the phase in which it is located. Thus, for the interfacial equilibrium $\text{Cu}^{2+} + 2e^- \rightleftharpoons \text{Cu}$ at interface d in scheme (1), Eq. (6) becomes Eq. (7).

$$\bar{\mu}_{\text{Cu}^{2+}}^{\text{soln(2)}} + 2\bar{\mu}_{e^-}^{\text{Cu}} = \bar{\mu}_{\text{Cu}}^{\text{Cu}} \quad (7)$$

Similar equations can be used to quantify the potential differences across each interface in

scheme (1), assuming that the potential difference across interface c is zero (a good approximation when the separator is equally permeable to all the ions in the two adjacent phases). For interfaces a, b, and d of scheme (1), Eqs. (8) can be written.

Interface a (partitioning of electrons between Cu' and Zn phases):

$$\bar{\mu}_{e^-}^{\text{Cu}'} = \bar{\mu}_{e^-}^{\text{Zn}} \quad (8a)$$

Interface b (equilibrium: $\text{Zn}^{2+} + 2e^- \rightleftharpoons \text{Zn}$):

$$\bar{\mu}_{\text{Zn}^{2+}}^{\text{soln}(1)} + 2\bar{\mu}_{e^-}^{\text{Zn}} = \bar{\mu}_{\text{Zn}}^{\text{Zn}} \quad (8b)$$

Interface d (equilibrium: $\text{Cu}^{2+} + 2e^- \rightleftharpoons \text{Cu}$)

$$\bar{\mu}_{\text{Cu}^{2+}}^{\text{soln}(2)} + 2\bar{\mu}_{e^-}^{\text{Cu}} = \bar{\mu}_{\text{Cu}}^{\text{Cu}} \quad (8c)$$

Rearranging Eqs. (8) gives Eq. (9).

$$\begin{aligned} \bar{\mu}_{e^-}^{\text{Cu}} - \bar{\mu}_{e^-}^{\text{Cu}'} &= \bar{\mu}_{e^-}^{\text{Cu}} - \bar{\mu}_{e^-}^{\text{Zn}} = \\ &= \frac{1}{2} [\bar{\mu}_{\text{Cu}}^{\text{Cu}} - \bar{\mu}_{\text{Cu}^{2+}}^{\text{soln}(2)}] - \frac{1}{2} [\bar{\mu}_{\text{Zn}}^{\text{Zn}} - \bar{\mu}_{\text{Zn}^{2+}}^{\text{soln}(1)}] \quad (9) \end{aligned}$$

Since $\mu_{e^-}^{\text{Cu}} = \mu_{e^-}^{\text{Cu}'}$ and for an electron $z_{e^-} = -1$, Eq. (3) can be used to write Eq. (10).

$$\bar{\mu}_{e^-}^{\text{Cu}} - \bar{\mu}_{e^-}^{\text{Cu}'} = -F(\phi^{\text{Cu}} - \phi^{\text{Cu}'} \quad (10)$$

The cell potential, E_{cell} , is simply the sum of all the equilibrium interfacial potential differences, as in Eq. (11), where it has been assumed that $\phi^{\text{soln}(2)} - \phi^{\text{soln}(1)} = 0$. Combining Eqs. (3), (10) and (11) (and

$$\begin{aligned} E_{\text{cell}} &= (\phi^{\text{Cu}} - \phi^{\text{soln}(2)}) + (\phi^{\text{soln}(1)} - \phi^{\text{Zn}}) \\ &+ (\phi^{\text{Zn}} - \phi^{\text{Cu}'} = (\phi^{\text{Cu}} - \phi^{\text{Cu}'} \quad (11) \end{aligned}$$

noting that the copper and zinc ions have a charge of 2+ and that the activity of the metal phases is unity) gives Eq. (12),

$$\begin{aligned} E_{\text{cell}} &= (\phi^{\text{Cu}} - \phi^{\text{Cu}'} = -\frac{1}{F} [\bar{\mu}_{e^-}^{\text{Cu}} - \bar{\mu}_{e^-}^{\text{Cu}'}] \\ &= \frac{1}{2F} [\mu_{\text{Cu}^{2+}}^{\text{soln}(2)} - \mu_{\text{Cu}}^{\circ} + \mu_{\text{Zn}}^{\circ} - \mu_{\text{Zn}^{2+}}^{\text{soln}(1)}] \\ &+ \frac{RT}{2F} \ln a_{\text{Cu}^{2+}}^{\text{soln}(2)} - \frac{RT}{2F} \ln a_{\text{Zn}^{2+}}^{\text{soln}(1)} \quad (12) \end{aligned}$$

where E_{cell} is the potential difference between the two copper leads [scheme (1)] measured at open circuit (that is, when no current passes through the cell) and is therefore the emf of the cell. Dropping the superscripts denoting the solutions, Eq. (12) can be rewritten as Eq. (13),

$$\begin{aligned} E_{\text{cell}} &= \left[E_{\text{Cu}^{2+}|\text{Cu}}^{\circ} + \frac{RT}{2F} \ln a_{\text{Cu}^{2+}} \right] \\ &- \left[E_{\text{Zn}^{2+}|\text{Zn}}^{\circ} + \frac{RT}{2F} \ln a_{\text{Zn}^{2+}} \right] \quad (13) \end{aligned}$$

where $E_{\text{Cu}^{2+}|\text{Cu}}^{\circ}$ and $E_{\text{Zn}^{2+}|\text{Zn}}^{\circ}$ are termed the standard electrode potentials of the copper and zinc half cells. Equation (13) is an example of a Nernst equation, which relates cell emf to the activities of the cell constituents. See CHEMICAL EQUILIBRIUM; CHEMICAL THERMODYNAMICS; ELECTROMOTIVE FORCE (EMF); THERMODYNAMIC PRINCIPLES.

LIBRIUM; CHEMICAL THERMODYNAMICS; ELECTROMOTIVE FORCE (EMF); THERMODYNAMIC PRINCIPLES.

Half cells. It is convenient to describe any electrochemical cell in terms of half cells. A half cell consists of an oxidant (Ox) and reductant (Red) such that $\text{Ox} + ne^- \rightleftharpoons \text{Red}$; species Ox and Red are commonly referred to as a redox couple. The Daniell cell, for example, comprises the two half cells $\text{Cu}^{2+} + 2e^- \rightleftharpoons \text{Cu}$ (redox couple is $\text{Cu}^{2+}|\text{Cu}$) and $\text{Zn}^{2+} + 2e^- \rightleftharpoons \text{Zn}$ (redox couple is $\text{Zn}^{2+}|\text{Zn}$) with the half-cell potentials $E_{\text{Cu}^{2+}|\text{Cu}}^{\circ}$ and $E_{\text{Zn}^{2+}|\text{Zn}}^{\circ}$ given by Eqs. (14) and (15). Such expressions are useful only if the E°

$$E_{\text{Cu}^{2+}|\text{Cu}} = E_{\text{Cu}^{2+}|\text{Cu}}^{\circ} + \frac{RT}{2F} \ln a_{\text{Cu}^{2+}} \quad (14)$$

$$E_{\text{Zn}^{2+}|\text{Zn}} = E_{\text{Zn}^{2+}|\text{Zn}}^{\circ} + \frac{RT}{2F} \ln a_{\text{Zn}^{2+}} \quad (15)$$

value for each half cell is known. Values of E° can be assigned to any given half cell by arbitrarily specifying that $E_{\text{H}^+|\text{H}_2}^{\circ}$ [the E° for the standard hydrogen electrode (SHE) half cell, $\text{H}^+(\text{aq}, a = 1) + e^- \rightleftharpoons \frac{1}{2}\text{H}_2(\text{g}, 1 \text{ atm})$] is zero. The temperature dependence $dE_{\text{H}^+|\text{H}_2}^{\circ}/dT$ is also specified as zero. E° values versus the SHE for selected half cells are given in the **Table 1**. In principle, any E° (and its temperature dependence) can be measured directly versus the SHE, or another half cell whose electrode potential and temperature dependence have been determined. Such an electrode is termed a reference electrode. The reference electrode is a half cell designed so that its potential is stable and reproducible, and it neither contaminates nor is contaminated by the medium in which it is immersed. Two convenient reference electrodes commonly used in aqueous systems are the saturated calomel and the silver|silver chloride electrodes (see table). See ELECTROCHEMISTRY; ELECTRODE POTENTIAL; OXIDATION-REDUCTION; REFERENCE ELECTRODE.

E° values and thermodynamics. The E° values associated with the two half cells that make up a cell provide fundamental thermodynamic information about the chemical reaction between the redox couples of the two half cells. Returning to the example of the Daniell cell and combining Eqs. (12) and (13) yields

Selected standard electrode potentials at 25°C

Electrode reaction	E°/V
$\text{Li}^+(\text{aq}) + e^- \rightleftharpoons \text{Li}(\text{s})$	-3.045
$\text{Zn}^{2+}(\text{aq}) + 2e^- \rightleftharpoons \text{Zn}(\text{s})$	-0.763
$2\text{H}^+(\text{aq}) + 2e^- \rightleftharpoons \text{H}_2(\text{g})$	0
$\text{Hg}_2\text{Cl}_2(\text{s}) + 2e^- \rightleftharpoons 2\text{Hg}(\text{s}) + 2\text{Cl}^- (\text{sat aq KCl})$	0.246
$\text{AgCl}(\text{s}) + e^- \rightleftharpoons \text{Ag}(\text{s}) + \text{Cl}^-(\text{aq})$	0.222
$\text{Cu}^{2+}(\text{aq}) + 2e^- \rightleftharpoons \text{Cu}(\text{s})$	0.337
$\text{Fe}(\text{CN})_6^{3-}(\text{aq}) + e^- \rightleftharpoons \text{Fe}(\text{CN})_6^{4-}(\text{aq})$	0.69
$\text{O}_2(\text{g}) + 4\text{H}^+(\text{aq}) + 4e^- \rightleftharpoons 2\text{H}_2\text{O}$	1.223
$\text{F}_2(\text{g}) + 2\text{H}^+ + 2e^- \rightleftharpoons 2\text{HF}(\text{aq})$	3.06

Eq. (16). The term on the left-hand side is the stan-

$$E_{\text{Cu}^{2+}|\text{Cu}}^{\circ} - E_{\text{Zn}^{2+}|\text{Zn}}^{\circ} = \frac{1}{2F} [\mu_{\text{Cu}^{2+}}^{\circ \text{soln}(2)} - \mu_{\text{Cu}}^{\circ} + \mu_{\text{Zn}}^{\circ} - \mu_{\text{Zn}^{2+}}^{\circ \text{soln}(1)}] \quad (16)$$

dard cell emf E_{cell}° (that is, the value of E_{cell} when all the cell constituents are in their standard states of unit activity), and the bracketed term on the right-hand side is $-\Delta G^{\circ}$, where ΔG° is the standard Gibbs energy change of the cell reaction $\text{Cu}^{2+} + \text{Zn} \rightarrow \text{Cu} + \text{Zn}^{2+}$. Hence Eq. (16) can be written as Eq. (17).

$$E_{\text{cell}}^{\circ} = -\frac{\Delta G^{\circ}}{2F} \quad (17)$$

More generally, for the two half-cell reactions $\text{Ox}_1 + n_1 e^- \rightleftharpoons \text{Red}_1$ and $\text{Ox}_2 + n_2 e^- \rightleftharpoons \text{Red}_2$, E_{cell}° is defined by Eq. (18). The cell reaction (19) involves the trans-

$$E_{\text{cell}}^{\circ} = E_{\text{Ox}_1|\text{Red}_1}^{\circ} - E_{\text{Ox}_2|\text{Red}_2}^{\circ} \quad (18)$$



fer of $n_1 n_2$ electrons, and Eq. (17) takes the general form (20), where $n = n_1 n_2$. If the cell is not in its standard state, the relation becomes Eq. (21).

$$E_{\text{cell}}^{\circ} = -\frac{\Delta G^{\circ}}{n_1 n_2 F} = -\frac{\Delta G^{\circ}}{nF} \quad (20)$$

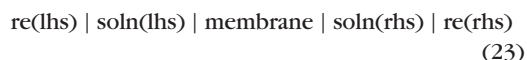
$$E_{\text{cell}} = -\frac{\Delta G}{nF} \quad (21)$$

Using the general thermodynamic relation $\Delta G^{\circ} = -RT \ln K_{\text{eq}}$, it is also possible to find the equilibrium constant K_{eq} of the cell reaction from the relation (22).

$$K_{\text{eq}} = \frac{[\text{Ox}_2]^{n_1} [\text{Red}_1]^{n_2}}{[\text{Ox}_1]^{n_2} [\text{Red}_2]^{n_1}} = \exp \left[-\frac{\Delta G^{\circ}}{RT} \right] = \exp \left[\frac{n_1 n_2 F E_{\text{cell}}^{\circ}}{RT} \right] \quad (22)$$

See FREE ENERGY.

Membrane potentials. When two ionic solutions of different composition are separated by a membrane, a potential difference can develop across the membrane. The complete cell requires electrical contacts with the solutions on each side of the membrane, accomplished with reference electrodes, as in scheme (23),



where re(lhs) and re(rhs) are the reference electrodes in the left and right solutions, respectively. The cell emf will be as in Eq. (24).

$$E_{\text{cell}} = E_{\text{re(rhs)}} + \phi^{\text{soln(rhs)}} - \phi^{\text{soln(lhs)}} - E_{\text{re(lhs)}} \quad (24)$$

The potential difference across the membrane, $\phi^{\text{soln(rhs)}} - \phi^{\text{soln(lhs)}}$, is a function of the membrane properties as well as the composition of the adjacent solutions. The simplest example occurs when only

one particular ion (such as the hydrogen ion, H^+) is transported across the membrane. The transported ion will tend to move from the side with high concentration to the side with low concentration; however, in the process of doing so, charge is separated across the membrane, producing a potential difference that exactly counters the ion transport. The membrane potential difference is given by Eq. (25),

$$\begin{aligned} \phi^{\text{soln(rhs)}} - \phi^{\text{soln(lhs)}} &= \frac{RT}{z_j F} \ln \frac{a_j^{\text{soln(lhs)}}}{a_j^{\text{soln(rhs)}}} \\ &\approx \frac{RT}{z_j F} \ln \frac{c_j^{\text{soln(lhs)}}}{c_j^{\text{soln(rhs)}}} \end{aligned} \quad (25)$$

where a_j and c_j are the activity and concentration of the transported ion. The cell emf is given by Eq. (26).

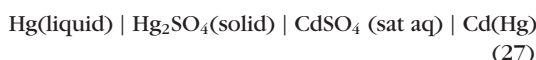
$$E_{\text{cell}} = E_{\text{re(rhs)}} + \frac{RT}{z_j F} \ln \frac{c_j^{\text{soln(lhs)}}}{c_j^{\text{soln(rhs)}}} - E_{\text{re(lhs)}} \quad (26)$$

If the concentration on one side, for example, $c_j^{\text{soln(lhs)}}$ is kept constant, E_{cell} will reflect any variation of the other concentration, $c_j^{\text{soln(rhs)}}$, and the system becomes a sensor for that ion.

As long as the membrane transports only a single ionic species, the potential difference is thermodynamically determined. As soon as the membrane transports more than one ionic species, irreversible thermodynamics come into play; the membrane potential difference will exhibit a complicated dependence on the concentrations and mobilities of the ions within the membrane as well as on their concentrations in the adjacent solutions. In certain circumstances, the response is still selective for one ion and the membrane electrode has analytical utility. The glass electrode used for measurement of pH, in which a thin glass membrane responds selectively to the H^+ ion, is an electrode of this type.

Potential differences known as liquid junction potentials arise where two different ionic solutions make contact through a permeable separator such as a glass frit. These interfere with thermodynamic measurements and can be minimized by interposing a salt bridge, that is, a concentrated immobilized solution of potassium chloride (KCl), or other salt in which the cation and anion have nearly identical mobilities, between the two solutions. Numerous theories based on irreversible thermodynamics have evolved to estimate the magnitude of liquid junction potentials. See BIOPOTENTIALS AND IONIC CURRENTS; ION-SELECTIVE MEMBRANES AND ELECTRODES.

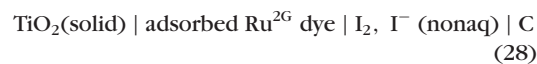
Uses. The emf produced by a single cell or a set of cells in series (a battery) is used as a dc power source for a wide array of applications, ranging from powering wristwatches to emergency power supplies. The emf of the saturated Weston cell (scheme 27;



sat aq = saturated aqueous solution) still serves as a high-level voltage reference for the National Institute of Standards and Technology. However, Josephson

arrays are now considered the most precise voltage references, and Zener diodes are used to produce reference voltages for many laboratory and field voltage measurement devices.

There are also a few photoelectrochemical cells with possible application in solar energy conversion, such as the Grätzel cell (scheme 28; nonaq = non-



aqueous solution), where the cell emf is produced by absorption of visible light.

The emf of a cell can also be used as an indicator of chemical composition. Devices that depend on the measurement of an open circuit-cell potential work well when the device is sensitive to a single analyte, but they are notoriously sensitive to interferences. An example is the so-called alkaline error that occurs when glass electrodes are used to measure very high pH values. For complex systems containing several species that can undergo electrochemical reaction, individual E° values can be determined using electroanalytical techniques that involve controlling the potential applied to a cell with measurement of the resultant current. Techniques such as polarography and cyclic voltammetry, for example, involve changing the potential of an indicator electrode and observing a wave or peak in the current at the redox potential of the species in solution; the height of the wave or peak indicates the concentration of that species. See BATTERY; ELECTROCHEMICAL TECHNIQUES; FUEL CELL; POLAROGRAPHIC ANALYSIS.

Cell voltages when current passes. E_{cell} values, which are thermodynamic quantities, do not give any information about the kinetics of the electrode reactions. When current passes through a cell, the cell voltage (that is, the potential difference between the terminals) will differ from the emf because of resistive losses within the cell (a function of cell design) and because the kinetics of the electrode reactions are not fast enough to sustain thermodynamic (nernstian) behavior at high rates of reaction. These kinetics are very dependent on temperature and the nature of the electrode material. When a cell is discharged through an external load, the cell voltage will be less than the emf; if the direction of current flow is reversed (using an external power supply), the cell voltage will be greater than the emf. In both cases the greater the current, the greater the deviation of the cell voltage from the emf. In some cases the electrode reactions are so slow that even the open-circuit cell voltage may not be a reliable measure of the emf.

For example, a cell comprising the two half cells $2\text{H}^+(\text{aq}) + 2\text{e}^- \rightleftharpoons \text{H}_2(\text{g})$ and $\text{O}_2(\text{g}) + 4\text{H}^+(\text{aq}) + 4\text{e}^- \rightleftharpoons 2\text{H}_2\text{O}$ (where g = gas), both under standard conditions, has a cell emf of 1.223 V at 25°C (see table). A highly sophisticated version of this cell, the hydrogen/oxygen fuel cell, is designed to minimize internal cell resistance, and has electrodes tailored to maximize the rates of the electrode re-

actions; nevertheless, the voltage under load is usually less than 0.8 V, with most of the voltage loss occurring at the oxygen electrode. When an external power supply is used to electrolyze water, the cell voltage required to produce H_2 and O_2 at a reasonable rate is ≥ 1.6 V; and this again depends critically on the cell design and the choice of electrodes.

Mary D. Archer; Stephen W. Feldberg
Bibliography. R. G. Compton and G. H. W. Sanders, *Electrode Potentials*, Oxford University Press, 1996; D. R. Crow, *Principles and Application of Electrochemistry*, 4th ed., Blackie Academic & Professional, 1994; A. R. Denaro, *Elementary Electrochemistry*, 2d ed., 1971; D. J. G. Ives and G. J. Janz, *Reference Electrodes: Theory and Practice*, Academic Press, 1961; J. Koryta, *Ion-Selective Electrodes*, Cambridge University Press, 1975; W. J. Moore, *Physical Chemistry*, 5th ed., Longman, 1974.

Electromotive force (emf)

A measure of the strength of a source of electrical energy. The term is often shortened to emf. It is not, of course, a force in the usual mechanical sense (and for this reason has sometimes been called electromotance), but it is a conveniently descriptive term for the agency which drives current through an electric circuit. In the simple case of a direct current I (measured in amperes) flowing through a resistor R (in ohms), Ohm's law states that there will be a voltage drop (or potential difference) of $V = IR$ (in volts) across the resistor. To cause this current to flow requires a source with emf (also measured in volts) $E = V$. More generally, Kirchhoff's voltage law states that the sum of the source emf's taken around any closed path in an electric circuit is equal to the sum of the voltage drops. This is equivalent to the statement that the total emf in a closed circuit is equal to the line integral of the electric field strength around the circuit. See ELECTRIC CURRENT; ELECTRIC FIELD; ELECTRICAL RESISTANCE; OHM'S LAW.

A source of emf E delivering a current I supplies electrical energy at a rate of EI (in watts). Therefore, the emf in a circuit is numerically equal to the work done in carrying a charge of 1 coulomb around the circuit.

Open-circuit emf. A source such as a battery has an internal resistance which causes an internal voltage drop when a current is drawn from it. Consequently the emf at the external terminals varies with the current. By extrapolating to zero current, the open-circuit emf may be derived. It is the open-circuit emf which is usually said to be the emf of a source; together with the value of the internal resistance, it enables the performance of the source in any circuit to be calculated.

Types of emf. An emf may be steady (direct), as for a battery, or time-varying, as for a charged capacitor discharging through a resistor. A class of time-varying sources of very great practical importance is that of generators of alternating (sinusoidal) emf. Emf's may be generated by a variety of physical, chemical, and

biological processes. Some of the more important are:

1. Electrochemical reactions, as used in direct-current (dc) batteries, in which the emf results from the reactions between electrolyte and electrodes. See BATTERY; ELECTROCHEMISTRY; ELECTROMOTIVE FORCE (CELLS).

2. Electromagnetic induction, in which the emf results from a change in the magnetic flux linking the circuit. This finds application in alternating-current rotary generators and transformers, providing the basis for the electricity supply industry. See ALTERNATING-CURRENT GENERATOR; ELECTROMAGNETIC INDUCTION; FARADAY'S LAW OF INDUCTION; TRANSFORMER.

3. Thermoelectric effects, in which a temperature difference between different parts of a circuit produces an emf. The main use is for the measurement of temperature by means of thermocouples; there are some applications to electric power generation. See THERMOCOUPLE; THERMOELECTRIC POWER GENERATOR; THERMOELECTRICITY.

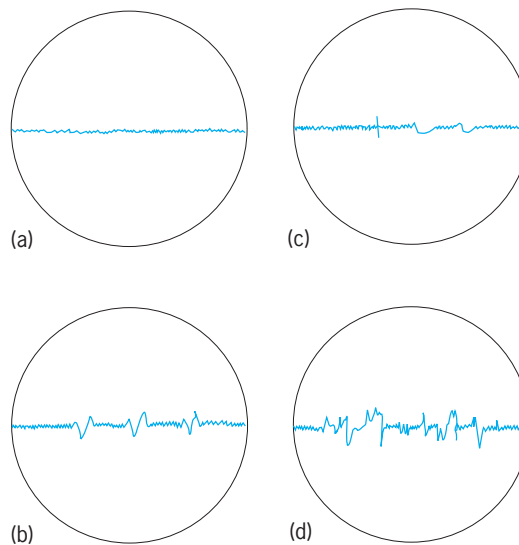
4. The photovoltaic effect, in which the absorption of light (or, more generally, electromagnetic radiation) in a semiconductor produces an emf. This is widely used for scientific purposes in radiation detectors and also, increasingly, for the generation of electric power from the Sun's radiation. See PHOTOVOLTAIC EFFECT; RADIOMETRY; SOLAR CELL.

5. The piezoelectric effect, in which the application of mechanical stress to certain types of crystal generates an emf. There are applications in sound recording, in ultrasonics, and in various types of measurement transducer. See MICROPHONE; PIEZOELECTRICITY; TRANSDUCER; ULTRASONICS. A. Earle Bailey

Electromyography

The detection and recording of electrical activity generated by muscle fibers. The basic elements of motor control in the body are the motor units which comprise motor neurons in the brainstem or spinal cord, their axons, and from ten to several hundred muscle fibers supplied by each motor neuron. Motor units vary in the size and properties of their motoneurons, the sizes and conduction velocities of their axons, the morphology of their nerve muscle junctions, and the structure and physiological properties of the muscle fibers supplied by each motor neuron. Voluntary activity recruits motor units in an orderly manner, beginning with smaller, twitch units with more slowly conducting axons recruited in weak contractions and progressing to much-harder-to-recruit, larger and faster-twitch motor units with more rapidly conducting axons.

Impulses originating in single motoneurons in response to various command signals from the central nervous system conduct to the periphery of the unit, normally causing all the muscle fibers in the unit to discharge. The electrical activity generated by the more or less synchronous discharges of all the mus-



Record of synchronous and asynchronous muscle discharge. (a) Normal resting muscle. (b) Partial contraction. (c) Degenerating muscle. (d) Dystrophic muscle.

cle fibers in the unit may be detected by recording electrodes on the skin surface or by needles inserted into the muscle (see *illus.*). Such potentials are called motor-unit action potentials and reflect the electrical activity generated by the whole motor unit.

Diseases affecting motor neurons are sometimes accompanied by spontaneous discharges of the axons (fasciculations) which, if near the skin, may be visible as twitches in the muscle. Additionally, degeneration of motor axons may leave some muscle fibers deprived of their normal innervation, some of which spontaneously fire. Such single muscle-fiber discharges are called fibrillations and are readily detected for diagnostic purposes by needle electrodes inserted into the muscle.

Electromyography may also be used to study primary muscle diseases such as the muscular dystrophies, and a wide variety of other metabolic inflammatory and congenital myopathies affecting the muscle fibers rather than motor neurons or their axons. See BIOPOTENTIALS AND IONIC CURRENTS; ELECTRODIAGNOSIS. William F. Brown

Electron

The negatively charged constituent of ordinary matter, responsible for the chemical properties of matter, and the carrier of electricity. The electron is the lightest known particle that possesses an electric charge. Its rest mass m_e is 9.1×10^{-31} kg (2.0×10^{-30} lb), about 1/1836 of the mass of the proton or neutron, which are, respectively, the positively charged and neutral constituents of ordinary matter. The rest mass of an electron can also be expressed as $0.511 \text{ MeV}/c^2$, where c is the speed of light. Electrons were discovered in 1895 by J. J. Thomson

in the form of cathode rays. The electron was the first elementary particle to be identified. *See* ATOMIC STRUCTURE AND SPECTRA; CATHODE RAYS; ELECTRIC CHARGE; ELECTRONVOLT; ELEMENTARY PARTICLE; NUCLEAR STRUCTURE.

Charge. The charge of the electron is $e = 1.602 \times 10^{-19}$ coulomb. The sign of the electron's charge is negative by convention, and that of the equally charged proton is positive. This is a somewhat unfortunate convention, because the flow of electrons in a metallic conductor is thus opposite to the conventional direction of the current. *See* CONDUCTION (ELECTRICITY).

The historical method of measuring the charge of the electron is the celebrated oil drop experiment of R. A. Millikan in 1909, in which the charges of droplets of oil in air are measured by finding the electric field which balances each drop against its weight. The weight of each drop is determined by observing its rate of free fall through the air, and using Stokes' formula for the viscous drag on a slowly moving sphere. The charges thus measured are integral multiples of e . For more precise values of e and m_e , *see* FUNDAMENTAL CONSTANTS.

Electrons and matter. Electrons are emitted in radioactivity (as beta rays) and in many other decay processes; for instance, the ultimate decay products of all mesons are electrons, neutrinos, and photons, the meson's charge being carried away by the electrons. The electron itself is completely stable, according to all available evidence. Electrons contribute the volume to ordinary matter; the volume of an atom is nearly all occupied by the cloud of electrons surrounding the nucleus, which occupies only about 10^{-13} of the atom's volume. Virtually all the mass of an atom, however, is in the nucleus. The chemical properties of ordinary matter are determined by the electron cloud. *See* BETA PARTICLES; CHEMICAL BONDING; QUANTUM CHEMISTRY; RADIOACTIVITY; VALENCE.

The electron obeys the Fermi-Dirac statistics, and for this reason is often called a fermion. One of the primary attributes of matter, impenetrability, results from the fact that the electron, being a fermion, obeys the Pauli exclusion principle; the world would be completely different if the lightest charged particle were a boson, that is, a particle that obeys Bose-Einstein statistics. *See* BOSE-EINSTEIN STATISTICS; EXCLUSION PRINCIPLE; FERMI-DIRAC STATISTICS.

Spin. Every elementary particle possesses an intrinsic angular momentum called its spin. The spin of the electron has the magnitude $\frac{1}{2}\hbar$, where \hbar is Planck's constant h divided by 2π . An electron thus has two spin states: spin up and spin down. To describe this, the nonrelativistic wave function is a two-component function, that is, a vector in a two-dimensional spin space; the two linearly independent vectors represent the two possible spin states. In 1928 P. A. M. Dirac derived the corresponding relativistic wave equation (Dirac equation). Here, the electron wave function must have four components; correspondingly, for a wave of given momen-

tum, there are four internal states. In addition to the two-valued spin coordinate, there is an energy coordinate; that is, for a momentum p , the energy can be $\pm\sqrt{(m_e c^2)^2 + (pc)^2}$, where c is the velocity of light. *See* ANGULAR MOMENTUM; ELECTRON SPIN; NONRELATIVISTIC QUANTUM THEORY; RELATIVISTIC QUANTUM THEORY; SPIN (QUANTUM MECHANICS).

Positron. The negative energy states were at first an embarrassment, for they extend downward indefinitely; an electron would cascade indefinitely downward in energy, radiating photons. Electrons obey the exclusion principle, however, and this conclusion can be avoided by assuming that in empty space all the negative energy states are already occupied, and so exclude any more electrons. A new process is possible now; such a negative energy electron, by absorbing energy, can go to a positive energy state. This leaves behind a hole in the sea of negative-energy electrons; the hole has a positive energy, because it represents a missing negative energy. In fact, such a hole has all the properties of an electron except that it appears to have a positive charge (because it represents a missing negative charge). This particle is the positron, discovered in 1932 by C. D. Anderson in a cloud-chamber study of cosmic radiation. *See* POSITRON.

The electron and the positron are on an equal footing; if one started with a Dirac wave equation for the positron, identifying electrons with holes in the negative-energy positron states, one would get an equivalent theory. The apparent dissymmetry inherent in the construction of the hole theory disappears from the results when the total charge, energy, and momentum of empty space is defined to be zero; actually, the dissymmetry can be avoided at all stages in the formalism of quantum field theory. *See* ANTIMATTER; QUANTUM FIELD THEORY; SYMMETRY LAWS (PHYSICS).

Magnetic moment. The electron has magnetic properties by virtue of its orbital motion about the nucleus of its atom and its rotation (spin) about its own axis. The intrinsic magnetic moment of the electron is predicted by the Dirac equation to be the value (Bohr magneton) shown in the equation below. The actual moment μ differs from μ_D

$$\mu_D = \frac{e\hbar}{2m_e}$$

by a small amount (anomalous magnetic moment) due to electromagnetic radiative corrections: $\mu = 1.0011 \mu_D$. This theoretical value, calculated using renormalized quantum field theory, agrees with the experimental value. *See* MAGNETON; QUANTUM ELECTRODYNAMICS.

Other leptons. The electron is the lightest of a family of elementary particles, the leptons. The other known charged leptons are the muon and the tau. These three particles differ only (as far as known) in mass; they have the same spin, charge, strong interactions (namely, none), and weak interactions. In a weak interaction a charged lepton is either unchanged (neutral weak current reaction) or changed

(charged weak current reaction) into an uncharged lepton, that is, a neutrino. In the latter case, each charged lepton (electron, muon, or tau) is seen to change only into the corresponding neutrino (ν_e , ν_μ , or ν_τ). See LEPTON; NEUTRINO; WEAK NUCLEAR INTERACTIONS.

Electrons are responsible for many phenomena of importance. Electrons moving through a conducting material constitute electricity. There is a magnetic field surrounding a current-carrying conductor, which makes possible electric motors and generators. Moving electrons are generally responsible for the production of light and electromagnetic radiation; examples include radio and television waves emitted by electrons moving in an antenna, microwaves generated in a magnetron in a microwave oven, and x-rays emitted by electrons passing through matter or electrons deflected in a magnetic field. Electron spin is also responsible for the behavior of magnetic materials. See ANTENNA (ELECTROMAGNETISM); BREMSSTRAHLUNG; ELECTRIC CURRENT; ELECTRICITY; ELECTROMAGNETIC RADIATION; GENERATOR; LIGHT; MAGNETIC MATERIALS; MAGNETISM; MAGNETRON; MOTOR; SYNCHROTRON RADIATION.

Electrons are responsible for chemical bonding in molecules. Additionally, the number of protons or electrons in an atom determines the elemental species of the atom, and thus the ranking of elements by atomic number (the number of electrons or protons in the atom) in the periodic table. See ATOMIC NUMBER; CHEMICAL BONDING; ELEMENT (CHEMISTRY); PERIODIC TABLE.

Electrons are a valuable probe of matter, as in an electron microscope. Electrons at high energies have a very short wavelength, and microscopes with atomic resolution can be constructed. See ELECTRON DIFFRACTION; ELECTRON MICROSCOPE.

Other information. It would be difficult to list all the articles wherein the electron forms an integral part of the discussion. The following articles and those referenced in the preceding discussion are intended to be merely a representative sample. See BAND THEORY OF SOLIDS; CHARGED PARTICLE OPTICS; COMPTON EFFECT; ELECTRON CAPTURE; ELECTRON CONFIGURATION; ELECTRON EMISSION; ELECTRON MOTION IN VACUUM; ELECTRONICS; EXCHANGE INTERACTION; FREE-ELECTRON THEORY OF METALS; MAGNETISM; PARTICLE ACCELERATOR; QUANTUM MECHANICS; RELATIVISTIC ELECTRODYNAMICS.

Charles J. Goebel; Alfred S. Schlachter

Electron affinity

The amount of energy release when an electron e at rest is captured by a species M , producing the negative ion M^- . The electron affinity of a species M can also be thought of as the ionization potential of the negative ion M^- . In these definitions, it is assumed that both M and M^- are in the lowest electronic, vibrational, and rotational state.

If the electron affinity of M is negative, the M^- ion is unstable with respect to decomposition into $M + e$. Most atoms have positive electron affinities, even though there is no net Coulomb attraction between the electron and the atom until the electron is close enough to be "a part of the atom." The simple rules of chemical valence provide a qualitative guide to the magnitude of electron affinities. Thus the noble gases, which have a filled outer electronic shell and are chemically inert, are not capable of binding an additional electron to form a negative ion. The largest electron affinities are possessed by the halogens, atoms which require only one additional electron to fill the valence shell. See HALOGEN ELEMENTS.

The major exception to this concept is that multiply charged negative ions—for example, O^{2-} , one of many such ions which are stable in solution—have not been observed in the gas phase. The ability to place more than one additional electron in the valence shell of a neutral atom or molecule arises from many particle electrostatic interactions with the ion's environment: either the solvent shell surrounding the ion in liquid solutions or the amorphous or crystalline region surrounding the ion in solids. See VALENCE.

Experimental methods. Accurate ionization potentials of the elements were known for a number of years before comparable data for electron affinities of the elements became available. In order to determine the ionization potential of an element, the element is vaporized, placed in an optical spectrometer, and observed for the onset of photoabsorption corresponding to the photoionization process, $h\nu + M \rightarrow M^+ + e$. The photon energy corresponding to the threshold wavelength for this process is the ionization energy of the species M . The analogous method for determination of an electron affinity is through observation of the threshold of the very similar photodetachment reaction, $h\nu + M^- \rightarrow M + e$. Again, the threshold wavelength for this process corresponds to the electron affinity of the species M . Unfortunately, it has not proved possible to produce the sufficiently large densities of negative ions required to utilize photoabsorption spectrometers to observe directly the threshold for the photodetachment process. Consequently, determination of accurate electron affinities lagged far behind the determination of accurate ionization potentials. However, an ion velocity modulation scheme has been devised that permits the direct observation of vibrational photoabsorption in ions; it is not generally applicable to the electron detachment process required to determine electron affinities.

The two major experimental advances that have enabled accurate electron affinity determinations have been the development of ion-beam techniques affording the detection of photoabsorption by the monitoring of photodetachment products and the availability of intense light sources in the form of lasers.

In experiments to determine electron affinities, negative ions are formed in an electrical

1 H 0.754201								2 He <0
3 Li 0.6180	4 Be <0	5 B 0.277	6 C 1.2629	7 N ≈0	8 O 1.4611215	9 F 3.399	10 Ne <0	
11 Na 0.54793	12 Mg <0	13 Al 0.441	14 Si 1.385	15 P 0.7465	16 S 2.077120	17 Cl 3.61272	18 Ar <0	
19 K 0.50147	20 Ca 0.043	31 Ga 0.30	32 Ge 1.233	33 As 0.81	34 Se 2.02069	35 Br 3.365	36 Kr <0	
37 Rb 0.48592	38 Sr 0.10	49 In 0.30	50 Sn 1.112	51 Sb 1.07	52 Te 1.9708	53 I 3.0591	54 Xe <0	
55 Cs 0.47163	56 Ba 0.18	81 Tl 0.2	82 Pb 0.364	83 Bi 0.946	84 Po 1.9	85 At 2.8	86 Rn <0	

21 Sc 0.188	22 Ti 0.079	23 V 0.525	24 Cr 0.677	25 Mn <0	26 Fe 0.151	27 Co 0.662	28 Ni 1.156	29 Cu 1.235	30 Zn <0
39 Y 0.307	40 Zr 0.426	41 Nb 0.893	42 Mo 0.746	43 Tc 0.55	44 Ru 1.05	45 Rh 1.137	46 Pd 0.557	47 Ag 1.302	48 Cd <0
57 La 0.5	72 Hf ≈0	73 Ta 0.322	74 W 0.815	75 Re 0.15	76 Os 1.1	77 Ir 1.565	78 Pt 2.128	79 Au 2.30863	80 Hg <0

Periodic table showing the best values for the electron affinities of the elements. All values are reported in electronvolts. The value <0 implies that the negative ion is unstable with respect to decomposition to an electron and a neutral atom. In general, the last decimal place in each value is uncertain. (After H. Hotop and W. C. Lineberger, *Binding energies in atomic negative ions: II, J. Phys. Chem. Ref. Data*, 14:731-750, 1985)

discharge, extracted through an aperture into a high-vacuum region, formed into a negative-ion beam, mass-analyzed, and intersected by an intense laser beam. The laser-beam-negative-ion-beam intersection takes place in a high-vacuum region where no collisions are likely. The occurrence of a photodetachment event is determined by detection of the photodetached electron.

Two experimental methods evolved which produce accurate electron affinities. In the first method the laser is tunable, and a search is made for the wavelength corresponding to the threshold for the photodetachment process. In this case the electron affinity is given directly by the threshold wavelength, and experimental methods have been developed that provide accuracies of 10^{-8} eV. In the second type of experiment, called photoelectron spectroscopy, a fixed-frequency laser (of known photon energy) is employed, and the electrostatic fields are used to determine the kinetic energy of the ejected electron. From simple energy conservation arguments, the electron affinity is then given by the photon energy less the kinetic energy of the ejected electron. This latter technique is quite general, but is limited in accuracy by the resolution of the electron energy analyzer (approximately 3×10^{-3} eV for state-of-the-art analyzers). See LASER; LASER SPECTROSCOPY; MASS SPECTROMETRY.

Periodic trends. These laser photodetachment studies dramatically improved knowledge of the electron affinities of the elements. The **illustration** is a periodic table showing values of the electron affinities of the elements. Most of the data shown here were obtained by using laser photodetachment methods. The periodic trends in electron affinities and the qualitative effects described earlier are immediately apparent. In addition, a number of more subtle trends are observable. For example, while it might be expected that filled-shell species such as the rare gases would not be capable of binding an additional electron, the illustration shows that half-filled shells (for example, N and P) also exhibit small or negative electron affinities. Again, this effect is the result of the fact that a half-filled valence shell is spherically symmetric and behaves somewhat as though it were a filled shell. A similar situation is also seen for half-filled *d* shells, as is the case for Mn, Tc, and Re in the transition metals.

These same techniques have provided a number of accurate electron affinity determinations for molecules and free radicals, and new insight into the structural and chemical properties of ions in the gas phase. See IONIZATION POTENTIAL; PERIODIC TABLE.

W. C. Lineberger

Bibliography. M. T. Bowers (ed.), *Gas Phase Ion Chemistry* 3, 1984; H. Hotop and W. C. Lineberger, Binding energies in atomic negative ions: II, *J. Phys. Chem. Ref. Data*, 14:731-750, 1985; I. N. Levine, *Physical Chemistry*, 4th ed., 1994; R. J. Saykally, Infrared laser spectroscopy of molecular ions, *Science*, 239:157-161, 1988.

Electron capture

A term that refers to two processes, one atomic and the other nuclear, in both of which an electron is captured: by a charged particle (an ion) passing through matter, when an electron is captured into an atomic orbital (in the atomic case); or by a nucleus (in the nuclear case), which results in radioactive decay of the nucleus.

Atomic process. The capture and loss of electrons is an important mechanism for slowing down fast positive ions (atoms with fewer electrons than their nuclear charge) passing through matter. Electron capture is significant generally when the velocity of the ion is similar to that of the target electrons to be captured. Fast ions passing through matter thus generally lose much of their energy near the end of their range passing through matter, as their velocity becomes similar to that of the bound target electrons.

An example of electron capture is shown in **Fig. 1**, which shows the cross section (the likelihood) for electron capture by a proton passing through lithium vapor. The reaction is $H^+ + Li \rightarrow H^0 + Li^+$ (or possibly $Li^{2+} + e^-$). The cross section (the likelihood of capturing an electron) becomes very small when the proton energy (which is proportional to the square

of the velocity) becomes large. See CHARGED PARTICLE BEAMS.

Nuclear process. The nuclear capture of electrons (K capture) occurs by a process quite different from atomic capture and leads to a nuclear reaction. An electron is captured by the nucleus; the electron is taken from the K shell (the electron shell nearest the nucleus) of the atom. The probability of electron capture by the nucleus depends on the amount of time the electrons spend near the nucleus, that is, on the size of the electron wave function at the nuclear center. Since to a very good approximation only electrons with zero orbital angular momentum have a wave function that is finite at the center, K capture generally occurs for capture of a K-shell electron. However, L capture is also possible. See RADIOACTIVITY.

Electron capture (K capture) generally occurs for a nucleus with too many protons (or too few neutrons). A K-shell electron from the atom is captured by the nucleus, creating a neutron and a neutrino, resulting in a more stable nucleus. The result of essentially transforming a proton into a neutron in the nucleus is to reduce the atomic number by one, resulting in a different element (since the number of electrons in an atom is equal to the number of protons in the nucleus, and the number of electrons or protons determines the element). Atomic mass is essentially unchanged by K capture, as the mass of a proton and neutron are essentially the same. The atom is missing a K-shell electron as a result of K capture; it stabilizes by emission of an x-ray or an Auger electron: an electron from an outer shell fills the vacancy in the K shell of the atom. Thus, K capture is a simultaneous atomic and nuclear event. See AUGER EFFECT.

K capture is favored when there is not enough energy for the nucleus to emit a positron, which requires 1.022 MeV. An example of a nucleus without sufficient energy to emit a positron is rubidium-

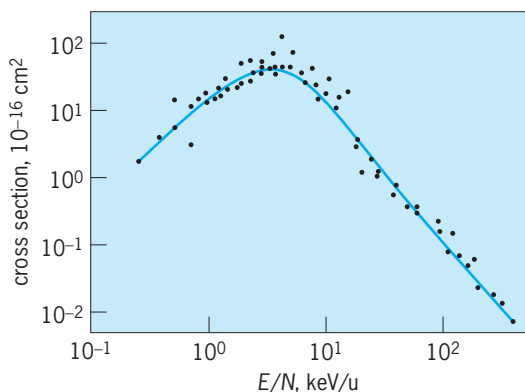


Fig. 1. Cross section (units of 10^{-16} cm^2) for capture of an electron by a proton passing through lithium vapor. Energy per nucleon (E/N) is in units of kiloelectronvolts per nucleon (keV/u). Data points from experiments and a curve fitted to the data are shown. (After T. J. Morgan et al., *Charge transfer of hydrogen ions and atoms in metal vapors*, *J. Phys. Chem. Ref. Data*, 14:971-1040, 1985).

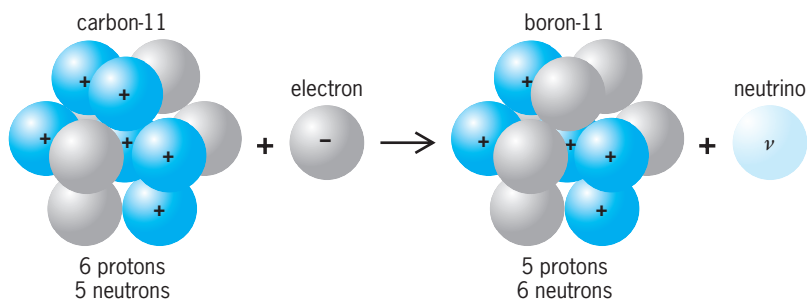


Fig. 2. Example of K capture: The reaction carbon-11 (6 protons and 5 neutrons) + electron \rightarrow boron-11 (5 protons and 6 neutrons) + neutrino. (Jefferson Lab, <http://education.jlab.org/glossary/electroncapture.html>)

83, which decays to krypton-83 by K capture. See POSITRON.

An example of K capture is shown in Fig. 2. Both carbon-11 and boron-11 have the same atomic mass (11); however, boron (atomic number 5) is one element before carbon (atomic number 6) in the periodic table.

Alfred S. Schlachter

Electron configuration

The orbital arrangement of an atom's electrons. Negatively charged electrons are attracted to a positively charged nucleus to form an atom or ion. Although such bound electrons exhibit a high degree of quantum-mechanical wavelike behavior, there still remain particle aspects to their motion. Bound electrons occupy orbitals that are somewhat concentrated in spatial shells lying at different distances from the nucleus. As the set of electron energies allowed by quantum mechanics is discrete, so is the set of mean shell radii. Both these quantized physical quantities are primarily specified by integral values of the principal, or total, quantum number n . The full electron configuration of an atom is correlated with a set of values for all the quantum numbers of each and every electron. In addition to n , another important quantum number is l , an integer representing the orbital angular momentum of an electron in units of $h/2\pi$ where h is Planck's constant. The values 1, 2, 3, 4, 5, 6, 7 for n and 0, 1, 2, 3, for l together suffice to describe the electron configurations of all known normal atoms and ions, that is, those that have their lowest possible values of total electronic energy. The first seven shells are also given the letter designations K, L, M, N, O, P, and Q, respectively. Electrons with l equal to 0, 1, 2, and 3 are designated s , p , d , and f , respectively. See QUANTUM MECHANICS; QUANTUM NUMBERS.

In any configuration the number of equivalent electrons (same n and l) is indicated by an integral exponent (not a quantum number) attached to the letters s , p , d , and f . According to the Pauli exclusion principle the maximum is s^2 , p^6 , d^{10} , and f^{14} . For example, the configuration $1s^2 2s^2 2p^6 3s$ of the

Distribution of electrons in the atoms*															Ground term	Ionization potential, eV	
Element and atomic number	K		L		M			N				O					
	1, 0 1s	2, 0 2s	2, 1 2p	3, 0 3s	3, 1 3p	3, 2 3d	4, 0 4s	4, 1 4p	4, 2 4d	4, 3 4f	5, 0 5s	5, 1 5p	5, 2 5d	5, 3 5f			
H	1	1	—	—	—	—	—	—	—	—	—	—	—	—	$2S_{1/2}$	13.5984	
He	2	2	—	—	—	—	—	—	—	—	—	—	—	—	$1S_0$	24.5874	
Li	3	2	1	—	—	—	—	—	—	—	—	—	—	—	$2S_{1/2}$	5.3917	
Be	4	2	2	—	—	—	—	—	—	—	—	—	—	—	$1S_0$	9.3227	
B	5	2	2	1	—	—	—	—	—	—	—	—	—	—	$2P_{1/2}$	8.2980	
C	6	2	2	2	—	—	—	—	—	—	—	—	—	—	$3P_0$	11.2603	
N	7	2	2	3	—	—	—	—	—	—	—	—	—	—	$4S_{3/2}$	14.5341	
O	8	2	2	4	—	—	—	—	—	—	—	—	—	—	$3P_2$	13.6181	
F	9	2	2	5	—	—	—	—	—	—	—	—	—	—	$2P_{3/2}$	17.4228	
Ne	10	2	2	6	—	—	—	—	—	—	—	—	—	—	$1S_0$	21.5645	
Na	11	Neon configuration			1	—	—	—	—	—	—	—	—	—	$2S_{1/2}$	5.1391	
Mg	12	Neon configuration			2	—	—	—	—	—	—	—	—	—	$1S_0$	7.6462	
Al	13	Neon configuration			2	1	—	—	—	—	—	—	—	—	$2P_{1/2}$	5.9858	
Si	14	Neon configuration			2	2	—	—	—	—	—	—	—	—	$3P_0$	8.1517	
P	15	Neon configuration			2	3	—	—	—	—	—	—	—	—	$4S_{3/2}$	10.4867	
S	16	Neon configuration			2	4	—	—	—	—	—	—	—	—	$3P_2$	10.3600	
Cl	17	Neon configuration			2	5	—	—	—	—	—	—	—	—	$2P_{3/2}$	12.9676	
Ar	18	Neon configuration			2	6	—	—	—	—	—	—	—	—	$1S_0$	15.7596	
K	19	Argon configuration				—	1	—	—	—	—	—	—	—	$2S_{1/2}$	4.3407	
Ca	20	Argon configuration				—	2	—	—	—	—	—	—	—	$1S_0$	6.1132	
Sc	21	Argon configuration				1	2	—	—	—	—	—	—	—	$2D_{3/2}$	6.5615	
Ti	22	Argon configuration				2	2	—	—	—	—	—	—	—	$3F_2$	6.8281	
V	23	Argon configuration				3	2	—	—	—	—	—	—	—	$4F_{3/2}$	6.7462	
Cr	24	Argon configuration				5	1	—	—	—	—	—	—	—	$7S_3$	6.7665	
Mn	25	Argon configuration				5	2	—	—	—	—	—	—	—	$6S_{5/2}$	7.4340	
Fe	26	Argon configuration				6	2	—	—	—	—	—	—	—	$5D_4$	7.9024	
Co	27	Argon configuration				7	2	—	—	—	—	—	—	—	$4F_{9/2}$	7.8810	
Ni	28	Argon configuration				8	2	—	—	—	—	—	—	—	$3F_4$	7.6398	
Cu	29	Argon configuration				10	1	—	—	—	—	—	—	—	$2S_{1/2}$	7.7264	
Zn	30	Argon configuration				10	2	—	—	—	—	—	—	—	$1S_0$	9.3942	
Ga	31	Argon configuration				10	2	1	—	—	—	—	—	—	$2P_{1/2}$	5.9993	
Ge	32	Argon configuration				10	2	2	—	—	—	—	—	—	$3P_0$	7.8994	
As	33	Argon configuration				10	2	3	—	—	—	—	—	—	$4S_{3/2}$	9.7886	
Se	34	Argon configuration				10	2	4	—	—	—	—	—	—	$3P_2$	9.7524	
Br	35	Argon configuration				10	2	5	—	—	—	—	—	—	$2P_{3/2}$	11.8138	
Kr	36	Argon configuration				10	2	6	—	—	—	—	—	—	$1S_0$	13.9996	
Rb	37	Krypton configuration							—	—	1	—	—	—	—	$2S_{1/2}$	4.1771
Sr	38	Krypton configuration							—	—	2	—	—	—	—	$1S_0$	5.6949
Y	39	Krypton configuration							1	—	2	—	—	—	—	$2D_{3/2}$	6.2173
Zr	40	Krypton configuration							2	—	2	—	—	—	—	$3F_2$	6.6339
Nb	41	Krypton configuration							4	—	1	—	—	—	—	$6D_{1/2}$	6.7589
Mo	42	Krypton configuration							5	—	1	—	—	—	—	$7S_3$	7.0924
Tc	43	Krypton configuration							5	—	2	—	—	—	—	$6S_{5/2}$	7.28
Ru	44	Krypton configuration							7	—	1	—	—	—	—	$5F_5$	7.3605
Rh	45	Krypton configuration							8	—	1	—	—	—	—	$4F_{9/2}$	7.4589
Pd	46	Krypton configuration							10	—	—	—	—	—	—	$1S_0$	8.3369

ground or lowest energy state of sodium means that there are two electrons with $n = 1, l = 0$; two with $n = 2, l = 0$; six with $n = 2, l = 1$; and one with $n = 3, l = 0$. Higher values of n and l can be achieved by excitation of normal atoms, either through photon absorption or by particle impact, temporarily moving one or more electrons to unfilled shells of larger radii and increasing the atom's electronic energy. Such excited electronic configurations are inherently unstable. Excited electrons ultimately fall back down to their normal locations closer to the nucleus, a process accompanied by the emission of one or more quanta of electromagnetic radiation (photons). While the atom is excited, the partially depleted shells are said to possess vacancies. *See* EXCLUSION PRINCIPLE.

An electron configuration is categorized as having even or odd parity, according to whether the sum of p and f electrons is even or odd. Strong spectral lines result only from transitions between configurations of unlike parity. *See* PARITY (QUANTUM MECHANICS).

Insofar as they are known from spectroscopic investigations, the electron configurations characteristic of the normal or ground states of the first 102 chemical elements are shown in the **table**. In the next-to-last column of the table, the spectral term of the energy level with lowest total electronic energy is shown. For all elements except protactinium (Pa), uranium (U), and neptunium (Np), an LS -coupling term is given. Here, the main part of the term symbol is a capital letter, S, P, D, F , and so on, that

Distribution of electrons in the atoms* (cont.)												
Element and atomic number	Configuration of inner shells	N		O			P			Q	Ground term	Ionization potential, eV
		4, 3 4f	5, 0 5s	5, 1 5p	5, 2 5d	5, 3 5f	6, 0 6s	6, 1 6p	6, 2 6d	7, 0 7s		
Ag 47	Palladium configuration	—	1	—	—	—	—	—	—	—	$2S_{1/2}$	7.5762
Cd 48		—	2	—	—	—	—	—	—	—	$1S_0$	8.9938
In 49		—	2	1	—	—	—	—	—	—	$2P_0^{\circ}$	5.7864
Sn 50		—	2	2	—	—	—	—	—	—	$3P_0$	7.3439
Sb 51		—	2	3	—	—	—	—	—	—	$4S_{3/2}^{\circ}$	8.6084
Te 52		—	2	4	—	—	—	—	—	—	$3P_2^{\circ}$	9.0096
I 53		—	2	5	—	—	—	—	—	—	$2P_{3/2}^{\circ}$	10.4513
Xe 54		—	2	6	—	—	—	—	—	—	$1S_0$	12.1298
Cs 55	The shells 1s to 4d contain 46 electrons	—	—	—	—	—	1	—	—	—	$2S_{1/2}$	3.8939
Ba 56		—	—	—	—	—	2	—	—	—	$1S_0$	5.2117
La 57		—	—	—	1	—	2	—	—	—	$2D_{3/2}$	5.5769
Ce 58		1	—	—	1	—	2	—	—	—	$1G_4^{\circ}$	5.5387
Pr 59		3	—	—	—	—	2	—	—	—	$4f_{9/2}^{\circ}$	5.473
Nd 60		4	—	—	—	—	2	—	—	—	$5f_4$	5.5250
Pm 61		5	—	—	—	—	2	—	—	—	$6H_{5/2}^{\circ}$	5.582
Sm 62		6	—	—	—	—	2	—	—	—	$7F_0$	5.6437
Eu 63		7	—	—	—	—	2	—	—	—	$8S_{7/2}^{\circ}$	5.6704
Gd 64		7	—	—	1	—	2	—	—	—	$9D^{\circ}$	6.1498
Tb 65		9	—	—	—	—	2	—	—	—	$6H_{15/2}^{\circ}$	5.8638
Dy 66		10	—	—	—	—	2	—	—	—	$5f_8$	5.9389
Ho 67		11	—	—	—	—	2	—	—	—	$4f_{15/2}^{\circ}$	6.0215
Er 68		12	—	—	—	—	2	—	—	—	$3H_6$	6.1077
Tm 69		13	—	—	—	—	2	—	—	—	$2F_{7/2}^{\circ}$	6.1843
Yb 70	14	—	—	—	—	2	—	—	—	$1S_0$	6.2542	
Lu 71	14	—	—	—	1	—	2	—	—	$2D_{3/2}$	5.4259	
Hf 72	The shells 1s to 5p contain 68 electrons	—	—	—	2	—	2	—	—	—	$3F_2$	6.8251
Ta 73		—	—	—	3	—	2	—	—	—	$4F_{3/2}$	7.5496
W 74		—	—	—	4	—	2	—	—	—	$5D_0$	7.8640
Re 75		—	—	—	5	—	2	—	—	—	$6S_{5/2}$	7.8335
Os 76		—	—	—	6	—	2	—	—	—	$5D_4$	8.28
Ir 77		—	—	—	7	—	2	—	—	—	$4F_{9/2}$	9.02
Pt 78	—	—	—	9	—	1	—	—	—	$3D_3$	8.9588	
Au 79	The shells 1s to 5d contain 78 electrons	—	—	—	—	—	1	—	—	—	$2S_{1/2}$	9.2255
Hg 80		—	—	—	—	—	2	—	—	—	$1S_0$	10.4375
Tl 81		—	—	—	—	—	2	1	—	—	$2P_{1/2}^{\circ}$	6.1082
Pb 82		—	—	—	—	—	2	2	—	—	$3P_0$	7.4167
Bi 83		—	—	—	—	—	2	3	—	—	$4S_{3/2}^{\circ}$	7.2855
Po 84		—	—	—	—	—	2	4	—	—	$3P_2$	8.414
At 85		—	—	—	—	—	2	5	—	—	$2P_{3/2}^{\circ}$	—
Rn 86		—	—	—	—	—	2	6	—	—	$1S_0$	10.7485
Fr 87		—	—	—	—	—	2	6	—	1	$2S_{1/2}$	4.0727
Ra 88		—	—	—	—	—	2	6	—	2	$1S_0$	5.2784
Ac 89		—	—	—	—	—	2	6	1	2	$2D_{3/2}$	5.17
Th 90		—	—	—	—	—	2	6	2	2	$3F_2$	6.3067
Pa 91		—	—	—	—	$2(^3H_4)$	2	6	1	2	$(4, 3/2)_{11/2}$	5.89
U 92		—	—	—	—	$3(4f_{9/2}^{\circ})$	2	6	1	2	$(9/2, 3/2)_6$	6.1941
Np 93		—	—	—	—	$4(f_{14}^{\circ})$	2	6	1	2	$(4, 3/2)_{11/2}$	6.2657
Pu 94		—	—	—	—	6	2	6	—	2	$7F_0$	6.0260
Am 95		—	—	—	—	7	2	6	—	2	$8S_{7/2}^{\circ}$	5.9738
Cm 96	—	—	—	—	7	2	6	1	2	$9D_2^{\circ}$	5.9914	
Bk 97	—	—	—	—	9	2	6	0	2	$6H_{5/2}^{\circ}$	6.1979	
Cf 98	—	—	—	—	10	2	6	0	2	$5f_8$	6.2817	
Es 99	—	—	—	—	11	2	6	0	2	$4f_{15/2}^{\circ}$	6.42	
Fm 100	—	—	—	—	12	2	6	0	2	$3H_6$	6.50	
Md 101	—	—	—	—	13	2	6	0	2	$2F_{7/2}^{\circ}$	6.58	
No 102	—	—	—	—	14	2	6	0	2	$1S_0$	6.65	

*Updated and revised with data from G. W. F. Drake (ed.), *Springer Handbook of Atomic, Molecular, and Optical Physics*, pp. 182–183, Springer, 2006.

represents the total electronic orbital angular momentum. Attached to this is a superior prefix, 1, 2, 3, 4, and so on, that indicates the multiplicity, and an anterior suffix, 0, $1/2$, 1, $3/2$, 2, $5/2$, and so on, that shows the total angular momentum, or J value, of the atom in the given state. The sign $^{\circ}$ above the J value signifies that the spectral term and electron configuration have odd parity. (For the elements Pa, U, and Np, a jj coupling term is given. Here, parentheses

enclose two numbers, the first of which is the total angular momentum of the electrons in the $5f$ shell, while the second is the total angular momentum of the electron in the $6d$ shell. The suffixes following the parentheses are the same as in LS coupling. For these elements the LS coupling of the electrons in the $5f$ shell is also given in the column that lists the number of such electrons.)

The last column of the table presents the first

ionization potential of the atom when this has been derived from spectroscopic observations. In any atomic spectrum, two or more spectral lines with certain similar properties may form a series such that the reciprocal wavelengths $1/\lambda$ (number of waves per centimeter = σ) can be closely represented by a formula of the Rydberg type, $\sigma = L - R/(n + \mu)^2$, in which L is the limit of the series. R is called the Rydberg constant, and the principal quantum number n has successive integral values to which a constant fractional part μ is added. The second term vanishes when n approaches infinity, and the series limit is thus evaluated. This limit is usually coincident with the ground state of the ion, and is thus a measure (in wave-number units) of the energy required to remove from an atom its least firmly bound electron and transform a neutral atom into a singly charged ion. The energy required to ionize an atom is usually expressed in electronvolts (1 eV = 8065.48 wave numbers) and is called its first ionization potential. See ATOMIC STRUCTURE AND SPECTRA; IONIZATION POTENTIAL; RYDBERG ATOM; RYDBERG CONSTANT.

Molecules consist of two or more atoms that at least partially share one or more electrons with one another. Sets of quantum numbers specify molecular electron configurations in much the same way as for atoms. However, the physical meaning of some molecular quantum numbers is different than that for the atomic case, as is the term notation. See MOLECULAR STRUCTURE AND SPECTRA. James E. Bayfield

Electron diffraction

The phenomenon associated with interference processes that occur when electrons are scattered by atoms to form diffraction patterns. The wave character of electrons is shown most strikingly, and doubtless most conclusively, by the phenomena of interference. For this reason, the diffraction of electrons presents the most obvious confirmation of quantum mechanics. Because of the dependence of the diffraction pattern on the distances between the atoms, electron diffraction is also an important tool for the study of the structure of crystals and of free molecules, analogous to the use of x-rays for these purposes. See X-RAY CRYSTALLOGRAPHY; X-RAY DIFFRACTION.

According to quantum theory, any particle moving with momentum $m\nu$ has a wavelength $\lambda = h/m\nu$, where h is Planck's constant. If the particle is an electron, and its velocity is the result of acceleration by a the potential difference V , this formula becomes Eq. (1), where m_0 and e are the rest mass

$$\lambda = \frac{h}{(2m_0eV)^{1/2}(1 + eV/2m_0c^2)^{1/2}} \quad (1)$$

and charge of the electron, and c is the velocity of light. The last factor in the denominator represents the relativity correction, which is negligible at low

voltages and amounts to only 5% at 100,000 V. For V in volts and λ in nanometers, Eq. (2) is a good

$$\lambda = \left(\frac{1.5}{V}\right)^{1/2} \quad (2)$$

approximation to Eq. (1) at nonrelativistic energies. According to energy $E = eV$, two major techniques of structure analysis with electron beams are distinguished: low-energy electron diffraction (LEED) [$E \simeq 5\text{--}500$ eV] and high-energy electron diffraction (HEED) [$E \simeq 5\text{--}500$ keV]; medium-energy electron diffraction (MEED) [$E \simeq 500$ eV–5 keV] is of little importance. In addition, electrons generated in condensed matter by incident electrons or x-ray photons are diffracted (in Auger electron diffraction and photoelectron diffraction). Unlike neutrons and x-rays, electrons penetrate matter only for a very short distance before they lose energy (by inelastic scattering) or are scattered elastically (diffracted). Due to the energy loss, the wavelength of an electron changes according to Eq. (2) so that it can no longer interfere with the other incident electrons. This loss of coherence occurs in condensed matter typically after mean free paths for inelastic scattering ranging from several tenths of nanometers at low energies ($E \simeq 50$ eV) to several tens of nanometers at high energies, which determines the application range of LEED and HEED. See COHERENCE; DE BROGLIE WAVELENGTH; DIFFRACTION; INTERFERENCE OF WAVES; MEAN FREE PATH; NONRELATIVISTIC QUANTUM THEORY; QUANTUM MECHANICS.

Low-Energy Electron Diffraction

LEED is used mainly for the study of the structure of single-crystal surfaces and of processes on such surfaces that are associated with changes in the lateral periodicity of the surface. A monochromatic, nearly parallel electron beam, of 10^{-4} to 10^{-3} m (4×10^{-3} to 4×10^{-2} in.) in diameter, strikes the surface, usually at normal incidence. The elastically backscattered electrons are separated from all other electrons by a retarding field and detected with a suitable movable collector or, more frequently—after acceleration to about 5 keV energy—on a hemispherical fluorescent screen with the crystal in its center (Fig. 1). The intensity of the diffraction spots can be measured as a function of the energy $E = eV$ of the incident electrons to obtain so-called $I(V)$ curves (Fig. 2).

To obtain such results, the surface must be carefully cleaned and kept in ultrahigh vacuum (UHV, pressure $\simeq 10^{-10}$ torr $\simeq 10^{-8}$ pascal) during the experiment and must be exposed to well-defined amounts of gases or vapors at low pressures if adsorption, condensation, or corrosion is to be studied. This requirement, as well as the fact that many surfaces are changed by the electron beam—for example, due to dissociation in the case of ionic crystals or due to desorption of adsorbed gases—limits the applicability of LEED. Nevertheless, the number of surfaces and surface processes that can be studied

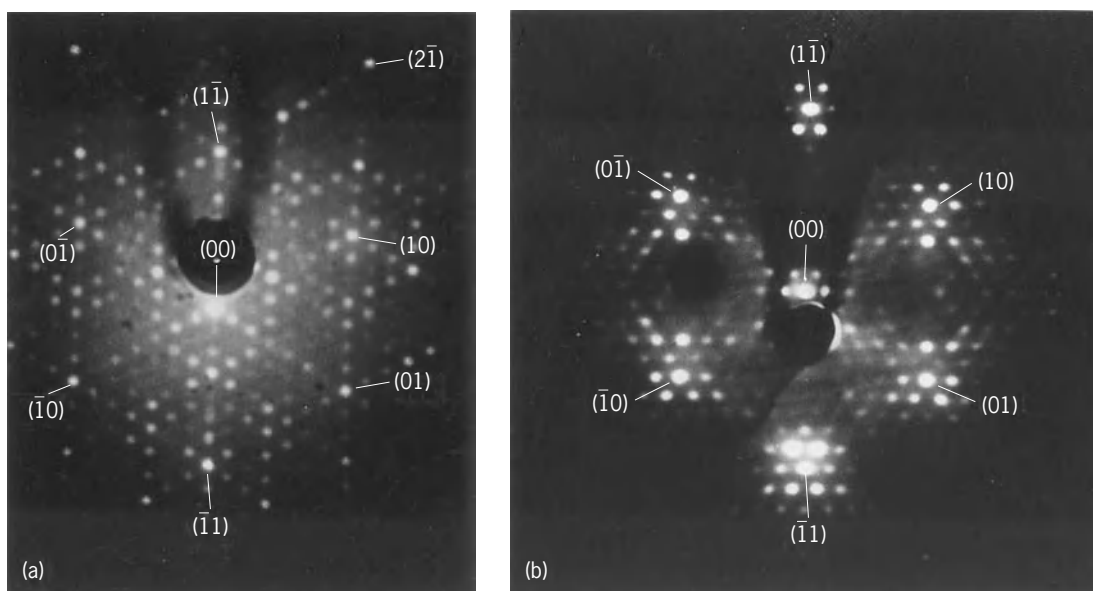


Fig. 1. LEED patterns from single-crystal surfaces. Numbers are indices (h, k) labeling diffraction spots of the surfaces. (a) Clean silicon (111) surface (7×7 structure). $E = 85$ eV. (b) Ordered carbon monoxide adsorption layer on a tungsten (110) surface. $E = 120$ eV.

is very large, and the main limitation of LEED is not caused by difficulties in obtaining LEED patterns but rather in evaluating them.

Interpretation of patterns. The difficulties in interpreting LEED patterns are caused by the strong elastic scattering of slow electrons by atoms. As a consequence of this, the electrons observed in the LEED pattern have been scattered several times on the average, in contrast to x-ray and neutron diffraction, in which single scattering is a good approximation. Thus evaluation of the spot intensities of LEED patterns, with the additional complications arising from the elastic scattering at the surface barrier and from inelastic scattering, requires large-scale computers. As a consequence, only the structures of surfaces whose unit meshes (two-dimensional unit cells) are small have been determined by LEED. For this reason, alternative evaluation methods have been developed in which multiple-scattering effects are reduced by proper averaging procedures.

The interpretation difficulties are considerably reduced if only the lateral periodicity of the surface, but not the location of the atoms in the unit mesh, is to be determined. The determination of the lateral periodicity of the surface structure and of the size and shape of domains with a given structure requires only the evaluation of the geometry of the pattern, that is, of the position, size, and shape of the diffraction spots but not of their intensity, and this evaluation is relatively easy. Thus the most important application of LEED is the study of surface phenomena associated with changes in lateral periodicity. A large amount of information on lateral disorder may be extracted from spot shapes and intensities. Special systems (spot profile analysis LEED, or SPALeED) have been developed for this purpose.

Surface structure. Clean surfaces of crystals frequently have the lateral periodicity that is expected from their bulk crystal structure. The analysis of $I(V)$ curves of such surfaces (Fig. 2) in general reveals that the distribution of the atoms normal to the surface in the topmost layers differs somewhat from that in the bulk. These relaxations, that is, displacements of atoms from their positions in the bulk, are only a few percent on densely packed surfaces. On more open surfaces they can amount to 10% or more and are frequently oscillatory. There are, however, also surfaces which show complicated superstructures (known as reconstructed surfaces), that is, structures whose lateral periodicity is larger than that of the substrate (Fig. 1a). Usually, such surfaces are analyzed by LEED

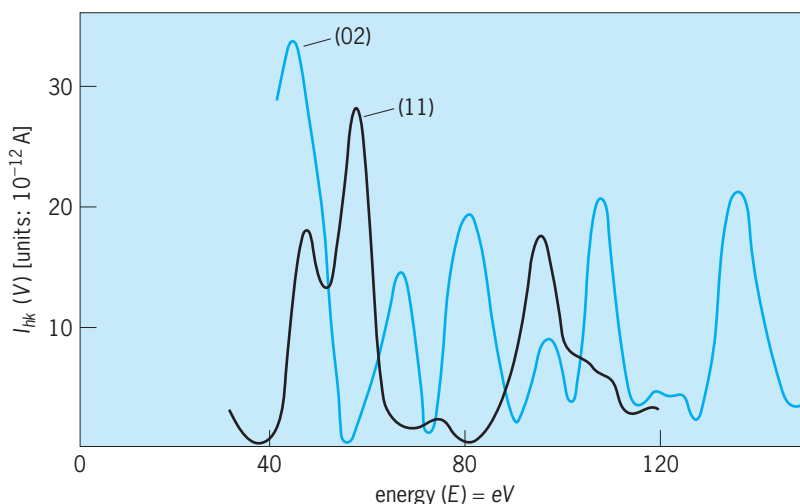


Fig. 2. $I(V)$ curves of spots (hk) in the LEED pattern from a clean tungsten (110) surface. (After G. A. Somorjai, *The Structure and Chemistry of Solid Surfaces*, Wiley, 1969)

only after structural models have been proposed on the basis of other structure studies. On ideal surfaces, the decrease of the spot intensities with temperature has shown that the amplitudes of the thermal vibrations of surface atoms normal to the surface can be much larger than in the bulk. See CRYSTAL STRUCTURE; LATTICE VIBRATIONS.

Adsorption. Foreign atoms can interact with a clean surface in a variety of ways. The most important contribution of LEED is to the understanding of chemisorption, which precedes corrosion and, in many cases, epitaxy. Here, not only the structure of many adsorption systems, mainly of gases on metals, or metals on other metals and semiconductors, has been studied, but also the kinetics of the adsorption and desorption process as well as changes in the adsorption layer upon heating. The combination of LEED with Auger electron spectroscopy (AES) and with work-function measurements has proven particularly powerful in these studies, because such methods give the coverage and information on the location of the adsorbed atoms normal to the surface. Combining LEED with other complementary techniques such as ion scattering spectroscopy, electron energy loss spectroscopy, or photoelectron spectroscopy has become increasingly popular and can enable the elimination of ambiguities in the interpretation of many LEED results. See ADSORPTION; AUGER EFFECT; CORROSION; ELECTRON SPECTROSCOPY; SURFACE AND INTERFACIAL CHEMISTRY; SURFACE PHYSICS.

SPLEED. This is LEED with spin-polarized electrons. It gives information on the magnetic structure of surfaces and thin films due to the spin dependence of the diffraction of slow electrons in magnetic materials. This dependence is of particular importance in the spin-polarized version of low-energy electron microscopy (LEEM), called SPLEEM, which allows the magnetic microstructures of surfaces and thin films to be studied.

LEEM. Low-energy diffracted electrons can be used to image the surface of a specimen if it is made the cathode in an immersion objective lens of an ultrahigh-vacuum electron microscope. Because of the strong elastic backscattering of very slow electrons, the resulting LEEM images (Fig. 3) can be recorded at video rates. This allows the study of surface processes such as phase transitions or film growth in real time. See ELECTRON MICROSCOPE.

High-Energy Electron Diffraction

HEED is used mainly for the study of the structure of thin foils, films, and small particles (thickness or diameter of 10^{-9} to 10^{-6} m or 4×10^{-8} to 4×10^{-5} in.), of molecules, and also of the surfaces of crystalline materials. A monochromatic, usually nearly parallel, electron beam with a diameter of 10^{-3} to 10^{-8} m (4×10^{-2} to 4×10^{-7} in.) is incident on the target. The forward-scattered electrons (backscattering is negligible) are detected by means of a fluorescent screen, a photoplate, or some other current-sensitive

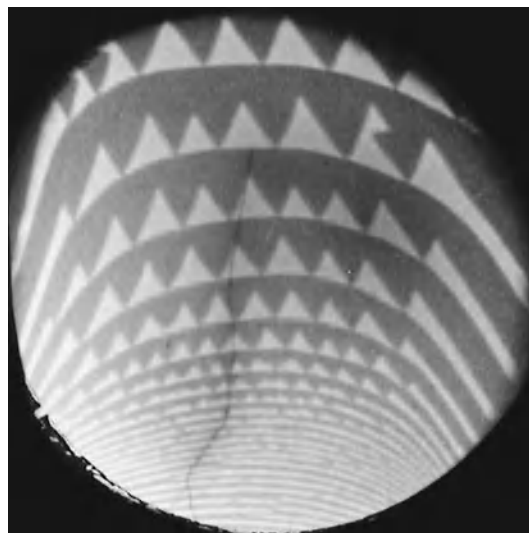


Fig. 3. LEEM image of a partially reconstructed silicon (111) surface. The bright triangles have (7×7) structure, the dark regions (1×1) .

detector, usually without the inelastically scattered electrons being eliminated.

Reflection HEED. The availability of UHV instruments has made it possible to study the structure of surfaces under clean conditions by reflection HEED (RHEED). Many experiments have confirmed the theoretical expectation that RHEED has a sensitivity comparable to that of LEED. Therefore, similar to LEED, RHEED can be used for the determination of the lateral arrangement of the atoms in the topmost layers of the surface (Fig. 4), including the structure of adsorbed layers. Although it is more convenient to deduce the periodicity of the atomic arrangement parallel to the surface from LEED patterns than from RHEED patterns, LEED frequently becomes inapplicable when the surface is rough. This usually occurs in the later stages of corrosion or in precipitation, for example, of silicon carbide on silicon, when small crystals grow on the surface. In such investigations RHEED is far superior to LEED because the fast electrons can penetrate the asperities and produce a transmission HEED (THEED) pattern. RHEED has become particularly important for thin-film growth monitoring via the specular beam intensity oscillations caused by monolayer-by-monolayer growth (Fig. 5). See CRYSTAL GROWTH.

Like LEED, RHEED gives little information on the chemical nature of the atoms that produce the diffraction pattern. Therefore, RHEED has been combined with analytical tools, such as Auger electron spectroscopy or x-ray emission spectroscopy. A second limitation inherent to both RHEED and LEED is the difficulty of determining the atomic positions normal to the surface. This is a consequence of dynamical effects and of absorption. Nevertheless, both techniques have given valuable information on corrosion, precipitation, adsorption, condensation, film growth, and other surface and interface phenomena that could not be obtained by other methods.

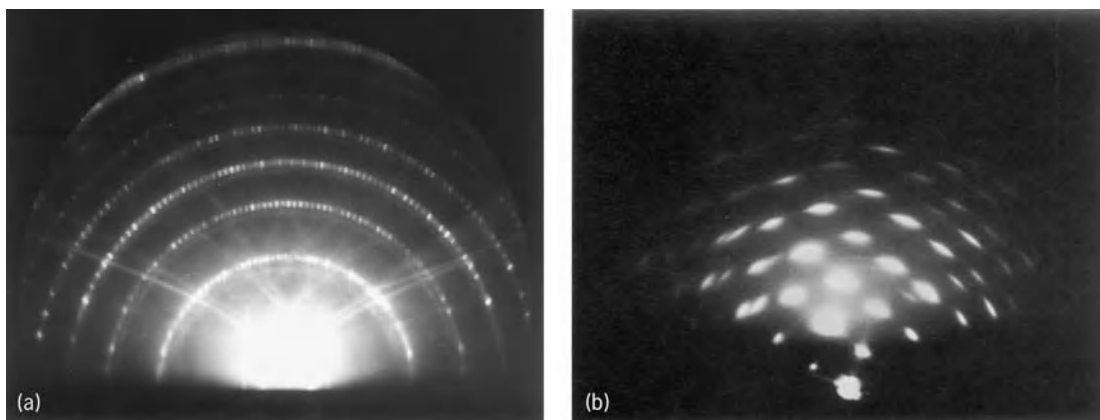


Fig. 4. RHEED pattern from (a) flat (111) surface of a silicon single crystal, and (b) rough surface of a polycrystalline evaporated film of calcium fluoride (CaF_2) with strong fiber texture, from which no LEED pattern can be obtained.

Scanning HEED. In scanning HEED (SHEED) the diffracted electrons are not recorded on photographic film but are directly measured electronically with sensitive detectors. By moving the detector across the diffraction pattern or by deflecting the diffracted electrons across a stationary detector (scanning), the intensity distribution in the diffraction pattern can be displayed quantitatively on an XY recorder. If an energy filter is put in front of the detector, the inelastically scattered electrons, which are usually not taken into account in the quantitative intensity evaluation, may be filtered out so that only the elastically scattered electrons are measured. This technique is particularly useful in transmission through polycrystalline samples, which produce a ring pattern (Debye-Scherrer diagram). It has also been used in single-crystal samples producing a spot

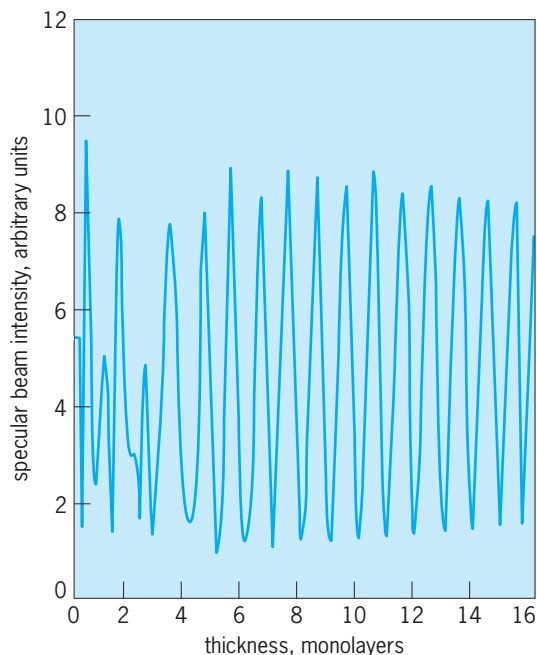


Fig. 5. RHEED specular beam intensity oscillations during molecular beam epitaxy of lead on a silicon (111) surface.

pattern (Laue diagram) and even in reflection diffraction from surfaces. The main application of SHEED is in the study of processes which are accompanied by changes of the intensity distribution, such as the growth of thin films and annealing and corrosion processes.

Transmission HEED. The technological importance of thin film and interface devices has led to an upsurge of thin film growth studies by conventional transmission HEED (THEED), usually combined with transmission microscopy. Information obtained this way has been mainly on the orientation of the crystallites composing the film (Fig. 6).

Diffraction electron microscopy. In transmission electron microscopy of crystalline specimens the main contrast mechanism is the diffraction contrast. To understand it, knowledge of the diffraction pattern and of diffraction theory is necessary. High-voltage electron microscopes (with electron energies of about 1 MeV) have become available. As a consequence, samples with thickness of the order 10^{-6} m (4×10^{-5} in.) can now be studied, as compared to the former upper thickness limit of about 10^{-7} m (4×10^{-6} in.). This has significance for metallurgy and the physics of metals and semiconductors. It is now possible to study the structure, distribution, and behavior of imperfections in metals, alloys, and semiconductors in bulk with little influence from the boundaries of the material. Other electron-microscope techniques that rely heavily on diffraction are the use of illumination systems producing tilted or conical illumination of the specimen and of objective apertures which transmit only diffracted electrons. With these techniques it has become possible to produce, by diffraction contrast, images of crystal planes 1.5×10^{-10} m (6×10^{-9} in.) apart. RHEED beams can be used for reflection electron microscopy (REM), though with considerable foreshortening due to the grazing incidence on the surface which is necessary for surface sensitivity and intensity.

Photoelectron diffraction. In photoelectron diffraction (PED), the electron wave emitted by one atom

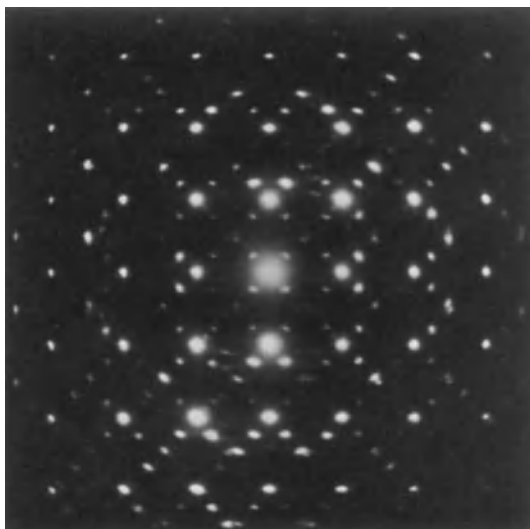


Fig. 6. THEED pattern from epitaxial gold particles on a lead sulfide (PbS) single-crystal film. (From A. K. Green, J. Dancy, and E. Bauer, *Growth of Au on PbS, PbSe, PbTe and SnTe thin film substrates*, *J. Vac. Sci. Technol.*, 10:494-502, 1973)

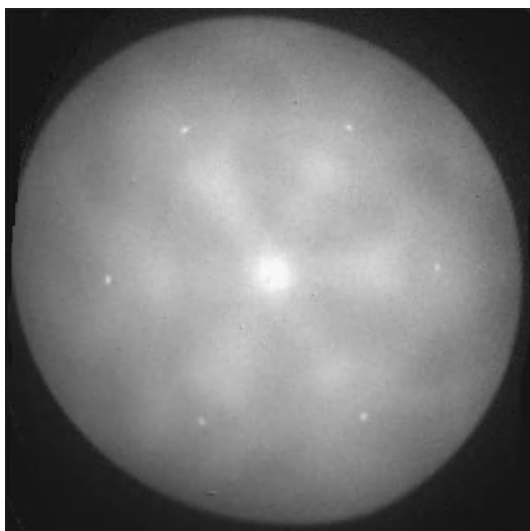


Fig. 7. PED pattern (diffuse features) with superimposed LEED pattern (sharp spots) of a lead film, five atomic layers thick, on a silicon (111) surface. Photon energy is 62 eV; electron energy, 39 eV. (Courtesy of T. Schmidt)

upon photoionization is diffracted by the surrounding atoms. The resulting diffraction pattern (Fig. 7) allows local structure analysis around the selected atom. At low energies (below several hundred electronvolts) the interpretation requires an effort comparable to that needed in LEED; at higher energies (above about 1000 eV) strong forward scattering frequently permits a simple qualitative analysis. Similar to LEEM, REM, and transmission electron microscopy (TEM), the photoemitted electron wave may also be used for imaging, a method known as x-ray photoemission electron microscopy (XPEEM) with synchrotron radiation. In this case the imaging is atom-specific, just as is PED.

E. Bauer

Diffraction in Gases and Liquids

Electron diffraction in gases and liquids is similar in principle to that in solids; the differences arise from the lack in gases and liquids of any highly regular arrangement of the component atoms. In gases the low density makes it possible to study diffraction by individual atoms and molecules. The results obtained from monatomic gases represent the density of electronic charge in the atom as a function of the distance from the nucleus. The results from gaseous polyatomic molecules represent the equilibrium distances between the atomic nuclei and the average amplitudes of vibration associated with these distances. Liquids have been studied much less thoroughly, both in theory and in practice, than have gases.

It should be remembered that the structures of molecules in the gaseous state may be different from those in the liquid or solid states; for example, hydrofluoric acid in the gas forms an $(\text{HF})_6$ ring, while in solution or liquid it forms an infinite chain ($-\text{FHFHFHFHFHFHF}-$); biphenyl in the gas is twisted with conformational equilibria, whereas in the solid state it is coplanar in one frozen conformer.

Applications. Typical questions of molecular structure studied by electron diffraction include those of configuration and size in many molecules, with special interest in the variation of chemical bond distances in related molecules (such as the 6.8-picometers or 2.7×10^{-10} in. decrease in C—F distance in the series CH_3F , CH_2F_2 , CHF_3 , CF_4), the variation in the angles between chemical bonds, the distinction between geometric isomers, the degree of restricted rotation around chemical bonds, and in general the relations between the geometry of molecules and their energy and chemical behavior. See CHEMICAL BONDING; MOLECULAR ISOMERISM; MOLECULAR STRUCTURE AND SPECTRA; STRUCTURAL CHEMISTRY.

The application of electron diffraction in liquids has been restricted because of the less precise information obtainable from molecules which lie in close contact but in irregular arrangements.

The liquids studied include certain metals (mercury, tin, tin-bismuth, and tin-aluminum alloys), paraffin oils, and thin layers of water. For monatomic molecules in the liquid state the distances between nearest neighbors are fairly uniform, whereas the separation between more distant neighbors varies considerably. Evaluation of the closest distance of approach can be made to about 10 pm (4×10^{-10} in.); major developments in the theory are still required.

The most interesting application of electron diffraction to a liquid or amorphous condensed phase has been in the examination of the surface layers on polished solids. The diffuse diffraction patterns often obtained may be the result of an amorphous arrangement in the surface or of the poor resolving power of very tiny crystalline particles. It is probable that both these states are produced in the polishing of different solids.

Theory and techniques. Structural information about gaseous molecules is obtained by having a fine beam of electrons pass through the gas and strike a photographic plate. The electrons in the beam interact with the charged particles in the atoms (electrons and nuclei) and are bent away from the original direction through varying angles, as registered in the pattern on the plate. The observed variation in the number of scattered electrons with increasing angle is interpreted as an interference effect; that is, electron scattering can be described in the language of diffracted waves for which the resultant is the sum of the component wavelets whose amplitudes and phases are influenced by the wavelength of the incident radiation and the relative positions of the scattering centers. The equivalent wavelength of the electrons, λ , is determined by their energy and is given by Eq. (1) or (2). The values of λ commonly used in diffraction by gases are in the range $7\text{--}5 \times 10^{-12}$ m ($3\text{--}2 \times 10^{-10}$ in.).

The intensity I of electrons scattered by a spherically symmetrical atom is computed by the Schrödinger wave equation; the result is given by Eq. (3), where k is a constant and f is a complex func-

$$I = k |f|^2 \quad (3)$$

tion called the atomic scattering factor. If the atomic number Z is not too large, Eq. (4) holds, where a is the Bohr radius and F is given by Eq. (5). Here $\rho(r)$ is

$$|f|^2 = [2(Z - F)/as^2]^2 \quad (4)$$

$$F = 4\pi \int_0^\infty r^2 \rho(r) [(\sin sr)/sr] dr \quad (5)$$

density of electronic charge at distance r from the nucleus, $s = 4\pi(\sin \theta)/\lambda$ with 2θ equal to the angle between the scattered electron and the original beam, and λ is the electron wavelength. Experimental observations on I as a function of s in monatomic gases lead to the determination of $\rho(r)$.

The intensity of electrons scattered by a collection of independent molecules having all possible orientations is given by Eq. (6), where the summations are

$$I(s) = k \sum_i [|f_i|^2 + (2/as^2)^2 S_i] + \sum_i \sum_j f_i^* f_j \int_0^\infty P_{ij}(r) [(\sin sr)/sr] dr \quad (6)$$

taken over all the atoms in the molecule and $(2as^2)^2 S_i$ represents the intensity of inelastic scattering by the i th atom. The double summation has a term for each pair of atoms i and j . The relative probability of finding the distance between the atom i and j at various values is represented by $P_{ij}(r)$. In principle the use of the observed intensity to determine P_{ij} for each pair of atoms in the molecule constitutes a structure determination for the molecule. The expression is simpler when the atomic motions are nearly harmonic, as in the case of CCl_4 . For this molecule the double summation has only two distinct terms. The first

is given by expression (7); the second for Cl—Cl is similar.

$$8f_{\text{Cl}} f_{\text{Cl}} \exp(-l_{\text{CCl}}^2 s^2) (\sin sr_{\text{CCl}}) / sr_{\text{CCl}} \quad (7)$$

The four parameters which describe the structure are the equilibrium distances, r_{CCl} and r_{ClCl} , and the average amplitudes of vibration, l_{CCl} and l_{ClCl} . Molecular parameters such as these can be determined with high precision when a proper treatment of the scattering factors is included, and when accurate experimental intensities are obtained.

Special equipment is required to meet the conditions assumed in the scattering theory. The electron beam is accelerated with a steady voltage between 30,000 and 70,000 V, which should be constant within about 0.01%. The beam is focused by electrostatic and magnetic lenses so that it has a diameter of no more than 0.1 mm (0.004 in.) at the photographic plate. The whole path of the beam is enclosed in a high-vacuum chamber (pressure of about 10^{-5} mm Hg or 10^{-3} Pa or 10^{-7} lbf/in.) so that no appreciable electron scattering will occur in the residual air. The gas specimen is introduced in a fine jet so that the volume in which the electrons meet the gas is no more than 1.0 mm (0.04 in.) in diameter; high-speed pumping is required to remove the gas as rapidly as possible. In front of the photographic plate a rotating sector is mounted to modify the intensity of electrons reaching the plate so that the normally rapid decrease of intensity with increasing angle of scattering is leveled off in a known way and the emulsion is able to register the incident electrons over a wide range of angle.

Interpretation of results. The pattern observed is a set of light and dark circular bands whose spacings and relative intensities depend on the composition and structure of the specimen molecules. The pattern is scanned by a recording microphotometer which yields (after calibration of the photographic emulsion) a tracing or a digital output of the experimental data of I as a function of s . These data are interpreted with the aid of a Fourier transform of the intensity expression of Eq. (6) to give the P_{ij} functions for the pairs of atoms. From the positions and widths of the peaks for CCl_4 in Fig. 8, it is determined that Eqs. (8) hold (1 pm =

$$\begin{aligned} r_{\text{CCl}} &= 176.6 \pm 0.3 \text{ pm} \\ l_{\text{CCl}} &= 6.0 \pm 0.5 \text{ pm} \\ r_{\text{ClCl}} &= 288.7 \pm 0.4 \text{ pm} \\ l_{\text{ClCl}} &= 6.8 \pm 0.3 \text{ pm} \end{aligned} \quad (8)$$

10^{-12} m = 3.937×10^{-11} in.). See FOURIER SERIES AND TRANSFORMS.

This relatively high precision in determining the distances in gas molecules has been achieved since about 1950 with the development of new instrumentation and the application of high-speed computing methods. The relative intensities of scattered electrons registered on the photographic plate can be measured to better than 0.1%; the time saved by the

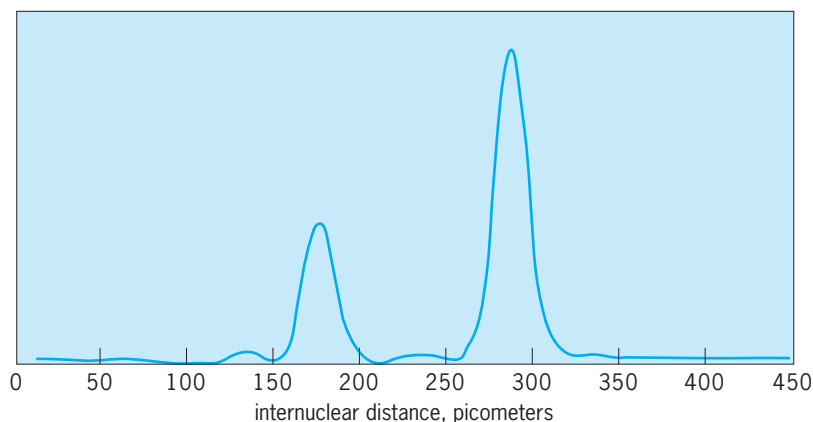


Fig. 8. Experimental distribution of internuclear distances in CCl_4 . The two prominent peaks represent the C-Cl and Cl-Cl distances.

use of electronic computing methods in interpreting the data permits a more rigid application of the criteria for satisfactory agreement between experiment and theory.

Over a thousand gases have been studied by the methods giving internuclear distances within 0.4 pm (1.6×10^{-11} in.). In favorable cases precisions exceeding 0.1 pm (4×10^{-12} in.) have been obtained; when these are compared with spectroscopic results, good agreement is found if allowance is made for the difference in the nature of the distances measured. The molecular structures of some 500 gases have been reported by earlier electron diffraction procedures with uncertainties 10 times larger or more.

The applicability of the gas diffraction method alone is limited to simple molecules. Only when the three-dimensional structure can be derived uniquely from a one-dimensional spectrum of internuclear distances blurred by thermal motion is a complete structure determination possible from electron diffraction data alone; light atoms in the presence of heavy ones in the same molecule are less precisely located; distances as close as 3 pm (1.2×10^{-10} in.) can barely be resolved. The range of the method is greatly increased, however, when some structural features in the molecule can be assumed from the results of other methods of investigation. For example, if a six-membered ring of atoms is known to have trigonal symmetry, the number of structural parameters is decreased from 12 to 2. See NEUTRON DIFFRACTION; SCATTERING EXPERIMENTS (ATOMS AND MOLECULES); SCATTERING EXPERIMENTS (NUCLEI).

Lawrence O. Brockway
Bibliography. W. Braun, *Applied RHEED: Reflection High-Energy Electron Diffraction During Crystal Growth*, Springer-Verlag, 1999; J. M. Cowley, *Diffraction Physics*, 3d ed., 1995; J. M. Cowley (ed.), *Electron Diffraction Techniques*, 2 vols. 1992, 1993; D. L. Dorset, *Structural Electron Crystallography*, 1995; C. S. Fadley, The study of structures by photoelectron diffraction and Auger electron diffraction, in R. Z. Bachrach (ed.), *Synchrotron Radiation*

Research: Advances in Surface and Interface Science, 1992; A. Howie and U. Valdre (eds.), *Surface and Interface Characterization by Electron Optical Methods*, 1988; P. K. Larsen and P. J. Dobson (eds.), *Reflection High-Energy Electron Diffraction and Reflection Electron Imaging of Surfaces*, 1988; D. L. Misell and E. B. Brown, *Electron Diffraction: An Introduction for Biologists*, 1988; M. A. van Hove, W. H. Weinberg, and C. M. Chan, *Low Energy Electron Diffraction*, 1987.

Electron emission

The liberation of electrons from a substance into vacuum. Since all substances are built up of atoms and since all atoms contain electrons, any substance may emit electrons; usually, however, the term refers to emission of electrons from the surface of a solid.

The process of electron emission is analogous to that of ionization of a free atom, in which the latter parts with one or more electrons. The energy of the electrons in an atom is lower than that of an electron at rest in vacuum; consequently, in order to ionize an atom, energy must be supplied to the electrons in some way or other. By the same token, a substance does not emit electrons spontaneously, but only if some of the electrons have energies equal to, or larger than, that of an electron at rest in vacuum. This may be achieved by various means. If a substance is heated, the atoms begin to vibrate with larger amplitudes, and electrons may absorb sufficient energy from these vibrations to be emitted in the process known as thermionic emission. Electrons may also be liberated upon irradiation of the substance with light (photoemission). Electron emission from a substance may be induced by bombardment with charged particles such as electrons or ions in the phenomenon called secondary emission. Field emission, or cold emission, is the emission of electrons under influence of a strong electric field. Electrons may also be emitted from one solid into another, but this is usually referred to as electron injection. For example, a metal may inject electrons into an insulator under certain circumstances. See FIELD EMISSION; PHOTOEMISSION; SECONDARY EMISSION; THERMIONIC EMISSION.

Adrianus J. Dekker

Bibliography. R. Gomer, *Field Emission and Field Ionization*, 1961, reprint 1993; G. Hohler, E. A. Nieksch, and J. Treusch (eds.), *Particle Induced Electron Emission*, 1991; A. Modinos, *Field, Thermionic and Secondary Electron Emission Spectroscopy*, 1984.

Electron-hole recombination

The process in which an electron, which has been excited from the valence band to the conduction band of a semiconductor, falls back into an empty state in the valence band, which is known as a hole. See HOLE STATES IN SOLIDS.

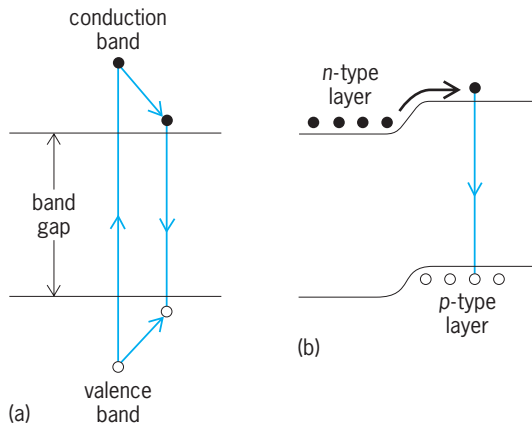


Fig. 1. Recombination of electrons and holes generated by (a) optical absorption and (b) a forward-biased pn junction.

When atoms come together to form a crystal, the discrete atomic energy levels broaden into a level continuum. The periodic potential of the crystal produces gaps in the continuum, separating it into a number of energy bands. In a pure semiconductor crystal, bonding electrons fill the valence band, and above the valence band in energy is an empty conduction band. The two bands are separated by the band gap, which has no levels, and the extent of this band gap is a characteristic of each semiconductor. The band gap is of order 1 eV for common semiconductor materials such as silicon (Si) or gallium arsenide (GaAs). See BAND THEORY OF SOLIDS.

Light with photon energies greater than the band gap can be absorbed by the crystal, exciting electrons from the filled valence band to the empty conduction band (Fig. 1a). The state in which an electron is removed from the filled valence band is known as a hole. It is analogous to a bubble in a liquid. The hole can be thought of as being mobile and having positive charge. The excited electrons and holes rapidly lose energy (in about 10^{-12} s) by the excitation of lattice phonons (vibrational quanta). The excited electrons fall to near the bottom of the conduction band, and the holes rise to near the top of the valence band, and then on a much longer time scale (of 10^{-9} to 10^{-6} s) the electron drops across the energy gap into the empty state represented by the hole. This is known as electron-hole recombination. An energy approximately equal to the band gap is released in the process. Electron-hole recombination is radiative if the released energy is light and nonradiative if it is heat. See PHONON.

Recombination at pn junctions. Electron-hole recombination requires an excited semiconductor in which both electrons and holes occupy the same volume of the crystal. This state can be produced by purely electrical means by forward-biasing a pn junction. The pn junction is formed by adjacent p -type and n -type layers containing excess holes and excess electrons, respectively. In the n -type layer, donor impurities replace some of the host atoms. Each donor impurity has one more electron in its outer shell than

the host atom, and these extra electrons go into the conduction band. Likewise, in the p -type layer, acceptor impurities replace some of the host atoms. The acceptor atoms contain one less electron in their outer shell than the host atoms. They complete their valence bonds by taking a valence electron from elsewhere in the crystal, thereby creating holes in the valence band.

The pn junction forms a barrier that keeps the electrons and holes separated. By contacting these layers and applying a forward-bias voltage to the pn junction, the barrier height can be reduced, enabling electrons to flow into the p -region or holes into the n -region, thus allowing electron-hole recombination to take place. The current passing through a pn diode in electrons per second equals the rate of electron-hole recombination (Fig. 1b). A major application of this phenomenon is the light-emitting diode (LED), which generates light by means of radiative electron-hole recombination resulting from forward-biasing the pn junction. See LIGHT-EMITTING DIODE; LUMINESCENCE; SEMICONDUCTOR DIODE.

Radiative recombination. Semiconductor materials may be divided into two types: direct-band-gap semiconductors and indirect-band-gap semiconductors. Efficient radiative recombination between free electrons and holes takes place only in direct-band-gap semiconductors. The continuous energy levels of electrons in the conduction and valence bands are characterized by a definite value of momentum (or more precisely, wave vector). During an optical transition, momentum is conserved, and since the photon carries away negligible momentum, transitions take place only between conduction-band and valence-band states having the same momentum. This is easily satisfied in direct-band-gap semiconductors, because electrons and holes collect at the conduction band at minimum and the valence band at maximum, and both extrema have the same momentum (Fig. 2a). However, for indirect-band-gap semiconductors (Fig. 2b), the conduction-band minimum and valence-band maximum have very

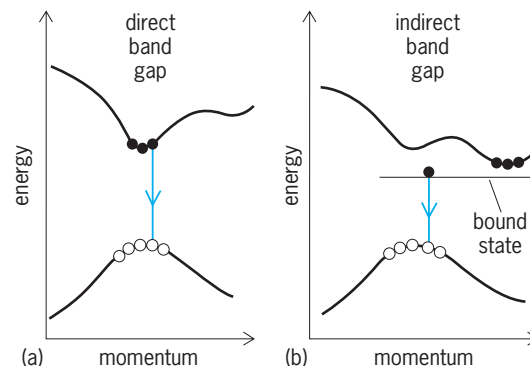


Fig. 2. Energy versus momentum for the conduction-band and valence-band states of (a) a direct-band-gap semiconductor and (b) an indirect-band-gap semiconductor. Optical transitions in indirect-band-gap semiconductors take place through bound states that have a momentum distribution as indicated.

different momenta, and consequently optical transitions between free electrons and holes are forbidden.

Radiative electron-hole recombination is possible in indirect-band-gap semiconductors when the transition is assisted by lattice phonons and impurities. Weak optical transitions can take place in which a momentum-conserving phonon is emitted along with the photon. Electrons can be trapped at impurities having energy levels within the energy gap. Unlike the free electron, the bound electron has a spread of momentum as a consequence of its localization (the uncertainty principle), with some components of momentum at the valence-band maxima (Fig. 2*b*). The isoelectronic impurity nitrogen substituted for phosphorus in gallium phosphide (GaP) and gallium arsenide phosphide (GaAsP) is used for this purpose. When heavily doped with nitrogen, these indirect semiconductors have moderate luminescence efficiency and are used as light-emitting diodes to generate green light. Bismuth impurities can be used in the same way to produce yellow light. It would be more desirable to use direct-band-gap semiconductors for this purpose, but currently there are no direct-band-gap semiconductors capable of forming *pn* junctions which emit these colors at room temperature. See CRYSTAL DEFECTS.

Apart from its application in light-emitting diodes and laser operation, radiative recombination, especially at low temperatures (approximately 2 K or -456°F), has been a very important tool for studying the interaction of electrons and holes in semiconductor crystals. At low temperatures, many weakly bound states such as excitons are able to form stably. These states interact only weakly with lattice vibrations and consequently emit much of their recombination radiation in the form of sharp (zero phonon) spectral lines. For example, at low temperatures, electrons become trapped at donor sites and holes at acceptor sites. The spectrum of donor-acceptor recombination contains hundreds of resolvable lines, one for each crystallographically inequivalent donor-acceptor pair. See EXCITON.

Nonradiative recombination. Competing with radiative recombination are the nonradiative recombination processes of multiphonon emission and Auger recombination. Nonradiative recombination of electrons and holes by multiphonon emission occurs in two steps. First, an electron is captured into a deep level near the middle of the energy gap, and subsequently the trapped electron state captures a hole from the valence band (Fig. 3). Multiphonon emission requires atoms very different from the host, or defects such as vacancy-impurity complexes, dislocations, or semiconductor-oxide interface states. Electron capture alters the positions of the surrounding lattice, in a process known as lattice relaxation that lowers the electronic energy of the bound electron (Fig. 3). The vibration of the lattice causes the electronic level to move within the energy gap and even cross into the bands. Capture takes place during the crossing.

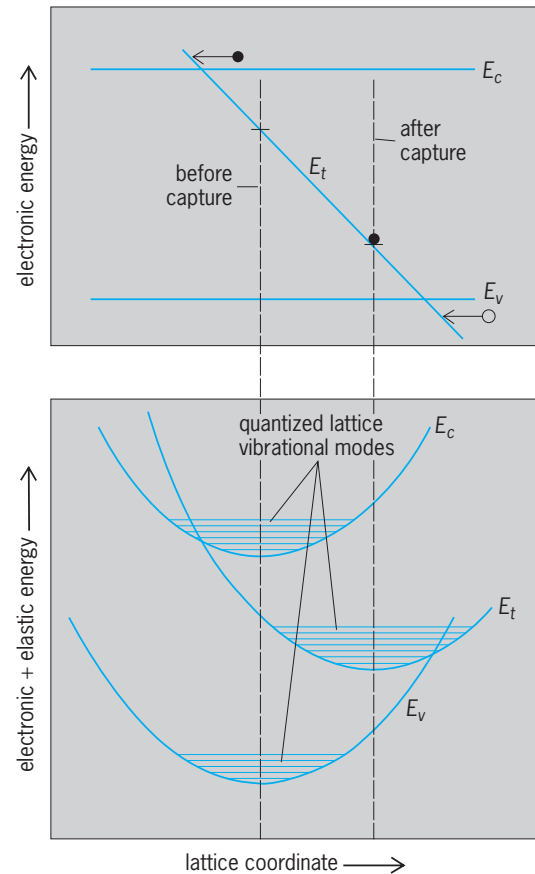


Fig. 3. Electron and hole capture into a deep level by multiphonon emission. After capture the lattice relaxes, lowering the electronic energy of the bound state E_t . E_c and E_v are the minima of the conduction band and maxima of the valence band, respectively.

In multiphonon emission, the energy lost during capture highly excites the vibrations of one or several atoms related to the defect site. This can greatly enhance the probability for an atom to jump to a different site, which also requires a highly excited atomic vibration. Such recombination-enhanced motion has been observed in detail for defects formed by ionizing radiation in GaAs and in Si. It is suspected that nonradiative recombination by multiphonon emission drives the movement of atoms at room temperature that are responsible for device degradation phenomena such as the climb of dislocations found in GaAs light-emitting diodes and lasers. See DIFFUSION; RADIATION DAMAGE TO MATERIALS.

Another source of nonradiative recombination is the Auger effect, in which an electron or hole takes up the energy released in the recombination of a pair of other charge carriers. Auger recombination has been observed for excitons bound to neutral donors in Si and GaP, a complex of two electrons and a hole. In these cases, the Auger recombination rate is several orders of magnitude greater than the radiative rate. Auger recombination has been shown to limit the performance of long-wavelength (1.3–1.6 micrometer) lasers and light-emitting diodes used in optical communication systems. The Auger effect reduces light-emitting diode

efficiency and greatly enhances the increase of laser threshold current with temperature. See AUGER EFFECT; LASER; OPTICAL COMMUNICATIONS; SEMICONDUCTOR.

Charles H. Henry

Bibliography. V. N. Aba Kumov, V. I. Perel, and I. N. Yassievich, *Nonradiative Recombination in Semiconductors*, 1991; G. Gillessen and W. Schairer, *Light-Emitting Diodes: An Introduction*, 1987; P. T. Landsberg, *Recombination in Semiconductors*, 1992; J. Singh, *Semiconductor Devices An Introduction*, 1994; S. M. Sze, *Semiconductor Devices*, 1985.

Electron lens

An electric or magnetic field, or a combination thereof, which acts upon an electron beam in a manner analogous to that in which an optical lens acts upon a light beam. Electron lenses find application for the formation of sharply focused electron beams, as in cathode-ray tubes, and for the formation of electron images, as in infrared converter tubes, various types of television camera tubes, and electron microscopes.

Any electric or magnetic field which is symmetrical about an axis is capable of forming either a real or a virtual electron image of an object on the axis which either emits electrons or transmits electrons from another electron source. Hence, an axially symmetric electric or magnetic field is analogous to a spherical optical lens.

The lens action of an electric and magnetic field of appropriate symmetry can be derived from the fact that it is possible to define an index of refraction n for electron paths in such fields. This index depends on the field distribution and the velocity and direction of the electrons. It is given by the equivalencies shown below. Here e is the charge of the electron,

$$n = \sqrt{\phi + \frac{2e\phi^2}{mc^2}} - \sqrt{\frac{e}{2m}} A \cos \chi$$

$$= \sqrt{\phi + 0.978 \cdot 10^{-6}\phi^2} - 0.297A \cos \chi$$

m its rest mass, ϕ the potential of the point in space under consideration (so normalized that the kinetic energy of the electron vanishes for $\phi = 0$), c the velocity of light, A the magnetic vector potential, and χ the angle formed by the path of the electron with the direction of the magnetic vector potential. For an axially symmetric field the magnetic vector potential is perpendicular to the plane passing through the axis of symmetry and the reference point. Its magnitude is equal to the magnetic flux through the circle about the axis through the reference point divided by the circumference of that circle. The numerical coefficients in the final expression for n apply if ϕ is measured in volts and A in gauss-centimeters. See ELECTRON MICROSCOPE.

Electron lenses differ from optical lenses both in the fact that the index of refraction is continuously variable within them and that it covers an enormous range. Furthermore, in the presence of a magnetic

field, n depends both on the position of the electron in space and on its direction of motion. It is not possible to shape electron lenses arbitrarily. See ELECTROSTATIC LENS; MAGNETIC LENS.

Edward G. Ramberg

Bibliography. D. A. DeWolf, *Basics of Electron Optics*, 1990; E. M. Slayter and H. S. Slayter, *Light and Electron Microscopy*, 1993.

Electron microscope

A device for forming greatly magnified images of objects by means of electrons. Electron microscopes serve primarily two purposes: the visual examination of structures too fine to be resolved with ordinary, or light, microscopes, and the study of surfaces that emit electrons. The first function made transmission electron microscopes essential research tools in biology, chemistry, and metallurgy. Beginning in the 1960s the scanning electron microscope came to play an increasingly important role in the study of the surfaces of solid objects at more moderate magnifications. Various emission electron microscopes serve more specialized research purposes.

Instruments

A transmission electron microscope, shown in close analogy with a light microscope in **Fig. 1**, consists in its simplest form of a source supplying a beam of electrons of uniform velocity, a condenser lens for concentrating the electrons on the specimen, a specimen stage for displacing the specimen which transmits the electron beam, an objective lens, a projector lens, and a fluorescent screen on which the final image is observed. For permanent record of the

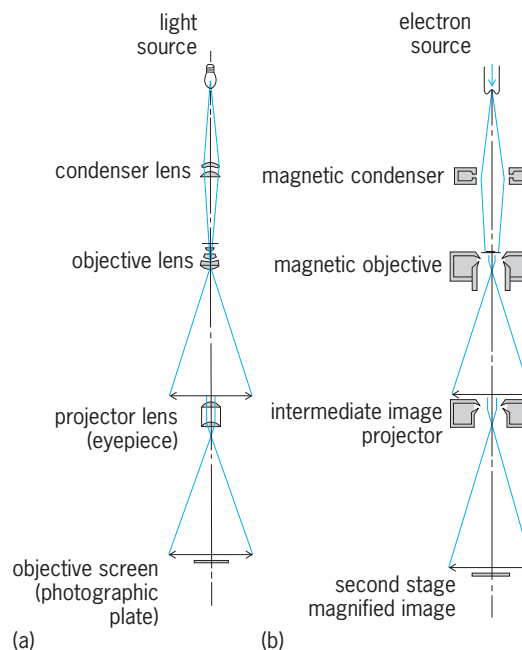


Fig. 1. Comparison between (a) light microscope and (b) magnetic electron microscope. (After V. K. Zworykin et al., *Electron Optics and the Electron Microscope*, Wiley, New York, 1945)

image, the fluorescent screen is replaced by a photographic plate or film.

Electrons are strongly scattered by all forms of matter including air. Hence the entire instrument must be evacuated to about 10^{-4} mmHg (10^{-7} atm or 10^{-2} pascal). Furthermore, the lenses cannot be material in nature. Instead, they are electric or magnetic fields, symmetrical about the axis of the instrument, that have the property of bending the electron paths toward the axis, just as converging glass lenses bend light rays toward their axis. The lenses in the instruments shown in Fig. 1 are magnetic fields formed at narrow gaps in the iron casings surrounding coils traversed by electric current. Lens strength is varied by varying the current. Most electron microscopes employ magnetic lenses of this type. These have yielded the highest resolution and magnification attained.

However, good results have also been obtained with electron microscopes employing unipotential electrostatic lenses and magnetic lenses excited by permanent magnets. See ELECTRON LENS; ELECTROSTATIC LENS; MAGNETIC LENS.

Resolution. A microscope can, at best, permit the discrimination of two point objects greater than $0.7\lambda/\sin \theta$ apart. Here λ is the wavelength of the illuminating radiation, and θ is the aperture angle of the cone of radiation that participates in forming the image (Fig. 2). For green light $\lambda = 500$ nanometers. Even for ultraviolet radiation in an immersion medium of refractive index 1.5, λ is no less than 170 nm.

Since light and ultraviolet microscope objectives can be designed to utilize practically all the radiation passing through the specimen, $\sin \theta \cong 1$ and the least resolvable distance for the ultraviolet microscope is about 100 nm. For 50- to 100-kV electrons, such as are commonly employed in electron microscopes,

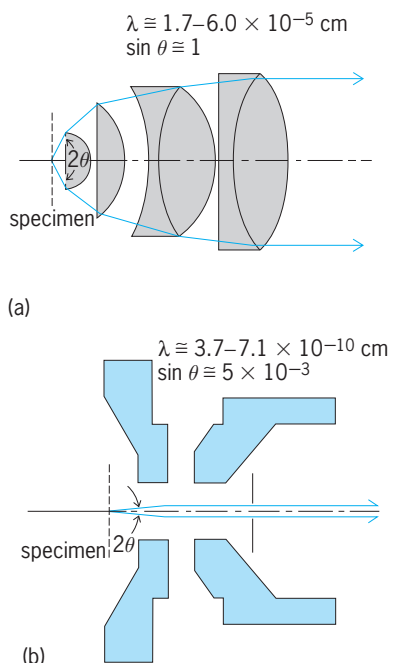


Fig. 2. Effective objective aperture in (a) light and (b) electron microscopes.

the wavelength range is 0.0053–0.0037 nm. Hence, even though a cone of radiation with an aperture angle less than 0.01 radian contributes to an image of optimum sharpness, object separations smaller than 0.3 nm have been resolved with the electron microscope.

Thus the electron microscope has several hundred times the resolving power of the light microscope. Similarly, whereas the maximum useful magnification of the light microscope is about 2000, that of the electron microscope may approach 1,000,000. The maximum useful magnification is the least magnification of the image that reveals to the observer all the specimen detail that the microscope is capable of conveying.

Transmission electron microscope. Electrons are commonly emitted from the tip of a fine tungsten-wire hairpin filament or, to further reduce the size of the effective electron source, from a sharply pointed segment of wire welded to the filament tip (Fig. 3). The filament is maintained at a carefully stabilized negative potential of 50–100 kV with respect to the remainder of the instrument. Electrons enter the instrument through an anode aperture. The intensity and convergence of the electron beam that is falling on the specimen are adjusted by varying the coil current of the condenser lens.

Image contrasts are formed by the scattering of electrons out of the narrow cone that contributes to the formation of the image; denser or thicker portions of the specimen scatter more electrons and hence appear darker in the image. The sharpness of the image observed on the screen is adjusted by varying the objective coil current, and its magnification by varying the projector coil current. Both currents must be carefully stabilized to yield high resolution. In modern instruments the function of the projector is performed by two magnetic lenses in tandem. In the instrument shown in Fig. 3 the normal magnification range from 1400 to 200,000 is realized by varying the coil currents. The magnification range can be extended downward to 500 with the aid of a special specimen holder and projector polepiece. A binocular viewer provides an additional magnification by 10, needed for the fine adjustment of the objective current to yield images of maximum sharpness. Comparable magnifications of the final image are obtained by enlarging the fine-grain photographic plates on which the electron image is recorded.

The image intensity required for precise visual focusing and specimen exploration is considerably higher than that needed for photographic recording. This intensity, and the corresponding specimen damage by the electron beam, can be greatly reduced through the use of an image intensifier optically coupled to the image on the fluorescent screen of the electron microscope on the one hand and the target of a television camera tube on the other (Fig. 4). The picture signal generated by the camera tube can be employed to reproduce the microscope image instantaneously on the screen of a television receiver, can be recorded on tape by a video recorder for instant playback at a later time, or can be distributed

to a group of receivers for teaching purposes. Since the image can be photographed on the screen of the television receiver, this system makes it also possible to avoid contamination of the vacuum in the microscope column resulting from the introduction of photographic material. Furthermore, it permits arbitrary control of image contrasts.

In addition to the standard transmission microscopes operating at 50–100 kV, a number of very high-voltage instruments have been constructed. Instruments designed for operation up to 1500 kV (Fig. 5) are necessarily large, because clearances in the gun and the high-voltage supply are required to suppress field emission, and because the focal lengths of electron lenses increase as the accelerating voltage of the electrons is increased. The advantage of high-voltage electron microscopy does not lie in greater resolving power—the effect of a shorter electron wavelength (0.00087 nm at 1,000,000 eV) is more than outbalanced by the increased aberration of the lenses—but in increased penetration, which is particularly valuable in the direct study of metal sections prepared with a microtome.

Scanning electron microscope. A modern version of the scanning electron microscope is illustrated in Fig. 6. The microscope column, with the source at the top and the specimen chamber at the bottom, is at the left. The square cathode-ray tube display for visual observation with a long-persistence yellow phosphor is in the center, and the cathode-ray tube display for photographic recording with a short-persistence blue phosphor is at the right. Figure 7 is a schematic diagram of the microscope column. The electron gun is similar to that in the transmission microscope, but operates at much lower voltages (1–20 kV). Three magnetic lenses image the first crossover, or narrowest cross section of the beam issuing from the filament tip, at a greatly reduced scale on the specimen. A pair of deflection coils actually mounted within the bore of the final lens serves to deflect the electron probe formed by the lenses over the surface of the specimen, the specimen and the cathode-ray displays being scanned in synchronism. The secondary electrons generated at the point of impact of the electron beam on the specimen are accelerated by 12 kV and impinge on a plastic scintillator optically coupled by a plastic light pipe to a multiplier phototube. The output current of this phototube controls the beam current, and hence the brightness in the display. See SCINTILLATION COUNTER.

The whole scanning electron microscope can be regarded as a slow-scan television system, with the microscope column constituting a very high-resolution camera tube, the specimen serving as camera tube target. Resolution is determined by the electron probe size at the specimen and magnification by the ratio of the scan amplitude in the display to that at the specimen. Thus the magnification can be varied, for example, in a range of 20–100,000, without any refocusing. The effective probe size is determined not only by the electron lens properties but also by the scattering of the electrons within the

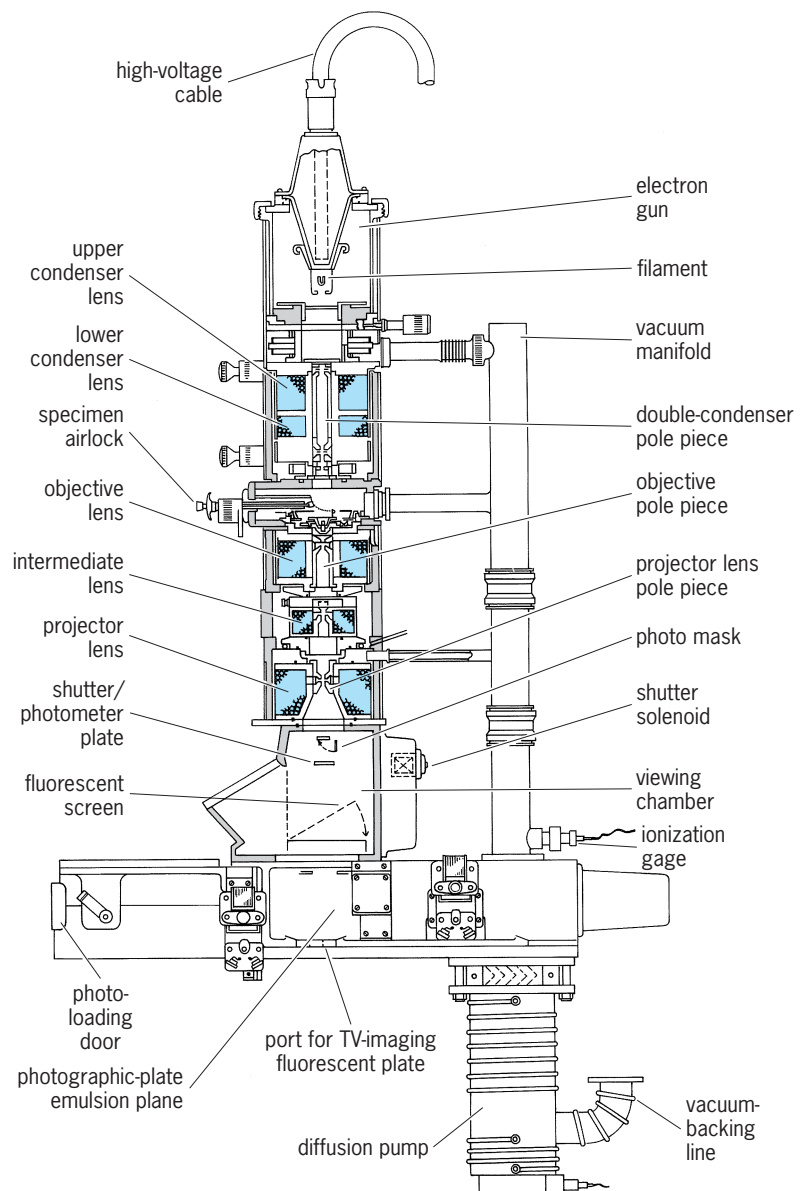


Fig. 3. Section through column of an electron microscope. This instrument has a normal magnification range from 1400 to 200,000. (Radio Corporation of America)

specimen. Since the range of electron diffusion decreases as the accelerating voltage of the electrons is reduced, relatively low operating voltages are desirable. The optimal resolution cited for the instrument in Fig. 6 is better than 20 nm, the normal operating resolution about 50 nm. At the same time, the great depth of focus (about 300 times that of a light microscope) endows images with considerable extent in depth with an extraordinarily lifelike, three-dimensional quality. This is lacking in light microscope images of the same magnification, since the image is sharp only at a rather sharply defined planar section through the object surface. See CATHODE-RAY TUBE; TELEVISION CAMERA TUBE.

Since the electron probe currents are very small (10^{-9} to 10^{-13}), the recording of noise-free high-resolution images demands relatively long scan periods, ranging from 2 to 1000 s; for visual observation

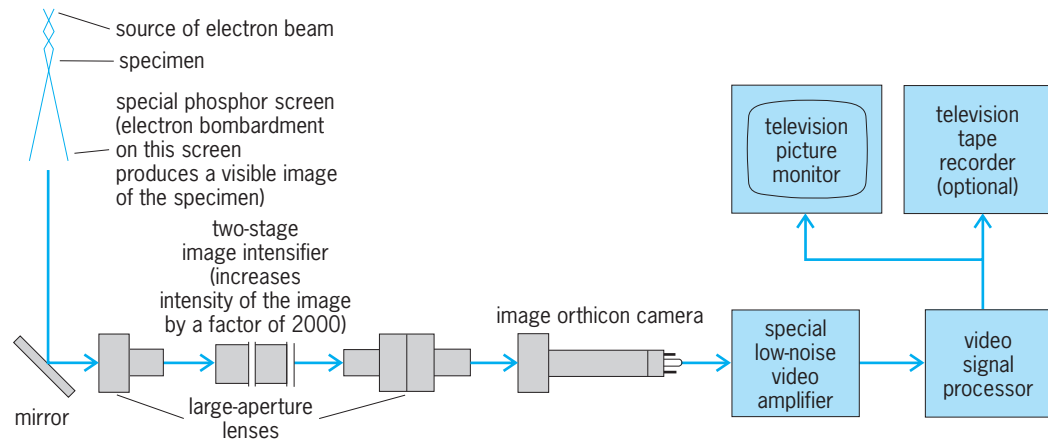


Fig. 4. Scanning electron microscope. (Cambridge Instrument Co., Ltd.)

on the long-persistence screen the frame period is reduced to 0.1 to 12 s.

The signal pickup system can be modified in various ways to show different properties of the object. In the examination of integrated circuit elements a comparison of the image obtained in the conductive mode (Fig. 8), the signal being derived from an appropriately chosen electrode on the specimen, with that obtained in the normal emissive mode is frequently useful in testing the operation of the system and locating defects. See INTEGRATED CIRCUITS; SCANNING ELECTRON MICROSCOPE.

Emission electron microscopes. These may take several forms (Fig. 9). In the immersion electron microscope, for which a rudimentary form is shown in Fig. 9a, the specimen is a flat conducting surface which may be heated, illuminated, or bombarded

by high-velocity electrons or ions so as to emit low-velocity thermionic, photo, or secondary electrons. These are accelerated to a high velocity in an immersion objective or cathode lens and imaged as in a transmission electron microscope.

Other emission microscopes employ simple projection without lenses. One example (Fig. 9b) is the cylindrical microscope used for examining the thermionic emission from a thin wire stretched along the axis of a cylindrical tube whose walls at high positive potential are coated with a fluorescent substance. Another (Fig. 9c) is the field-emission microscope, in which field electron currents are drawn from an extremely fine, rounded point on the end of a wire of a metal such as tungsten. The electrons drawn out of the point follow nearly straight lines toward a positive-potential screen. The display on the screen shows variations in the emission characteristics (as determined by the work function and surface nonuniformities) over the surface of the point. See FIELD-EMISSION MICROSCOPY; THERMIONIC EMISSION.

Low-temperature techniques. Apart from freeze-etching, low-temperature techniques find several other applications in electron microscopy.

One of these is employment of a superconducting objective lens, as demonstrated by H. Fernandez-Moran. The lens field is produced by an air-core solenoid of niobium-zirconium alloy, cooled by a liquid-helium bath and short-circuited after excitation so as to remain in the persistent-current mode. While the immersion of the microscope column in liquid helium presents a very serious inconvenience, such superconducting solenoids may generate higher magnetic fields and can thus provide, in principle, lenses of shorter focal length and lower aberration than conventional iron pole-piece objectives. Furthermore, they are characterized by extraordinary stability. See SUPERCONDUCTIVITY.

Liquid air-cooling of the specimen chamber also serves to suppress contamination deposited on the specimen through the interaction of the electron beam with traces of organic vapor; without special precautions, contamination buildup can quickly



Fig. 5. TV image intensification system for the electron microscope. (Radio Corporation of America)



Fig. 6. Megavolt electron microscope at U. S. Steel research laboratory at Monroeville, Pennsylvania. (*Radio Corporation of America*)

obscure fine specimen detail at very high magnifications.

Finally, freeze-drying and freeze substitution provide ways of preparing organic specimens for observation with minimal distortion of their cell structure. In the former instance the specimen may be cooled extremely rapidly by immersion in liquid air-cooled propane and then be warmed slowly in vacuum, so that the ice sublimates rather than melts. The infusion of the embedding material then takes place either in vacuum or a dry-nitrogen atmosphere. In freeze substitution ice is removed by diffusion into a liquid solvent instead of sublimation in vacuum.

Edward G. Ramberg

Electron Microscopy

The application of the electron microscope to examination and investigation of the ultramicrostructure of materials has become so extensive that there is hardly an area in biological and nonbiological research where electron microscopy does not play

a role. In biological and medical research the development of sectioning techniques has extended electron microscopy down to observations at the macromolecular level for delineation of the complex organization of cell components such as membranes, mitochondria, endoplasmic reticulum, and ribosomes. Other techniques have been developed for studies on the structure of virus and even individual proteins and nucleic acids. Nonbiological solid materials have also become objects of extensive and fruitful investigation. Diffraction microscopy, together with the development of adequate techniques of sectioning and thinning crystalline materials such as metals, now makes the observation of defect structure in solids an important aspect of electron microscopy.

The very different character of the interaction of electrons and light with matter produces the greatest practical differences between electron microscopy and optical microscopy. Contrast in the optical microscope image is produced largely by

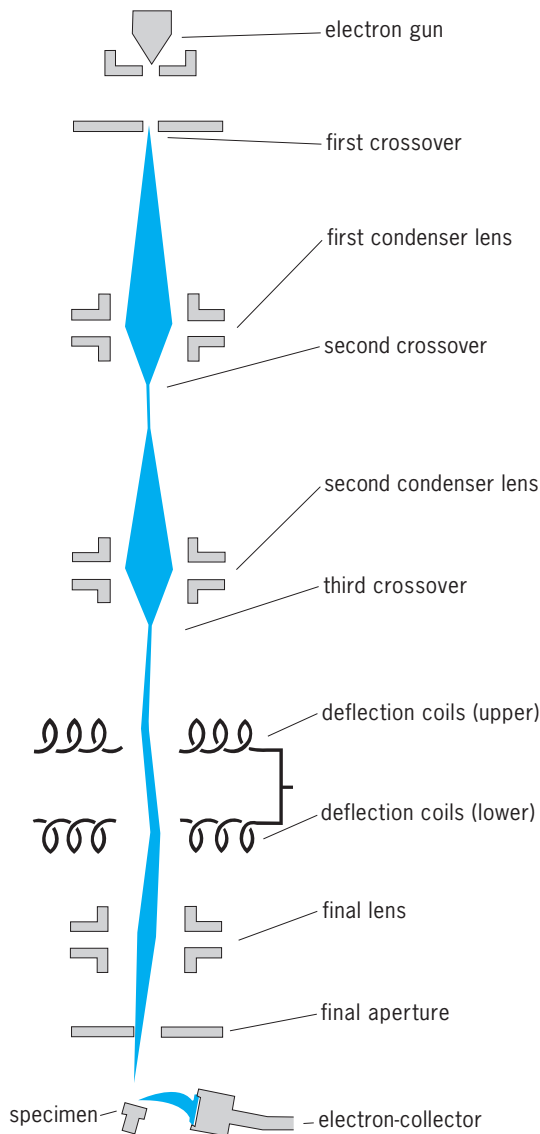


Fig. 7. Scanning electron microscope column. (Cambridge Instrument Co., Ltd.)

spectral adsorption of light passing through the specimen. Contrast in the electron microscope image is produced primarily by the scattering of electrons from the specimen. Electrons in the 50–100-kV range, used in conventional commercial electron microscopes, are strongly scattered by matter. It is this strong interaction that determines the requirements for specimen preparation and places serious limitations on the examination of materials. As indicated previously, to avoid scattering of electrons by ambient gas molecules, an electron beam must be confined to a column evacuated to a pressure of the order of 10^{-4} to 10^{-5} mmHg (10^{-7} to 10^{-8} atm or 10^{-2} to 10^{-3} Pa). Thus, specimens must be so prepared that they can be placed in a high vacuum. Dehydration is required, and this requirement immediately places a severe restraint on all biological systems. Moreover, the total amount of material in a specimen must be very small to prevent excessive heating and the specimen must also be extremely thin. Substrates for mounting specimens are typically of the order of 10 nm thick. Sections of solid materials should be no thicker than the order of 0.1 micrometer. Biological tissues are usually sectioned to thicknesses of about $0.05 \mu\text{m}$ to obtain high-resolution detail. Metal foils are usually thicker, ranging from 0.1 to $0.2 \mu\text{m}$ to obtain more information on their defect characteristics. As thickness is increased, multiple scattering of electrons by matter destroys the contrast in the image; finally, energy transferred from the beam by inelastic collisions becomes so large that the specimen may be destroyed.

High-voltage electron microscopes have been developed in which the electrons are accelerated to energies of 500 keV to over 1 MeV, giving them greater penetrating power. Some of the severe limitations indicated above are relaxed to a degree in these high-voltage instruments. Thicker sections can be examined; less damage occurs in the specimen, and small specimen wet chambers with hydrated biological materials can be utilized. Development of high-voltage electron microscopy began in the mid-1960s. The problems and complexities of high-voltage electron microscopes limit their application, and most investigations will continue to be done within the limitations set by the 50–100-kV instruments.

Biological applications. Techniques and applications of electron microscopy are so specialized that it is convenient to discuss them in separate categories of biological and nonbiological, even though there are basic techniques that overlap both fields such as substrate preparation and shadow-casting.

Substrates. The material to be examined in the electron microscope must, in general, be mounted on an extremely thin supporting membrane which has a mass density per unit area less than that of the object being observed. In practice a membrane of the order of 10 nm thick is mounted on a fine electrodeposited metal mesh with openings $50\text{--}100 \mu\text{m}$ on a side, and the specimen material is deposited on these mounted substrates. Cellulose nitrate films prepared by casting a film on water from a solution of cellulose nitrate in amyl acetate produces a suitable substrate.

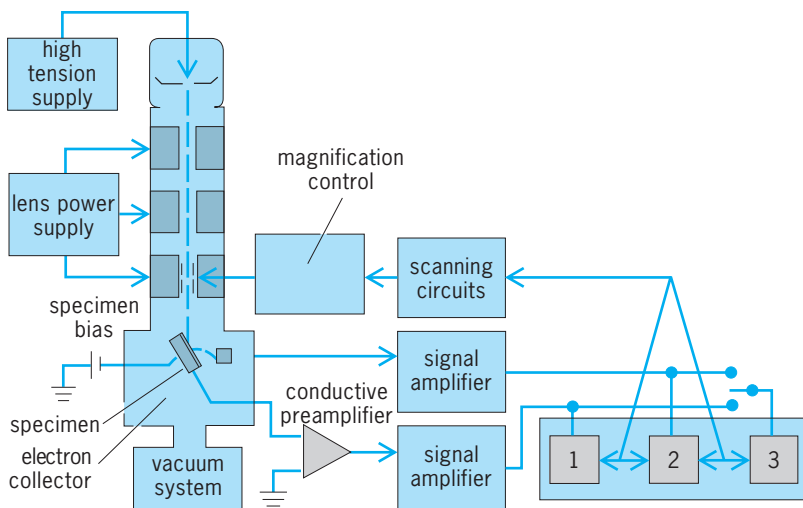


Fig. 8. Alternative operation of scanning electron microscope in conductive and emissive mode. (Cambridge Instrument Co., Ltd.)

Another substrate material is Formvar, a polyvinyl formal plastic, cast on glass from solutions of the plastic in dichloroethane (ethylene dichloride). More durable substrates can be prepared by evaporating carbon on a glass slide, stripping off the carbon film by floating it onto a water surface, and mounting the film on the standard mesh specimen screens. Carbon-film substrates as thin as 5 nm can be produced and handled in this fashion.

Shadow-casting. Materials of low mass density will not produce appreciably greater scattering than the substrate film upon which they must be mounted. Contrast in electron images from these specimens will therefore be low, and since this low differential of mass density applies especially to organic materials, it presents a large obstacle in imaging very small biological objects. The technique of shadow-casting developed by R. C. Williams and R. W. G. Wyckoff is a very powerful means of enhancing contrast in electron microscopy of small objects of low mass density. It is now possible to resolve globular macromolecular material with diameters as small as 2–3 nm and linear structures such as strands of nucleic acid less than 1 nm thick. Resolution in shadow-cast specimens is set by the granularity of the shadow-casting material.

The technique involves vacuum deposition of an extremely thin metal layer on the specimen at an oblique angle. Any projection above the surface receives a heavier deposit of the metal on the side facing the metal evaporating source and casts a shadow in which no metal is deposited (Fig. 10). The metal shadowing layer produces a topographical representation of the surface of the specimen and, as discussed below, is an essential step in all replication of surfaces for study by electron microscopy.

After shadow-casting, the metal layer is the specimen imaged by the electron microscope, for the scattering from the shadow-cast specimen is predominantly that of the deposited metal layer. The areas in the shadow where no metal is deposited will scatter relatively few electrons, and the respective areas of the image will appear dark on the photographic plate that has been exposed to the high electron density incident in these portions of the image. Projecting areas facing the evaporation source will have heavy layers of the metal and will scatter most electrons, thus appearing light on the photographic record of the image. An example of a shadow-cast specimen is shown in Fig. 11, in which the image is reproduced in the customary manner with the same relative light and dark as on the original electron micrograph negative. The effect is as though strong illumination were falling obliquely on the specimen.

The metals most generally used are uranium, platinum, or chromium. The procedures developed for producing thin films by metal evaporation are applied. A tungsten filament charged with the metal to be evaporated is heated in a vacuum of about 10^{-5} mmHg (10^{-8} atm or 10^{-3} Pa) to a temperature at which the metal has a vapor pressure of the order 10^{-2} mmHg (10^{-5} atm or 1 Pa), and a layer of mean

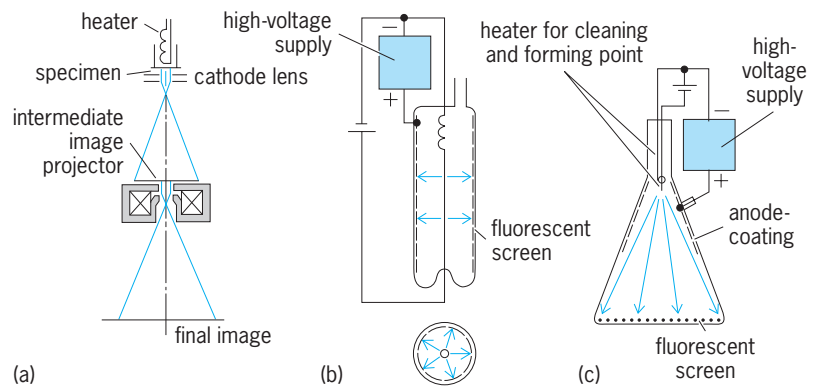


Fig. 9. Emission electron microscopes. (a) Immersion electron microscope. (b) Cylindrical thermionic-emission microscope. (c) Field-emission microscope.

thickness of about 0.5 nm is deposited on the specimen at the desired angle.

Particulate materials. Examination of particulate biological materials such as bacteria, virus, or mitochondria of cells requires proper dispersion of a very small sample of the material on the substrate. Methods must be utilized that prevent aggregation and surface tension effects of drying.

Particles already dispersed and suspended in an aqueous medium can be deposited directly on plastic or carbon substrate films. The effects of surface tension on drying, the aggregation of the particles, and the effects of changing salt concentration and pH on drying must be considered. A method of minimizing these effects is to spray the particle suspension on the substrate in very small droplets of the order of 5–10 μm in diameter. If buffers are needed, volatile ones such as ammonium acetate, ammonium carbonate, or veronal acetate must be used. A further advantage of the spray technique is the rapid drying of the small droplets; this minimizes the effects of changing concentration on the biological materials. By depositing the droplets on substrates cooled to liquid nitrogen temperatures and subliming the ice

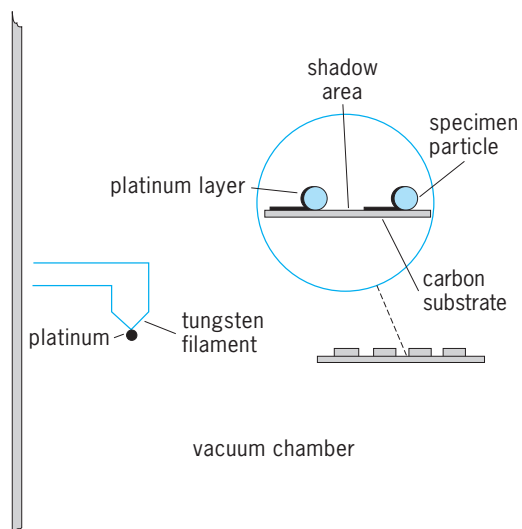


Fig. 10. Schematic of shadow-casting procedure.

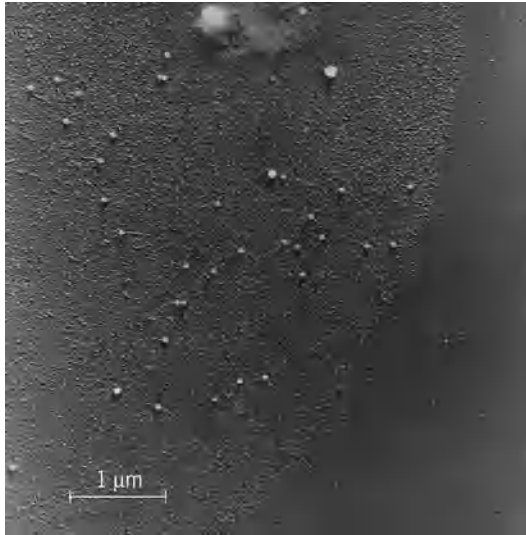


Fig. 11. Edge of spray droplet with bacteriophage.

under vacuum, the surface tension effects of drying can be further reduced.

The spray technique can be used for quantitative work, and assays can be made on suspensions of particles of unknown concentration in which a direct count is obtained on the number of particles per unit volume. Since each droplet is an aliquot sample of the total specimen, when the volume of

the droplet is known, a count of the particles in a drop pattern gives the assay. By mixing the suspension to be assayed with a given volume of a known concentration of monodispersed polystyrene latex particles and counting the number of polystyrene in each drop, the volume of the drop is obtained. When a highly purified material is available for assay, it is possible to take an aliquot portion to dryness, to weigh it, and from the electron microscope particle count to calculate the molecular weight of the particle. **Figure 12** shows tobacco mosaic virus (rods) being assayed using the polystyrene spheres.

The electron microscope is capable of resolving macromolecules of biological interest. Spraying a purified sample of macromolecular material onto a very smooth surface such as freshly cleaved mica produces a suitably dispersed specimen. This specimen is shadowed directly on the mica with platinum, is backed with an evaporated substrate of carbon, and is stripped off the mica by using methods described below.

Negative staining is another important technique used to enhance the image contrast obtained at very high magnifications from biological specimens of a particulate nature such as subcellular components, bacterial surface structure, virus particles, and protein molecules. The particulate material is treated much as described above but with a soluble heavy metal salt added to the aqueous suspension medium. A typical negative staining procedure uses a 2% solution of potassium phosphotungstate buffered and adjusted to a neutral pH of 6.8 to 7.4. A droplet of the suspension either sprayed or deposited on the substrate film dries with an embedding matrix of the heavy metal salt surrounding it. Examination in the electron microscope produces an image of reverse contrast. The surrounding matrix strongly scatters the electron beam, producing a dark ground, while particles such as viruses which should not react with the stain appear light. If the virus or other particle under examination has substructure, the heavy metal matrix can penetrate the interstices and delineate the subunits of the particle (**Fig. 13**).

Tissues and cells. Successful development of techniques for sectioning biological tissues for examination at high resolution provided a means of directly observing the complex ultrastructure of organisms. The understanding obtained of the organization of the cell and the structure of its organelles profoundly changed the fields of histology and cytology and is the most significant contribution of electron microscopy (**Figs. 14 and 15**). See CELL (BIOLOGY).

The first successful thin sectioning using a modified conventional microtome was reported in 1948 by D. C. Pease and R. F. Baker. After that techniques developed rapidly and the methods for study of the ultrastructure of cells and tissues became highly developed. Microtomes, using fractured glass or polished diamonds for knives, cut serial sections for reproduction down to thicknesses as low as a few hundred angstroms, or in special cases as thin as 5 nm. Techniques for fixation, dehydration, embedding,

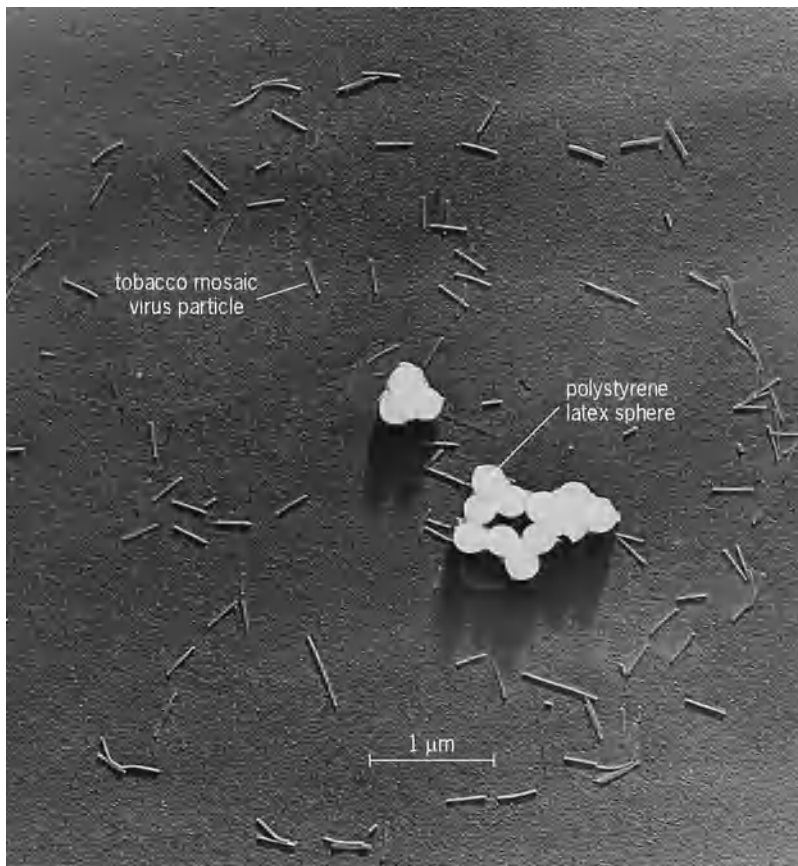


Fig. 12. Spray droplet pattern for assay of tobacco mosaic virus (rods) using polystyrene latex particles (spheres).

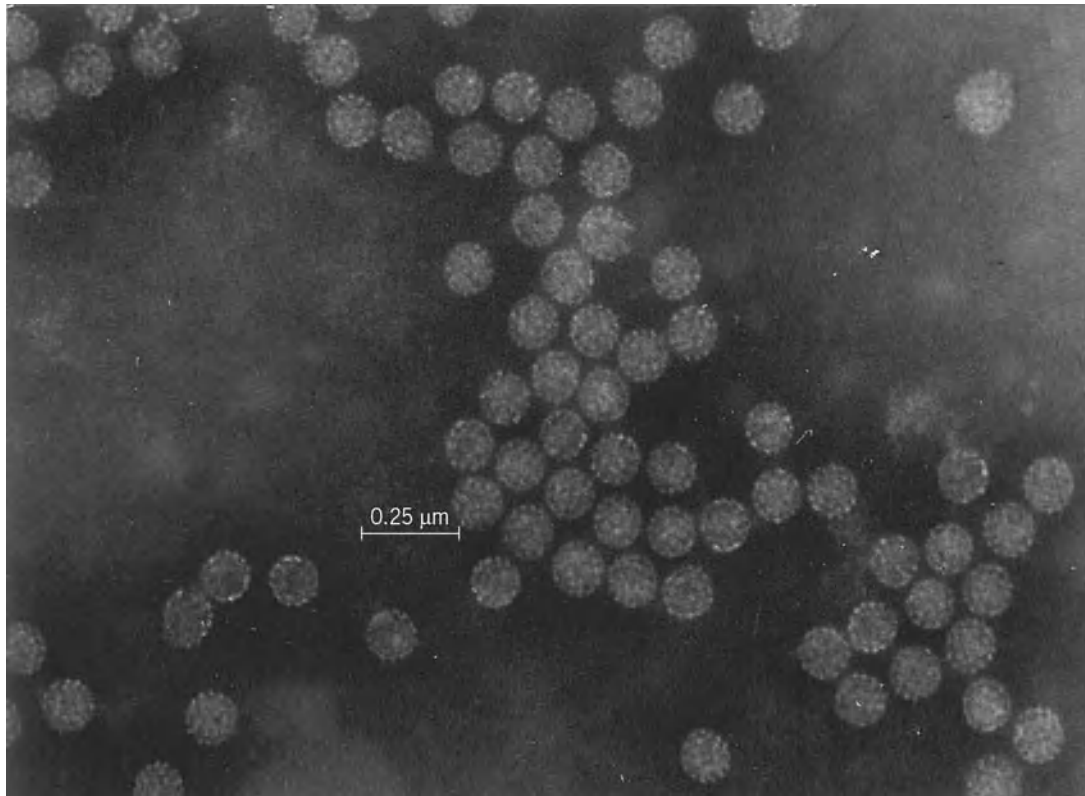


Fig. 13. Purified polyoma virus, negative stain 2% phosphotungstic acid buffered to pH 6. (Courtesy of S. S. Breese, Jr.)

and staining of tissues were developed which are adequate for a large range of observations on the complex detail of cell structure. The widely used fixatives are based on osmium tetroxide, potassium permanganate, or the aldehydes properly buffered to control pH. Embedding media were extended to a range of polymeric materials, the most successful being the epoxy resins and polyester resins. Staining of tissues either before or after sectioning can be very helpful to enhance contrast in very thin tissue sections. The general principle of staining for electron microscopy is to introduce heavy atoms such as lead or uranium into the tissue to make it more effective in scattering the electron illumination and thus to increase the visibility of those parts of the tissue to which the heavy atoms are attached. The specificity of stains available is limited, and the electron microscopist does not have available to him the range of specific staining materials with easily identifiable color contrasts that play such important roles in histological work with the light microscope. Among staining materials widely used are such heavy metal compounds as phosphotungstic acid, uranyl acetate, and lead hydroxide as well as the fixative osmium tetroxide.

There are special techniques for localization of specific entities in a cell. Immunochemical staining was developed by adapting the fluorescent antibody technique of light microscopy. Ferritin-labeled antibodies have been prepared that react with the homologous antigen and thus can be located in the cell by the increased contrast associated with this entity in the electron microscope image of the sectioned tis-

sue. See FLUORESCENCE MICROSCOPE; IMMUNOCHEMISTRY.

Another technique adapted from light microscopy and applied with increased resolution is autoradiography. A tissue is labeled with a radioactive tracer material by injecting the animal or the nutrient medium before the tissue is taken. Tritium incorporated in thymidine or similar constituents is the most widely used label. The thin section containing the radioactive element is covered with a very thin layer of photographic emulsion and allowed to stand in the dark for weeks or months. Ionizing radiation exposes the photographic emulsion producing a latent image in the silver halide grains of the emulsion at those sites at which radioactive decay has occurred. After development and fixation of the photographic emulsion and staining of the tissue section, the composite section and developed photographic emulsion is examined in the electron microscope. The developed silver grains locate the sites at which the label was attached in the section with resolutions of better than 100 nm (Fig. 15). See AUTORADIOGRAPHY.

Great depth of focus, a characteristic of electron microscopes in general and of the scanning microscope in particular, gives the scanning electron microscope a decided advantage over the light microscope for topographical observations on biological structures (Fig. 16).

Nonbiological materials. The most severe limitation in observation of nonbiological materials is their bulk nature and the fact that only very small and extremely thin specimens can be examined. Since solid material often loses its bulk properties when reduced to



Fig. 14. Electron micrograph of longitudinal section of a spider's leg just above tarsal-metatarsal joint cut through the vibration receptor (ear). (Courtesy of M. Salpeter)

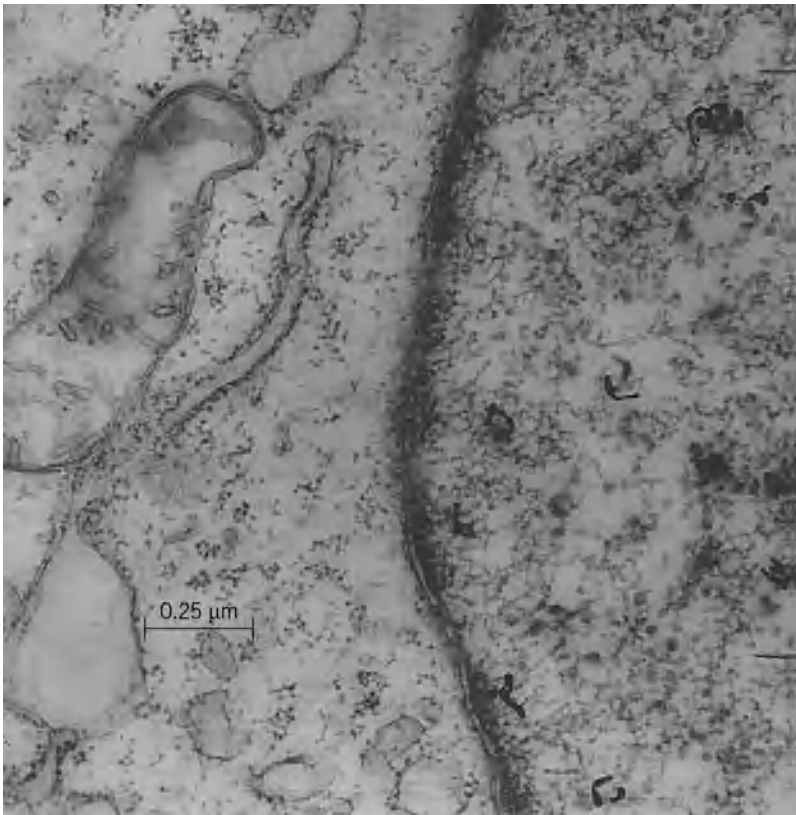


Fig. 15. Electron microscope autoradiogram of thin section of a mesenchymatous cell from the regenerating limb of adult newt (*Triturus*) labeled with tritiated thymidine. The dense black bodies over the nucleus on the right are developed silver grains, indicating radioactive decay in the tissue. (Courtesy of M. Salpeter)

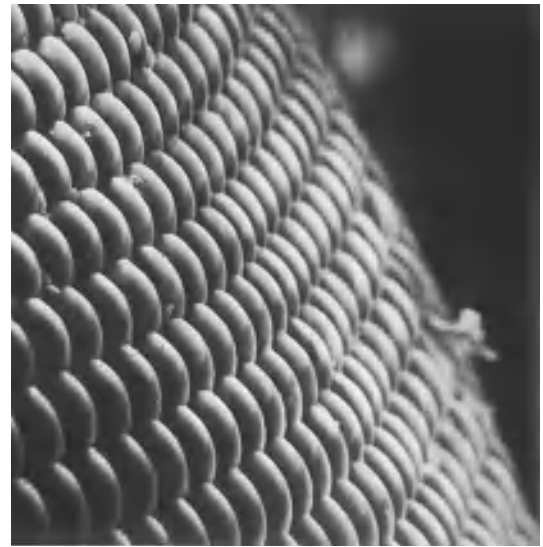


Fig. 16. Scanning electron micrograph of the surface of a fly's eye. (RCA Research Labs, Princeton, New Jersey)

very thin films or is too finely dispersed, high-voltage electron microscopes are used most actively for studies on nonbiological materials.

Particulate materials. Proper dispersion of the sample is the main consideration. If the material to be examined is a dry powder, a dispersion technique must be employed. One method used extensively is to mill the powder in a viscose solution of cellulose nitrate in amyl acetate until the aggregates are properly dispersed. The suspension is then diluted with amyl acetate to appropriate concentration for casting a suitably thin film with the particles suspended in it. Shadow-casting is usually employed to enhance the contrast or to give information on particle heights or shapes from the characteristics of their shadows.

Surfaces by replication. Surfaces of materials too thick to be placed directly in the electron microscope can be examined in the transmission electron microscope by preparing a thin film replica of the surface. The shadow-casting technique which produces contrast effects characteristic of the topography of the surface on which it is deposited greatly enhances the results obtainable by replication (Fig. 17).

A number of variations in method are employed in producing suitable replicas of surfaces. Each type of surface requires the proper choice from a variety of techniques and materials developed for application to different surfaces. The materials used for thin-film replicas include cellulose nitrate, Formvar, and evaporated carbon films. Carbon replicas are preferred for their greater strength and stability. A relatively smooth surface may be replicated directly by coating a film of Formvar on it by flooding the surface with a 0.5% solution of Formvar in ethylene dichloride. For greater ease of stripping, the replica may be backed by coating a cellulose nitrate film on the Formvar film. The backed replica can be wet-stripped, either by allowing the plastic film to float off on water if it will easily separate from



Fig. 17. Replica of a fractured surface of a magnesia ceramic. (Courtesy of T. T. Sheng)

the surface being replicated or the material being replicated can be dissolved away with a suitable solvent or acid to release replicating film. The replica film is mounted on specimen grids with the replication surface up, shadowed, and the backing film, if present, dissolved away in a solvent that does not attack the replica film. In the case of a Formvar replica backed with cellulose nitrate, amyl acetate is satisfactory. Dry stripping is accomplished by sandwiching the specimen grids between the film and Scotch tape and pulling off the film. Sometimes a stripping layer of a hydrophilic material is coated on the surface to be replicated. A detergent material such as Victawet can be evaporated onto the surface to be replicated, the metal shadowing deposited directly over the Victawet coat in the same vacuum chamber and, to complete the replica, a carbon substrate deposited by evaporation. This pre-shadowed replica is perhaps the most satisfactory and is used if the material to be replicated is suitable. On very rough surfaces, such as wood and fractured materials, it may be necessary to form a thick negative replica of a material such as polystyrene which can be molded to the surface by heat and pressure. This can then be replicated by shadowing, evaporating a carbon film on it, and releasing the replica by dissolving the polystyrene away with a suitable solvent. Materials containing water have been successfully replicated by using polyvinyl alcohol as the material for the negative replica. Purified crystalline virus in aqueous suspension and similar materials can be replicated by a rapid freezing technique. A frozen suspension is fractured to expose a fresh surface, placed in a vacuum chamber, and warmed to -110°F (-80°C) until the ice matrix is etched away slightly, leaving the material to be replicated in relief. This is cooled again to liquid nitrogen temperatures, the surface is

replicated by shadowing, and a carbon substrate is evaporated on.

Crystalline films. The electron microscope finds wide application in investigations of crystalline materials. Two techniques are used to obtain specimens of bulk materials thin enough for direct transmission microscopy. Bulk samples can be reduced to a foil of 0.004–0.008 in. (0.1–0.2 mm) thick by processes such as jet machining, spark machining, or cutting and grinding. A final thinning is carried out by chemical or electrochemical methods which process the specimen to thicknesses of 0.1–0.2 μm for direct observation in the electron microscope. The other method utilizes various methods of deposition to grow films of suitable thicknesses and characteristics (Fig. 18).

Crystalline films obtained by these methods will strongly diffract the electron beam, and image contrast is determined by the nature of these diffraction interactions. When any part of the crystalline specimen is oriented so that a strong Bragg reflection occurs, the diffracted beam is scattered outside of the

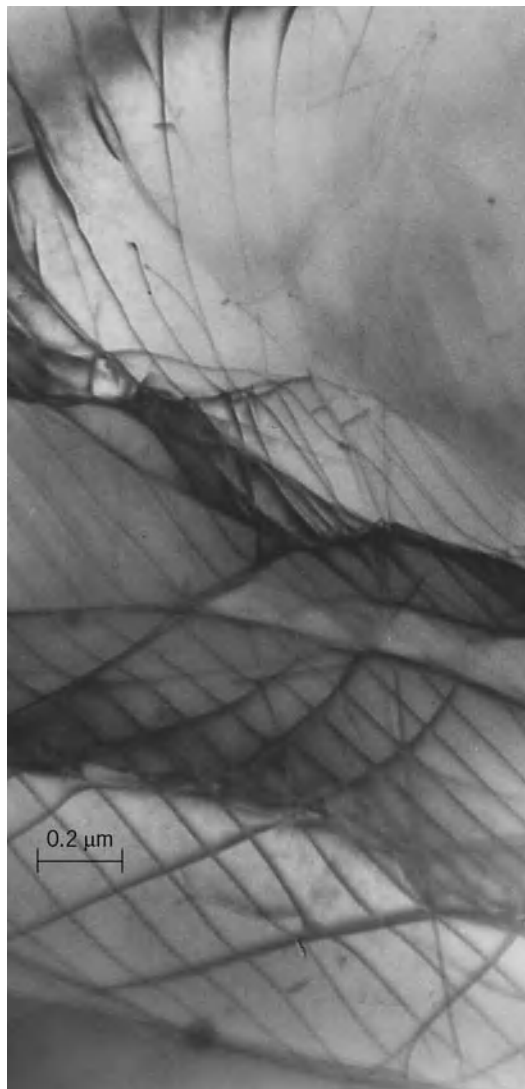


Fig. 18. Thin magnesia foil showing dislocation networks. (Courtesy of T. T. Sheng)

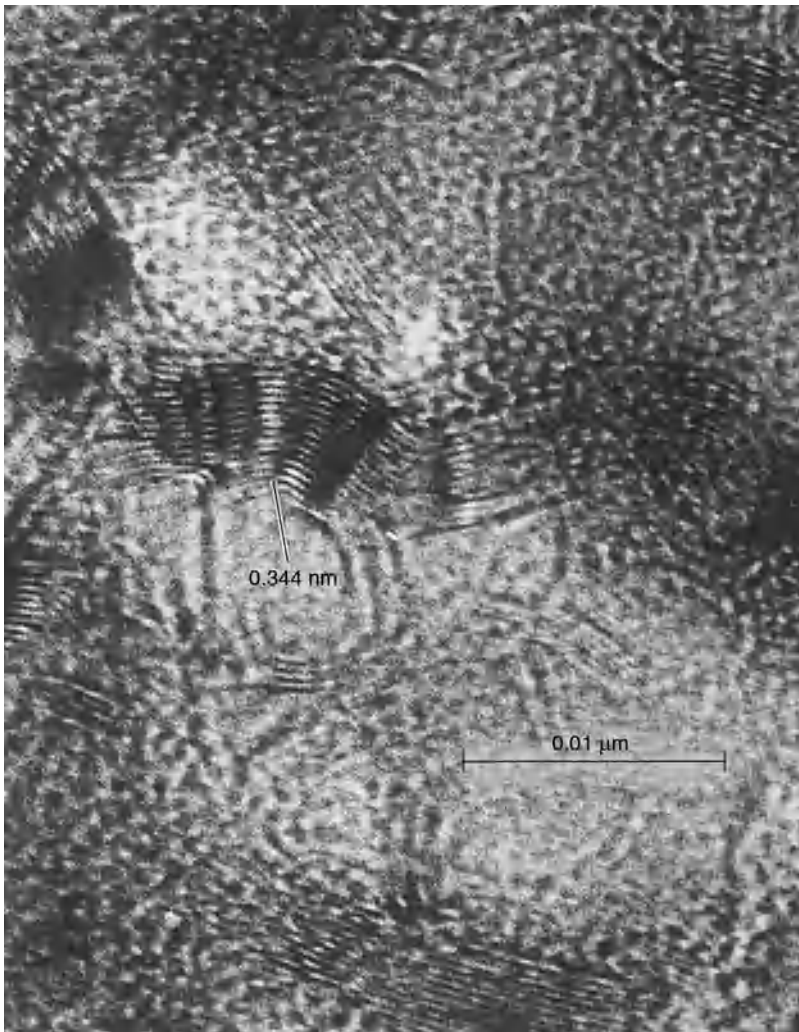


Fig. 19. Graphitized carbon showing lattice resolution of the 0.344-nm lattice spacing.

objective aperture and is no longer available in the imaging process. The image of this specimen then appears dark in areas corresponding to regions where Bragg reflection has occurred, so that high contrast is produced in the electron micrograph. Thus the crystal structure and its imperfections produce contrast effects in the imaging process that make this application of electron microscopy an important tool for study of defect structure in solids and the structure of thin films.

The electron optics of the electron microscope can be adjusted by a simple change of the focal length of one lens to produce the electron diffraction pattern from the same small area of specimen under observation with the normal electron optics for microscopy. See ELECTRON DIFFRACTION.

When the aperture stop of the objective lens is adjusted to allow one or two of the diffracted beams to contribute to the image formation in addition to the undeviated electron beam, a lattice resolution image is obtained that resolves lattice planes in crystalline material. Very thin deposited films or crystals which grow as thin platelets are required as specimens for these extremely high-resolution observations. By careful adjustment two or more lattice planes can be

very highly resolved to spacings of less than 0.2 nm. If there are imperfections in the lattice, these defects will also be observed under certain orientations of the crystal. Graphitized carbon can also produce lattice resolution images as shown in Fig. 19.

Limitations on observations. Specimen artifacts, instrumental defects, and image characteristics are factors in the limitations on observations.

Specimen artifacts. Artifacts introduced in biological specimen preparation are primarily caused by the dehydration process necessary to bring the material under the vacuum of the electron microscope. Where effects such as distortion from surface tension and chemical changes from changing ionic concentration are severe, observations must be supplemented with other methods and careful evaluation made of the extent to which preparation artifacts may have changed a living, dynamic system.

The specimen will also be altered when exposed to the electron beam. Careful procedures will eliminate changes in the micromorphology of the material. Irradiation of a specimen with a low-intensity beam will cause chemical changes in the material, so that, for example, cellulose nitrate film becomes insoluble in all solvents that normally dissolve it. Also, the polymerized butyl methacrylate matrix used to embed tissue sublimates out when thin sections are exposed to low-intensity beams, leaving the cytological material and enhancing the contrast observed in the electron micrograph. On the other hand, exposure to high-intensity electron beams can produce a rise in temperature of hundreds of degrees Celsius in the specimen material. If a thin section with a polymeric embedding medium is subjected directly to a high-intensity beam, the polymer will melt, causing considerable distortion in the section. Thermal effects can cause considerable change in some materials, and care must be exercised in avoiding irradiation with too-high-intensity beams. Examination at high electron optical magnifications ($\times 1,000,000$) requires very intense beams, and special illuminating systems such as those employing a double condenser are utilized to minimize heating effects.

Contamination of the specimen also occurs under exposure to the electron beam. A layer of low-density, inert, carbonaceous material is laid down, sometimes at rates as high as several angstroms per second. The rate of deposition is a function of beam intensity, specimen temperature, and the amount of hydrocarbons present from pump oils and gasket grease. By surrounding the specimen with a chamber kept at -110°F (-80°C) the deposition can be reduced to a minimum. Special chambers designed to achieve this effect are essential for all high-resolution work where detail and contrast would soon be obliterated by deposition of a contamination layer.

Instrumental defects. High-resolution electron microscopes can be adjusted to attain an instrumental resolving power of better than 0.5 nm. Resolution obtained is usually limited by the resolution inherent in the specimen, especially for biological materials. Achievement of such a high level of instrumental performance requires a high degree of skill in the

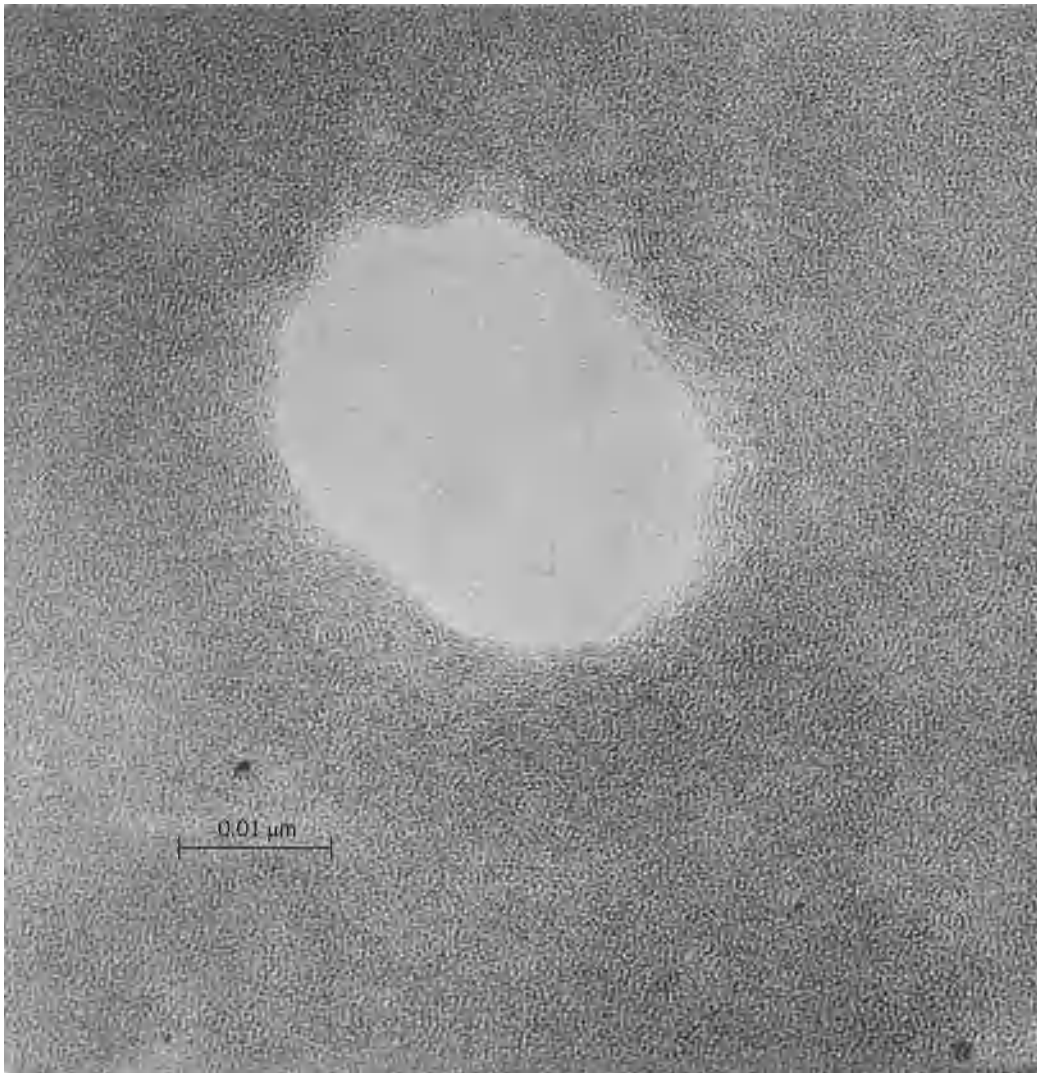


Fig. 20. Evaporated carbon film with small hole. Very slight underfocus shows Fresnel fringe inside hole. Phase contrast produces structural detail in film.

adjustment by the electron microscopist. The effect which first limits the resolving power even in a well-adjusted instrument is a deviation from axial symmetry in the objective lens producing an astigmatic image. Correction is made by introducing a compensating cylindrical field at right angles to the existing one, either by shims of ferromagnetic material oriented at the correct azimuth in the gap of the pole piece or by an electrostatic or magnetic field of correct orientation and magnitude. Externally adjustable stigmators, of both magnetic and electrostatic types, are available which enable the operator to compensate the electron microscope and obtain resolutions of 0.5 nm or better.

Contrast in the electron image can be enhanced by using a stop in the objective lens limiting the angular apertures that are necessary when relatively thick specimens such as sections and replicas are examined. However, an objective aperture can cause deterioration of resolving power, especially if asymmetric contamination forms on the aperture and causes increased astigmatism in the image. Good practice

requires careful checking on instrumental performance using test specimens of a noncrystalline material such as carbon film with small holes to check asymmetry or astigmatism and the phase contrast in the out-of-focus image.

Image characteristics. One of the most prominent characteristics of the electron microscope is the contour effect. An instrument with a good, small effective source of illumination produces an image where Fresnel diffraction is clearly observable when the objective lens is out of true focus. In practice only the first-order fringe is strongly in evidence outlining any sharp boundaries in the image; hence the so-called contour effect. When the objective-lens current is too low, so that the focal length of the lens is too long, the image is said to be underfocused. At an edge where there is a sharp change in contrast, the outside of the edge is bordered by a dark contour. Where focal length is too short, the image is said to be overfocused and the contour now appears inside the edge. At true focus the Fresnel fringe vanishes and a soft image is obtained. Most electron

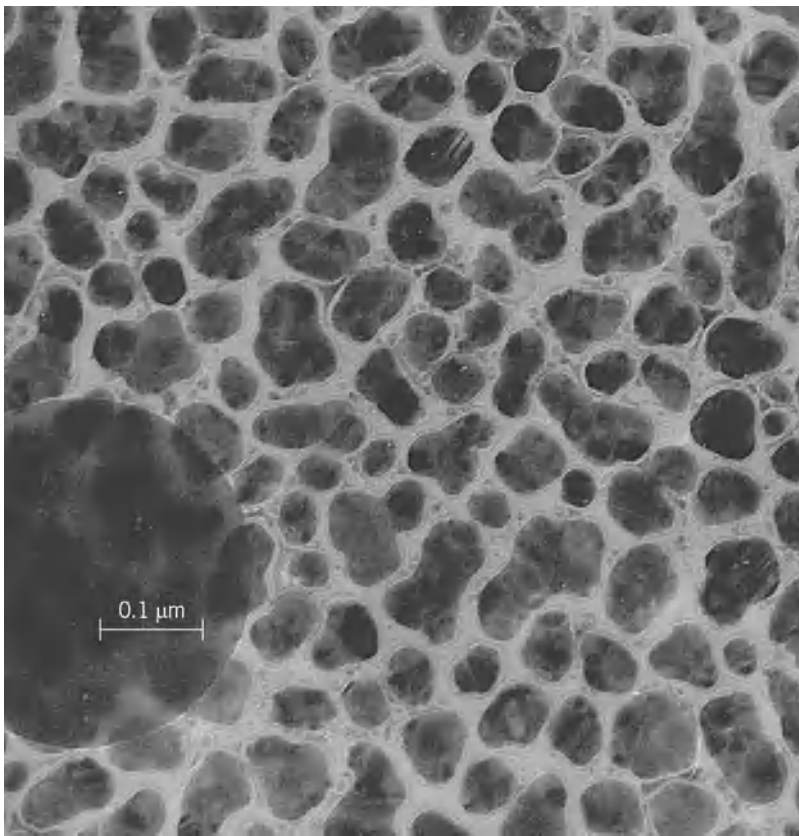


Fig. 21. Transmission electron micrograph of evaporated copper deposited on mica and transferred to a carbon substrate. Contrast effects are caused by different orientations and lattice imperfections.

micrographs are taken at a slightly underfocus setting of the objective lens, but for critical high-resolution work a focal series should be taken. Again, a carbon film with holes is a good test specimen (Fig. 20).

Crystalline material can produce anomalous contrast effects. If a crystal is oriented so that the incident electron beam is at an angle satisfying the Bragg condition for diffraction, the crystal can scatter the total incident beam into the Bragg angle while the same crystal when not oriented at the Bragg angle will scatter a fraction of the incident beam, depending on its thickness and its density.

The images obtained under the two different conditions of orientation can be of very different photographic density, and care must be taken in interpreting contrast effects of images of crystalline materials. **Figure 21** shows the strong contrast effects which can be produced by the electron interference interactions when scattered from a crystal lattice. The particles are of approximately uniform thickness, and the observed contrast effects are caused by different orientations and lattice imperfections.

The extension of the working resolving power of the electron microscope shows the significant role that diffraction phenomena play in the imaging characteristics of the instrument. Proper understanding of image contrast can be achieved only with a wave optical treatment of the imaging process. At levels of resolution of less than 0.5 nm, phase contrast plays a most important role. Phase contrast occurs with a

defocusing of the image to adjust for the phase shift of the scattered electron wave by the spherical aberration of the objective lens. Phase contrast produces the contrast observed in the image of the carbon film shown in Fig. 19. Crystalline structures such as thin graphite layers show strong phase-contrast effects. The lattice resolution method is essentially a phase-contrast image. Biological macromolecular detail as small as 0.4–0.5 nm can be resolved by phase-contrast imaging when the specimen is suspended on the edge of a fiber support.

Benjamin M. Siegel
Bibliography. P. B. Bell, Jr., *Scanning Electron Microscopy of Cells in Culture*, 1984; H. Bethge and J. Heydenreich (eds.), *Electron Microscopy in Solid State Physics*, 1987; R. Dickersin, *Diagnostic Electron Microscopy: A Text-Atlas*, 2d ed., 2000; V. A. Drits, *Electron Diffraction and High-Resolution Electron Microscopy of Mineral Structures*, 1987; H. Fujita (ed.), *History of Electron Microscopy*, 1987; P. Goodhew and F. Humphreys (eds.), *Electron Microscopy and Analysis*, 3d ed., 2000.

Electron motion in vacuum

Motion of electrons in a space freed sufficiently from matter so that collisions with other particles play a negligible role. The motion of electrons in vacuum is controlled by electric and magnetic fields whose force on the electrons is proportional to their magnitude. Electric and magnetic fields may arise from the presence of electrodes, currents, and magnets surrounding the evacuated space in which a particular electron moves, as well as from the presence of other charged particles within this space. This article deals with the nonrelativistic motion of electrons in static electric and magnetic fields, the effect of space charge on the electron paths, and motion in the time-varying fields that are encountered, for example, in the cathode-ray oscilloscope. For a discussion of the motion of electrons and other charged particles in cases where the velocities approach that of light see RELATIVISTIC ELECTRODYNAMICS

Static electric fields. An electron moving in a plane of symmetry of an electric field will remain indefinitely in that plane because the electrical forces acting on the electron lie within the plane. A plane of symmetry is here defined as one for which the mirror image of the potential distribution in front of the plane coincides with the potential distribution in back of the plane. Newton's second law of motion for the electron moving in the plane with rectangular coordinates x, y takes the form of Eqs. (1) and (2). Here $-e$ and m are the charge and mass of the

$$m \frac{d^2x}{dt^2} = -eE_x = e \frac{\partial \phi}{\partial x} \quad (1)$$

$$m \frac{d^2y}{dt^2} = -eE_y = e \frac{\partial \phi}{\partial y} \quad (2)$$

electron, t is time, E_x and E_y are the x and y components of the electric field, and ϕ is the electric

potential, normalized as stated by Eq. (3). Newton's

$$\frac{m}{2} \left[\left(\frac{dx}{dt} \right)^2 + \left(\frac{dy}{dt} \right)^2 \right] = e\phi \quad (3)$$

law as stated in Eqs. (1) and (2) implies that the speed of the electron is small enough so that its mass can be regarded as constant. For an electron having an energy of 10 kilovolts, the increase in mass is about 1%.

Path equation. Elimination of time from Eqs. (1) and (2) leads to the path equation (4). If potential distri-

$$\frac{d^2y}{dx^2} = \frac{1}{2\phi} \left[1 + \left(\frac{dy}{dx} \right)^2 \right] \left(\frac{\partial\phi}{\partial y} - \frac{dy}{dx} \frac{\partial\phi}{\partial x} \right) \quad (4)$$

bution $\phi(x,y)$ is known, and if position and velocity of the electron at one point within the field are also known, the electron's path can be determined by integrating Eq. (4). For simple electrode structures, the potential distribution can be determined analytically by solving Laplace's equation in the form of Eq. (5). More generally, it can be found by construct-

$$\nabla^2\phi = \frac{\partial^2\phi}{\partial x^2} + \frac{\partial^2\phi}{\partial y^2} + \frac{\partial^2\phi}{\partial z^2} = 0 \quad (5)$$

ing a large-scale model of the electrode structure and immersing it in an electrolytic tank so that the surface of the liquid (usually slightly acidified tap water) coincides with the plane of symmetry of interest. With potentials proportional to the actual potentials applied to the model electrodes, an equipotential line on the surface can be found by determining the points at which a probe at the potential in question draws no current.

The path equation (4) leads to Eq. (6) for the ra-

$$R = \frac{2\phi}{\partial\phi/\partial n} \quad (6)$$

dius of curvature R of the paths. Here $-\partial\phi/\partial n$ is the component of the electric field normal to the electron path. If an equipotential plot has been prepared, this relation permits graphical plotting of an electron path (Fig. 1). The path is approximated by a series of circular arcs, the radius of curvature between successive equipotential lines being computed from the preceding relation for R .

Paraxial-ray equation. Electrostatic fields having not only a plane of symmetry but an axis of symmetry, which represents the intersection of an infinite family of planes of symmetry, have particular practical importance. Equation (4) still applies here, provided that y is identified with r , the distance from the axis, and x with z , the distance measured along the axis. The Laplace equation in the new coordinates z, r takes the form of Eq. (7). Equation (7) is solved quite

$$\frac{\partial^2\phi}{\partial r^2} + \frac{1}{r} \frac{\partial\phi}{\partial r} + \frac{\partial^2\phi}{\partial z^2} = 0 \quad (7)$$

generally by the series shown as Eq. (8), where $\Phi(z)$

$$\phi(z, r) = \Phi(z) - \frac{r^2}{4} \frac{d^2\Phi}{dz^2} + \frac{r^4}{64} \frac{d^4\Phi}{dz^4} \dots \quad (8)$$

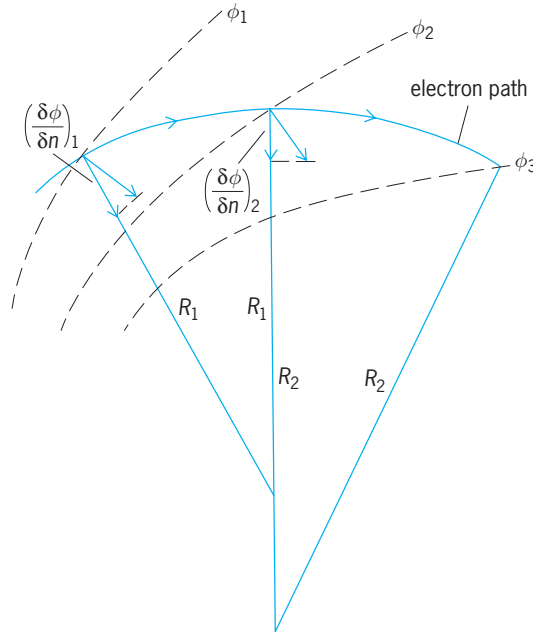


Fig. 1. Path plotting in an electrostatic field. Broken lines $\phi_1, \phi_2,$ and ϕ_3 are equipotential lines. R_1 and R_2 are radii or curvature of the electron path.

is the potential on the axis of symmetry. Thus, the potential everywhere within the axially symmetric electrode structure is fully determined by the potential variation along the axis. Substitution of ϕ from Eq. (8) into Eq. (4) with retention of terms of the first order only in r and dr/dt leads to the paraxial ray equation (9). Equation (9) applies to electrons

$$\frac{d^2r}{dz^2} + \frac{1}{2} \frac{d\Phi}{dz} \frac{dr}{dz} + \frac{1}{4} \frac{d^2\Phi}{dz^2} r = 0 \quad (9)$$

whose paths depart relatively little, both in slope and in distance, from the axis of the field.

Equation (9) is linear in r . Thus, if one electron path intersects the axis in two points, all electron paths passing through one of the points also pass through the other. In brief, the electric field images the one point into the other. It can be shown further that this imaging property is not limited to the axis but applies to extended areas about the axis, so that the axially symmetric electric field acts on the paths of electrons in the same manner as glass lenses act on light rays. Departures of the exact path from the paraxial equation result in image defects or

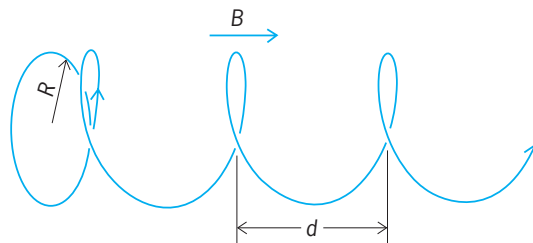


Fig. 2. Motion of an electron in a uniform magnetic field. Its path is in general a helix with pitch d , radius R , and axis parallel to the field.

aberrations similar in character to those observed for glass lenses. See ELECTROSTATIC LENS.

Magnetic fields. A magnetic field exerts on an electron of velocity \mathbf{v} a force \mathbf{F} perpendicular to both the direction of motion and the direction of the field. In vector notation this Lorentz force is given by Eq. (10). Here, \mathbf{b} is the magnetic induction.

$$\mathbf{F} = -e(\mathbf{v} \times \mathbf{b}) \quad (10)$$

The components of the Lorentz force are given by Eqs. (11)-(13). Because this force is perpendicular

$$F_x = e \left(b_y \frac{dz}{dt} - b_z \frac{dy}{dt} \right) \quad (11)$$

$$F_y = e \left(b_z \frac{dx}{dt} - b_x \frac{dz}{dt} \right) \quad (12)$$

$$F_z = e \left(b_x \frac{dy}{dt} - b_y \frac{dx}{dt} \right) \quad (13)$$

to the direction of motion, it does no work on the electron, whose velocity consequently remains unchanged in magnitude. A uniform magnetic field parallel to the z axis is described by Eqs. (14). Newton's

$$b_z = B \quad b_x = b_y = 0 \quad (14)$$

second law leads to a constant z component of the velocity. The magnitude of the velocity component v_{xy} in the xy plane is similarly constant, the square of the total velocity being equal to the sum of the squares of the components. For the motion projected on the xy plane, Newton's second law thus takes the form of Eq. (15). Here R is the radius of curvature of the

$$\frac{mv_{xy}^2}{R} = ev_{xy}B \quad (15)$$

projected path. R is seen to be a constant, so that

the projected path is a circle with radius given by Eq. (16). Here α is the angle which the electron path

$$R = \frac{mv_{xy}}{eB} = \frac{\sin \alpha}{B} \left(\frac{2m\phi}{e} \right)^{1/2} = \frac{3.37\phi^{1/2}}{B} \sin \alpha \quad (16)$$

makes with the field direction, and ϕ is the accelerating potential of the electrons. If B is measured in gauss, and ϕ in volts, R is in centimeters. The frequency with which the circle is traversed by the electron is given by Eq. (17). This frequency, in s^{-1} ,

$$f = \frac{v_{xy}}{2\pi R} = \frac{eB}{2\pi m} = 2.8 \times 10^6 B \quad (17)$$

the cyclotron frequency, thus depends only on the magnetic field strength.

The complete motion of the electron (Fig. 2) is thus a helix about a magnetic line of force, with a pitch d in centimeters given by Eq. (18). All elec-

$$d = \frac{v_z}{f} = \frac{2\pi \cos \alpha}{B} \left(\frac{2m\phi}{e} \right)^{1/2} = 21.08 \frac{\phi^{1/2}}{B} \cos \alpha \quad (18)$$

trons passing through a point with equal axial velocity components pass through a series of points separated by d on the same magnetic field line. An initially divergent electron beam is held together by a uniform magnetic field, because an electron path which intersects a particular field line can never depart from it by more than twice the radius R of the helix. Uniform magnetic fields are widely used for keeping electron beams from spreading, for example, in traveling-wave tubes and klystrons. Electron beams will also follow magnetic field lines, if these are gently curved. This property is utilized in the magnetic deflection of beams in certain television camera tubes, such as the image orthicon and the vidicon. In these tubes, a weak transverse

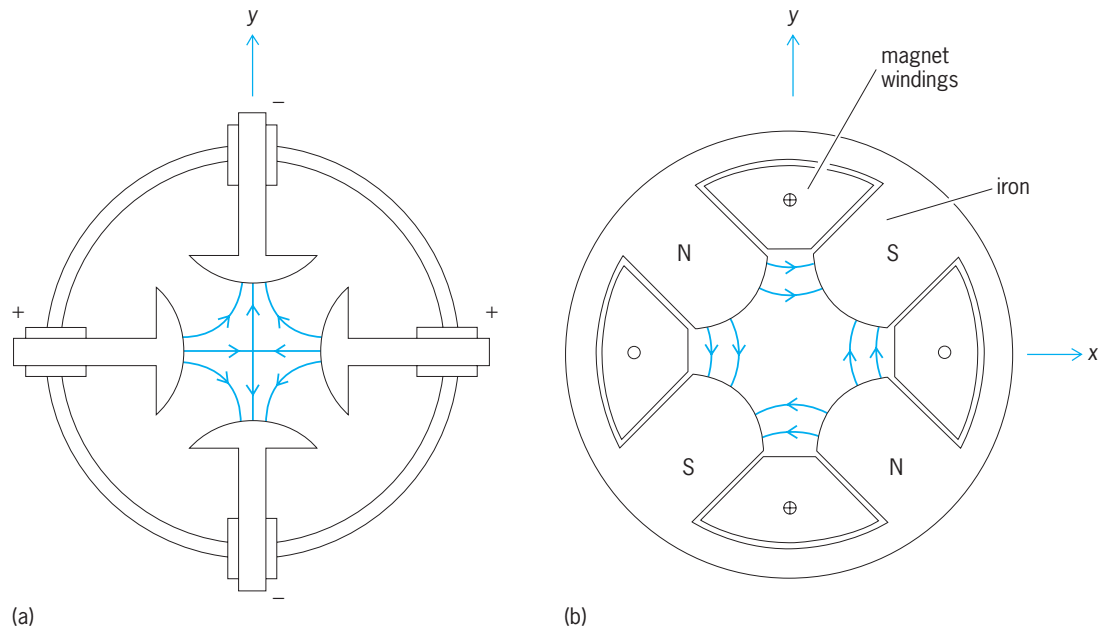


Fig. 3. Two types of quadrupole lens. (a) Electrostatic quadrupole. (b) Magnetic quadrupole.

magnetic deflection field is superposed on a strong longitudinal magnetic focusing field. See TELEVISION CAMERA TUBE; VAN ALLEN RADIATION.

Motion in nonuniform magnetic fields with axial symmetry is conveniently treated as a special case of motion in combined electric and magnetic fields.

Combined fields with axial symmetry. Motion is now expressed by Eqs. (19)–(21). Coordinates z , r , and θ

$$m \frac{d^2 z}{dt^2} = e \left(\frac{\partial \phi}{\partial z} + b_z r \frac{d\theta}{dt} \right) \quad (19)$$

$$m \left[\frac{d^2 r}{dt^2} - r \left(\frac{d\theta}{dt} \right)^2 \right] = e \left(\frac{\partial \phi}{\partial r} - b_z r \frac{d\theta}{dt} \right) \quad (20)$$

$$m \frac{1}{r} \frac{d}{dt} \left(r^2 \frac{d\theta}{dt} \right) = e \left(b_z \frac{dr}{dt} - b_r \frac{dz}{dt} \right) \quad (21)$$

represent distance along the axis, perpendicular distance from the axis, and azimuthal angle about the axis. Terms b_z and b_r are the axial and radial components of the magnetic induction. From Eqs. (19)–(21), Eq. (22a), a path equation expressing the varia-

$$\frac{d^2 r}{dz^2} = \frac{1}{2\phi^*} \left[1 + \left(\frac{dr}{dz} \right)^2 \right] \left(\frac{\partial \phi^*}{\partial r} - \frac{dr}{dz} \frac{\partial \phi^*}{\partial z} \right) \quad (22a)$$

$$\phi^* = \phi (1 - D^2)$$

$$= \phi - \left[\frac{C}{r} + \left(\frac{e}{2m} \right)^{1/2} A \right]^2 \quad (22b)$$

$$C = r^2 \frac{d\theta}{dz} \phi^{1/2} \left[\left(\frac{dr}{dz} \right)^2 + r^2 \left(\frac{d\theta}{dz} \right)^2 + 1 \right]^{-1/2} - \left(\frac{e}{2M} \right)^{1/2} r A \quad (22c)$$

tion of the radial distance r with the axial distance z , is derived, with ϕ^* defined by Eq. (22b). Here ϕD^2 is a shorthand symbol for the last term in Eq. (22b) and C is, except for a universal multiplying constant, the angular momentum of the electron about the axis at a point where the magnetic field vanishes, given by Eq. (22c). Here A is the magnetic vector potential, which is numerically equal to the magnetic flux through a circle about the axis through the reference point divided by the circumference of that circle.

At the same time, the azimuth θ of the electron changes according to Eq. (23). Equations (22) can

$$\theta = \phi_0 + \int_{z_0}^z \frac{D}{r(1-D^2)^{1/2}} \times \left[1 + \left(\frac{dr}{dz} \right)^2 \right]^{1/2} dz \quad (23)$$

be solved by graphical and numerical methods useful for determining electron paths in electrostatic fields. The general paraxial equation is obtained by substituting the expansion of Eq. (24) into Eq. (22b).

$$A = \frac{r}{2} B(z) - \frac{r^3}{16} \frac{d^2 B(z)}{dz^2} \dots \quad (24)$$

Here $B(z)$ is the magnetic induction along the axis. Substitution of Eq. (24) and that for the electrostatic potential and retention of terms of the first order in r and dr/dz only lead to Eqs. (25). With B in

$$\frac{d^2 r}{dz^2} + \frac{1}{2\Phi} \frac{d\Phi}{dz} \frac{dr}{dz} + \left(\frac{1}{4\Phi} \frac{d^2 \Phi}{dz^2} + \frac{eB^2}{8m\Phi} - \frac{C^2}{\Phi r^4} \right) r = 0 \quad (25a)$$

$$\theta = \theta_0 + \int_{z_0}^z \left[\frac{C}{r^2 \Phi^{1/2}} + \left(\frac{e}{8m\Phi} \right)^{1/2} B \right] dz \quad (25b)$$

$$\frac{C}{r_0^2 \Phi_0^{1/2}} = \left(\frac{d\theta}{dz} \right)_0 - \left(\frac{e}{8m\Phi_0} \right)^{1/2} B(z_0) \quad (25c)$$

gauss, ϕ in volts, and z in centimeters, $e/8m$ equals $0.022 \text{ volt}/(\text{gauss}\cdot\text{cm})^2$.

Quadrupole fields. The principal use of axially symmetric electrostatic and magnetic fields is to converge electron pencils; thus, in the cathode-ray tube, electrons diverging from a point in front of the emitter, called the crossover, are converged to a small spot on the viewing screen. Similarly, in an electron-imaging device such as the electron microscope, electrons diverging from a point on an object are brought to focus at a corresponding point of the

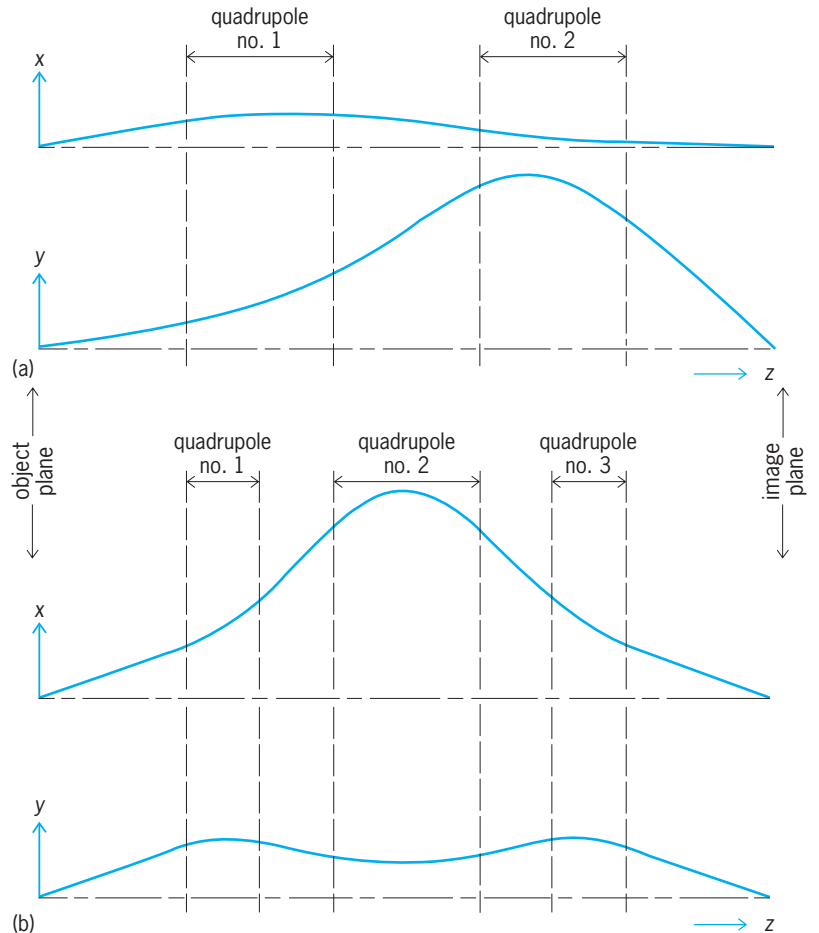


Fig. 4. Diagrammatic representation of electron paths in principal planes of image-forming symmetrical (a) doublet and (b) triplet.

image. However, this converging effect is only secondary; the primary effect of an axially symmetric electrostatic field is to accelerate or decelerate the electrons parallel to the axis, and the primary effect of an axially symmetric magnetic field is to give electrons diverging from a point on the axis a rotation about the axis. Correspondingly, the converging action or refractive power of a conventional electrostatic lens is proportional not to the first power but to the square of the ratio of the electrostatic field to the accelerating voltage, whereas the converging action of a conventional magnetic lens is proportional to the ratio of the square of the magnetic field to the accelerating potential at low electron energies (<0.5 MeV) and to the square of the ratio of the field to the accelerating potential at high electron energies (>0.5 MeV). Because breakdown phenomena place practical limits on the magnitude of electrostatic fields, and because pole-piece saturation limits (static) magnetic fields, conventional electron lenses are relatively ineffective in converging high-energy particles; they can be properly characterized as weak-focusing lenses.

Quadrupole fields, on the other hand, constitute strong-focusing lenses, insofar as their converging action is directly proportional to the electrostatic or magnetic field respectively. Consequently they (particularly magnetic quadrupole fields) have assumed great importance in the development of high-energy particle accelerators. In this application they permit the narrow confinement of the particle beam and consequently greatly reduce the cost of construction. Apart from the concentration and focus-

ing of high-energy particle beams, strong-focusing lenses have found application as projector lenses (but not as objectives) of electron microscopes. Weak quadrupole fields (stigmators) also are used to compensate residual asymmetries in conventional electron lenses. See ELECTRON MICROSCOPE; PARTICLE ACCELERATOR.

Figure 3 shows sections normal to the lens axis of an electrostatic and a magnetic quadrupole lens. In the first instance the electrostatic potential is symmetric with respect to the two principal (xz and yz) planes and antisymmetric with respect to the planes through the axis forming an angle of 45° with respect to the principal planes. In the second instance the magnetostatic potential is antisymmetric with respect to the principal planes and symmetric with respect to the two 45° planes. Fulfillment of Laplace's equation (5) now demands the form of Eq. (26) for

$$\phi(x, y, z) = \frac{a}{2}(x^2 - y^2) - \frac{d^2 a}{dz^2} 2(x^4 - y^4) \dots \quad (26)$$

the electrostatic potential. Magnetic field B , on the other hand, may be written as Eqs. (27)-(29). In Eqs. (26)-(29), a and A are functions of z . Neglecting the

$$B_x = Ay - \frac{1}{6} \frac{d^2 A}{dz^2} y(3x^2 + y^2) \dots \quad (27)$$

$$B_y = Ax - \frac{1}{6} \frac{d^2 A}{dz^2} x(x^2 + 3y^2) \dots \quad (28)$$

$$B_z = 0 - \frac{1}{6} \frac{d^3 A}{dz^3} xy(x^2 + y^2) \dots \quad (29)$$

higher-order terms, which become negligible for a long lens, one finds that the force exerted by the field on the electron is given by Eqs. (30)-(32).

<i>Electrostatic lens</i>		<i>Magnetic lens</i>	
$F_x = -eax$	(30a)	$F_x = -eAv_2x$	(30b)
$F_y = eay$	(31a)	$F_y = eAv_2y$	(31b)
$F_z = 0$	(32a)	$F_z = 0$	(32b)

The force is thus, to a first approximation, entirely in a plane normal to the axis. In one principal plane it is such as to produce convergence (direction of force opposite to that of the displacement from the axis), and in the other, such as to produce divergence (direction of force equal to that of the displacement). Two successive quadrupole fields of opposite polarity (a doublet) are required to form a sharp, real image of an object on the axis. This image will have different magnifications in the two principal directions. If the image is to be undistorted, a triplet of three successive quadrupole fields is required. These electron paths are illustrated in Fig. 4.

The shaping of the poles influences only terms of the fifth and higher orders in the x and y coordinates in the field expressions. For long lenses all terms except the linear terms vanish if the poles are rectangular hyperbolic cylinders. In practice, circular-cylinder pole caps are more easily realized and give practically the same results. Like axially symmetric lenses, quadrupole lens systems exhibit aberrations, although these are more complex in character.

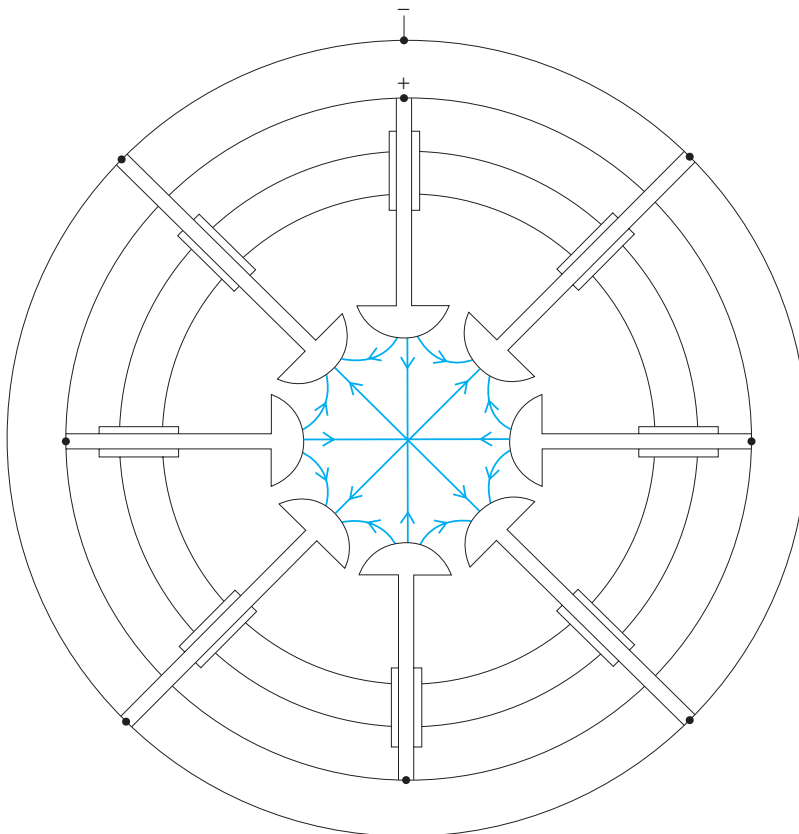


Fig. 5. Electrostatic octupole.

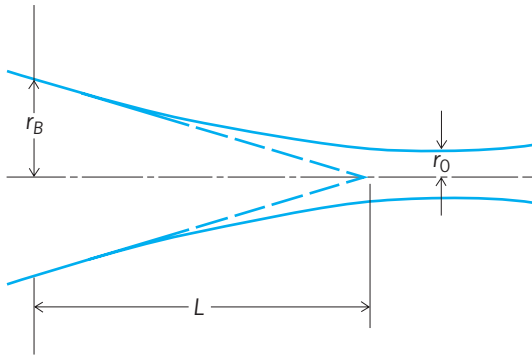


Fig. 6. Widening of electron beam by space-charge repulsion as it traverses from left to right.

Octupole fields (Fig. 5), which produce transverse forces proportional to the third power of the displacement from the axis, in the principal planes and in the 45° planes, are commonly employed to correct aberrations of quadrupole lens systems.

Effect of space charge. Space charge of either positive or negative sign can influence the paths of electrons. Space charge of positive sign is formed by electron beams passing through an imperfectly evacuated space. The beam electrons collide with gas atoms and ionize them. The heavy ions remain in the path of the electron beams for some time and prevent it from spreading. The luminous nodular or thread beams which are produced in this way are favorite objects for demonstration. See SPACE CHARGE.

Electron beams in high vacuum, on the other hand, are subject only to the mutually repulsive forces between the electrons themselves. The repulsion is reduced, but never canceled, by the action of the magnetic fields that surround charges in motion; for two electrons moving with the same velocity v parallel to each other, the ratio of the magnetic attractive force to the electrostatic repulsive force is v^2/c^2 , where c is the velocity of light. Hence, the magnetic force is significant only for electrons of very high energy.

The action of the remainder of the electrons in the beam upon any one electron can be approximated adequately by that of a continuous charge distribution equal to the average space-charge distribution. The behavior of the edge ray of a uniform circular beam of current I aimed at a point of convergence a distance L from the initial cross section of radius r_B may serve as an example (Fig. 6). If the variation of the potential along the axis of the beam is neglected (that is, if ϕ is assumed to be constant), and the charge density ρ is regarded as uniform within any beam cross section, ρ is given by Eq. (33) and the

$$\rho = \frac{I}{\pi r^2} \left(\frac{m}{2e\phi} \right)^{1/2} \quad (33)$$

path equation becomes Eq. (34). Here ϵ is the dielec-

$$\frac{d^2 r}{dz^2} = \frac{\pi \rho}{\epsilon \phi} r = \left(\frac{m}{2e} \right)^{1/2} \frac{I}{\phi^{3/2} \epsilon r} \quad (34)$$

tric constant of vacuum. As the result of the repulsive force of space charge, the ray under consideration

does not cross the axis, but reaches a minimum separation r_0 from the axis and diverges from this point on. Integration of the differential equation, Eq. (34), gives the radius as Eq. (35). For example, if $r_B =$

$$\begin{aligned} r_0 &= r_B \exp \left[-\epsilon \left(\frac{e}{2m} \right)^{1/2} \frac{r_B^2 \phi^{3/2}}{L^2 I} \right] \\ &= r_B \exp \left(-3.3 \times 10^{-5} \frac{r_B^2 \phi^{3/2}}{L^2 I} \right) \quad (35) \end{aligned}$$

1 mm, $\phi = 10,000$ volts, $L = 10$ cm, and $I = 0.001$ ampere, then $r_0 = 0.037$ mm.

Time-varying fields. In the preceding discussion, it was assumed that the electric and magnetic fields traversed by the electrons were constant in time. The total energy of the electron, or the sum of the kinetic energy and the potential energy, is then a constant. Because the potential energy is a function of position only, so is the kinetic energy. This is no longer true if the fields change appreciably in a period corresponding to the transit time of the electrons. See KLYSTRON; MAGNETRON; PARTICLE ACCELERATOR; TRAVELING-WAVE TUBE.

Beam deflection in the cathode-ray oscilloscope ceases to be proportional to the potential difference V applied to the deflection plates if V changes appreciably during the passage of the electron beam through the deflection field. If $V = V_0 \cos 2\pi ft$, an integration of the transverse impulse impressed on the electron passing between two parallel plates of length l and separation d expresses the deflection angle α in the form of Eq. (36). Here ϕ is

$$\tan \alpha = \frac{V_0 l}{2\phi d} \frac{\sin u}{u} \cos(2\pi ft) \quad (36)$$

the accelerating potential of the beam, and t is the time the electron passes through the center of the

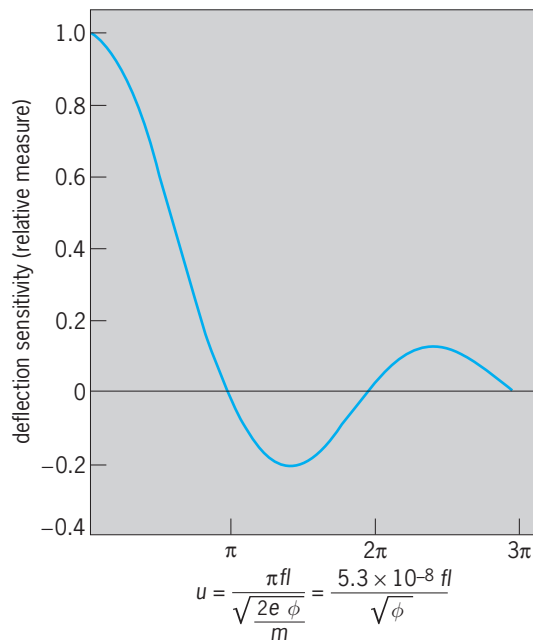


Fig. 7. Deflection sensitivity of cathode-ray oscilloscope as a function of frequency f .

deflection field. If f (the frequency of the applied voltage) is measured in s^{-1} , l in centimeters, and ϕ in volts, then Eq. (37) is applicable. The quantity

$$u = \pi f l / (2e\phi/m)^{1/2} = 5.3 \times 10^{-8} f l / \phi^{1/2} \quad (37)$$

$(\sin u)/u$ represents the ratio of the deflection sensitivity at frequency f to that at low frequencies ($f \rightarrow 0$) (Fig. 7). Thus, for a 10-kV beam and deflection plates 1 cm in length, the response is 95.4% at 1000 MHz (10^9 s^{-1}), 82.3% at 2000 MHz, 40.2% at 4000 MHz, and it drops to zero at about 5930 MHz. In this discussion, the deflection field is assumed to be a sharply cutoff uniform field, with the effects of fringe fields neglected.

Edward G. Ramberg

Bibliography. D. A. DeWolf, *Basics of Electron Optics*, 1990; P. W. Hawkes and E. Kaspar (eds.), *Principles of Electron Optics*, 1996; M. Szilagy, *Electron and Ion Optics*, 1988.

Electron paramagnetic resonance (EPR) spectroscopy

The study of the resonant response to microwave- or radio-frequency radiation of paramagnetic materials placed in a magnetic field. It is sometimes referred to as electron spin resonance (ESR). Paramagnetic substances normally have an odd number of electrons or unpaired electrons, but sometimes electron paramagnetic resonance (EPR) is observed for ions or biradicals with an even number of electrons. EPR spectra are normally presented as plots of the first derivative of the energy absorbed from an oscillating magnetic field at a fixed microwave frequency versus the magnetic field strength. The dispersion may also be detected.

To overcome the intrinsic low sensitivity of the magnetic dipole transitions responsible for EPR, samples are placed in resonant cavities. Routine experiments are carried out in the steady state at a fixed microwave frequency of approximately 9 gigahertz by slowly sweeping the magnetic field through resonance. Free electrons resonate in a magnetic field of 3250 gauss (325 millitesla) at the microwave frequency of 9.1081 GHz, whereas organic free radicals resonate at slightly different magnetic fields characteristic of each particular molecule. See ELECTRON SPIN; MAGNETIC RESONANCE; PARAMAGNETISM.

Applications. EPR spectroscopy is used to determine the electronic structure of free radicals as well as transition-metal and rare-earth ions in a variety of substances, to study interactions between molecules, and to measure nuclear spins and magnetic moments. It is applied in the fields of physics, chemistry, biology, archeology, geology, and mineralogy. It is also used in the investigation of radiation-damaged materials and in radiation dosimetry.

The basic physics of transition-metal ions and rare-earth ions present in low concentrations in diamagnetic host crystals has provided a theoretical basis for how electronic structure is modified by the surrounding atoms. Particular applications include

probing phase transitions in solids and studies of pairs and triads of magnetically interacting ions.

Applications of EPR in chemistry include characterization of free radicals, studies of organic reactions, and investigations of the electronic properties of paramagnetic inorganic molecules. Information obtained is used in the investigation of molecular structure. EPR is used widely in biology in the study of metal proteins, for nitroxide spin labeling, and in the investigation of radicals produced during reaction processes in proteins and other biomacromolecules. EPR has proved to be an important technique for interdisciplinary investigations of photosynthetic systems. By means of EPR, more than 20 proteins that function in the mitochondrial respiratory chains of mammals have been identified, and details regarding their electron transfer processes have been elucidated.

Sensitivity and resolution. EPR spectra due to pure transition-metal compounds or crystals show broad lines due to interactions between electron magnetic moments. In contrast, in nuclear magnetic resonance (NMR) the interactions between the much smaller nuclear magnetic moments are about 1 million times less intense. Small linewidths and usable resolution require use of dilute solutions in diamagnetic solvents or low concentrations of paramagnetic ions in isomorphous diamagnetic crystal lattices. Concentrations of typically less than 10^{16} molecules per cubic centimeter usually suffice for organic free radicals, but concentrations 10,000 times greater may be needed for inorganic molecules and for transition-metal and rare-earth ions in solids because of the greater linewidths, even for dilute systems. The number of spins required depends on the detectable limit of commercial spectrometers of 10 billion spins with a linewidth of 1 gauss (0.1 tesla). The larger the linewidths, the larger the number of spins required.

The observation of EPR spectra depends on spin-lattice relaxation, which is the exchange of magnetic energy with the thermal motion of the crystal or molecule. For transition-metal ions and rare-earth ions, experiments often require operation at or near liquid helium temperature (4 K; -269°C ; -452°F). Organic free radicals can usually be studied successfully at room temperature. See MAGNETIC RELAXATION.

Solids. EPR spectra from single crystals clearly provide the greatest amount of information. These include crystals containing small concentrations of paramagnetic ions substituting for the regular ions in the crystal or, for organic molecules, small fractions of free radicals produced by ionizing radiation. Spectra which may contain up to several hundred lines are often highly anisotropic; that is, they change with the orientation of the magnetic field direction in the crystal. Transition-metal-ion and rare-earth-ion EPR spectra in crystals are generally much more anisotropic than free radicals due to the intrinsic anisotropy of the electron magnetic moments, and of other effects that are important when there is more than one unpaired electron.

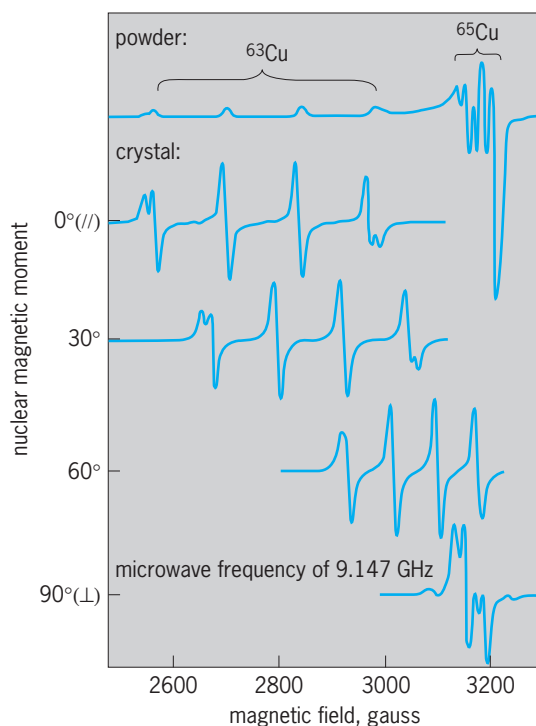


Fig. 1. Powder and single-crystal EPR spectra due to Cu^{2+} ions in calcium cadmium acetate hexahydrate at room temperature. The single-crystal spectra, with the magnetic field making angles of 0° , 30° , 60° , and 90° with the symmetry axis ($//$), illustrate anisotropy of the electron magnetic moment and the four-line copper hyperfine structure. The presence of the two copper isotopes, ^{63}Cu and ^{65}Cu , is revealed by the partly resolved features. The outer lines belong to ^{65}Cu , which has the larger nuclear magnetic moment. The lack of a characteristic four-line pattern at 90° is well understood. The powder spectra show absorptionlike features associated with 0° crystal spectra and evidence for the second copper isotope (^{65}Cu). (After J. R. Pilbrow, *Transition Ion Paramagnetic Resonance*, Oxford University Press, 1990, © J. R. Pilbrow)

The occurrence of many lines is due to interactions of the orbital motions of electrons with the electric potential of the local surrounding atoms, and to hyperfine interactions between the paramagnetic electrons and nuclear magnetic moments of the paramagnetic ion and surrounding atoms. In the case of free radicals, symmetric or nearly symmetric characteristic hyperfine patterns are observed. From knowledge of hyperfine interactions with nuclei whose spins and magnetic moments are known, the electron distribution throughout a molecule may be determined. Since hyperfine interactions vary as the reciprocal of the cube of the distance between the center of the free radical and the nucleus, structural information may be obtained in addition to electron densities.

The relative values of nuclear moments of different isotopes of a given element often can be determined from the ratio of the hyperfine splittings produced in the same chemical environment (Fig. 1).

Liquids and motional averaging. Spectra in solution due to free radicals are often quite simple as a result of motional averaging, and this clearly gives less information than would be obtained from a single-crystal investigation. Linewidths are very narrow (approx-

mately 0.1 gauss or smaller). By varying the temperature above and below room temperature, EPR spectra range from the frozen solution at low temperatures, with a powderlike spectrum, to rapid motional averaging at room temperature where anisotropies are averaged out.

The intermediate region can provide information about slow molecular motions, which is especially important for nitroxide spin labels selectively attached to different parts of macromolecules such as the components of natural and synthetic phospholipid membranes, liquid crystals, and proteins. Such measurements have revealed important structural and functional information. Since the nitrogen hyperfine interaction in nitroxide free radicals has characteristic anisotropy, molecular reorientation times in the range 10^{-9} to 10^{-5} second can be determined. A typical free-radical spectrum in solution shows hyperfine splittings that split into two lines (Fig. 2). These splittings are due to two inequivalent nitrogen-14 (^{14}N) nuclei, a single phosphorus-31 (^{31}P) nucleus, and a weak coupling to a single hydrogen (^1H) nucleus (or proton). Other commonly occurring nuclear hyperfine splittings in free-radical EPR spectra may be due to deuterium (^2H), boron-11 (^{11}B), carbon-13 (^{13}C), or nitrogen-15 (^{15}N). When there are many equivalent protons, the splitting pattern can be calculated from binomial coefficients.

Electron transfer between molecules. Studies of chemical reaction processes involving electron transfer between molecules take advantage of stable free radicals with resolved hyperfine structure, where each line is associated with a particular precession frequency of the unpaired electron spin coupled to a particular configuration of nuclear spins. If an electron jumps to a different molecule nearby with a different arrangement of nuclear spins, its electron magnetic moment will then precess at a slightly different frequency. If the jump rate is small compared with the difference in precession frequencies, the observed linewidths will increase by the jump rate (in hertz). When the jump rates are fast compared with the difference between the precession frequencies, spectra characteristic of each molecule are observed. This latter situation also applies to inorganic transition-metal ion dimer complexes in crystals and frozen solutions where averaging is not observed and the anisotropic spectra provide useful structural clues. See ELECTRON-TRANSFER REACTION.

Electron nuclear double resonance (ENDOR). This multiple-resonance method involves simultaneous irradiation at both microwave frequencies and radio frequencies, the latter usually being swept over a range covering the nuclear frequencies. It can be understood as the detection of nuclear magnetic resonance (NMR) through a change in the EPR signal as the nuclear frequencies are swept through one or more resonances. ENDOR patterns for equivalent nuclei are much simpler than in a normal EPR spectrum. For example, for ^1H , the ENDOR pattern consists of two lines for any number of equivalent nuclei, whereas the normal EPR spectrum would

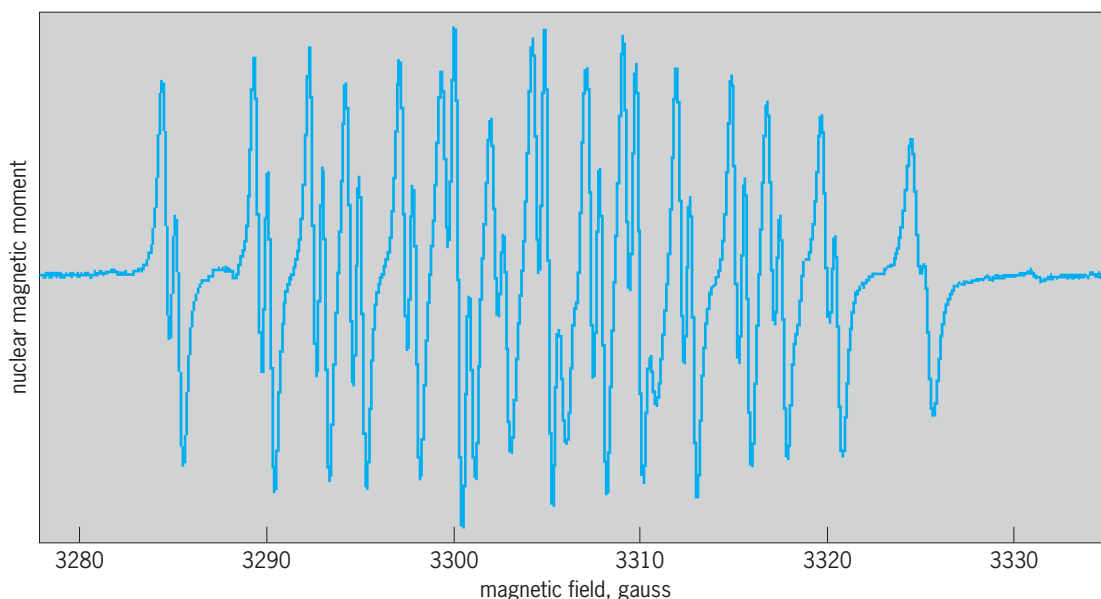


Fig. 2. Room-temperature solution spectrum of the betaine free-radical cation in toluene at a microwave frequency of 9.28877 GHz. The spectrum would be more symmetric at higher microwave frequencies. (Courtesy of Dr. G. R. Hanson)

yield multiline spectra with intensities determined by the binomial coefficients. Linewidths are typically as low as 10 kHz compared with typical values in EPR of approximately 1 MHz for free radicals and 10 MHz for transition-metal ions. ENDOR has been successfully applied to determine properties of hydrogen atoms trapped in alkali halide and calcium fluoride crystals; to locate paramagnetic defects in diamonds, semiconductors, ferroelectrics; and to study organic free radicals.

Electron-electron double resonance (ELDOR). This method involves irradiation at two different microwave frequencies. The changes in the spectrum resulting from either sweeping the second frequency or the magnetic field are monitored through detection at the first frequency. This has been important for studies of slow motions using nitroxide spin labels.

Optical detection of magnetic resonance (ODMR). This method takes advantage of the much higher sensitivity of electric dipole transitions. In ODMR, paramagnetic states are optically excited, and the EPR signal is detected through changes in the optical absorption as the magnetic field is swept through one or more resonances. It has been particularly important for studying spin-triplet states that have non-magnetic ground states and defects in solids, particularly in semiconductors. See NUCLEAR MAGNETIC RESONANCE (NMR).

Pulsed methods. It is often of interest to study the EPR spectra of many chemical or biological systems as a function of time after the microwave power is suddenly switched off. In practice, this is usually accomplished with controlled short high-power pulses of microwave radiation. Pulses as short as 2 billionths of a second have been used. The free induction decay (FID) contains all of the spectral information, and this is mathematically transformed

into the normal spectrum using the Fourier transform. The microwave circuit in EPR spectrometers places a limit on measurements, typically until at least 70 billionths of a second following the microwave pulse. While the free induction decay can be observed for many free radicals, it is normally difficult to observe for transition-metal and rare-earth ions. Use of the spin echo is resorted to in these cases.

Electron spin echoes occur whenever two or more short, intense pulses of microwave radiation are applied to a sample with inhomogeneously broadened lines due to many overlapping narrow lines resulting from resonances corresponding to different arrangements of nuclear spin splittings. In the simplest (two-pulse) echo, during the interval τ between the two pulses, the magnetic moments become dephased and lose their coherence. The second pulse is arranged so as to refocus the magnetic moments, thus producing an echo at a later time τ . This means that, provided relaxation times permit, the experimenter has considerable control as to when echoes will occur. If τ is increased in small steps in a two-pulse echo, or suitable time intervals are incremented for multiple-pulse echoes, the echo intensity as a function of time often shows one or more periodic modulations known as electron spin echo envelope modulation (ESEEM). These modulated echo patterns may be mathematically transformed to yield the spectrum of nuclear frequencies. Two-dimensional (2D) forms of spectroscopy are also possible by increasing two different pulse intervals. The resulting 2D plots enable precise correlation of spectra to structural and electronic properties of the active paramagnetic species.

ESEEM and ENDOR provide complementary information. In samples where the unpaired electron in a molecule, or from an ion in a crystal, has high

symmetry, ENDOR spectra are strongest near or along principal, or symmetry, directions; whereas ESEEM has its greatest intensities away from principal directions.

ESEEM has been used effectively to probe weak hyperfine interactions in solids, in assisting with the determination of molecular structure, and in illuminating the local environment of paramagnetic ions at active sites in proteins.

J. R. Pilbrow

Bibliography. L. Kevan and M. K. Bowman (eds.), *Pulsed and Continuous Wave ESR*, 1990; J. R. Pilbrow, *Transition Ion Electron Paramagnetic Resonance*, 1990; C. P. Poole and H. A. Farach, *Handbook of Electron Spin Resonance*, 1994; J. A. Weil, J. R. Bolton, and J. E. Wertz, *Electron Paramagnetic Resonance: Elementary Theory and Practical Applications*, 1993.

Electron-positron pair production

A process in which an electron and a positron are simultaneously created in the vicinity of a nucleus or subatomic particle. Electron-positron pair production is an example of the materialization of energy predicted by special relativity and is accurately described by quantum electrodynamics. Pair production usually refers to external pair production, in which the positron (positively charged antielectron) and electron are created from a high-energy gamma ray as it passes through matter. Electron-positron pairs are also produced from internal pair conversions in nuclei, decays of unstable subatomic particles, and collisions between charged particles. See QUANTUM ELECTRODYNAMICS.

External pair production. In external conversion, the energy of an incoming gamma ray (a high-energy electromagnetic photon) is directly converted into the mass of the electron-positron pair (Fig. 1). The photon energy $h\nu$ (where h is Planck's constant and ν is the photon frequency) must therefore exceed twice the rest mass of the electron $2m_0c^2$, equal to 1.022 MeV (m_0 is the electron mass, c the velocity of light). In order to conserve both energy and momentum in this process, the pair must be created near a nucleus, which recoils to balance the momentum of the incoming photon with the momenta of the created electron and positron. Because the nucleus is so much heavier than the elec-

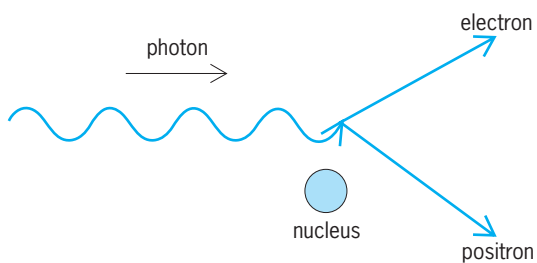


Fig. 1. External electron-positron pair production.

tron, it carries away almost no energy from the pair, and the energy of the photon in excess of $2m_0c^2$ is shared unequally as kinetic energy by the positron and electron. Individually, the electron and positron each exhibit a distribution of kinetic energies ranging from zero to the maximum available energy, $E_{\max} = h\nu - 2m_0c^2$, correlated with one another so that their sum is equal to E_{\max} . Similarly, the positron and electron are emitted over a broad range of angles, although they exhibit a tendency to move in the same direction, which reflects the momentum of the incoming photon. For incident photon energies above 5 MeV, external pair production is the dominant mechanism by which gamma rays are absorbed in matter. See GAMMA-RAY DETECTORS; GAMMA RAYS; PHOTON.

When a highly energetic charged particle collides with matter or with another particle, positron-electron pairs may be copiously produced by the external conversion of bremsstrahlung photons emitted as the incident particle is decelerated. This process is repeated as the created particles in turn produce pairs, until the average available energy is below the pair production threshold, giving rise to an electromagnetic cascade or shower. Showers are commonly observed in cosmic-ray interactions and are often used to detect high-energy electrons and photons produced by particle accelerators. See COSMIC RAYS.

Internal pair creation. Internal pair creation differs from external conversion in that the positron and electron are created directly from energy liberated by the deexcitation of an excited nucleus (produced, for example, in radioactive decay or nuclear collisions) to a state of lower energy, if the transition energy exceeds the pair mass threshold of $2m_0c^2$. Internal pair creation usually occurs only 10^{-3} times as often as deexcitation by gamma-ray emission, although the exact pair creation probability, as well as the angular correlation between the emitted positron and electron, depends on the nuclear charge and upon the energy and multipolarity of the nuclear transition. See RADIOACTIVITY.

Particle decay. Many unstable subatomic particles, such as the neutral Z boson and J/ψ meson, decay into a positron-electron pair alone or with other particles. Since the decaying parent particle is massive, momentum is conserved without the presence of an additional nucleus as is required in external conversion. Decay into a single pair alone creates a positron and electron with equal and opposite momenta. They share the available kinetic energy equally, that is, the equation below is satisfied, where

$$E_{e^+} = E_{e^-} = \frac{Mc^2 - 2m_0c^2}{2}$$

M is the mass of the parent, and they are emitted back to back in the rest frame of the parent particle. See ELEMENTARY PARTICLE; MESON.

Superheavy nuclear collisions. In the 1980s, a new mode of positron-electron pair production was discovered in collisions of very heavy nuclei at high energies (Fig. 2). Narrow positron and electron peaks

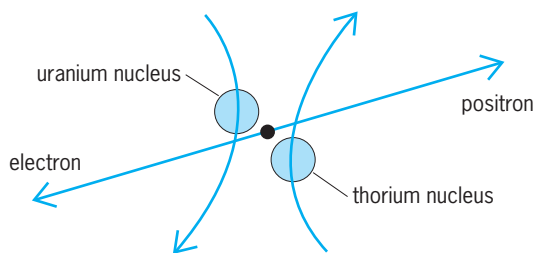


Fig. 2. Anomalous positron-electron pair creation in superheavy quasiautomic collisions.

were detected in coincidence from 1.4-GeV uranium-on-thorium quasiautomic collisions; three sets of correlated structures were found. In each case, the corresponding peak in the distribution of the sum of the positron and electron energies is narrower than the individual Doppler-broadened lines. This suggests the correlated cancellation of the positron and electron laboratory Doppler shifts, which is characteristic of an 180° back-to-back emission-angle correlation in the rest frame of their source. Explanations for these structures involving instrumental backgrounds, uncorrelated emission of separate positron and electron lines, external or internal pair production, and multibody final states are unable to describe the observed narrow width of the sum-energy peak. The data appear to reflect the production and decay of at least three previously undetected neutral objects with masses of 1.63, 1.77, and 1.84 MeV/c^2 . Attempts to explain these data in terms of one or several new elementary particles have been unsuccessful so far. See POSITRON; QUASIAUTOM.

Thomas E. Cowan

Bibliography. M. L. Burns and A. K. Harding (eds.), *Positron-Electron Pairs in Astrophysics*, 1983; R. D. Evans, *The Atomic Nucleus*, 1955, reprint 1982; W. Greiner (ed.), *Physics of Strong Fields*, 1987; K. Siegbahn (ed.), *α -, β -, and γ -Ray Spectroscopy*, 1968; J. Thompson, *Soft Lepton Pair and Photon Production*, 1991.

Electron-probe microanalysis

A method used for determining the elemental composition of materials, based on the x-rays emitted by different elements when bombarded with high-energy electrons. It is a micro method that can detect x-ray photons emitted by the atoms within a small volume excited by an electron beam focused to 10 nanometers diameter or less. In biology electron-probe microanalysis can be used to determine the composition of cell organelles without isolating them and therefore altering the distribution of diffusible elements. Studies with electron-probe analysis have thrown light on the major problems of excitation-contraction coupling in muscle, visual transduction in retinal rods, and the functions of the endoplasmic reticulum and mitochondria in the regulation of cytoplasmic calcium.

High-energy electrons can ionize atoms, ejecting an inner-shell electron. To fill the resultant vacancy,

an outer-shell electron falls into the ionized shell (Fig. 1); the atom remains in the higher-energy excited state. The emission of a characteristic x-ray by the ionized atom is one of the mechanisms for releasing its excess energy. Another mechanism is the emission of Auger electrons; the probability of ejecting an x-ray instead of an Auger electron is the fluorescence yield. In the case of electron-probe microanalysis the source of the exciting electrons is the electron gun, which, in modern electron microscopes equipped with field-emission guns, can produce a focused beam narrower than 1 nm. The x-rays emitted as the result of atomic ionization are called characteristic x-rays, because their energy is characteristic of the core shell of the ionized element and of the shell from which the electron relaxed into the vacancy. X-rays excited by ionization of the innermost (K) shells are called K x-rays, those due to L-shell ionization, L x-rays. Typical energies of characteristic K x-rays of the biologically important elements are given in the table. See AUGER EFFECT; SURFACE PHYSICS; X-RAY FLUORESCENCE ANALYSIS.

X-ray detection and data processing. At the great sensitivity required for biological applications at high spatial resolution, x-rays are detected with energy-dispersive detectors. They are solid-state devices in which incident x-rays create electron-hole pairs. In those most commonly used, lithium-drifted silicon (Si-Li) detectors, each 3.8 eV of energy discharged by an x-ray photon creates one electron-hole pair. Hence, following the arrival of a high-energy x-ray, the charge deposited and then swept off the detector is a measure of the energy of the incident x-ray. The charge after pulse and digital processing, can be stored in one of the channels of a multichannel analyzer. This procedure builds up a spectrum of the number of emitted x-rays as a function of their

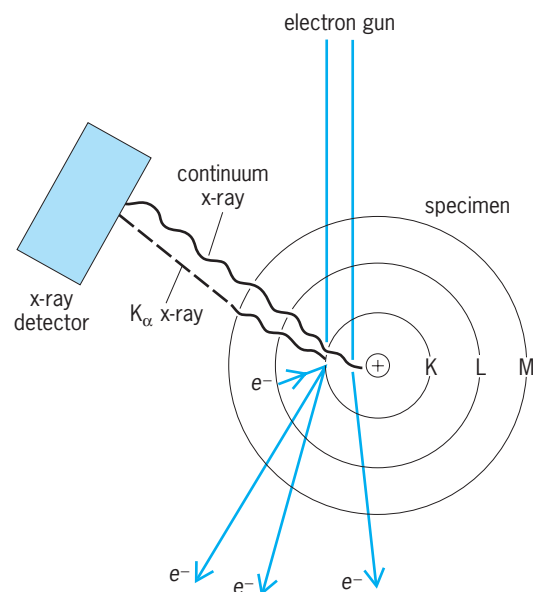


Fig. 1. Ionization of an atom in a specimen and production and detection of the characteristic and continuum x-rays generated by the incident electrons. Unscattered and elastically or inelastically scattered electrons are shown leaving the specimen plane.

Energies of K x-rays	
Element	Energy, keV
Sodium	1.041
Magnesium	1.253
Phosphorus	2.013
Sulfur	2.306
Potassium	3.310
Calcium	3.691

energies (Fig. 2). The energy resolution of modern x-ray detectors, usually specified at the $K\alpha$ manganese line, is approximately 150 eV. This represents significant broadening of the intrinsic (4 eV or less) width of the x-ray lines by electronic noise and other factors, and can cause overlap of adjacent x-ray peaks. The analysis of overlapping peaks, required for accurate determinations, is performed by computer.

Quantitative analysis. The concentrations of elements are determined from the relationship, given in the equation below, between the number of atoms

$$I_x = \omega_x \cdot q_x \cdot N_x \cdot T_x$$

(N_x) in an irradiated space and the number of characteristic x-rays (I_x) emitted. Here ω_x is the fluorescence yield, q_x the ionization cross section (the probability of an atom being ionized), N_x the number of atoms of element X in the space irradiated, I the probe current, and T_x the detector transmission efficiency. Thus, the number of x-rays generated in a given space is directly proportional to the number of atoms and the number of exciting electrons. The two major influences on the performance of the detector are the geometric efficiency with which the emitted x-rays reach it and the extent of x-ray absorp-

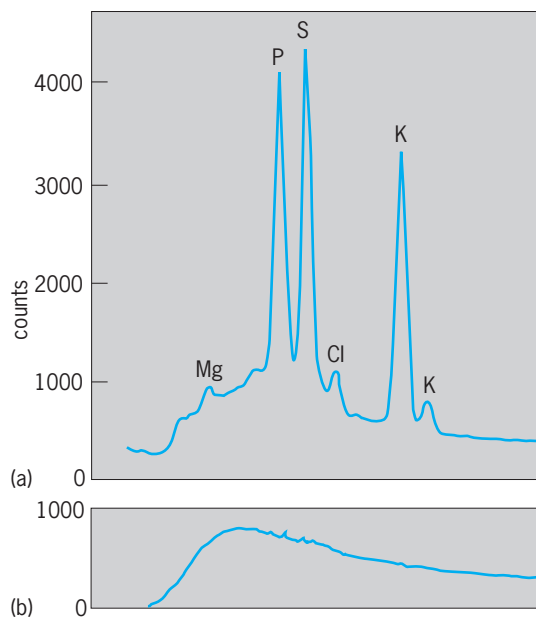


Fig. 2. Typical cellular x-ray spectrum, indicating high concentration of potassium: (a) a series of characteristic elemental peaks superimposed on (b) a background of x-ray continuum, shown with the characteristic peaks removed by computer processing.

tion by its window. The geometry of electron microscope lenses is such that less than 1% of emitted x-rays reaches an energy-dispersive detector. Low-energy x-rays are absorbed by the detector window: elements below sodium in atomic number are not detected, and the detection efficiency for sodium and magnesium is less than that for, say, potassium. Elements of low atomic number can be analyzed with windowless x-ray detectors or electron energy-loss spectroscopy. See ELECTRON MICROSCOPE; ELECTRON SPECTROSCOPY.

X-ray photons are also absorbed by the material being analyzed, complicating the analysis of thick and bulky specimens. Spatial resolution is also reduced in thick specimens through broadening of the probe by electron scattering. Biological studies at the high spatial resolution required for analyzing cell organelles are performed on sections less than 200 nm thick.

The concentration of elements is also best determined in ultrathin specimens, because excitation of an atom with electrons also creates continuum x-rays, which result from inelastic scattering by atomic nuclei. Continuum x-rays are distributed over a broad energy band (Fig. 2a); their number is proportional to the total mass within the irradiated space. The same energy-dispersive x-ray detectors simultaneously detect characteristic and continuum x-rays, and the ratio of the two is a measure of elemental concentration, usually expressed as mmol/kg, dry weight. With modern instrumentation, concentrations of calcium as low as 0.3 mmol/kg can be detected within a 100-nm³ volume or less.

X-ray mapping. X-ray mapping, a second mode of electron-probe microanalysis application, yields a two-dimensional map of elemental distribution. X-ray maps are obtained by scanning an area of a specimen with the electron beam and collecting the x-ray spectrum at each picture point (pixel) as the beam is rastered. In this manner, the concentration of different elements at each picture point can be determined, and the resultant values displayed as relative intensities on x-ray maps of, for example, phosphorus or calcium. Powerful electron sources are required with the small probes that give high spatial resolution, to obtain low-noise x-ray maps of biological specimens. Field-emission guns are the strongest available and provide the highest current density in a focused beam. The distribution of several elements in specific granules of the atrium of a rat heart, mapped with a field-emission gun, is shown in Fig. 3. See X-RAY MICROSCOPE.

Specimen preparation by rapid freezing. Extreme care must be taken to prevent the loss or translocation of diffusible elements during specimen preparation. Freezing at rates of almost 10,000°C/s (18,000°F/s) is used almost universally, followed by vacuum dehydration. Valuable information, albeit at lower resolution, has been obtained by analysis of frozen hydrated (ice) specimens of secretory epithelia.

Biological applications. Important biological structures have been analyzed with electron-probe

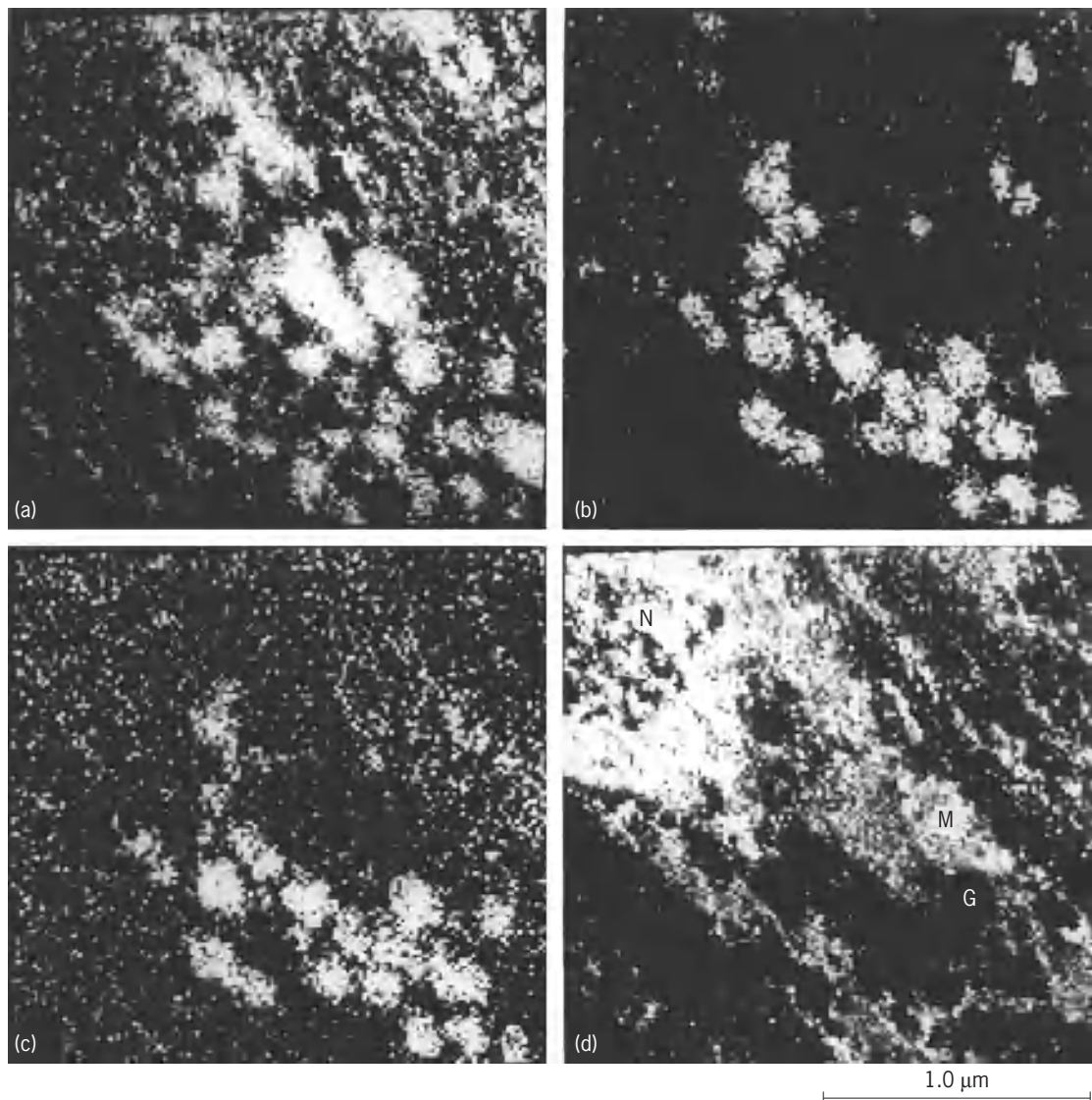


Fig. 3. X-ray maps of a rat heart cryosection, showing the distribution of (a) sulfur, (b) calcium, (c) chlorine, and (d) phosphorus in a $2\text{-}\mu\text{m}^2$ region of the nuclear pole, containing mitochondria (M), atrial granules (G), and condensed chromatin inside the nuclear (N) envelope. The intensity of gray indicates distribution; black represents the fewest and white the most numerous atoms of each element.

microanalysis, including sarcoplasmic and endoplasmic reticulum, mitochondria, and nuclei. Largely as a result, mitochondria are now known to accumulate significant amounts of calcium only under pathological conditions when cytoplasmic calcium rises to abnormally high levels; in normal cells the endoplasmic (in muscle the sarcoplasmic) reticulum, in conjunction with plasma membrane calcium pumps, is the principal regulator of cytoplasmic calcium. Electron-probe microanalysis has revealed subcellular movements of magnesium in and out of mitochondria and sarcoplasmic reticulum that were not anticipated from studies of isolated organelles. Electron-probe microanalysis of bacteria and bacterial spores has revealed the distribution of magnesium, calcium, and other elements, as well as an increase in cell calcium during bacterial division. The application of this technique to abnormal red blood cells in sickle cell anemia has shown that their high calcium content is

concentrated in small vacuoles. *See* MASS SPECTROMETRY; SECONDARY ION MASS SPECTROMETRY (SIMS); SPECTROSCOPY.

Andrew Paul Somlyo

Bibliography. A. Le Furgey, M. Bond, and P. Ingram, *Frontiers in electron probe microanalysis: Application to cell physiology*, *Ultramicroscopy*, 24:185-219, 1988; A. P. Somlyo, Cell calcium measurement with electron probe and electron energy loss analysis, *Cell Calcium*, 6:197-212, 1985; A. P. Somlyo, Compositional mapping in biology: X-rays and electrons, *J. Ultrastruc. Res.*, 83:135-142, 1985.

Electron spectroscopy

A form of spectroscopy which deals with the emission and recording of the electrons which constitute matter—solids, liquids, or gases. The usual form of spectroscopy concerns the emission or absorption

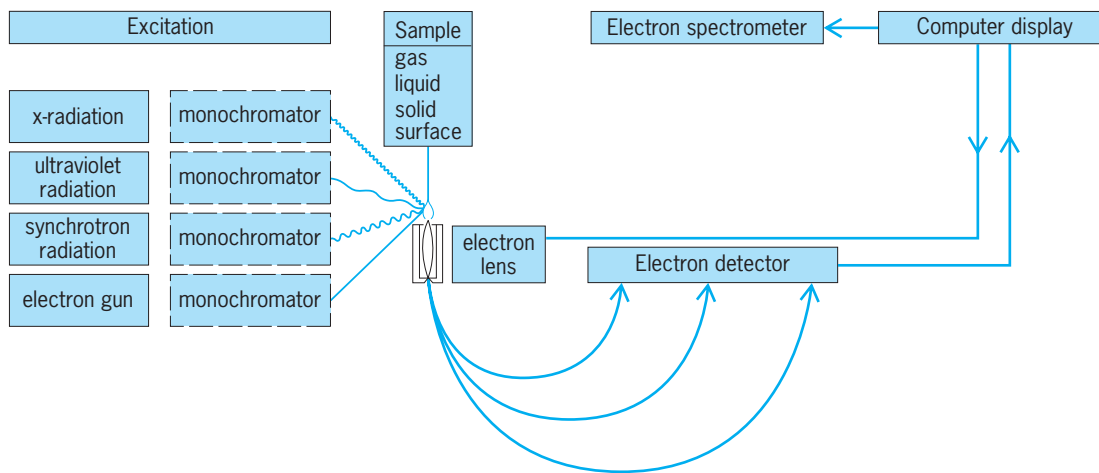


Fig. 1. Excitation of electron spectra recorded with high-resolution instruments. (After H. Siegbahn and L. Karlsson, *Photoelectron spectroscopy, Handbuch der Physik*, vol. 31, 1982)

of photons [x-rays, ultraviolet (uv) rays, visible or microwave wavelengths, and so on]. Electron spectra can be excited by x-rays, which is the basis for electron spectroscopy for chemical analysis (ESCA), or by uv photons, or by ions (electrons; Fig. 1). By means of x-ray or uv photons with energy E_{bv} , photoelectron spectra (PES) are produced when electrons with binding energies E_b are emitted with energy

E_{kinetic} from bound molecular states, according to the equation below. E_r is a usually negligible small recoil

$$E_{\text{kinetic}} = E_{bv} - E_b - E_r - \phi$$

energy, and ϕ a small work function correction which can be attributed to contact potentials. For insulating solid materials, precautions are taken for

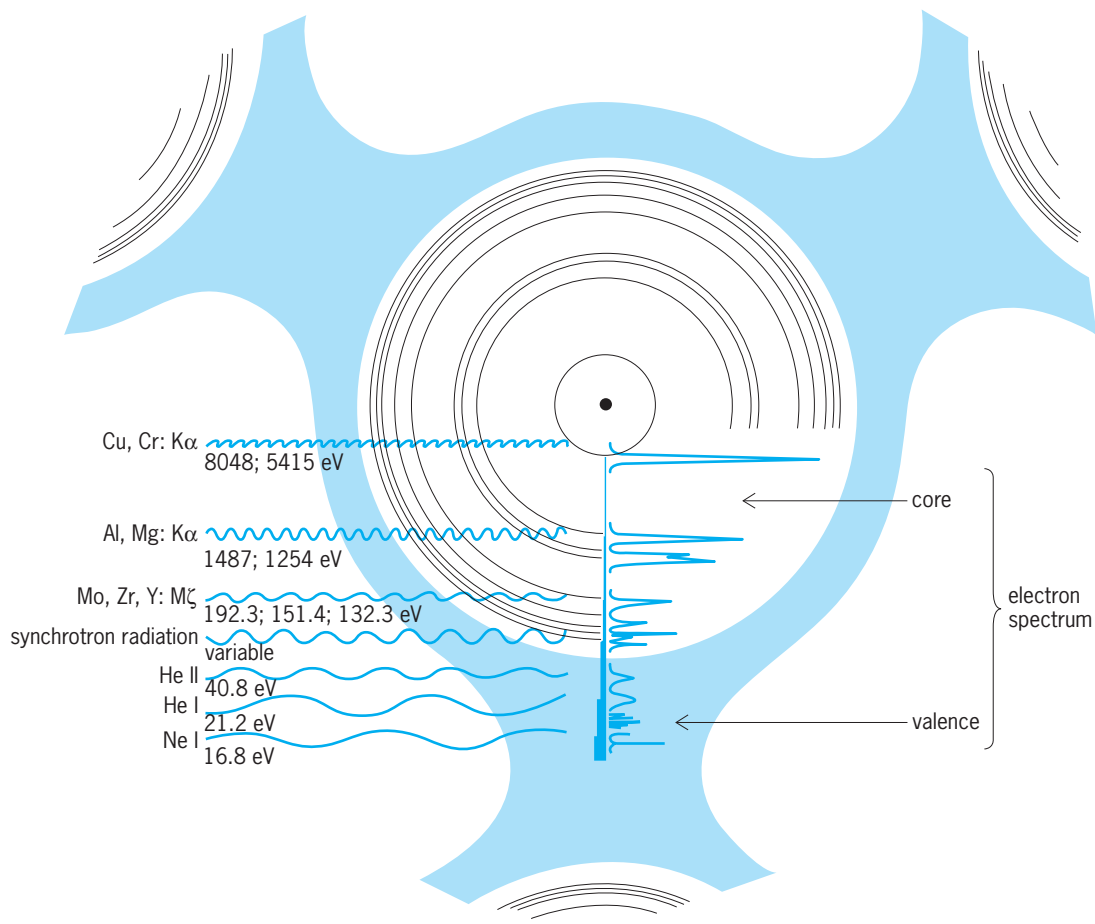


Fig. 2. The system of levels of an atom in a molecule can be divided into a valence electron region and an atomic-core region. Excitation of electron lines from the various regions can be made at different photon energies.

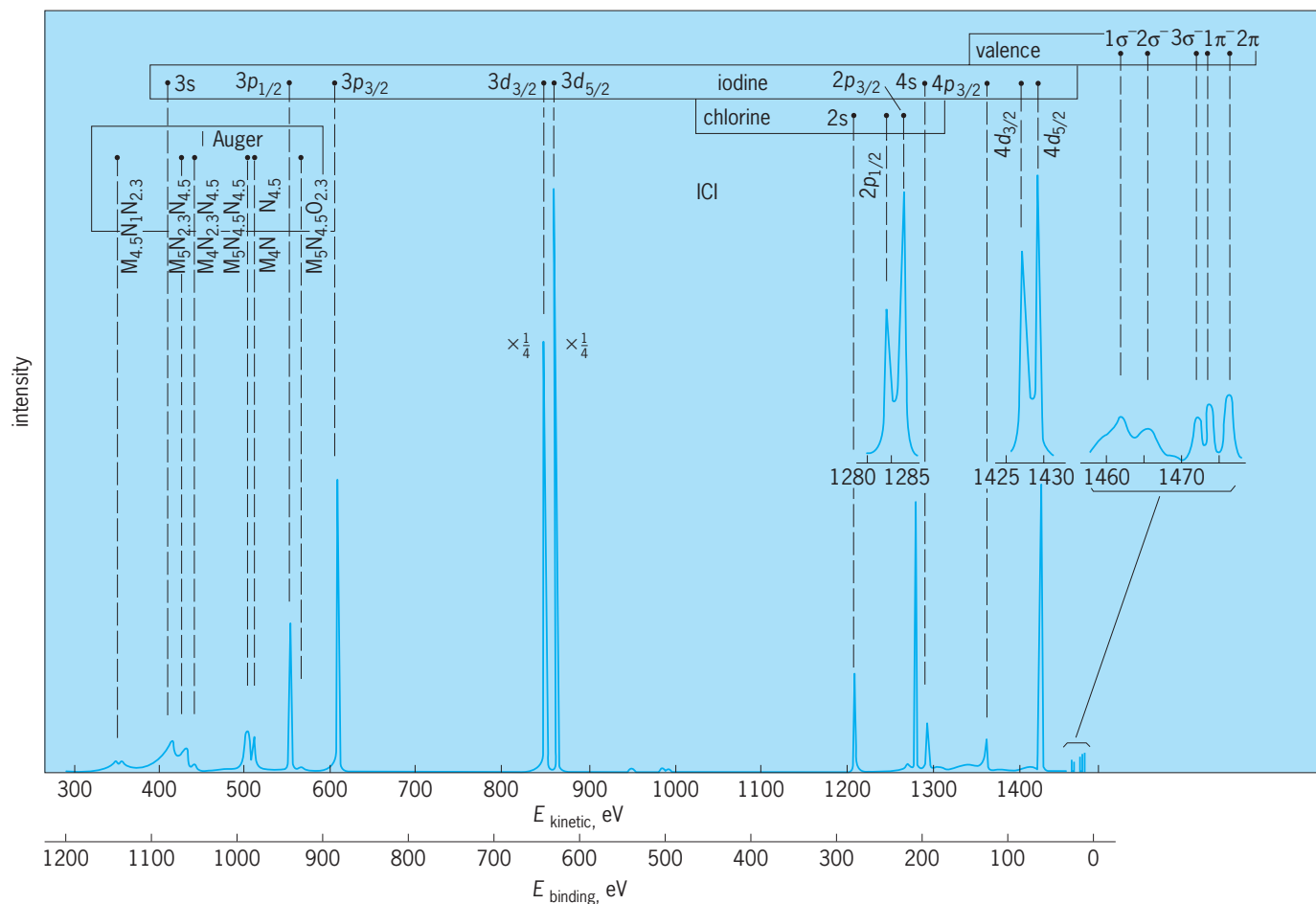


Fig. 3. Electron spectrum of the vapor of iodine chloride at a pressure of 0.04 torr (5 Pa), excited monochromatic Al $K\alpha$ radiation with the ESCA instrument (scan spectrum). The M- and N-shell levels are excited in iodine, and the L levels in chlorine. At low binding energies (high kinetic energies), the valence spectrum of the molecule is recorded. The MNN Auger electron lines of iodine can also be seen.

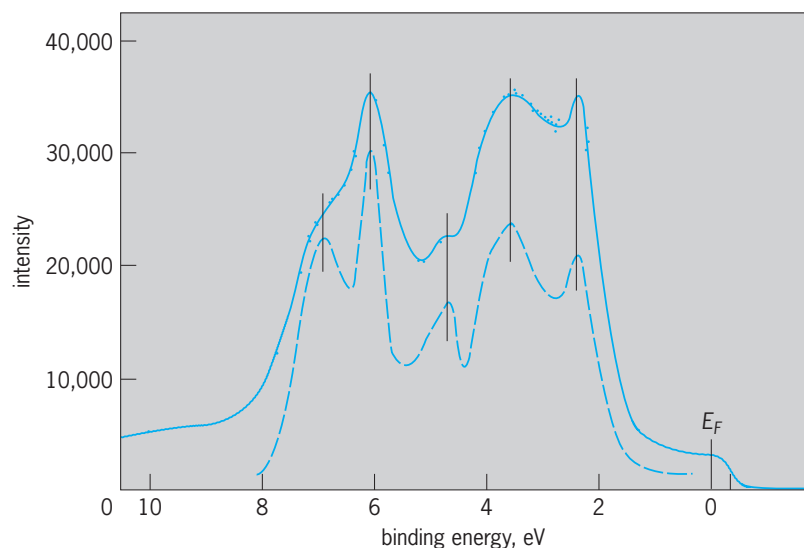


Fig. 4. Valence-band electron spectrum of gold excited by monochromatized Al $K\alpha$ radiation (solid curve). Comparison is made with a theoretically calculated density of states (broken curve) which has been folded with a 0.25-eV gaussian corresponding to the experimental resolution. The energy scale of the theoretical density-of-states curve has also been expanded by 7.5%, which leads to a near-perfect correspondence between theoretical and experimental band profiles. E_F is the Fermi edge. $h\nu = 1486.7$ eV. Foil, $\theta = 45^\circ$. (After U. Gelius et al., *Experimental spectrum*, Nucl. Instr. Meth., B1:85 1984; and N. E. Christensen and B. O. Seraphin, *Theoretical density of states*, Phys. Rev., B4:3321, 1971)

stabilizing the surface potential. For solids, the binding energies are preferably referred to the Fermi level, and for gases, to the vacuum level.

Modes of excitation. By means of ESCA, complete sets of photoelectron lines can be excited from the internal (core) levels as well as from the external (valence) region (Fig. 2). Also, complete sequences of the Auger electron lines are automatically obtained in this mode. A convenient source of excitation is the x-radiation from Al at 1487 eV. With the use of spherically bent quartz crystals this radiation can be further monochromatized to an inherent width of 0.2 eV. The commonly used light source for uv excitation of photoelectron spectra within the valence electron region is a helium lamp that produces highly monochromatic radiation at 21.2 eV (internal width less than 10 meV) and (with less intensity) at 40.8 eV. Intermediate in energy for excitation are ultrasoft x-rays (Y M ζ at 132 eV). Another source of excitation in the intermediate region is provided by synchrotron radiation, which can be continually varied by means of a suitable monochromator. Excitation by means of an electron beam is an alternative mode to obtain Auger electron lines that has the advantage of ease of production. Ordinary photoelectron

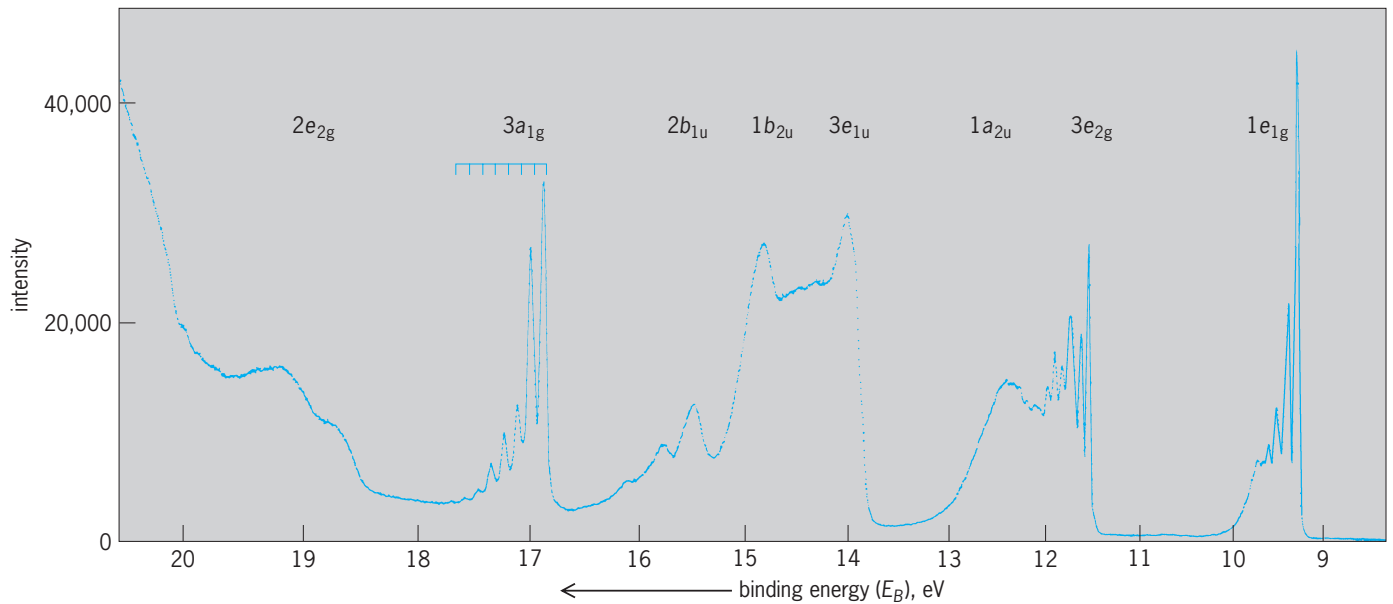


Fig. 5. Benzene valence electron spectrum excited by HeI radiation at 21.22 eV, showing vibrational structures.

lines are not produced by that radiation. In order to compensate for the low-signal-to-background ratio when Auger electron spectra are excited by an electron beam impinging on a solid surface, the spectral distribution is frequently differentiated. See AUGER EFFECT.

Applications. In ESCA, only electrons which are expelled from a surface layer of less than 50 nanometers of a solid material contribute to the electron line with the kinetic energy given above. Electrons from the interior of the material are scattered out from the line, and form a low background which does not interfere with the line character of the ESCA spectrum. The

electron lines are extremely sharp and well suited for precision measurements. With a high-resolving ESCA spectrometer which has a magnetic or electrostatic focusing dispersive system, the electron lines have widths which are set by the limit caused by the uncertainty principle (the "inherent" widths of atomic levels). With a suitable choice of radiation, electron spectroscopy reproduces directly the electronic level structure from the innermost shells (core electrons) to the atomic surface (valence or conduction band; Fig. 3). Furthermore, all elements from hydrogen to the heaviest ones can be studied even if the element occurs together with several other

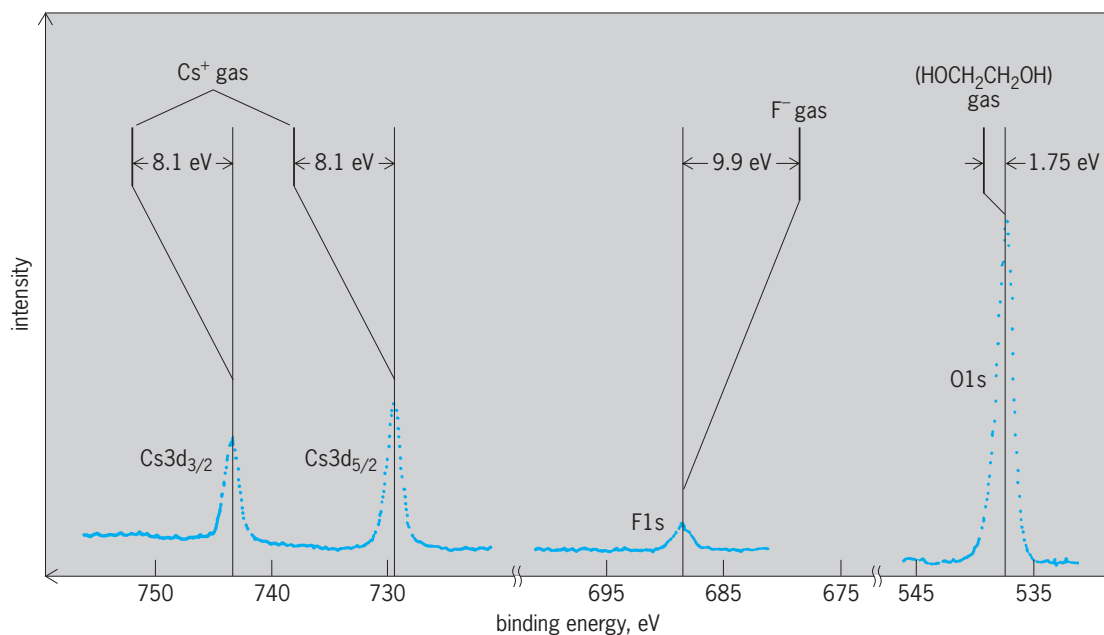


Fig. 6. ESCA spectrum of cesium fluoride (CsF) in glycol ($\text{HOCH}_2\text{CH}_2\text{OH}$) solution. Comparison is made with the corresponding electron-binding energies in the gas phase. (After H. Siegbahn et al., *Phys. Scripta*, 27:431, 1983)

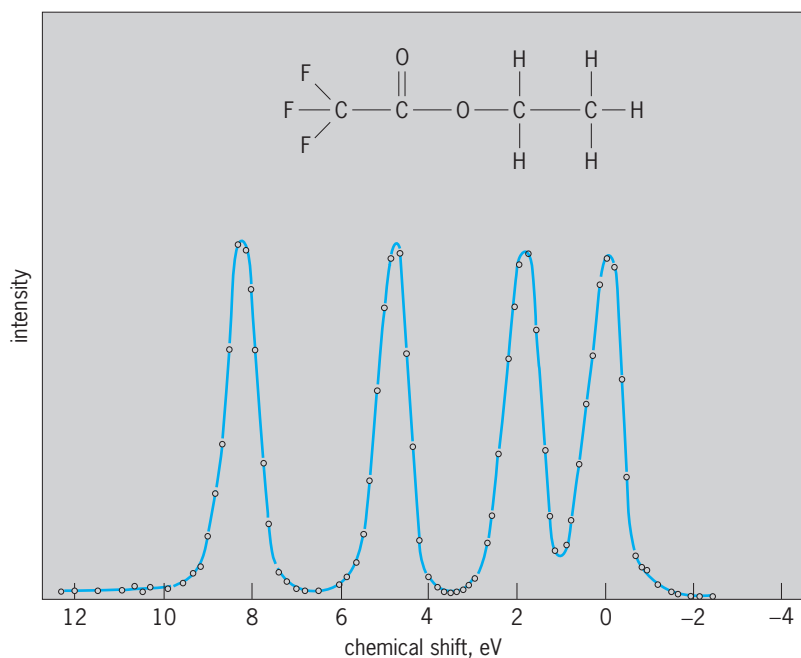


Fig. 7. Electron spectrum from the 1s level of carbon in ethyl trifluoroacetate; binding energy = 291.2 eV.

elements and even if the element represents only a small part of the chemical compound. See LINE SPECTRUM; UNCERTAINTY PRINCIPLE.

When applied to solid materials, ESCA is a typical surface spectroscopy with applications to problems such as chemical surface reactions, for example, corrosion or heterogeneous catalysis. In such cases, ultrahigh vacuum conditions are usually required, with a vacuum less than 10^{-9} torr (10^{-7} pascal). ESCA also reproduces bulk matter properties such as valence electron band structures (Fig. 4). With uv excitation, the resolution is sufficient to resolve vibrational structures in valence electron spectra for gases (Fig. 5). In fact, electron spectroscopy can supply a detailed knowledge of the valence orbital structure for all molecules which can be brought into gaseous form with pressures of 10^{-5} torr (10^{-3} Pa) or more. A convenient pressure region is 10^{-2} torr (10 Pa). Liquids and solutions of various compositions can also be studied by ESCA techniques (Fig. 6).

ESCA chemical shifts. When atoms are brought close together to form a molecule, the electronic orbitals of each atom are perturbed. Inner orbitals, that is, those with higher binding energies, may still be regarded as atomic and belonging to specified atoms within the molecule, whereas the outer orbitals combine to form the valence-level system of the molecule. These orbitals play a more or less active part in the chemical properties. The chemical bonds affect the charge distribution so that the original atoms can be regarded as charged to various degrees; a neutral molecule has a net charge of zero. The individual atoms in the molecule can be regarded as spheres with different charges set up by the transfer of certain small charges from one atomic sphere to the neighboring atoms taking part in the chemical bond. Inside each charged sphere the atomic potential is constant, in accordance with classical electro-

static theory. The result of this atomic potential is to shift the whole system of inner levels in an atom by a small amount, the same amount for each level. Levels belonging to different atoms in the molecule are shifted differently, however, and if the ESCA chemical shifts for individual atoms in the molecule are measured, a mapping can be made of the distribution of charge or potential in the molecule. This is a reflection of the chemical bondings between the atoms, which in turn can be described by the orbitals in the valence-level system.

A unique feature of ESCA is that, if the exact position of the electron lines characteristic of the various elements in the molecule is measured, the area of inspection can be moved from one atomic species to another in the molecular structure. Figure 7 shows the electron spectrum from the 1s level of the carbon in ethyl trifluoroacetate. All four carbon atoms in this molecule are distinguished in the spectrum. The lines appear in the same order from left to right, as do the corresponding carbon atoms in the structure that is shown in the figure. If the structure of the molecule is known, the charge distribution can be estimated in a simple way by using, for example, the electronegativity concept and assuming certain resonance structures. More sophisticated quantum-chemical treatments can also be applied. Conversely, if, by means of ESCA, the approximate charge distribution is known, conclusions about the structure of the molecule can be drawn.

Experimental evidence obtained so far for various elements in a large number of molecules indicates strong correlations between chemical shifts and calculated atomic charges. A typical correlation curve for the 2p level in sulfur obtained from more than 100 compounds containing this element is shown in Fig. 8. Similar curves are obtained for other elements such as carbon, nitrogen, oxygen, and phosphorus. Chemical shifts are also observed in the electron lines

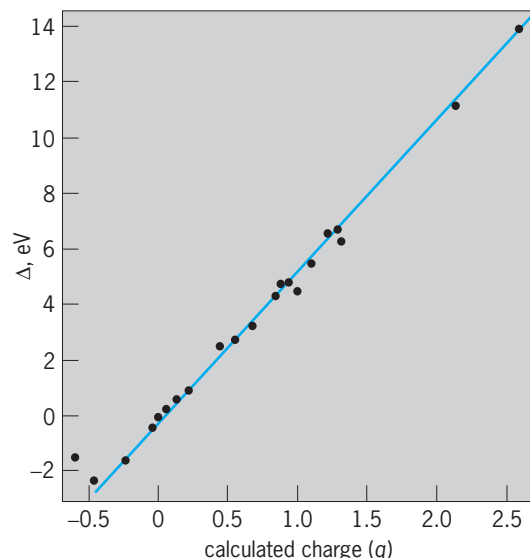


Fig. 8. Binding-energy shifts for the sulfur 2p electrons versus the calculated charge. The points indicate averages from more than 100 compounds.

due to the Auger effect. Second-order chemical shifts from groups situated farther away in the molecule (inductive effects) are also observed.

The theoretical importance of the chemical shift effect lies in the study of molecular electronic structure. In other contexts, its usefulness lies precisely in its ability to specify chemical composition, in particular, at surfaces. For example, a metal and its oxide give two distinctly different lines because of the chemical shift effect. It is often easy to follow the rate of oxidation at a surface, or the adsorption of a gas layer onto it. The chemical composition of a surface or the changes due to various chemical or physical treatments are amenable to detailed examination. See ATOMIC STRUCTURE AND SPECTRA; CHEMICAL BONDING; ELECTRON CONFIGURATION; ELECTRONEGATIVITY; MOLECULAR ORBITAL THEORY; SPECTROSCOPY.

Kai Siegbahn; Hans Siegbahn

Bibliography. T. L. Barr, *Electron Spectroscopy-Chemical Analysis (ESCA)*, 1994; C. R. Brundle and A. D. Baker (eds.), *Electron Spectroscopy, Theory, Techniques and Applications*, vols. 1-5, 1977-1984; M. Cardona and L. Ley (eds.), *Photoemission in Solids*, vols. 1 and 2, 1978, 1979; K. Siegbahn et al., *ESCA Applied to Free Molecules*, 1969; K. Siegbahn et al., *ESCA: Atomic, Molecular and Solid State Structure Studies by Means of Electron Spectroscopy*, 1967; G. C. Smith, *Surface Analysis by Electron Spectroscopy*, 1994; E. L. Wolf, *Principles of Electron Tunneling Spectroscopy*, 1985, reprint 1989.

Electron spin

That property of an electron which gives rise to its angular momentum about an axis within the electron. Spin is one of the permanent and basic properties of the electron. Both the spin and the associated magnetic dipole moment of the electron were postulated by G. E. Uhlenbeck and S. Goudsmit in 1925 as necessary to allow the interpretation of many observed effects, among them the so-called anomalous Zeeman effect, the existence of doublets (pairs of closely spaced lines) in the spectra of the alkali atoms, and certain features of x-ray spectra. See SPIN (QUANTUM MECHANICS).

All theory that concerns itself with electronic, nuclear, atomic, and molecular phenomena includes the electron spin in its formulation to obtain a theoretical structure consistent with experimental observation. The electron thus possesses the intrinsic property of spin angular momentum (rotational motion about an axis), in addition to the intrinsic properties of charge and mass. See ELECTRON.

The spin quantum number is s , which is always $1/2$. This means that the component of spin angular momentum along a preferred direction, such as the direction of a magnetic field, is $\pm 1/2\hbar$, where $\hbar = h/2\pi$ and h is Planck's constant. The spin angular momentum of the electron is not to be confused with the orbital angular momentum of the electron associated with its motion about the nucleus. In the

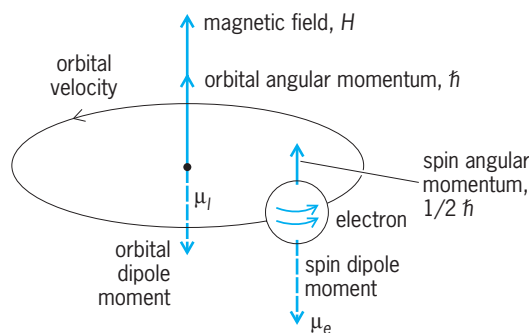


Fig. 1. Diagram of an electron in circular motion in a magnetic field, showing spin and orbital angular momenta and resulting magnetic moments.

latter case the maximum component of angular momentum along a preferred direction is $l\hbar$, where l is the angular momentum quantum number and may be any positive integer or zero. The total orbital angular momentum is equal to $\sqrt{l(l+1)}\hbar$. In this discussion the terms angular momentum or magnetic dipole moment describe the maximum component of these quantities along a field direction. See ANGULAR MOMENTUM; QUANTUM NUMBERS.

Electron magnetic moment. The electron has a magnetic dipole moment by virtue of its spin. The approximate value of the dipole moment is the Bohr magneton μ_0 which is equal, in SI units, to $eb/4\pi m = 9.27 \times 10^{-24}$ joule/tesla, where e is the electron charge measured in coulombs, m is the mass of the electron, and c is the velocity of light. The orbital motion of the electron also gives rise to a magnetic dipole moment μ_l that is equal to μ_0 when $l = 1$. In Fig. 1 is shown the simple case of an electron in circular motion in a magnetic field. The electron is shown with a spin angular momentum parallel to the orbital angular momentum. (Physically it could equally well be in the opposite direction.) Since the electron has a negative charge, the directions of the orbital dipole moment and spin dipole moment (both are vectors) are opposite to that of the orbital angular momentum. For a positron (a positively charged particle having the same mass and magnitude of charge as the negatively charged electron), the magnetic moments are positive, that is, in the same direction as the angular momentum. See MAGNETON; POSITRON.

The orbital magnetic moment of an electron can readily be deduced with the use of the classical statements of electromagnetic theory in quantum-mechanical theory; the simple classical analog of a current flowing in a loop of wire describes the magnetic effects of an electron moving in an orbit. The spin of an electron and the magnetic properties associated with it are, however, not possible to understand from a classical point of view. The classical radius of the electron is $e^2/(8\pi\epsilon_0 mc^2) = 1.141 \times 10^{-13}$ cm, where ϵ_0 is the permittivity of free space; a reasonable distribution of mass and electric charge within this radius which would lead to a magnetic moment μ_0 leads to the calculation of a peripheral velocity of the electron far greater than the velocity

of light. This is, of course, wholly precluded by the special theory of relativity. No theory of the structure of the electron has been formulated which makes the spin of the electron amenable to simple pictorial understanding. Nevertheless the interpretation of the spectra of atoms and molecules, the magnetic properties of materials, and other phenomena on an atomic scale unambiguously require that the electron have the property of spin. To the extent that physical theory assumes these properties and does not concern itself with questions about the structure of the electron, it is quite adequate to deal with a large range of physical phenomena in a highly quantitative way. See RELATIVITY.

In the Landé g factor, g is defined as the negative ratio of the magnetic moment, in units of μ_0 , to the angular momentum, in units of \hbar . For the orbital motion of an electron, $g_l = 1$. For the spin of the electron the appropriate g value is $g_s \simeq 2$; that is, unit spin angular momentum produces twice the magnetic moment that unit orbital angular momentum produces.

The following discussion is limited to atoms which have a single electron outside of closed electron shells. Both the orbital and spin angular momenta of the electrons within closed shells add up in such a way that their net angular momentum is zero. The single electron outside closed shells may have $l = 0, 1, 2, \dots$. By the usual rules developed from quantum mechanics the total angular momentum quantum number of the electron, which is called j , is $l \pm s$ (Fig. 1) except when $l = 0$, in which case the total angular momentum is s . For instance, when $l = 1$, $j = 1/2$ or $3/2$. These relations may be represented by vector diagrams, as shown in Fig. 2. Since the revolution of the electron about the nucleus produces a magnetic field at the electron, and since the electron has a magnetic dipole moment, the energy of the atom is different when the vector s is parallel to l ($j = l + 1/2$) and when s is antiparallel (180°) to l ($j = l - 1/2$). Thus the spin of the electron causes a doubling of the energy levels in all single-electron atoms except when $l = 0$, in which case level is single. In the case of sodium, the familiar yellow lines (the D lines) comprise a closely spaced doublet that

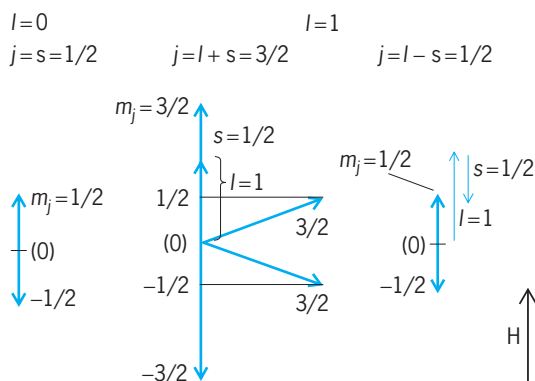


Fig. 2. Diagram illustrating addition of l and s vectors into a resultant j vector. The allowed values of m_j are also shown. In no case does the (0) give an allowed value. H = magnetic field.

arises from a transition from a state of $l = 1$ to one of $l = 0$. The doubling of the lines is a direct consequence of the existence of the electron spin.

Energy-level splitting. The total electronic magnetic moment of an atom depends on the state of coupling between the orbital and spin angular momenta of the electron (Fig. 2). In the single-electron case, an atom in the state for which $l = 0$ has only the magnetic moment associated with the spin, and this moment may be oriented in either of two directions with respect to an externally applied magnetic field. However, when the atom is in a state for which $l = 1$, j can be $1/2$ or $3/2$ and the Landé g factors for these two states are $2/3$ and $4/3$, respectively. That is, the magnetic moment per unit angular momentum is equal neither to that which characterizes the orbital motion nor to that which characterizes the spin motion. In a magnetic field, a single energy level characterized by l and j is split into several energy levels, each described by the component of j , m_j , along a magnetic field, where $m_j = j, j - 1, \dots, -(j - 1), -j$. The energy of the level is the zero field energy plus the term $m_j g_j \mu_0 H$. In a transition between two energy levels that gives rise to a spectral line, the energy of the emitted or absorbed photon is the zero field energy difference between the two levels plus the difference between two magnetic energy terms as given previously. The resultant splitting of the line, which is single at zero magnetic field, into a line of two or more components is called the Zeeman effect. The Zeeman effect has been called anomalous when it is not explicable purely in terms of the orbital motion of the electron. The introduction of the electron spin allows the interpretation of all observed Zeeman effects. See ZEEMAN EFFECT.

Atomic beam measurements. Prior to 1940, spectroscopic measurements had been made only on optical spectra, and Zeeman effects were small effects superimposed on lines of considerable natural width. Exact measurements of the splitting of lines in a magnetic field, therefore, could not be made on such lines, and all data on the Zeeman effect were consistent with the statement that $g_s = 2g_l$. With the development of spectroscopy by the atomic beam method, a new order of precision in the measurement of the frequencies of spectral lines became possible. By using the atomic-beam techniques, it became possible to measure g_s/g_l directly, with the result $g_s/g_l = 2$ (1.001168 ± 0.000005). The magnetic moment of the electron therefore is not μ_0 but $1.001168\mu_0$, or equivalently the g factor of the electron departs from 2 by the so-called g factor anomaly defined as $a = (g_s - 2)/2$ so that $\mu = (1 + a)\mu_0$. Thus the first molecular beam work gave $a = 0.001168$. See MOLECULAR BEAMS.

Calculation of g -factor anomaly. It is not possible to give a qualitative description of the effects which give rise to the g -factor anomaly of the electron. The detailed theoretical calculation of the quantity is in the domain of quantum electrodynamics, and involves the interaction of the zero-point oscillation of the electromagnetic field with the electron.

Comparison of theoretical determination of a with its experimental measurement constitutes the most accurate and direct existing test of the theory of quantum electrodynamics.

Theoretical work on the g -factor anomaly, based on the principles of quantum electrodynamics, began simultaneously with its experimental discovery. The initial prediction of the value of the anomaly was made by J. Schwinger in 1948, who showed that the anomaly could be written as $a = 0.5(\alpha/\pi)$, where α , the fine-structure constant, is given by $\alpha = e^2/\hbar c \simeq 1/137$. Shortly thereafter, it was shown that a could be expressed as a power series in α (more customarily, α/π), that is, one can write $a = A(\alpha/\pi) + B(\alpha/\pi)^2 + C(\alpha/\pi)^3 + D(\alpha/\pi)^4 + \dots$. Calculation of B , C , and D from quantum electrodynamics has proven to be very difficult, with errors occurring at various stages of the work. However, it is now generally agreed that $A = 0.500$, $B = -0.328478965\dots$, $C = 1.181241456\dots$, and $D = -1.7283 \pm 0.0035$. In order to find a theoretical value for a , an accurate value of α is required. There are several values of α , obtained from different sources, which may be used to determine a . The most accurate value of α is based upon results obtained by a measurement of the recoil velocity of rubidium atoms in an optical lattice and is given by $\alpha^{-1}_{\text{Rb}} = 137.03599878 \pm 0.00000091$. The result of substituting this value of α into the expression for a gives the most recent theoretical value of a as $a(\text{theoretical}) = 0.0011596521887 \pm 0.000000000077$. The uncertainty quoted for this number is due primarily to error in α . Finally, the weak and strong interactions have been calculated to add about 0.000000000017 to $a(\text{theoretical})$. See ELECTRICAL UNITS AND STANDARDS; FUNDAMENTAL CONSTANTS; HALL EFFECT; QUANTUM ELECTRODYNAMICS; STRONG NUCLEAR INTERACTIONS; WEAK NUCLEAR INTERACTIONS.

Measurement of g -factor anomaly. The initial molecular-beam method for measuring a was based on measurement of g_s/g_b , which yields $1 + a$ rather than a itself. Such experiments have been superseded by two types of experiments which have the major advantages of measuring a directly, and of measuring a for electrons trapped in electric and magnetic fields but not bound to atoms. The second feature means that various corrections due to the presence of an atom are no longer necessary.

The new experiments may be readily understood, qualitatively, by noting that if an electron (spin angular momentum $\hbar/2$) moves at a low velocity $V(V/c \ll 1)$ perpendicular to a magnetic field B , the particle will rotate in the field at a frequency (the cyclotron frequency) given by $f_c = (1/2\pi)(eB/m)$, while its magnetic moment (and spin) will precess about the field at a frequency (the spin precession frequency) given by $f_s = (1 + a)f_c$. The difference between these frequencies, often called the $g - 2$ or difference frequency, is given by $f_D = f_s - f_c = af_c$. Thus f_D represents the rate at which the spin of the particle rotates relative to the particle's velocity.

One technique for measuring a consists of trapping very low-energy electrons in a magnetic field

and measuring f_D and f_c simultaneously, with a being given by the relation $a = f_D/f_c$.

A slightly different version of the same technique consists of trapping electrons whose velocity is about half the speed of light (electrons accelerated to several hundred thousand volts) in a magnetic field and determining the same ratio. Relativistic effects manifest themselves in both the expressions for f_s and f_c , but cancel exactly when the ratio f_D/f_c is formed. The high-energy technique constitutes an excellent check of the separate measurements, as well as one of the most precise tests extant of the special theory of relativity, which predicts that the frequency ratio should be independent of velocity. See RELATIVISTIC ELECTRODYNAMICS.

The most accurate experimental result has been obtained by using the low-energy technique. This result is $a(\text{experimental}) = 0.00115965218085 \pm 0.0000000000076$, and it is in agreement, within the limits of error, with $a(\text{theoretical})$. The comparison of $a(\text{experimental})$ with $a(\text{theoretical})$ presented above constitutes the most precise confrontation of any experiment with a theoretical prediction in the history of science, and it provides the value of the fine-structure constant, $\alpha^{-1} = 137.035999710 \pm 0.000000096$, with a relative precision of 0.7 part per billion (10^9), which is an order of magnitude more precise than any other measurement available at present. See ATOMIC STRUCTURE AND SPECTRA; GYROMAGNETIC RATIO; QUANTUM MECHANICS.

Arthur Rich; Toichiro Kinoshita

Bibliography. R. Eisberg and R. Resnick, *Quantum Physics of Atoms, Molecules, Solids, Nuclei, and Particles*, 2d ed., 1985; G. Gabrielse et al., New determination of the fine structure constant from the electron g value and QED, *Phys. Rev. Lett.*, 97:030802, 2006; P. J. Mohr and B. N. Taylor, CODATA recommended values of the fundamental physical constants: 2002, *Rev. Mod. Phys.*, 77:1-107, 2005; B. Odom et al., New measurement of the electron magnetic moment using a one-electron quantum cyclotron, *Phys. Rev. Lett.*, 97:030801, 2006; H. C. Ohanian, *Modern Physics*, 2d ed., 1995; J. W. Rohl, *Modern Physics from α to Z^0* , 1994.

Electron tomography

A powerful technique for determining three-dimensional (3D) structure from a series of two-dimensional (2D) images recorded in a transmission electron microscope. Transmission electron microscopy uses high-energy electrons to illuminate specimens at resolutions that in favorable cases surpass 1.0 nanometer (1/25,000,000 of an inch) and can resolve atoms. Resolution is not as great with electron tomography (4–8 nm) but is still much higher than can be achieved by light microscopy or medical imaging methods. However, the weak penetrating power of electrons requires operation under high vacuum and generally limits the technique to specimens less than about 1 micrometer in diameter. The most common application is

in biological research to study the 3D architecture of subcellular components such as mitochondria, nerve synapses, the Golgi apparatus, muscle cells during contraction, the arrangement of molecular components in bacterial cells, chromosome fibers, the kinetochore, and other structures. Such studies provide important data for understanding how cellular components carry out their vital functions. Electron tomography is also being used to study the structure of certain inorganic crystals used in the semiconductor industry.

Principles of operation. Transmission electron microscopy provides a translucent interior view of the specimen, analogous to what is seen in an x-ray image, rather than the surface view that we are accustomed to seeing in everyday life. Although this enables the microscopist to see inside the specimen, it has the disadvantage that the image formed is a 2D projection of the 3D object with features from different depths superimposed (Fig. 1a). The solution is to reconstruct a 3D image of the original specimen from a series of 2D projections (Fig. 1b). This approach is also employed in medical imaging methods such as computerized axial tomography (CAT scanning) and magnetic resonance imaging (MRI).

The series of 2D images required for tomographic reconstructions is obtained by tilting the specimen in the electron beam. This is in contrast to medical imaging in which the patient is kept stationary

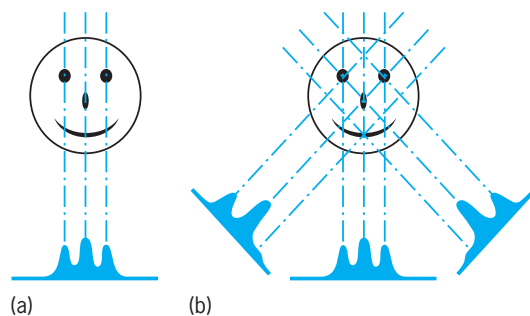


Fig. 1. Projection imaging and tomographic reconstruction. (a) Formation of one-dimensional (1D) projection of a two-dimensional (2D) smiley face. To form a true projection, the amount of penetrating radiation that strikes a unit length of the recording medium must be a measure of the amount of mass traversed by the beam in its path through the specimen. This is illustrated by the image density graph superimposed upon the 1D projection, and the broken lines that represent the radiation path through areas of maximum density in the specimen (that is, through the eyes or nose plus the mouth of the face). In transmission electron microscopy, true projections are formed as long as the specimen is not thicker than 0.2–0.5 micrometer. (b) Tomographic reconstruction of the 2D smiley face. The first step in tomography is to systematically record a series of 1D images from different viewing directions. The density in each 1D image is then projected back into a 2D area by simply distributing the density evenly along the original projection direction. When these “back-projection rays” from the different 1D images are summed, they intersect and reinforce one another only at the points where there was significant mass in the original object. In this way the 2D image is reconstructed from the 1D images. The same strategy is used to reconstruct a 3D image from 2D projection images. (Adapted with permission from B. F. McEwen and M. Marko, *Three-dimensional transmission electron microscopy and its application to mitosis research*, *Meth. Cell Biol.*, 61:82–111, 1999)

while the beam source and detector are rotated to provide the necessary views. An important limitation of electron tomography is the inability to collect tilt data over the full 180° range required for complete angular coverage. (A 360° range is not necessary because, in projection imaging, views from the opposite side of the specimen provide the same information.) The difficulty in electron tomography is that typical electron microscope specimens are very thin (50–500 nm) and mounted on a support film that is more than 1000 times wider than the specimen. As a result, images can be recorded only over a limited tilt range of ±60 or 70° (a total range of 120–140°) before the specimen becomes too thick to view. This gives rise to a range of missing angular coverage that produces a well-characterized distortion that must be taken into account when interpreting electron tomographic reconstructions.

Electron tomographic data sets are generally collected with 1 or 2° tilt angle intervals, because the resolution obtained in the resulting reconstruction is dependent upon how finely the angular range is sampled. Thus, a typical tilt series contains 60–140 images. Recording so many images was time-consuming and tedious with older microscopes, but recent improvements to mechanical design, along with incorporation of direct digital recording of the images and computer control of the microscope, have enabled newer microscopes to automatically record a tilt series of 140 images in less than 30 min. These technical developments in electron microscopy, coupled with the rapidly increasing computational power and the disk storage capacity of modern computers, are transforming electron tomography into a high-throughput tool that can be used ever more routinely in research and diagnostic settings.

Volume analysis and sample application. Perhaps the most surprising aspect of electron tomography is that it generally takes longer to analyze an electron tomographic reconstruction than it does to collect the data and compute the reconstruction. The difficulty is that the investigator is faced with displaying a three-dimensional tomographic reconstruction on a two-dimensional computer screen. However, in contrast to the original imaging situation, one now has a 3D image in the computer that can be taken apart and examined from any direction. Usually the quickest way to locate components of the structure is to view the reconstruction volume as a sequential stack of thin slices. However, understanding the 3D relationships of component parts and communicating the finding to a larger audience generally require constructing a 3D representation. For complex structures the most common approach is to segment select components away from the rest of the tomographic reconstruction by manual tracing and assignment of a color code.

Figure 2 illustrates the application of electron tomography to the study of the mammalian kinetochore. The kinetochore is a specialized structure that forms at a specific location on chromosomes during cell division. It functions to attach chromosomes to a bipolar spindle, generate chromosome motion, and exert control over the events of cell

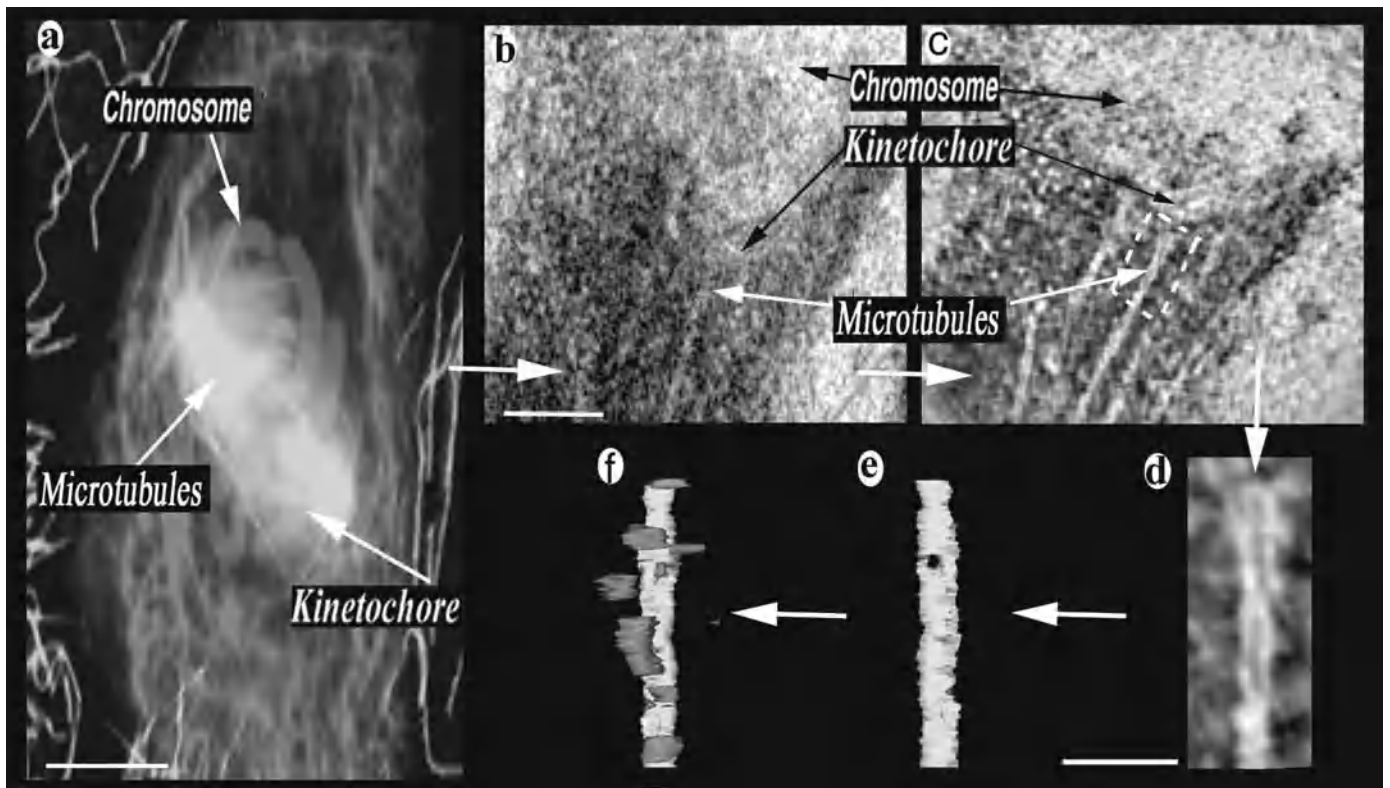


Fig. 2. Using electron tomography to analyze cell division. (a) An immunofluorescent light microscopy image of a mammalian cell during cell division. (b) Electron microscopy image of a kinetochore from the same type of cell. The kinetochore is the region of the chromosome that attaches to the mitotic spindle. (c) Single 2D slice from the 3D electron tomographic reconstruction of the kinetochore in b. Note the sharper microtubules in the tomographic reconstruction. (d) Single microtubule extracted from the tomographic reconstruction. (e) 3D model created by tracing and color-coding (light gray) the microtubule in d. (f) Same model as in e showing kinetochore components (dark gray) that are attached to the microtubule.

division. These functions are essential for proper cell division and the viability of the organism. Figure 2a is a gray-scale version of an immunofluorescence light microscopy image of a dividing cell, with arrows indicating a chromosome, location of a kinetochore, and microtubules. Microtubules are the chief structural components of the mitotic spindle. In the actual immunofluorescent image, microtubules are stained green, chromosomes blue, and keratin fibers, which form a cage around the spindle, red. One of the crucial yet poorly understood aspects of chromosome alignment is the interaction between the kinetochore and spindle microtubules. It is clear from Fig. 2a that the resolution of light microscopy is inadequate to study this interaction. Although single microtubules are visible in an electron microscopy projection image (Fig. 2b), details of the microtubule ends become visible only in single slices from the tomographic reconstruction (Fig. 2c, d). Manual tracings from trimmed slices such as in Fig. 2d were stacked into a 3D image, which are shown as a surface view in Fig. 2e. Kinetochore connections to the microtubules were also traced, and are shown in dark gray along with the microtubule in Fig. 2f. These and similar 3D models are being used to test current hypotheses on how the kinetochore interacts with microtubule ends.

Future directions and conclusions. Recent technical advances and the inherent versatility of electron tomography have accelerated its use and expanded the

range of its applications. In particular, electron tomography is now being used for high-quality frozen-hydrated preparations of biological specimens. These preparations are made by rapidly freezing the specimen to prevent the formation of normal crystalline ice, which severely damages cellular structure. Such preparations are a snapshot of the specimen in its native hydrated environment and, thereby, superior to conventional preparations made by solvent extraction, heavy-metal staining, and plastic embedding. Previously the sensitivity of frozen-hydrated specimens to damage from electron irradiation had precluded such applications.

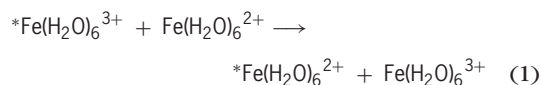
Electron tomography has been established as an important method for determining the structure of subcellular components and organelles, and the number and range of applications are likely to continue growing rapidly in the near future. See COMPUTERIZED TOMOGRAPHY; ELECTRON MICROSCOPE; FLUORESCENCE MICROSCOPE; MEDICAL IMAGING; OPTICAL MICROSCOPE.

Bruce F. McEwen

Bibliography. W. Baumeister, R. Grimm, and J. Walz, Electron tomography of molecules and cells, *Trends Cell Biol.*, 9:81–85, 1999; B. F. McEwen and M. Marko, The emergence of electron tomography as an important tool for investigating cellular ultrastructure, *J. Histochem. Cytochem.*, 49:553–563, 2001; B. F. McEwen and M. Marko, Three-dimensional transmission electron microscopy and its application to mitosis research, *Meth. Cell Biol.*, 61:82–111, 1999.

Electron-transfer reaction

A reaction in which one electron is transferred from one molecule or ion to another molecule or ion. Electron-transfer reactions are ubiquitous in nature. Some are deceptively simple [for example, reaction (1), where the asterisk is used to identify a



specific isotope]; others look very complicated (for example, the long-range electron transfers found in biology). The widespread occurrence of electron-transfer reactions has stimulated much theoretical and experimental work. *See* BIOLOGICAL OXIDATION; CHEMOSMOSIS; PHOTOSYNTHESIS.

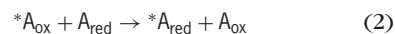
The development of electron-transfer theory, as well as the measurement of the rates of oxidation-reduction (redox) reactions, began in the late 1940s. A great deal of experimental work has been carried out by inorganic chemists, and during the 1960s and 1970s the reactivity patterns of a large number of complexes were uncovered. Using this body of knowledge as a reference point, bioinorganic chemists studying the rates and mechanisms of biological electron-transfer reactions are seeking parallels with the redox behavior of less complicated metal complexes. *See* BIOINORGANIC CHEMISTRY; COORDINATION COMPLEXES.

Metalloproteins are more than just metal ions in disguise. For example, the redox behavior as well as virtually every property of a protein (excluding its amino acid sequence) are dependent upon the solution pH; redox proteins are very large polyelectrolytes whose prosthetic groups are typically found buried in the protein interior. Thus pH-dependent behavior in a metalloprotein may be much more complex than in a simple metal complex. Another important distinction between the redox reactions of proteins and the electron transfers of small metal complexes is the magnitude of the electron donor-to-acceptor distance. The relevant distance for small molecules is generally taken to be van der Waals contact, which is roughly 0.3 nanometer. In contrast, several experiments have shown that electrons can jump at significant rates across distances of 1 nm or more in protein interiors. *See* INTERMOLECULAR FORCES; PH; PROTEIN.

Experimental investigation of the factors that control the rates of biological redox reactions has not come as far as the study of the electron transfers of metal complexes, because of the large number of variables that must be dealt with, and problems associated with protein purification and stability. Many experimental approaches have been pursued, including the covalent attachment of redox reagents to the surfaces of metalloproteins and the use of mutant proteins (obtained by gene cloning). *See* BIOELECTRONICS.

Self-exchange and cross reactions. The simplest reactions in solution chemistry are electron self-

exchange reaction (2), in which the reactants and



products are the same (the asterisk is used to identify a specific isotope). The only way to determine chemically that a reaction has taken place is to introduce an isotopic label. There is no change in the free energy ($\Delta G^\circ = 0$) for this type of reaction.

Much more common are cross reaction (3), where



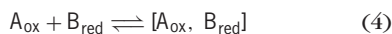
A_{ox} is the oxidized reactant, B_{red} is the reduced reactant, A_{red} is the reduced product, and B_{ox} is the oxidized product. For these reactions, $\Delta G^\circ \neq 0$. The experimental determination of cross-reaction rates is generally more straightforward than the determination of self-exchange rates. Either the reactants are simply mixed together or a thermodynamically unstable system is generated rapidly (via pulse radiolysis, flash photolysis, or temperature-jump relaxation) to initiate the redox reaction. In almost all cases, electronic absorption spectroscopy has been used to monitor the progress of protein cross reactions. *See* SPECTROSCOPY.

Inner-sphere and outer-sphere mechanisms. Both types of electron-transfer reactions (self-exchange and cross reactions) can be classified broadly as inner sphere or outer sphere. In an inner-sphere reaction, a ligand is shared between the oxidant and reductant in the transition state. An outer-sphere reaction, on the other hand, is one in which the inner coordination shells of both the oxidant and reductant remain intact in the transition state. There is no bond breaking or bond making, and no shared ligands between redox centers. Long-range electron transfers in biology are all of the outer-sphere type.

Electron-transfer theory. The simplest electron transfer occurs in an outer-sphere reaction. The changes in oxidation states of the donor and acceptor centers result in a change in their equilibrium nuclear configurations. This process involves geometrical changes, the magnitudes of which vary from system to system. In addition, changes in the interactions of the donor and acceptor with the surrounding solvent molecules will occur. The Franck-Condon principle governs the coupling of the electron transfer to these changes in nuclear geometry: during an electronic transition, the electronic motion is so rapid that the nuclei (including metal ligands and solvent molecules) do not have time to move. Hence, electron transfer occurs at a fixed nuclear configuration. In a self-exchange reaction, the energies of the donor and acceptor orbitals (hence, the bond lengths and bond angles of the donor and acceptor) must be the same before efficient electron transfer can take place. *See* FRANCK-CONDON PRINCIPLE.

The incorporation of the Franck-Condon restriction leads to the partitioning of an electron-transfer reaction into reactant (precursor complex) and product (successor complex) configurations.

Steps (4)–(6) go from reactants to products: K is the



equilibrium constant for the formation of the precursor complex $[A_{\text{ox}}, B_{\text{red}}]$, and k_{et} is the forward electron-transfer rate to produce the successor complex $[A_{\text{red}}, B_{\text{ox}}]$.

Some degree of electronic interaction, or coupling between donor and acceptor, is required if the redox system is to pass from the precursor to the successor complex. (This passage involves an intermediate species known as a transition state or activated complex.) The magnitude of the electronic coupling determines the behavior of the reactants once the transition state is reached. Two cases can be distinguished. If the electronic coupling is very small, as in so-called nonadiabatic reactions, there is a very low probability that the reactants will undergo electron transfer, leading to very little product formation. If the electronic interaction is sufficiently large, as is the case for adiabatic reactions, the probability of electron transfer occurring when the reactants reach the transition state is unity. The degree of adiabaticity of a reaction is characterized by a transmission coefficient κ , whose value ranges from 0 to 1. For systems in which the electronic coupling is sufficiently large, κ equals 1. This situation occurs when the reacting centers are close together, the donor and acceptor orbitals overlap, and no substantial changes in geometry are involved. The transmission coefficient is much less than 1 for electron transfer over long donor-to-acceptor distances, as commonly found in protein electron-transfer reactions and in some electron transfers in large organic molecules.

Long-range electron transfer. The rate of long-range electron transfer between an electron donor (B_{red}) and an electron acceptor (A_{ox}) depends on both the electronic coupling between A_{ox} and B_{red} (which is a function of the intersite $A_{\text{ox}}/B_{\text{red}}$ distance d , the nature of the intervening medium, and the relative $A_{\text{ox}}/B_{\text{red}}$ orientation, where // represents the protein medium that separates the donor and the acceptor) and an activation energy term. A standard theoretical rate equation (7) expresses k_{et} in terms of these

$$k_{\text{et}} = \nu [\exp(-\beta d)] \{ \exp[-(\lambda + \Delta G^\circ)^2 / 4\lambda RT] \} \quad (7)$$

factors; here ν is a frequency factor, β is a medium- and orientation-dependent quantity, d is the intersite λ is the reorganization energy, ΔG° is the reaction free energy of the electron-transfer process, R is the universal gas constant, and T is the absolute temperature of the system. Experiments in several laboratories have been designed to estimate the values of λ and β in modified metalloproteins, rigid organic molecules, and protein-protein complexes.

Modified metalloproteins. Long-range intramolecular electron transfer is being studied in synthetically

modified metalloproteins. Ruthenium (Ru) amines have been chosen as modification agents primarily because they can be attached specifically to surface histidines, and the resulting complexes are kinetically stable in solution in the relevant oxidation states. X-ray crystallographic structure determinations have been completed for the native proteins used in these studies, and so the distance and intervening medium between the donor and acceptor sites in the modified protein are accurately known. [X-ray crystallographic studies of a modified protein show that only minimal structural perturbations result from covalent attachment of ruthenium to a surface histidine (His).] The driving force ΔG° for the electron transfer between donor and acceptor of a modified protein can be determined accurately through standard electrochemical experiments. Thus electron transfer can be studied in a system where two of the important variables in Eq. (7), d and δG° , are reliably known. See X-RAY CRYSTALLOGRAPHY.

Sperm whale myoglobin (Mb) has four surface histidines that can be modified by using ruthenium amine complexes. The pentaammine-ruthenium(III) $[a_5\text{Ru}]$ complex has been covalently attached to these histidines. After separation and purification of the reaction products, four singly labeled myoglobin species are obtained, with the closest edge-to-edge distances from the heme ranging from 1.27 nm (His-48) to 2.2 nm (His-12). Electron transfer from the His-48 ruthenium complex to the iron (Fe)-containing heme has been measured by selectively reducing the Ru center on the exterior of the protein using flash photolysis techniques. This is followed by electron transfer from Ru^{2+} to the heme Fe^{3+} in the interior of the protein. The forward electron-transfer rate k_{et} (Ru^{2+} to Fe^{3+}) in $a_5\text{Ru}$ -(histidine-48)-modified myoglobin $[a_5\text{Ru}(48)\text{MbFe}]$ was found to be 0.02 s^{-1} ; the reverse rate k_r (Fe^{2+} to Ru^{3+}) was uniquely determined to be 0.04 s^{-1} by employing a method developed to generate the $\text{Fe}^{2+}/\text{Ru}^{3+}$ mixed-valence species. These kinetic results clearly demonstrate that long-range electron transfer is reversible in this system.

Cytochrome *c* (from horse heart) also has been modified with the $a_5\text{Ru}$ reagent. In contrast to the myoglobin system, cytochrome *c* has only one surface histidine (His-33) that readily binds ruthenium. The edge-to-edge distance from this histidine to the heme is 1.16 nm. By using flash photolysis, the Ru^{2+} -to- Fe^{3+} electron-transfer rate was found to be approximately 30 s^{-1} . Pulse radiolysis has also been employed on the same system, with a measured electron-transfer rate of 50 s^{-1} . See CYTOCHROME.

Driving-force dependences have been addressed in two ways with modified proteins. First, the amine ligands on the ruthenium label have been varied to produce complexes with different reduction potentials. Second, the iron from the heme has been removed and replaced with zinc, magnesium, and palladium. After photoexcitation, the metalloporphyrin functions as an electron donor, leading to electron transfer to the surface-bonded ruthenium.

With these methods, a wide range of driving forces for intramolecular electron transfer in cytochrome *c* and myoglobin has been investigated. At the highest driving forces in cytochrome *c*, rates of roughly $3 \times 10^6 \text{ s}^{-1}$ were found. Values of the reorganization energy λ between 1.3 and 2 eV have been obtained from these experiments.

The distance dependence of k_{et} has been investigated in the $a_5\text{RuMbZn}$ system. Values of β in the 8–10-nm⁻¹ range have been extracted from these studies, thereby indicating that a protein is a good medium for electronic coupling of a donor and acceptor.

Rigid organic molecules. Long-range electron-transfer reactions have been investigated in molecules containing organic donors and acceptors separated by rigid hydrocarbon units. With hydrocarbon spacers of 1–1.2 nm, electron-transfer rates as high as 10^{10} s^{-1} have been observed. The rates in the organic molecules are 10,000 times faster than protein electron transfers at comparable donor-to-acceptor distances. In one study involving hydrocarbon spacers of lengths 0.7–1.4 nm, β was found to be 8.5 nm⁻¹.

Long-range electron-transfer reactions in rigid organic molecules have been studied only in organic solvents; electron-transfer theory suggests that the reorganization energy λ should be smaller in non-aqueous media than reorganization energies for comparable long-range electron transfers in aqueous solutions. Results obtained are consistent with theory; a reorganization energy of roughly 1 eV was found in a thorough investigation of driving-force effects on the rate of long-range electron transfer in a donor (hydrocarbon spacer)-acceptor system.

Protein-protein complexes. In biologically relevant partners, both redox centers are buried in a protein matrix; the entire electron transfer occurs through a protein medium and includes a protein-protein interface. The reorganization energy λ has been measured by determining both the temperature dependence and the free-energy (ΔG°) dependence of electron-transfer rates in these systems. These protein-protein complexes also are well suited for the study of the effects of the medium on biological electron-transfer rates. Two approaches have been used. One is to compare electron-transfer rates between homologous (same species) and heterologous (different species) protein redox pairs. The other is to use the technique of site-directed mutagenesis to change amino acids that lie between the two redox centers in the protein-protein complex.

The most commonly used experimental system involves metal substitution of zinc for iron in one of the porphyrin redox centers in the protein-protein complex. This system allows facile initiation of electron transfer through photoexcitation of the zinc porphyrin (ZnP). The excited zinc porphyrin (ZnP*) may then decay to the ground state or transfer an electron: the electron-transfer reaction is $\text{Fe}^{3+} // \text{ZnP}^* \rightarrow \text{Fe}^{2+} // \text{ZnP}^+$. The Fe^{2+} in $\text{Fe}^{2+} // \text{ZnP}^+$ will transfer an electron back to ZnP^+ , thereby regenerating the starting material. The forward rate ($\text{ZnP}^* \rightarrow \text{Fe}^{3+}$) is k_{et} , and the back rate ($\text{Fe}^{2+} \rightarrow \text{ZnP}^+$) is k_b . See PORPHYRIN.

Hybrid hemoglobins. Hybrid hemoglobin is the best characterized of the protein-protein complexes. In this system, the donor-acceptor separation distance and intervening medium are known from x-ray structure determinations. Both $\alpha(\text{Zn})\beta(\text{Fe})$ and $\alpha(\text{Fe})\beta(\text{Zn})$ hybrids have been studied. Electron transfer ($\text{ZnP}^* \rightarrow \text{Fe}^{3+}$) between the α_1 and β_2 subunits of hemoglobin has been observed: the metal-metal distance is 2.5 nm (edge-to-edge porphyrin distance is 2.0 nm). The ΔG° for the electron transfer is approximately 0.8 V, and at room temperature k_{et} is approximately 100 s^{-1} . The temperature dependence of k_{et} has been analyzed by both quantum-mechanical and semiclassical theories; values of λ in the 2–3-eV range have been obtained. See HEMOGLOBIN.

Effects of the medium. The effects of the medium on long-range electron transfer have been investigated in the cytochrome *c* peroxidase (CCP)/cytochrome *c* (cyt *c*) complex. By using cyt *c* $\text{Fe}^{3+} // \text{CCPZn}$, ΔG° for back electron transfer is approximately 0.85 V. When yeast cytochrome *c* peroxidase was used in this reaction, k_{et} for the complex with yeast cytochrome *c* was found to be 10-fold faster than that for complexes with tuna cytochrome *c* or horse cytochrome *c*. The rate k_b was found to be 1000-fold faster for yeast than for tuna cytochrome *c*. These data indicate that specific medium and interface contact effects may enhance the reaction in the homologous yeast pair. The fact that the homologous and heterologous pairs show such radically different rate enhancements for k_{et} and k_b suggests a conformational change after the forward electron transfer. See CONFORMATIONAL ANALYSIS.

By using site-directed mutagenesis, phenylalanine-87, which is located above the heme pocket, has been replaced by tyrosine and serine in yeast cytochrome *c*. Electron-transfer studies with CCPZn have demonstrated that k_{et} is not affected by these mutations. The back electron transfer, on the other hand, is unaffected for the tyrosine mutant but is decreased 10,000-fold for the serine mutant. These data suggest that aromatic amino acids act as mediators for electron transfer. Although reorganization energy effects caused by the mutation to serine cannot be ruled out, these experiments demonstrate that changes in the intervening protein medium can produce dramatic electron-transfer effects.

Driving-force dependence of k_{et} . The ΔG° dependence of the electron-transfer rate has also been studied in cyt *c* // CCP by using a series of substituted horse cytochromes *c*. These experiments yielded an estimate of 1.5–2.0 eV for the reorganization energy λ . The λ values for this protein-protein complex and for hybrid hemoglobin are in reasonable agreement with those found for $a_5\text{Ru}$ -modified proteins. The origin of the relatively large reorganization energies for protein electron transfers in aqueous solution is being actively investigated in several laboratories. See CHEMICAL THERMODYNAMICS; OXIDATION-REDUCTION.

Bruce E. Bowler; Walther R. Ellis, Jr.; Harry B. Gray; Thomas J. Meade

Bibliography. J. R. Bolton, N. Mataga, and G. M. McClendon (eds.), *Electron Transfer in Inorganic, Organic, and Biological Systems*, 1991; G. L. Closs and J. R. Miller, Intramolecular long-distance electron transfer in organic molecules, *Science*, 240:440–447, 1988; J. Mattay (ed.), *Electron Transfer I*, 1994; S. L. Mayo et al., Long-range electron transfer in heme proteins, *Science*, 233:948–952, 1986; R. A. Scott, A. G. Mauk, and H. B. Gray, Experimental approaches to studying biological electron transfer, *J. Chem. Educ.*, 62:932–938, 1985.

Electron tube

A device in which electrons can travel through a sealed chamber containing at least two electrodes and gas at a very low pressure. The gas pressure usually ranges from about 10^{-6} to 10^{-9} atm (10^{-1} to 10^{-4} pascal). At the low extreme of this pressure range, electron tubes are sometimes referred to as vacuum tubes, and at the high extreme as gas tubes. See GAS TUBE; VACUUM TUBE.

Electron emission. At least one of the electrodes must emit electrons, and at least one must collect electrons. The emitting electrode, the cathode, may emit electrons through one or more of four mechanisms: thermionic or primary emission, secondary emission, photoelectric emission, or field emission. Electrons must acquire more energy than they have in the conduction band of a metal in order to escape from the surface of a metal. They acquire this energy, respectively, in the four mechanisms listed above, from heat, electron or ion impact, a photon impact, or an external electric field. Photoelectric emission, first observed by H. Hertz (1887), is used in light-sensing devices, often in combination with secondary electron multiplication to amplify the current. Secondary emission, sometimes in combination with thermionic emission, plays an important role in magnetrons invented by A. W. Hull (1921) and in crossed-field amplifiers. Field emission is used in some experimental amplifiers, flat-panel display devices, and x-ray tubes, but by far the most common type of emitting electrode used in electron tubes, including cathode-ray tubes, is the thermionic cathode. See CATHODE-RAY TUBE; FIELD EMISSION; FLAT-PANEL DISPLAY DEVICE; MAGNETRON; PHOTOEMISSION; PHOTOMULTIPLIER; PHOTOTUBE; SECONDARY EMISSION; X-RAY TUBE.

Thermionic emission was first observed by T. A. Edison (1883) as “negative electricity” emitted from an incandescent filament (the Edison effect). J. J. Thompson (1897) identified the particle emitted from the heated filament as one of very small mass (an electron). O. W. Richardson (1902), M. von Laue (1918), and S. Dushman (1923) derived equations that describe the phenomenon. Pure metals such as tungsten and tantalum can reach temperatures sufficient to produce large currents (about 1 A/cm²) through thermionic emission. Monolayers of elements such as thorium can reduce the filament temperature required for thermionic emission to about

1650°C (3000°F) by reducing the electron energy required of an electron to escape the surface (the work function). Tungsten or nickel coated with mixtures of barium, calcium, and strontium oxides discovered by A. Wehnelt (1903) or porous tungsten matrices impregnated with barium-calcium aluminate (the Philips cathode) can emit large currents at even lower temperatures (800–1150°C or 1500–2100°F). See THERMIONIC EMISSION.

Diodes. A diode is a two-electrode tube, with a cathode and a collecting electrode. Edison must have built a diode in order to discover thermionic emission. A. Fleming (1904) developed the first thermionic diode using an oxide cathode. Because the collecting electrode is usually operated at a positive potential with respect to the cathode in order to collect much of the available electron current from the cathode, it is called an anode. Even so, because of the thermal energy of thermionic electrons, the anode can collect some electrons when it has a slightly negative potential. In a diode, with a good vacuum, the electric field of the electrons in the space between the cathode and the anode just cancels the electric field produced by the anode at the cathode surface. Thus the number of electrons in this space must increase in proportion to anode voltage while their kinetic energy also increases in proportion to the voltage. The electron velocity is proportional to the square root of the kinetic energy, and because current is proportional to density times velocity, the current is proportional to the 3/2 power of the voltage. This is referred to as the space-charge-limited current. The proportionality constant is an invariant of the diode design and is known as the perveance. In gas diodes, ion neutralization of the electron space charge permits the maximum cathode current (determined by the supply voltage and the load resistance) to flow even at very low anode voltages. See CHILD-LANGMUIR LAW.

Diodes were used in great numbers as rectifiers until about 1950 or 1960. By that time, most had been replaced by silicon rectifiers, often used in series, to replace even very high voltage vacuum rectifiers. I. Langmuir (1923) and K. Blodgett (1924) calculated the space-charge-limited current and the potentials between electrodes in diodes with circular cylindrical and spherical electrodes. This theory remains important because it is used as a starting point in the design of highly convergent electron guns for very high power microwave tubes. See DIODE; SEMICONDUCTOR RECTIFIER.

Triodes. L. DeForest (1906) added a third electrode to a diode in order to control the current flow from cathode to anode. This third electrode, the grid, took the form of a fairly open array or mesh made of wires with a diameter small compared to their spacing. In this geometry, much of the electric field from the anode terminates on the grid, and the field from the grid that terminates on the cathode exerts a primary influence on the space-charge current that flows to or through the grid. When the grid is at a negative potential with respect to the cathode, current flows due to the anode field that leaks through the grid, but

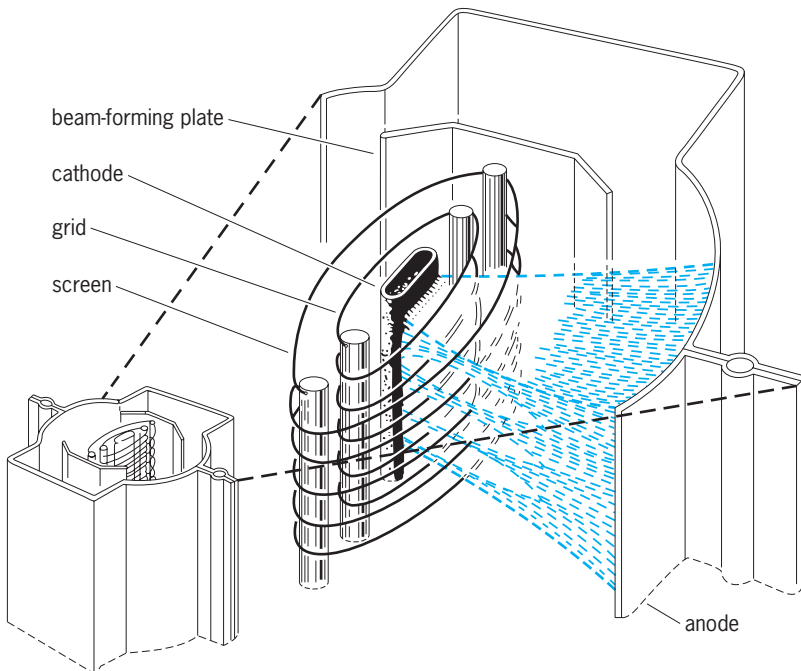


Fig. 1. Cutaway view of electrode arrangement in a beam-power tetrode. (*General Electric*)

the grid can collect no current. When the grid and anode are both positive, much more current flows and divides between the grid and anode.

By 1915, triode amplifiers and oscillators dominated the development of communications equipment. Gas-filled triodes, often called thyratrons and invented by Hull (1929), like gas-filled diodes, conduct much more current with much lower voltage drop than high-vacuum triodes because of space-charge neutralization by positive ions. The grid can start current flow and ionization, but cannot stop current flow until the ions are removed by reversing the polarity of the anode voltage.

Tetrodes and pentodes. Unfortunately, in triode amplifiers at high frequencies, the capacitance between the anode and grid electrodes, in combination with typical grid circuit reactances, can cause positive feedback, regeneration, or oscillation unless circuits that provide compensating negative feedback are used. For this reason W. Schottky (1919) invented the tetrode, which has a second or screen grid between the first or control grid and the anode. This grid was operated at a constant positive potential and effectively shielded the control grid from the anode. At large signal levels, it also created a problem by collecting secondary electrons emitted from the anode as a result of primary electron impacts when the instantaneous voltage on the anode was less than the screen grid voltage.

This problem was dealt with in two ways. G. Jobst and D. H. Tellegen (1926) introduced the pentode, which has a third very open suppressor grid between the screen grid and the anode. It was connected to the cathode. This created an electric field which returned secondary electrons to the anode. A more elegant solution to the secondary electron

problem was provided in the beam-power tetrode (**Fig. 1**). In these tetrodes the anode was placed far enough from the screen grid that the charge of the electrons traveling between the screen grid and anode actually depressed the potential in the space between the screen and anode enough to return secondary electrons to the anode. This solution was analyzed by C. E. Fay, A. L. Samuel, and W. Shockley (1938). Beam-power tetrodes were popular in high-power transmitters and audio amplifiers because, in such applications, screen-grid-current interception and consequent power dissipation are problems. They are still popular in amateur radio equipment and high-power amplifiers used by musicians.

Replacement with solid-state devices. During 1950–1970, advances in solid-state diode and transistor technology, together with the development of small- and large-scale integrated circuits, resulted in the replacement of most low-power tubes with solid-state devices. There was really no contest because solid-state devices operate at much lower power input and at a much lower impedance level (lower voltage and high current) and hence provide higher bandwidth as well as efficiency. See INTEGRATED CIRCUITS; SEMICONDUCTOR DIODE; TRANSISTOR.

High-power tubes. However, for applications in which power serves a useful purpose (for example, transmitters), particularly at high frequencies, the outcome has been quite different. Because no electronic amplifier is 100% efficient, a high-power amplifier must dissipate power as heat. If the amplifier can run hot, it will be small and lightweight. If the temperature is limited to a low value, as it is when using most transistors, the amplifier will be large and heavy. Most of this size and weight will be in the metal heat sinks attached to the transistors.

The tubes and transistors discussed so far act as valves that control the flow of a current to a load. The potential energy of the current is derived from a direct-current power source. There is another, rather different class of electron tubes, most of which are referred to as microwave tubes, in which electrons are accelerated to a velocity at which they have a kinetic energy that is equivalent to the full voltage of the power supply that was used to accelerate them. If these electrons are bunched periodically in time, they can be made to give up their energy to the electric field in a gap or gaps in a very high frequency or microwave circuit. Microwave tubes include the inductive output tube, the klystron, traveling-wave tubes, crossed-field devices, and cyclotron-resonance devices.

There are no solid-state analogs to such tubes because the current flow in solid-state devices is too collisional. In electron tubes the electron density is very low, about 10^9 – 10^{11} electrons per cubic centimeter. At these densities, electrons do not collide with one another. Instead they move coherently and acquire a kinetic energy equal to the time integral of the scalar product of the electric field force and the velocity. As a result, in higher-voltage devices these electrons travel very rapidly (0.3 the velocity of light at 30 kV). Gap dimensions can be

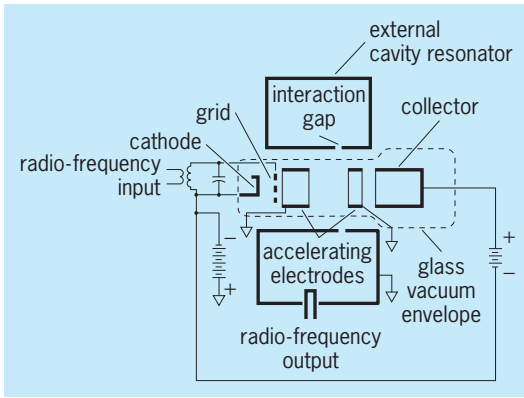


Fig. 2. Inductive output amplifier, showing separation of electron energy extraction from electron collection. (After R. S. Symons, *Tubes: Still vital after all these years, IEEE Spectrum*, 35(4):52-63, April 1998)

large for small transit-angle gaps. In a semiconductor, there are about 10^{19} carriers per cubic centimeter which have a great deal of thermal motion upon which is superimposed a low drift velocity proportional to the electric field. The carriers frequently collide with, and give up their energy to, the crystal lattice.

Inductive output tubes. In the inductive output tube (Fig. 2), invented by A. V. Haeff (1939), an electron beam is amplitude-modulated with a grid and then accelerated through a hole in the first accelerating electrode to form the high-velocity beam of electrons that passes through a gap in the center conductor of the coaxial external cavity resonator and into the collector. Inductive output tubes are used in many

television transmitters operating between 470 and 900 MHz.

Klystrons. The klystron (Fig. 3), invented by R. Varian and S. Varian (1939), has a similar output cavity and collector, but has a beam which is first accelerated in a diode electron gun and then velocity-modulated in another reentrant cavity gap. Fast electrons overtake slowed electrons and yield an intensity-modulated beam by the time the electrons reach the output cavity. Additional cavities may be interposed between the input and output cavities (Fig. 3) to provide very high gain (often as high as 60 dB). See KLYSTRON.

Traveling-wave tubes. In traveling-wave tubes (Fig. 4), invented by R. Kompfner (1946), a high-velocity electron beam is velocity-modulated by, and gives up its energy to, periodically loaded or helical waveguides which slow the electromagnetic wave to a velocity nearly equal to that of the electron beam. Again very high gain is possible. See TRAVELING-WAVE TUBE.

Properties of linear-beam tubes. In all these linear-beam tubes, the electron beams are usually formed by highly convergent electron guns designed using techniques first devised by J. R. Pierce (1940). Computer simulation is now used to refine such designs. Steady magnetic fields provided by either electromagnets or permanent magnets are used to guide the electrons through the radio-frequency circuits. When the electrons leave these focusing fields, they repel each other due to their electric charge, and in addition, they cross radial magnetic field lines which also deflect them outward. The electron collectors on linear-beam tubes may have collection areas as much

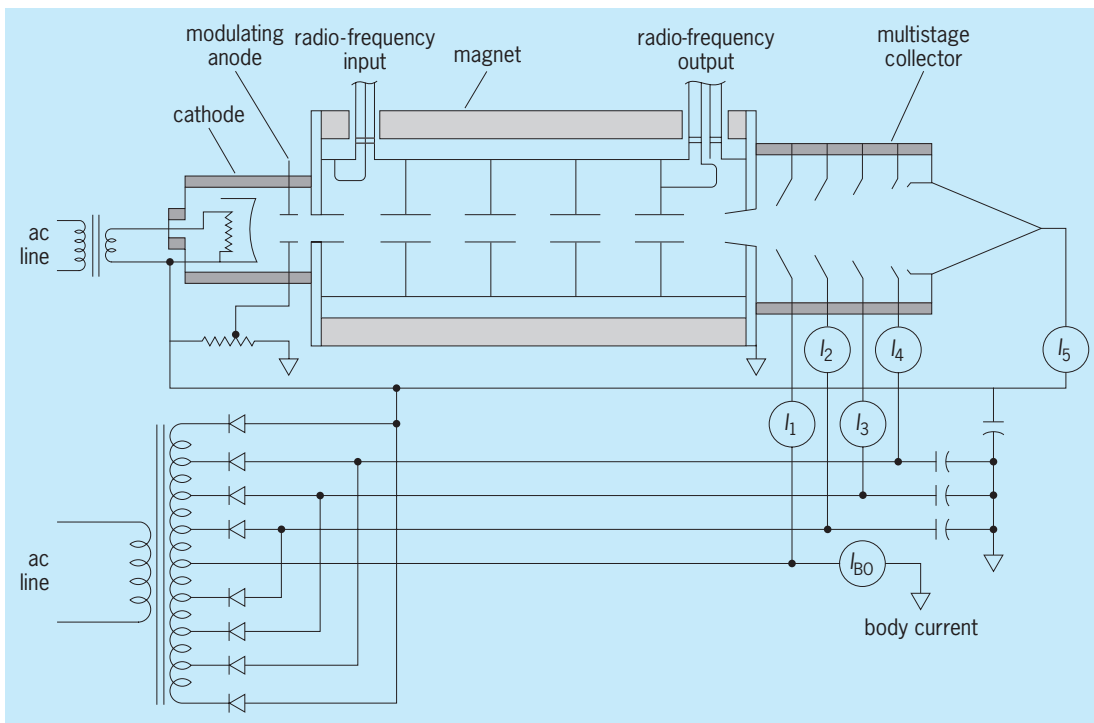


Fig. 3. High-efficiency, depressed-collector-klystron amplifier. (After D. Christiansen, ed., *Electronic Engineers' Handbook*, 4th ed., McGraw-Hill, 1996)

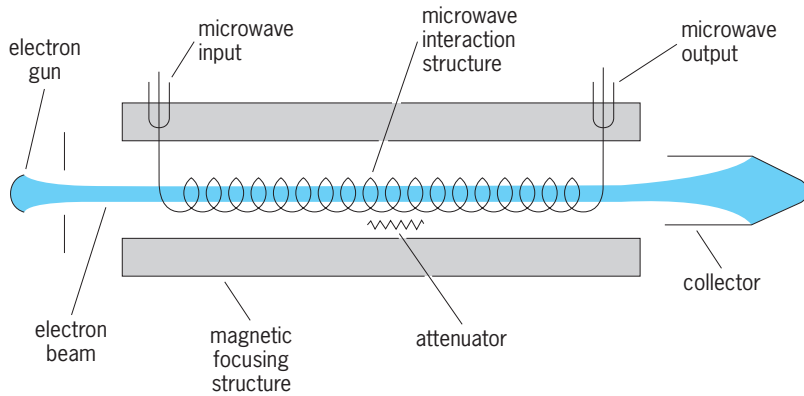


Fig. 4. Basic elements of a typical traveling-wave tube. (After D. Christiansen, ed., *Electronic Engineers' Handbook*, 4th ed., McGraw-Hill, 1996)

as 10,000 times the cross-sectional area of the electron beam, and continuous power as high as 1 MW has been achieved at frequencies as high as 10,000 MHz from such tubes. On klystrons and traveling-wave tubes, sometimes several collectors at different potentials are used (Fig. 3). These multistage depressed collectors, suggested by C. V. Litton (1940), can sort electrons by energy. Electrons that gave only small amounts of energy to the electric field in the radio-frequency circuit are collected on lower-potential electrodes. Tubes using these collectors are more efficient. R. Symons (1993) realized that an input-output tube with a multistage depressed collector could have the unique property of providing almost constant efficiency over a large range of signal amplitudes because both the current and the average collection voltage of the electrons could be made to go up and down with signal amplitude. Inductive-output tubes, klystrons, and traveling-wave tubes are used in television broadcasting, satellite communications systems, radar, scientific accelerators, medical accelerators used

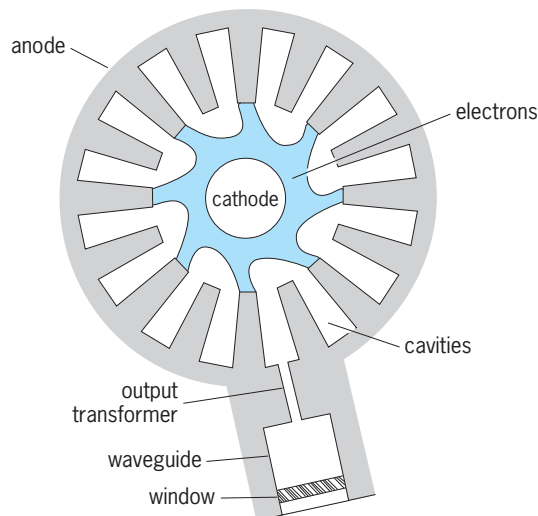


Fig. 5. Conventional microwave structure. (After D. Christiansen, ed., *Electronic Engineers' Handbook*, 4th ed., McGraw-Hill, 1996)

for cancer therapy, and military countermeasures equipment.

Magnetrons. In magnetrons (Fig. 5) and crossed-field amplifiers, electrons circulate about a cylindrical cathode in a radial static electric field and an axial magnetic field. Concentric with, and outside, the cathode is a periodically loaded transmission line that propagates a wave having circumferential electric field components that travel in synchronism with the rotating electron cloud. The electrons follow orbits that allow them to take energy from the radial static electric field and transfer it to the circumferential radio-frequency electric field of the wave on the circuit. Magnetrons are used in huge quantities in household microwave ovens. They and crossed-field amplifiers are also used in ground-based, shipboard, and airborne radars.

Cyclotron-resonance devices. Cyclotron-resonance devices including gyrotrons, gyroklystrons, and gyro-traveling-wave tubes again employ electrons that have been accelerated to the full energy provided by the electrical power supply. The beam is formed in a magnetic field so that it has a great deal of momentum perpendicular to the magnetic field, and the electrons follow helical paths. A radio-frequency electric field perpendicular to the axis of the electron trajectories will modulate the energy of the electrons and hence the relativistic mass and the cyclotron frequency. This azimuthal velocity modulation causes the electrons to draw into rodlike bunches that can give up their energy to a circuit supporting either the same alternating electric field that bunched them (in a gyrotron), or to an alternating electric field in another circuit (in a gyroklystron). Cyclotron-resonance devices are important because, up to a point, space-charge forces tend to create tighter bunches of electrons instead of driving electrons apart as they do in other linear beam tubes (inductive output tubes, klystrons, and traveling-wave tubes). For this reason, cyclotron-resonance devices can be built using very long circuits producing very weak electric fields, and as a result, having very low losses at very high frequencies. Efficient gyrotrons have been built at frequencies as high as several hundred gigahertz and have produced continuous power of hundreds of kilowatts. See GYROTRON; MICROWAVE TUBE.

Robert S. Symons

Bibliography. D. Christiansen (ed.), *Electronic Engineers' Handbook*, 4th ed., 1996; J. W. Gewartowski and H. A. Watson, *Principles of Electron Tubes*, 1965; A. S. Gilmour, Jr., *Microwave Tubes*, 1986; K. R. Spangenberg, *Vacuum Tubes*, 1948.

Electron wake

The pattern of electron density fluctuation and electromagnetic disturbance set up by the passage of a swift ion through condensed matter. In dense media that can sustain well-defined resonance oscillations at a frequency ω_0 , wakes of periodic character will form behind swift charged particles having speed v . The periodicity in space, λ , the distance between

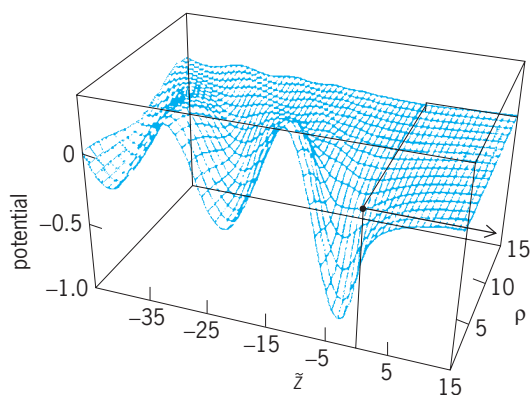


Fig. 1. Theoretical prediction of the scalar electric potential due to the electron wake of an ion penetrating a metallic medium having properties similar to those of aluminum. The ion has a speed three times the Bohr speed of 2.2×10^6 m/s. Coordinates \bar{z} and ρ are measured from the ion position in directions parallel and perpendicular to the ion velocity, respectively, in units of the Bohr radius, $a_0 = 0.0529$ nm. The potential, which does not include that of the bare ion charge, is scaled by the ion charge and is given in units of the potential at a distance of a_0 from a proton. (After P. M. Echenique, R. H. Ritchie, and W. Brandt, *Spatial excitation patterns induced by swift ions in condensed matter*, *Phys. Rev.*, B20:2567–2580, 1979)

troughs of the wake, is given by Eq (1). The oscillations (Fig. 1) trail behind the ion, move with the ion velocity, and have the frequency ω_0 . In addition, close collisions between the ion and electrons of the medium cause electrons to recoil to form the analog of a bow wave ahead of the ion. See RESONANCE (QUANTUM MECHANICS).

In metals such as aluminum, the resonant frequency ω_0 is just the plasma frequency, ω_p , of the conduction electrons. The plasma frequency is given by Eq. (2), in SI units, where n_0 , e , and m^* are the

$$\omega_p = \left(\frac{n_0 e^2}{\epsilon_0 m^*} \right)^{1/2} \quad (2)$$

equilibrium electron density and the charge and effective mass of the electron, respectively, and ϵ_0 is the permittivity of free space. Frictional effects, due to collisions with ions of the medium, give plasma oscillations a finite lifetime. Consequently, the wake is damped at distances of the order of 10λ , which is typically about 20 nanometers behind the ion. Dynamical screening effects limit the wake to lateral distances on the order of v/ω_0 from the track of the particle. The maximum amplitude of the electric potential is given approximately by $Ze\omega_0/(4\epsilon_0 v)$, in SI units, where Ze is the ion charge. Electron wakes may be generated in media such as amorphous carbon and in insulators as well. See ELECTRICAL UNITS AND STANDARDS; FREE-ELECTRON THEORY OF METALS; PLASMA (PHYSICS).

The wake at the position of the guiding ion is of special significance. The electric field there times the ion charge represents the reaction of the medium to the ion and yields the stopping power of the medium

for the ion, that is, the energy loss per unit path length.

Coulomb explosion. When molecular ions are injected into a solid with speeds greater than $v_0 = 2.2 \times 10^6$ m/s, the so-called Bohr speed (the speed of an electron in the ground state of hydrogen according to the Bohr model), the valence electrons are stripped, leaving atomic ions to propagate as clusters of correlated charged particles through the medium. A dicluster is composed of two atomic ions traveling close together at nearly the same velocity. A wake is formed given by a (generally nonlinear) superposition of wakes due to the individual ions of the cluster. The dynamically modified Coulomb repulsion between its constituents causes the cluster, in effect, to explode. A pair of ions traveling with the same initial velocity, and created exactly abreast of one another, will recede rapidly from one another because of the Coulomb force acting on them, until polarization of the medium screens out the repulsive force at separations on the order of v/ω_0 . At the same time, the center of mass of the pair will continue to move in the original direction. In most cases of interest, the initial speed is much greater than the relative speed of recession acquired by the pair. In typical experiments, the cluster-particle interaction probes the slope of the cluster wake potential near the origin, because the foils used in most experiments are so thin that the separation of the ions after traveling through the foil does not greatly exceed their initial separation. An important effect of the wake interaction is to cause the cluster to lose energy at a faster rate than would its isolated constituents traveling at the same speed, because the wake field of a given ion in a cluster acts, in most experiments, to retard the other ions of the cluster.

Measurements of the angular deflections and energy losses of protons resulting from diclusters formed from swift $(\text{HeH})^+$ or $(\text{OH})^+$ ions bombarding thin foils yield angular distributions (Fig. 2) that have a circular character due to the action of a Coulomb explosion. Such distributions generally have a large peaked region on the perimeter due to the trailing protons that are focused by the wake of the other ion in the Coulomb explosion. There is also a much smaller peak due to protons that lead the ion. Such experiments have been important in establishing the structure of molecular ions that were formerly not well known.

In other work with diclusters, the oscillatory character of the wake is vividly displayed in a two-foil experiment. A cluster enters the first foil and, after passing through a vacuum separating the carbon foils, enters the second one. The trailing ion experiences the wake force of the leading one. The dependence of the yield of secondary electrons from a final target on the distance between carbon foils shows the characteristic oscillatory behavior. See COULOMB EXPLOSION.

Nonlinear effects. It has been proposed that a moving positron (and possibly an electron) may become self-bound in the wake that it creates in matter. This composite entity might be termed a wakeon.

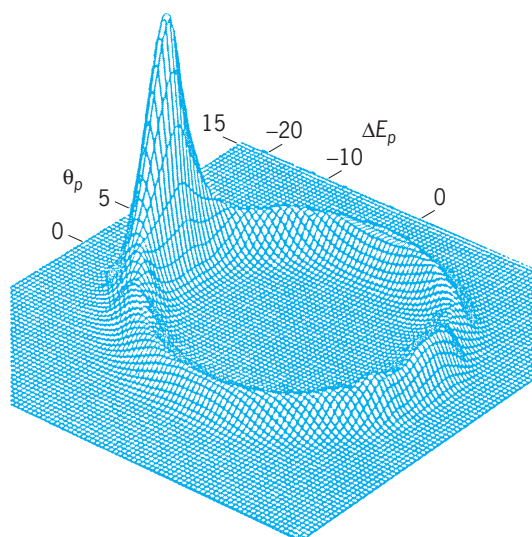


Fig. 2. Distribution of protons emerging from a carbon foil of 15.3 nm thickness that is bombarded by OH^+ molecular ions having an energy of 11.2 MeV. The distribution is given with respect to energy loss ΔE_p in kiloelectronvolts and angular deflection θ_p in milliradians measured with respect to the direction of the incident ions. The measurements were made in coincidence with the detection of the ion O^{6+} . (After A. Breskin et al., *Coulomb explosion of 11.2 OH^+ ion in carbon foils*, *Nucl. Instrum. Meth.*, 170:93–97, 1980)

A theoretical formulation, allowing for the nonlinear response of the medium in the framework of soliton theory, has been invoked to determine the effective mass of a positron in different metals. It is found that the effective mass is appreciably larger than the mass in vacuum and appears to agree with experimentally determined values. See CHARGED PARTICLE BEAMS; POSITRON; SOLITON.

Rufus H. Ritchie

Bibliography. P. M. Echenique, F. Flores, and R. H. Ritchie, Dynamic screening of ions in condensed matter, *Sol. State Phys.*, 43:229–308, 1990.

Electronegativity

Electronegativity, according to L. Pauling, is “the power of an atom in a molecule to attract electrons to itself.” With the concept of electronegativity, a vast number of observations of chemical and physical properties have been either correlated or predicted. Quantitative definitions and scales of electronegativity have been based not on electron distribution itself but on properties which were assumed to reflect electronegativity.

The electronegativity of an element depends upon its valence state and thus is not an invariant atomic property. As an example, the electron-withdrawing ability of an sp_n hybrid orbital centered on carbon and directed toward hydrogen increases as the percentage of s character in the orbital increases in the series ethane < ethylene < acetylene. Thus, according to this concept of orbital electronegativity, each element exhibits a range of electronegativity values. In the following paragraphs, a few scales of electronegativity will be discussed.

The original scale, proposed by Pauling in 1932, is

based upon the difference Δ between the energy of the A—B bond in the compound AB_n and the mean of the energies of the homopolar bonds A—A and B—B, as in Eq. (1). The A—B bond energy exceeds

$$\Delta = E(\text{A—B}) - \frac{E(\text{A—A}) + E(\text{B—B})}{2} \quad (1)$$

the arithmetic means of the A—A and B—B bonds to an increasing extent as the elements A and B diverge in electron-attracting ability. The difference in electronegativities of bonded atoms is proportional to the square root of the energy difference Δ , as in Eq. (2). The proportionality factor 0.208 converts

$$\chi_A - \chi_B = 0.208\sqrt{\Delta} \quad (2)$$

the units of energy from kilocalories to electronvolts. After the electronegativity of one element is arbitrarily assigned, other values of electronegativity can be calculated from thermochemical data. Values for selected elements in common oxidation states are presented in the **table**. Electronegativity increases with increasing oxidation state. For example, $\chi_{\text{Sn(II)}} = 1.80$ and $\chi_{\text{Sn(IV)}} = 1.96$.

R. S. Mulliken proposed that the electronegativity of an element is given by the average of the valence-state ionization potential and electron affinity: $\chi_M = (IP_V + EA_V)/2$. The quantities IP_V and EA_V are not observable properties of the ground state of an atom, but are energies for a hypothetical state of the isolated atom having the same electronic configuration (hybridization, electron-electron interaction, and so forth) as the atom in the molecule. The Mulliken approach has a sound theoretical basis, is consistent with Pauling's original definition, and gives orbital electronegativities, not invariant atomic electronegativities. Valence-state ionization potentials and electron affinities have been calculated from the equations $IP_V = IP_g + P^+ - P^0$ and $EA_V = EA_g + P^0 - P^-$, where IP_g and EA_g are ground-state potentials and

Average electronegativities from thermochemical data

Element	Value	Element	Value
H	2.20	Al	1.61
Li	0.98	Ga	1.81
Na	0.93	In	1.78
K	0.82	Tl	2.04
Rb	0.82	C	2.55
Cs	0.79	Si	1.90
Be	1.57	Ge	2.01
Mg	1.31	Sn	1.96
Ca	1.00	Pb	2.33
Sr	0.95	N	3.04
Ba	0.89	P	2.19
Sc	1.36	As	2.18
Ti	1.54	Sb	2.05
V	1.63	Bi	2.02
Cr	1.66	O	3.44
Mn	1.55	S	2.58
Fe	1.83	Se	2.55
Co	1.88	F	3.98
Ni	1.91	Cl	3.16
Cu	1.90	Br	2.96
Zn	1.65	I	2.66
B	2.04		

affinities, respectively, and P^+ , P^0 , and P^- are promotion energies of the positive ion, atom, and negative ion, respectively. The calculation of d -orbital electronegativity by the Mulliken method has not been accomplished for nontransition elements due to the lack of spectroscopic data. Since electronegativity is a sensitive function of d -orbital hybridization and since the extent of d -orbital participation generally cannot be ascertained quantitatively, the calculations of electronegativities for the heavier elements are limited.

The energy of an ion relative to the neutral atom can be expressed as a power series $E = aq + bq^2 + cq^3 + dq^4$, where q is the formal oxidation state or ionic charge for a particular state of ionization. Z. R. P. Iczkowaski and J. L. Margrave defined the electronegativity of a neutral atom as the derivative, $\chi_{IM} = (dE/dq)_{q=0}$. As a fairly good approximation, the last two terms in the above power series can be dropped, giving Eq. (3). The units of χ_{IM} are

$$\chi_{IM} = \frac{dE}{dq} = \frac{d(aq + bq^2)}{dq} = a + 2bq \quad (3)$$

energy/electron and the magnitudes are the same, in accordance with theory, as those of Mulliken if E is evaluated only from the electron affinity and the first ionization potential. The quantity a is the Mulliken electronegativity for a neutral atom, and electronegativity is shown by Eq. (3) to increase linearly with increasing positive charge. By using IP_V and EA_V values, H. H. Jaffé and coworkers calculated the electronegativities of certain vacant, singly occupied, and doubly occupied orbitals.

Electronegativity was defined by A. L. Allred and E. G. Rochow as the force of attraction between a nucleus and an electron from a bonded atom. The electrostatic force was calculated simply from the effective nuclear charge and the atomic radius, as in Eq. (4).

$$\chi = \frac{0.359Z_{\text{eff}}}{r^2} + 0.744 \quad (4)$$

A quantum-defect electronegativity scale has been developed from potentials based on atomic spectral data, and a nonempirical scale has been calculated by an ab initio method using floating gaussian orbitals.

Other methods for calculating electronegativities utilize such observables as bond-stretching force constants, electrostatic potentials, spectra, and covalent radii. The fact that the various scales of electronegativity have different dimensions (energy^{1/2}, energy/electron, force, potential, and so forth) or no dimension reflects the widespread results of differences in electron-attracting ability. The measurement of electronegativities involves observations of properties dependent upon electron distribution. Close agreement of electronegativity values obtained from measurements of several diverse properties lends confidence and utility to the concept. A. Louis Allred

Bibliography. F. A. Cotton and G. Wilkenson, *Advanced Inorganic Chemistry*, 6th ed., 1999; L.

Pauling, *Nature of the Chemical Bond*, 3d ed., 1960; W. W. Porterfield, *Inorganic Chemistry: A Unified Approach*, 2d ed., 1993; A. G. Sharpe, *Inorganic Chemistry*, 3d ed., 1992.

Electronic display

An electronic component used to convert electrical signals into visual imagery in real time suitable for direct interpretation by a human operator. It serves as the visual interface between human and machine. The visual imagery is processed, composed, and optimized for easy interpretation and minimum reading error. The electronic display is dynamic in that it presents information within a fraction of a second from the time received and continuously holds that information, using refresh or memory techniques, until new information is received. The image is created by electronically making a pattern from a visual contrast in luminance between (1) individual electrically alterable picture elements (pixels) in a matrix array of pixels in flat-panel displays (FPDs) or (2) electrically excited and nonexcited areas in a phosphor film in cathode-ray tubes (CRTs). High-information-content (HIC) displays are those displays that have a sufficient number of pixels (75,000 to 2,000,000) to show standard or high-definition television images, or comparable computer images.

The use of electronic displays for presentation of graphs, symbols, alphanumeric, and still and video pictures has doubled every several years, in parallel with the rapid expansion of microelectronics. Electronic displays have largely replaced traditional mechanical devices, counters, galvanometers, and, to a degree, hardcopy (paper) means for presenting information. This change is due to the increased use of computers, microprocessors, inexpensive large-scale integration (LSI) electronics, and digital mass memories. The success of the digital watch, handheld calculator, and personal computer is directly attributable to the availability of inexpensive LSI electronics and electronic displays. See CALCULATORS; COMPUTER; COMPUTER GRAPHICS; COMPUTER STORAGE TECHNOLOGY; INTEGRATED CIRCUITS; MICROCOMPUTER; MICROPROCESSOR; WATCH.

Electronic transducers and four-digit (or more) FPDs have been used to replace the galvanometer movement, thermometer scale, barometer movement, and other forms of scientific instrumentation. Large signs, arrival and departure announcements, and scoreboards also use electronic means to portray changing messages and data. The major electronic display application is in home color television.

The computer terminal using a CRT or FPD is the most important industrial application of electronic displays. The personal computer terminal with a microprocessor and mass memory serve to replace the office paper, typewriter, and file cabinet. See COMPUTER PERIPHERAL DEVICES; WORD PROCESSING.

With advances in high-information-content FPD technology, the electronic display industry went through a dramatic change in the early 1990s. This

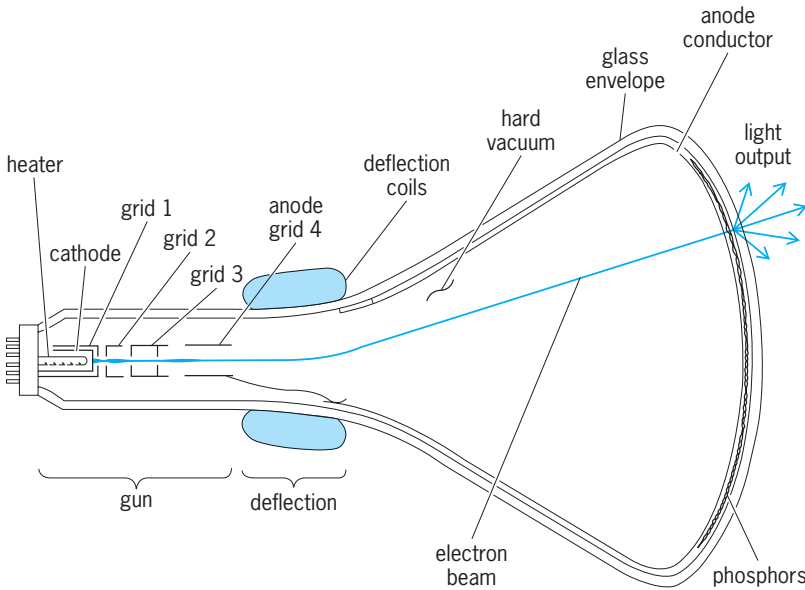


Fig. 1. Cathode-ray tube using electrostatic and magnetic deflection.

was primarily due to breakthroughs and manufacturing advances in liquid-crystal displays (LCDs). Color active-matrix LCDs (AMLCDs) have better performance than color CRTs for video and computer graphics and are thin and portable. However, the AMLCDs are several times more expensive than CRTs, so their use remains restricted to applications such as personal computers and television receivers where the CRT will not fit. See LIQUID CRYSTALS.

Cathode-ray tube. The primary applications of the CRT are in home entertainment television, computer monitors, scientific and electrical engineering oscilloscopes, radar display, and alphanumeric and graphic electronic displays. See OSCILLOSCOPE; RADAR; TELEVISION.

The CRT (Fig. 1) has a viewing screen coated with a phosphor which emits light when struck with a beam of high-energy electrons. The electrons are emitted from the cathode at the rear of the tube in a beam that is focused electrostatically (or magnetically) to a dot or spot on the phosphor screen and positioned in horizontal and vertical coordinates

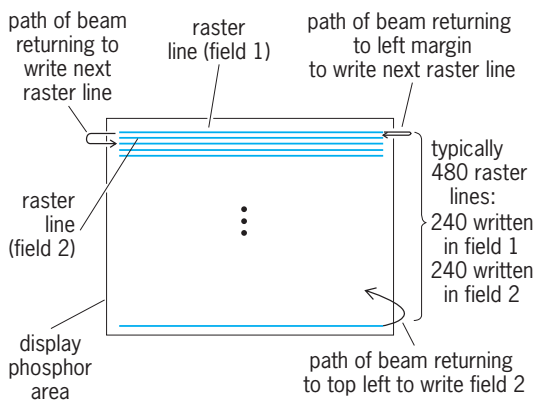


Fig. 2. Cathode-ray tube raster for television, using interlace.

by magnetic (or electrostatic) forces. The cathode grids and electron-focusing lenses are incorporated into a subassembly called the gun. The beam is accelerated toward the phosphor by a high voltage of approximately 20 kV between the cathode and the anode.

Imagery is created on the screen as the CRT raster (Fig. 2) is traced out. The video signal is applied, after amplification, directly to grid 2 of the gun, and controls the amplitude of the electron beam and thus the luminous output at the display surface. The deflection coils steer the beam to trace out the raster. The horizontal-scan deflection signal causes the beam to trace out the horizontal lines and then fly back for the next line. The vertical-scan deflection signal causes the beam to be stepped down the raster and then re-traced to the top left corner at the beginning of each frame. Two fields are interlaced, with the raster lines of one traced between those of the other. The purpose of the interlace is to minimize flicker in the picture and reduce video bandwidth. Computer CRTs do not use interlace, but continuous scan to prevent flicker. See CATHODE-RAY TUBE; PICTURE TUBE; TELEVISION RECEIVER.

Flat-panel displays. Because of the depth dimension of the CRT, there has been a concentrated effort to develop FPDs. A primary motivating factor has been to achieve a flat high-information-content display which could be hung on a wall or carried in a briefcase. Over the years, the electrical phenomena most extensively developed for FPDs have been gas discharge (plasma), electroluminescence, light-emitting diode, cathodoluminescence, and liquid crystallinity. The cost of FPDs is higher than CRTs on a character-per basis for HIC displays. See CATHODOLUMINESCENCE; ELECTROLUMINESCENCE; LIGHT-EMITTING DIODE.

Before the 1990s, cost and performance limitations restricted the use of FPDs to specialized applications. With the advent of two major breakthroughs in the 1980s, LCDs emerged as the leading FPD. These two breakthroughs were the invention of the supertwisted nematic (STN) LCD, often referred to as passive-matrix LCD (PMLCD), and the evolution of a manufacturable AMLCD. The STN inventions improve LCDs so that they can be made as a HIC FPD because of improvements in the matrix addressing capability. Both of these types are made in monochrome and full color. The AMLCD is fast enough for video and for multimedia displays, comparable to a CRT. The PMLCD is fast enough for computer use but not video; however, PMLCDs are about half the cost of AMLCDs.

Matrix addressing. Flat-panel displays are typically matrix-addressed. A row is enabled to accept display information in parallel via the column lines. The electronics commutate through the rows, serving the same purpose as the vertical deflection amplifier of the CRT. The column data are shifted into the column drivers, and at the proper time applied to the column lines.

The thickness of the flat-panel matrix addressing electrodes (Fig. 3) is approximately 0.1 in. (2.5 mm).

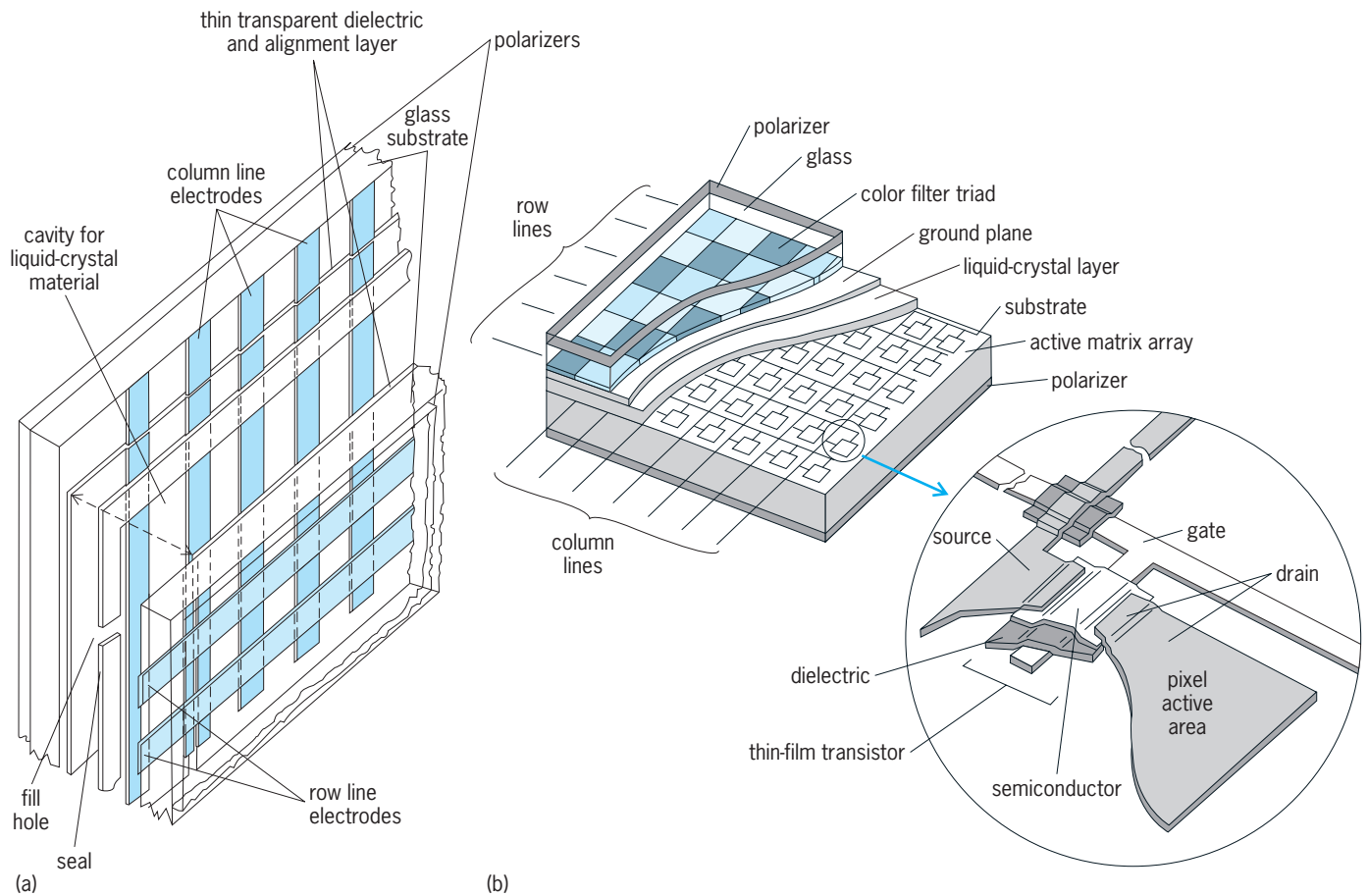


Fig. 3. Exploded sections of liquid-crystal displays (LCDs). (a) Passive-matrix LCD (PMLCD). Column-line and row-line electrodes are made of transparent indium–tin oxide. Liquid-crystal material aligns to the alignment layer (which has been mechanically rubbed with cotton cloth or similar means to impart microgrooves on the surface) during filling steps. Alignment rub directions are rotated 90° for twisted nematic LCDs and $180\text{--}270^\circ$ for supertwisted nematic LCDs. (b) Active-matrix LCD (AMLCD).

The row and column lines are spaced at 80 lines per inch (31.5 lines per centimeter). The intersection of each row line with each column line defines a pixel. A pixel is a picture element, and denotes the smallest addressable element for spatial information in an electronic display. To control the pixels, the row and column electronic drivers are attached at the edges of the glass, using matrix-addressing techniques.

Comparison of technologies. Each of the flat-panel technologies (Table 1) has unique advantages. The more important considerations include luminous efficiency, addressability, duty factor, gray scale, color, and cost. The capabilities in each area have evolved to the point of cost-effective applications to commercial and military products. For the vectorgraph and video display categories, liquid-crystal, electroluminescent, and gas-discharge or plasma displays have emerged as the most cost-effective approaches.

AMLCD technology. LCDs are manufactured in many configurations. The highest-performing and most widely used LCD is the AMLCD, where the active-matrix format is a thin-film (field-effect) transistor (TFT) [Fig. 3b]. The most common semiconductor used in the TFT is amorphous silicon; next is polysilicon; single-crystal silicon and cadmium selenide are also used. For medium-to-large displays, the semicon-

ductor is almost universally amorphous silicon with a mobility typically of 0.5–1.0. This mobility, although low by semiconductor standards, is fast enough for electronic displays and, more importantly, the TFT can be made on an inexpensive borosilicate glass substrate. The polysilicon TFT is typically annealed amorphous silicon on quartz and is used for small (less than 2 in. or 5 cm in diagonal) displays. Single-crystal silicon, which has the highest mobility, is also used in small displays.

LCD applications. Applications of LCDs include portable computing devices and memory aids; electronic games; moving map devices in conjunction with global positioning satellites (GPS) for navigation of boats, aircraft, and motored vehicles; educational devices and video for portable handheld devices; and displays for television, entertainment, education, and information dissemination. LCDs are also used in helmet displays and virtual-reality displays. See SATELLITE NAVIGATION SYSTEMS; VIDEO GAMES.

The smaller AMLCDs are used in camcorder viewfinders, helmet displays, and projectors. The medium-size displays are used in all the applications that are viewed directly, such as portable computers or television receivers.

PMLCDs are used with the STN configuration of

TABLE 1. Flat-panel technologies

Technology	Phenomena
Emissive displays	
Gas discharge (GD)	Cathode glow from conducting gaseous discharge
Plasma panel (PDP)	Alternating-current capacitively coupled gas discharge
Light-emitting diode (LED)	Electron injection in a forward-biased <i>pn</i> semiconductor junction
Vacuum fluorescence (VFD)	Electron bombardment of phosphor in hard vacuum under control of a grid
Electroluminescence (ELD)	Electron conduction in polycrystalline phosphors due to high electric field
Flat cathode-ray tube (FCRT)	Electron bombardment of phosphor in hard vacuum under control of a grid or cathode
Field-emitter display (FED)	Electron bombardment of phosphor in hard vacuum under control of row and column electrodes, where the cathode is a field-emitter tip
Organic light-emitting diode (OLED)	Complex organic semiconductor compounds that emit photons when conducting electrons
Nonemissive displays	
Liquid crystallinity	Electrostatic rotation of organic compounds which exhibit liquid crystallinity
Electrochromism	Charging and discharging chemical systems (battery) which exhibit a color change in accordance with Faraday's law
Colloidal suspensions	Electrostatic transport or rotation of light-absorbing particles in a colloidal suspension
Micromechanical display (MMD)	Mirror displays on a silicon chip, deflectable with electrostatic force

LCD in devices where some reduction in performance is acceptable, such as portable computer displays that do not require video imagery. PMLCDs in the basic twisted nematic (TN) mode are used in inexpensive toys, automotive panels, meter displays, and watches.

Display categories. Electronic displays can be categorized into four classifications; pseudoanalog, alphanumeric, vectorgraphic, and video (Table 2). Each classification is defined by natural technical boundaries and cost considerations. The categorization is useful in visualizing the extent to which electronic displays are used.

Special-purpose displays. The above categorization emphasizes direct-view-type electronic displays, which are of primary interest to industry. There are other special-purpose displays used in very sophisticated applications.

Projection display. In the projection display an image is generated on a high-brightness CRT or similar electronic image generator, and then optically projected

onto a larger screen. To illuminate screens larger than approximately 3×4 ft (1×1.3 m) and in color, multiple CRTs or light valves are used. The light valve is any direct-view display optimized for reflecting or transmitting the image, with an independent collimated light source for projection purposes. Light valves create images to control the reflection of light to be projected onto the screen. This permits powerful light sources such as xenon lamps to be used independent of the image-generating technique. Oil-film light valves, micromechanical displays, and liquid-crystal light valves are examples of devices used in large command and control and theater-size electronic display presentations.

Three-dimensional imagery. True three-dimensional imagery can be created electronically by several techniques. One technique requires goggles using PLZT (lead zirconate titanate modified with lanthanum) electrooptical ceramic eyepieces or liquid-crystal shutters over each eye. The eyepieces are electronically controllable shutters, with the ability to be reversibly switched from open to closed in microseconds. Two images from two television cameras placed to obtain the desired stereoscopic effects are electronically interlaced and displayed on one CRT monitor. Each image is sequentially displayed from each camera at television video rates. The goggles are synchronized to be opened and closed so that only the right eye sees the image from the right camera and the left eye sees the image from the left camera. The viewer sees true three-dimensional perspective while looking at the CRT monitor through the goggles. See ELECTROOPTICS.

Three-dimensional displays are made integral to and self-contained in a helmet and are called virtual-reality displays. A separate display is used for each eye to show the stereo pair of images. The display is typically either a miniature CRT or a miniature LCD, with appropriate optics to present the image at the proper focal length in front of the eye. In a virtual-reality display, the user cannot see beyond the displays as the intent is to completely control the images seen by the viewer. The viewer is visually engulfed in a synthetic visual environment. See VIRTUAL REALITY.

Helmet-mounted and heads-up displays. Helmet-mounted displays (sometimes called visually coupled displays) and heads-up displays are used in aircraft, automobiles, and toys. In both of these displays the image is projected, usually from a CRT onto a combining glass, and collimated to be in focus at infinity. The combining glass screen is designed to reflect the display imagery to the viewer, usually at selected wavelengths of light, while being sufficiently transmissive for the viewer to see the scene beyond. The primary application for the heads-up display is to present critical aircraft performance, such as speed and altitude, on a combining glass at the wind-screen for pilot monitoring while permitting the pilot to look out the window for other aircraft or the runway. The primary application for the helmet-mounted display is to present, on a combining glass within the visor of the helmet of a helicopter gunner, primary information for directing firepower. The angular direction of

TABLE 2. Electronic display spectrum of applications

Classification	Characteristics	Applications	Electronic technologies
Pseudoanalog	Dedicated arrangement of discrete pixels used to present analog or qualitative information	Meterlike presentations, go/no-go messages, legends and alerts, analog-like (watch) dial	Gas discharge, light-emitting diodes, liquid crystal, incandescent lamps
Alphanumeric	Dedicated alphanumeric pixel font of normally less than 480 characters; most common is 4- and 8-character numeric displays	Digital watches, calculators, digital multimeters, message terminals, games	Liquid crystal, light-emitting diodes, vacuum fluorescent, gas discharge, incandescent lamps
Vectorgraphic	Large orthogonal uniform array of pixels which are addressable at medium to high speeds; normally, monochromatic with no gray scale; may have memory; normally, over 480 characters and simple graphics	Computer terminals, arrivals and departures, scheduling terminals, weather radar, air-traffic control, games	Cathode-ray tube, plasma panels, gas discharge, vacuum fluorescent, electroluminescence, PMLCD, LED
Video	Large orthogonal array of pixels which are addressed at video rates (30 frames per second); monochromatic with gray scale or full color; standardized raster scan addressing interface, arrays of pixels approximately 240 rows by 320 columns and larger	Entertainment television, graphic arts, earth resources, video repeater, medical electronics, aircraft flight instruments, computer terminals, command and control, games	Cathode-ray tube, plasma panels, electroluminescence, AMLCD, LED

the helmet is sensed and used to control weapons or radar to point in the same direction in which the pilot gunner is looking. See AIRCRAFT INSTRUMENTATION.

Color. Color can be created on a CRT equipped with a shadow mask duplicating quite closely all colors that occur in nature. This is done in the CRT by using three different electron guns and three phosphors in a triad of red, green, and blue on the screen at each pixel. The shadow mask is a metal screen with one hole for each triad of phosphor dots (pixel). It is precisely located and aligned with the phosphor screen. Electron beams from each of three guns are constrained by each shadow mask hole to hit each respective phosphor dot. The gun, shadow-mask holes, and phosphor dots are aligned during manufacturing so that the three beams converge to pass through the single hole (or slit) in the shadow mask and then diverge as the beams emerge with sufficient separation to impact the three different phosphor dots. If all three guns are on simultaneously, the eye, upon close inspection, sees a red, a green, and a blue dot of light at each pixel. However, at a normal viewing distance, the three dots merge together when focused in the retina of the eye, and from the laws of additive color, the pixel appears white. The phosphor itself is whitish, but it emits different colors because of the dopants in the phosphor, and the image is created from the emitted light.

Penetration phosphors are also used to create color on CRT displays to eliminate the need for the shadow mask and extra guns. However, the color is limited and the brightness is low. Normally, two phosphors are placed on the screen in two layers or in microspheres of two layers. The gun and CRT anode are operated in two energy states to produce either a high-energy or a low-energy electron beam switchable in time. The high-energy beam penetrates

the first phosphor layer and is stopped at the second layer. It then excites the second layer to produce its characteristic color. The low-energy beam is stopped by the first phosphor layer and excites it to produce its characteristic color. The two phosphor colors most often used are red and green. Intermediate-energy beams make it possible to fractionally excite both layers to get color combinations of red and green such as yellow and orange. Practical considerations limit this color approach to these four. Full color would require at least three primary colors such as red, green, and blue. Three layers of phosphors are limited in brightness due to practical considerations and have not been commercially available.

Monochromatic color is readily produced by FPD technologies. One color of a plasma panel is normally orange, owing to the neon gas. Full-color plasma panels are available in medium resolution. Other monochromatic colors are feasible. Light-emitting diode (LED) luminance is normally red, yellow, or green. Blue LEDs are available and have high luminous efficiency. Electroluminescence is normally yellow or green, owing to the manganese or copper activator, respectively. Full-color electroluminescent displays have been demonstrated but are not cost-effective in production. In flat-panel as in CRT displays, full color is produced by using a triad of red, green, and blue for each pixel, or sequential flashing at 180 Hz of red, green, or blue at each pixel area.

Full color is very important for entertainment television displays. Most industrial electronic displays do not necessarily need full color, but full color has been used more frequently as color display costs have decreased. In these applications, the color display instrument is usually a CRT using the shadow-mask color technique.

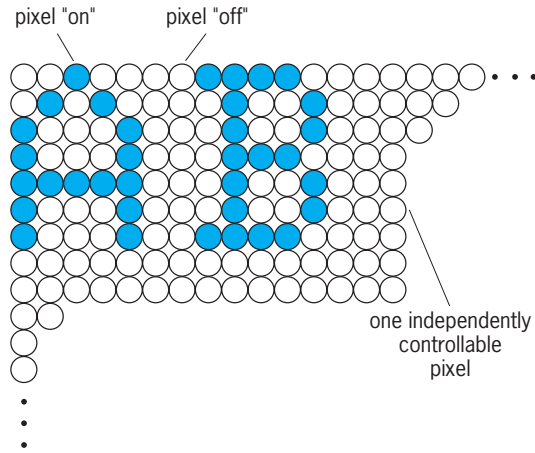


Fig. 4. Pixel array used for creating electronic display images.

Display technique. The essence of electronic displays is based upon the ability to turn on and off individual pixels (Fig. 4). A typical HIC display will have a quarter million pixels in an orthogonal array, each under individual control by the electronics. The pixel resolution is normally just at or below the resolving power of the eye at one minute of arc. Thus, a good-quality picture can be created from a pattern of activated pixels.

The pixel concept for electronic displays has evolved from the modern FPD technologies and digital electronics. It has been extended to the analog-raster-scan CRT in the following way: The electron beam from the gun is deflected magnetically (or electrostatically), so as to sweep across the phosphor and thereby cause a line to luminesce on the face of the CRT. In digitally modulated CRTs, the cathode is modulated by a sine wave as the beam is swept across the face of the CRT. Here, instead of a continuous line, a string of dots results. Each dot corresponds to a pixel. Each pixel is on when the beam density is high, and off when the beam is off. The beam is turned off between pixels. The pixels are refreshed from 30 to 60 times per second on a CRT.

Pixels are created in all the rows of the entire CRT raster in what is called a CRT-digital raster. This approach is commonly used in industrial applications and computer terminals, since it is easily interfaced with digital electronics. Home entertainment television uses an analog-raster-scan approach (Fig. 2).

There are some applications in which a nonraster approach is used to create alphanumeric characters and vectors on CRTs. The electron beam is deflected under control of the deflection amplifiers to stroke out each line of the image. When characters and vectors are generated this way, they are like Lissajous figures, as opposed to (digital or analog) raster characters and vectors. The Lissajous characters and vectors are best suited to large (25-in. or 63.5-cm diagonal) CRTs and where there are numerous vectors, straight lines, and curves. Vectors and curves drawn with the raster technique have stair steps. Lissajous vectors and curves are always smooth and continuous. Lissajous techniques have increasingly yielded

to the digital raster as the cost of digital electronics has improved. See LISSAJOUS FIGURES.

Font. With flat-panel and CRT digital display techniques, alphanumeric character fonts are created by turning on the appropriate pixels in an array. One standard size is a 5×7 array with two pixels between characters and two pixels between rows (Fig. 4). All the letters and numbers can be created on this common array format. Several other combinations of pixels may be used to create the letter A. The viewer soon becomes accustomed to the minor variation. Readers do not read pixels but read letters and words, and therefore the exact detail of the character pixel pattern is of a secondary consideration. In general, the more pixels available in the basic array, the more esthetically pleasing is the character, at the cost of additional electronics to control the extra pixels. The viewing distance is normally far enough so that the pixels blur together.

A very efficient and elegant array has evolved for portraying numeric characters and many letters of the alphabet, called the seven-bar font (Fig. 5). Each bar is a pixel by definition. This font was initially considered crude when compared to the Leroy font and other more esthetic printer fonts. It is now universally accepted. A similar 14 bar font is sometimes used for alphanumeric characters.

Display electronic addressing. The numeric or alphanumeric display electronic drive may be performed in a single LSI metal oxide semiconductor (MOS) chip mounted in a single dual in-line package (DIP) suitable for direct assembly on a printed circuit board. All the timing, logic, memory, resistors, and drivers are contained in a single chip. A four- or seven-bit character code is serially fed into the chip for display (Fig. 6).

A computer terminal (Fig. 7) will incorporate a microprocessor unit and a CRT controller to perform master control and CRT housekeeping tasks. A line of video data is loaded into the shift register for serial drive of the video amplifier at the time a display line starts. The shift register is loaded during flyback

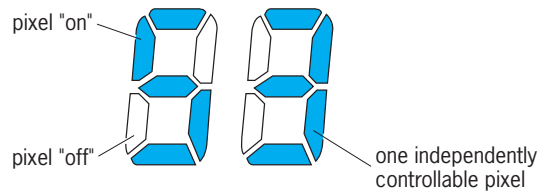


Fig. 5. Seven-bar numeric font.

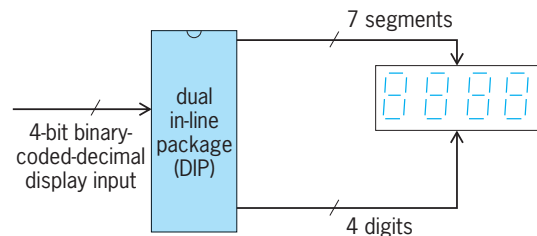


Fig. 6. Block diagram of LSI MOS electronic drive for numeric display.

time from the character generator. The character generator is a decoder transforming the alphanumeric information in coded form such as ASCII (American National Standard Code for Information Interchange) into pixel signals for each raster line. The ASCII code is stored in the display refresh memory. The entire cathode-ray-tube frame is stored in this memory, and is continuously displayed until changed under control of the microprocessor unit. New information can come from the mass memory, the keyboard, or other subsystems on the data bus and address bus, all under control of the microprocessor unit.

The CRT controller performs all the housekeeping tasks for proper CRT operation. These include the vertical- and horizontal-raster synchronization signals, the cursor control, blinking, blanking, interlace, paging, and scrolling. The CRT controller is a single LSI MOS chip, as are most of the other units in the terminal (Fig. 7).

Custom chip sets are available for controlling CRTs and flat-panel displays. The sets include processor, memory, and all programmable functions necessary for self-contained operation. The chip set operates off the main bus and formats the red, green, and blue signals and synchronization signals to the display. HIC FPDs typically come with their own electronics, with a red-green-blue plus sync interface similar to that used in CRTs.

Applications. Because of performance and cost issues, there now are only two major display technologies, CRTs and LCDs. Full-color displays are now expected by users, except for some niche areas such as medical imaging where resolution and gray shades are traditionally more important than color. Display technologies which cannot efficiently produce saturated red, green, and blue colors for electronically controlled additive mixing to produce full color can no longer find applications. Those technologies include electroluminescence, vacuum fluorescence, and field-emitter displays.

The CRT remains the most widely used display technology. It is made in the traditional 3:4 aspect ratio and a 9:16 aspect ratio for high-definition television (HDTV) with resolutions of over 2,000,000 pixels and sizes of up to 38 in. (1 m) in diagonal. For larger screen sizes, smaller projection monochrome red, green, and blue CRTs are used for consumer products up to 60 in. (1.5 m) and industrial products up to 100 in. (2.5 m) in diagonal screen size. The dominant position of the CRT is maintained because of its good performance and relatively low cost. The only limitation is the bulky size of the CRT which makes it impractical to ship and install when exceeding 40 in. (1 m) in diagonal size. CRT projectors using three or more smaller monochrome tubes are practical but more costly and have competition from other technologies in both price and performance.

If a CRT cannot be used due to its depth and volume, an LCD is probably the best alternative. LCDs can be made in almost any resolution and in sizes from 0.5 to 22 in. (1.2 to 55 cm) in diagonal, with possibilities of tiling four panels together. The LCD uses photolithography in manufacturing and is lim-

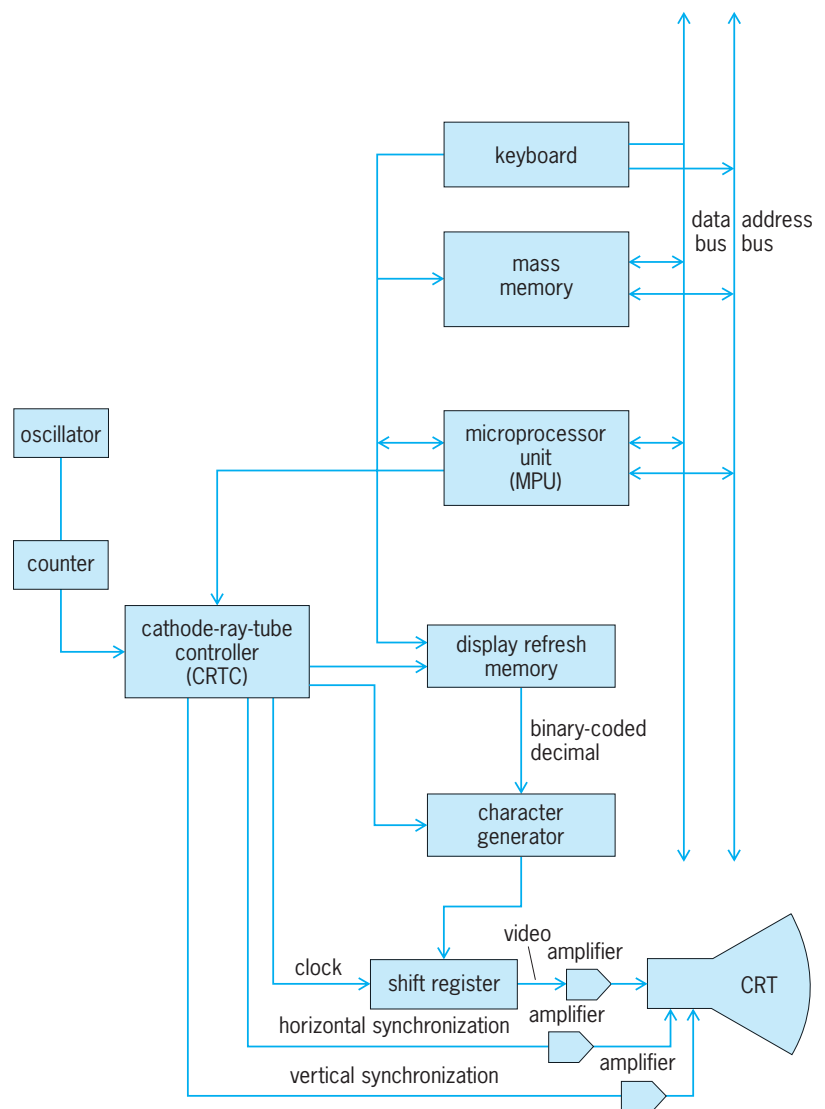


Fig. 7. Block diagram of LSI electronics for a smart computer terminal cathode-ray-tube display.

ited in size and resolution by the limitations in the photolithographic machinery and process. For larger sizes, LCDs are made as projectors.

For direct-view sizes from 20 to 60 in. (50 to 150 cm) in diagonal, the plasma display panel is the only viable technology. PDPs use screen printing manufacturing technology, and therefore the larger sizes at HDTV resolutions are possible when 50 lines/in. (2 lines/mm) resolution is appropriate. PDPs are ideally suited to the television-on-the-wall picture applications, but they are too expensive and consume too much power for home television. They are used mainly in marketing booths, public exhibits, and industrial applications.

For low-resolution displays of 20 lines/in. (1 line/mm) or less and large-size direct-view displays, the LED has become the preferred technology. With the development of an efficient blue LED in combination with red and green LEDs, this technology is ideally suited to marquees and very large (billboard-size) displays in such applications

as sports arenas and roadside advertising. The LED sign is made from individual lamps assembled into printed circuit board modules tiled together into large arrays. The LED can be made bright enough to be read in full daylight but may consume kilowatts of power in large sign applications. Because of the lamp assembly technique, LEDs are relegated to low resolution and are not suitable for consumer television or computer monitors. Lawrence E. Tannas, Jr.

Bibliography. B. Bahadur, *Liquid Crystals, Applications and Uses*, 3 vols., 1991-1993; P. A. Keller, *The Cathode-Ray Tube: Technology, History, and Applications*, Palisades Press, 1991; L. W. MacDonald and A. C. Lowe (eds.), *Display Systems: Design and Applications*, Wiley, 1997; T. J. Nelson and J. R. Wullert II, *Electronic Information Display Technologies*, 1997; W. C. O'Mara, *Liquid Crystal Flat Panel Displays: Manufacturing Science and Technology*, 1995; S. Sherr, *Electronic Displays*, 2d ed., 1993; L. E. Tannas, Jr., Color in electronic displays, *Phys. Today*, 45(12):52-57, December 1992; L. E. Tannas, Jr., et al., *Flat-Panel Display Technologies: Japan, Russia, Ukraine, and Belarus*, Noyes Publications, 1995; L. E. Tannas, Jr., *Flat-Panel Displays and CRTs*, 1985; J. Whitaker, *Electronic Displays: Technology, Design, and Applications*, 1994.

Electronic equipment grounding

The connecting of electronic equipment to an electromagnetic reference common to itself, its power source, its environment, and the environment of its users. Electronic equipment is grounded to protect users from shock, to protect the equipment from spurious currents or voltages, and especially to isolate it from noise that contaminates its environment.

Need for thorough grounding. Usually electronic equipment is powered from the electric service that supplies the building where the equipment is used. To provide security against electric shock or fire ignition, equipment frames and power-conductor enclosures are connected and grounded. See GROUNDING.

In addition, grounding (together with the shielding, isolation, compensation, and equalization) is applied to minimize the entrance of extraneous signals (noise) into the equipment. Basically, electronic equipment can be treated as a sensor or signal source, a signal circuit, and central equipment (Fig. 1). The sensor may be any of a wide variety of input devices, such as a tape reader, card reader, memory readout, medical sensor, industrial transducer, or microphone. Information from the sensor travels along the signal circuit in digital or analog form at low level. At the central equipment the signal is processed to be used for such purposes as recording or control. Reliability of the output use depends directly on the integrity of the signal received at the central equipment.

The sophistication of electronic equipment has advanced tremendously. The necessary intensity of electrical signals has been diminished generally to

below 1 V. The amount of time allocated for the transmission of one bit of information has been progressively diminished to less than 1 microsecond. These advances aggravate the likelihood that spurious error signals will be of sufficient magnitude to create an error response. The effect of an error response, if not identified and corrected for, can produce tragic results in critical areas such as manned space flight. The controlling parameter can be expressed as the signal-to-noise ratio. The problem is ultimately one of ensuring acceptably high levels of signal-to-noise ratio. See ELECTRICAL NOISE.

The most vulnerable circuits are typically those associated with the information-gathering function, because they generally operate at low signal level and are followed by high-gain amplifiers within the data-processing system. Such an input circuit is shown in Fig. 1. The signal-sensing element may be in the next room, on a different floor of the same building, in a different building, or even in a different city. Techniques that are effective in establishing a high signal-to-noise ratio in this critical circuit, and techniques needed to handle the less critical circuits, are reviewed here.

Noise coupling into signal circuit. Spurious unwanted electric impulses (noise) may be mixed with true information signals in numerous ways, the most common and important being electromagnetic induction, electrostatic induction, and conductive coupling (Fig. 2).

Magnetic induction. A changing magnetic field intensity in the space through which the signal circuits run (Fig. 2a) creates noise in the equipment. Such fields are present in the space surrounding any open conductor carrying substantial magnitudes of changing current. Intense fields exist around the anode leads of arc furnaces and the throat conductors of flash and spot welders. A helpful approach is to recognize that the inductive reactance of a current-carrying conductor, shown on the electrical engineer's one-line diagram as a property of the conductor, is in fact a property of the space-distributed magnetic field which surrounds the current-carrying conductor. Any electrical circuit which links a portion of such a space field will display a fractional part of the reactance voltage drop of the power conductor. See ELECTROMAGNETIC INDUCTION.

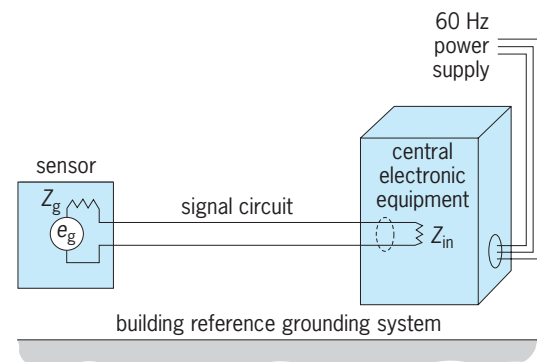


Fig. 1. The three essential components for communicating information electrically.

In Fig. 2a the fractional part of the changing space-distributed magnetic field ΔH_1 which threads between the two signal-circuit conductors will induce a normal-mode error signal in that circuit. The terminal sensor is often an electromagnetic device which, in the presence of a changing space-distributed magnetic field ΔH_3 in a particular direction, will display a normal-mode induced error voltage. A similar space field ΔH_2 threading between the signal circuit conductors and the nearby conducting structure which connects the central processing framework with the remote terminal structure will induce a common-mode potential which can appear on the signal circuit.

Electrostatic induction. A changing electrostatic space field (Fig. 2b) can induce a similar error voltage in the signal circuit by capacitive coupling. Electrostatic space-distributed fields will be created by unshielded, energized electrical conductors. Insulation over the conductor is not an attenuator. Intense fields can be expected in areas where ac high-potential testing is being done. Unshielded gaseous discharge lamps develop annoying space fields rich in harmonics and sharp step changes.

A changing-magnitude electrostatic field ΔE_1 , if coupled to one of the signal-circuit conductors through a capacitance C_1 , will impart to that conductor an error voltage. To the other signal conductor will be communicated a similar error voltage produced by space field ΔE_2 through coupling capacitance C_2 . Even though ΔE_1 equals ΔE_2 and C_1 equals C_2 , a normal-mode error signal voltage will be created if the signal circuit is not symmetrically balanced in its coupling to the central-equipment reference ground plane. Also, the terminal sensor may contain such parts as terminals, coils, and winding surfaces that are capacitively coupled to the ambient space field ΔE_3 . Unless coupling capacitance C_3 displays a symmetrical balance to the two signal conductors, a portion of the induced error voltage will be of the normal-mode type. Even an intentional grounding connection on the signal circuit at the central equipment, if not precisely balanced, will cause a fractional part of the common-mode, electrostatically induced voltage to appear as a normal-mode type.

The presence of large-magnitude, common-mode voltages on signal circuits may directly damage or even destroy electronic components and may create troublesome normal-mode error signals if only slight deviations from perfect symmetry occur on signal circuits.

Conduction between units. Normal-mode error voltages are also created by conductive coupling (Fig. 2c). For the most part this origin of noise is identical with that commonly referred to as ground loop problems. The error voltage owes its origin to the fact that the reference grounding conductor system is not at a common potential throughout. Consider that one signal-circuit conductor is grounded at the central equipment in location G_1 . This would be conventional practice when using coaxial signal circuits and is sometimes used on two-conductor wire circuits. The presence of a second ground connection

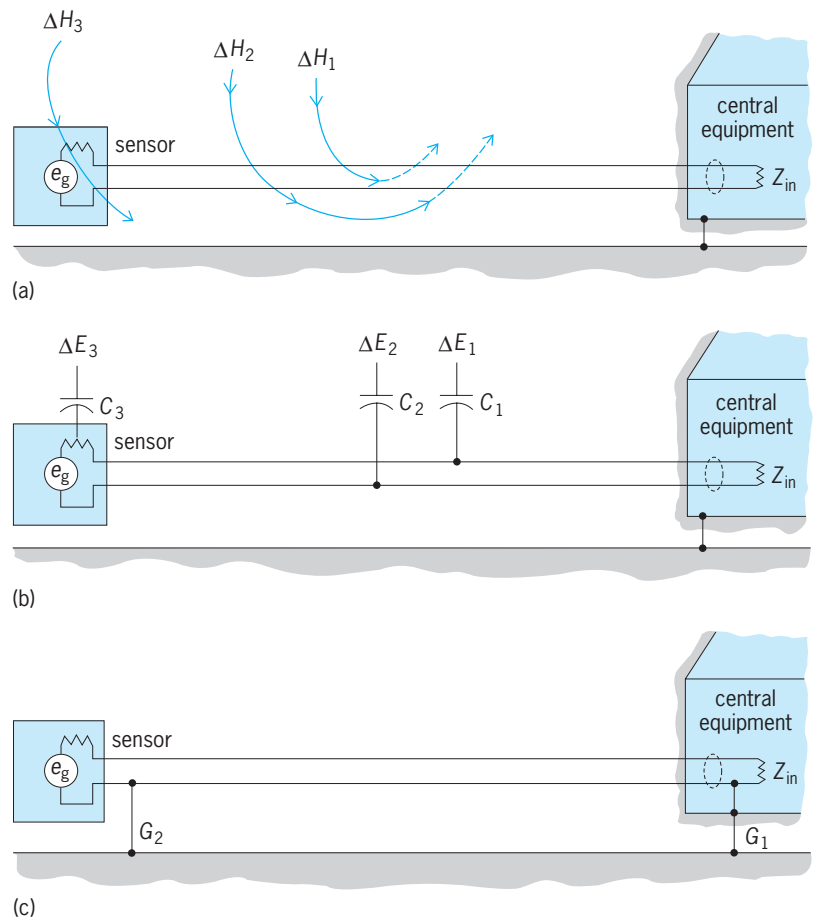


Fig. 2. Common sources of error signals. (a) Electromagnetic coupling. (b) Electrostatic coupling. (c) Conductive coupling.

on the same signal conductor at a different location G_2 will likely result in severe noise in low-level signal circuits. After making the first ground connection at G_1 , it will be unlikely that another spot can be found along that conductor with a zero-voltage difference to the adjacent reference ground. When a second ground connection is made along a signal conductor, it forces the voltage difference, which had previously existed here, to vanish and appear as an impedance voltage drop (IZ drop) around the loop circuit formed by the second ground connection. See ELECTRICAL IMPEDANCE.

By reviewing the character of the impedance elements forming this loop, it becomes quite evident that the signal-circuit conductor between grounding points G_1 and G_2 will principally account for the ground loop impedance. The remainder of the loop will be made up of such low-resistance conductors as heavy building structural members and large pipes. Thus most of the voltage difference which, prior to the second ground connection, had existed between the signal conductor and the adjacent reference ground G_2 has, after the second ground connection, become an impedance voltage drop along only one signal conductor between locations G_1 and G_2 . This voltage represents a normal-mode error voltage in the signal circuit.

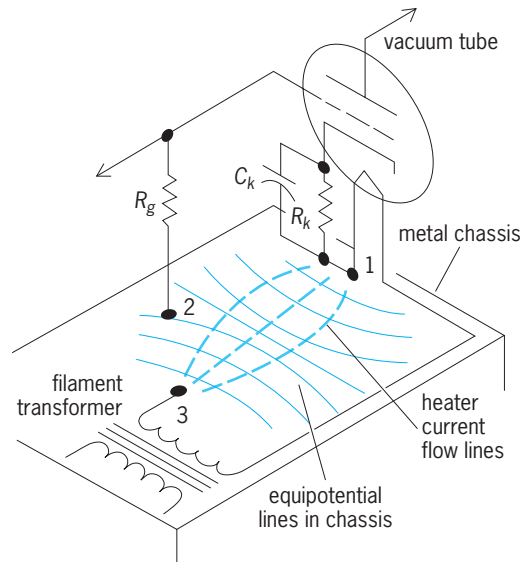


Fig. 3. Improper grounding of an individual amplifier stage to a metal chassis, so that the noise voltage is coupled into the signal circuit.

Conduction within one unit. At times the separation of only inches in the two ground connections can result in intolerable noise, as the following example illustrates. In a vacuum-tube amplifier (**Fig. 3**) the cathode and one heater lead are bonded to the chassis at point 1. One side of the filament transformer is bonded to the chassis at point 3; therefore, heater current flows through the chassis, as shown by the current flow lines. Normal to the current flow lines are equipotential lines, as shown. The grid resistor is bonded to the chassis at point 2. A noise signal e_{12} is thereby inserted into the grid circuit, resulting from the potential drop along the chassis between cathode and grid bonding points. This drop may amount to only a few microvolts, but if it enters an early stage of a high-gain amplifier, excessive noise may be noticed in the output.

To avoid such problems, all bonding connections for a given stage should be made to only one point on the chassis. Also, the flow of large currents through the chassis of any high-gain amplifier, whether vacuum-tube or solid-state, must be avoided. Where high-gain amplifiers are controlling high-power ac or dc circuits having grounded terminals, it is best to insulate the amplifier chassis from the frames of the rest of the machinery and to provide a single ground strap from chassis to frame. This will prevent large currents from passing through the amplifier chassis. See **CIRCUIT (ELECTRONICS)**.

In general, each stage of a high-gain amplifier should have its own chassis ground. Each individual amplifier should have its own ground lead to a common amplifier ground point. All machinery and other power devices should be bonded together, with a single lead to the amplifier ground point. Finally, if additional shielding is provided, such as a metal cabinet, it should have a single ground lead to the common ground point for the amplifier and the power equipment. This common ground point would then be

grounded to earth by bonding to a cold-water pipe or by other methods in conformance with the local electrical codes.

Ground gradients. The medium- to high-voltage substation may display large voltage differences between the potential on the local station grounding conductors and remote earth. Information-gathering (or sensing) circuits associated with electronic equipment extending into or close to such a substation area will generally require extreme design care to avoid dangerous common-mode potentials. It is not unreasonable to expect steady-state voltage levels of as much as 50 volts. During a line-to-ground fault on an outgoing overhead line, the potential of the substation ground mat might be elevated to several thousand volts, relative to remote mean earth potential by voltage gradients. Common-mode voltages from such causes can be compensated by insulating transformers or equalizing transformers or by other means of bridging large-magnitude, common-mode voltages.

The ordinary commercial building may present unexpected problems relative to common-mode voltage components due to voltage gradients in the building grounding conductor system. It is widely believed that the grounding requirements as spelled out in section NEC 250-23(a) of the National Electrical Code ensure that all electrical load current within the building will be returned to the service equipment (the point of electric service entrance) on power conductors, independent of the building structure.

An innocent-appearing exception can violate this concept. Section NEC 250-60 defines a grounding exception (rather generally used) for the frames of electric ranges, wall-mounted ovens, and counter-mounted cooking units. The code exception allows these appliance frames to be grounded by connection to the power system grounded conductor (white wire) if the electric service is a 240/120-V, single-phase, three-wire system or is taken from a 208Y/120-V, three-phase, four-wire system. It is unrealistic to assume that the appliance frames in question will not also be in contact with the building structure. The result is that the white load-current-carrying conductor becomes connected to the building frame through the heating appliance.

Thus a building devoted essentially to commercial activities, served with 208Y/120-V electric power, may have a snack bar installed on the sixth floor. Contained therein may be appliances grounded to the white wire as allowed by NEC 250-60. Instead of an electrically dry building frame, it may be found that between the sixth floor and the ground floor there is distributed the same electrical voltage drop as exists on the grounded power-system conductor between the snack bar connection and the service entrance.

When the service is three-phase and four-wire, all third-harmonic currents (and their multiples) in the entire array of line-to-neutral connected loads combine in an additive fashion in the neutral conductor (white wire) to aggravate the harmonic voltage drop

along the white wire. In both the three-phase and single-phase power-supply cases, the increasing use of time-modulated (SCR-controlled) current in line-to-neutral connected devices (such as fans and lighting units) makes for the presence of much step-front hash in the voltage drop along the neutral conductor, which becomes voltage gradient in the building structure if the NEC 250-60 exception is employed. See SEMICONDUCTOR RECTIFIER.

Efforts have been made to restrict the use of the NEC 250-60 exception to appliance circuits which originate as branch circuits at the service equipment. It must be recognized, however, that the ground reference potential on one floor of an office building may differ from that of another floor by as much as several volts (third harmonic) and contain substantial fast-front hash.

Section NEC 250-23 of the National Electrical Code clearly prescribes that there shall be no grounding connection to the grounded power conductor of the electrical system downstream of the service equipment for that establishment. When this rule is respected, all of the load-system power current is returned to the grounded conductor at the service via insulated conductors. There is no opportunity for electrical noise, as it exists on the grounded power conductor (the white wire) to be conductively transferred to the system grounding conductors. There are three exceptions in the NEC text which may permit that rule to be bypassed. The first one applies to an in-plant separately derived electrical system located remote from the service equipment of the establishment. It is common to apply a permanent grounding connection at the supply-machine neutral. This will create a cross bond between the grounded conductor (white wire) and the grounding conductor downstream of the main service equipment. The second one applies to the case of an electric-supply circuit run to a second independent building. The NEC requires a cross bond between the grounded and the grounding conductor at the point of entry to the second building, unless an independent grounding conductor has been included with the power conductors from the supply point in the main building. The third exception, pertaining to ranges, counter-mounted cooking units, wall-mounted ovens, and clothes dryers, can be troublesome. While the National Electrical Code says merely that the power-circuit grounded conductor (white wire) may be used as the grounding conductor, it is commonly found that the appliance frames so grounded by connection to the white wire are also in contact with building metal frame or metal piping systems. It is this "back door" connection which allows communication of the electrical noise to the establishment grounding conductor system.

Progress is being made in restricting the use of the white wire as a grounding conductor.

A few instances of a cross-bond between the grounded power conductor and the equipment grounding conductor downstream of the service equipment have been observed as a result of an

in-plant emergency or standby generator. These have been associated with 208Y/120-V and 480Y/277-V solidly grounded power systems serving hospitals, telephone exchanges, police centers, and other such important establishments. When operating on an emergency or standby basis, it is proper that a grounding connection exist at the local generator as prescribed for a separately derived electric power system in section NEC 250-26 of the National Electrical Code. It is common to find an intentional permanent grounding connection at the local-generator neutral junction. Unless the neutral line of the supplied load is switched, along with the phase conductors, the result is a permanent grounding connection on the "white wire" at the generator location, contrary to the planned design pattern. Location of the in-plant generator adjacent to the service equipment would eliminate the problem, as would also inclusion of the neutral conductor in the transfer switching operation.

There will occasionally occur brief intervals of much greater than normal voltage gradients along the building grounding conductor system. An insulation failure on a power conductor may permit a large-magnitude, ground-short-circuit current to flow through the building structure toward the service equipment location for the brief interval permitted by the overcurrent protector.

Importance of low-noise environment. The noise problems presented to the electronic equipment designer increase directly with increased ambient noise levels in the area. The quality level employed in grounding and shielding of power-system conductors has a marked effect on the ambient noise levels.

In an occupied building much can be learned of the general ambient noise levels by a study of the electric facility practices employed and by a general inspection of the building. Additional specific information may be obtained by test. It is important to look for unusual sources of high-energy radiation which may penetrate the building interior space. A nearby radio station antenna may create strong fields at one fixed frequency. High-power radar-scanning installations may be a source of annoying ambient fields of pulsed high frequency.

After an appraisal of the severity of the ambient electrical noise level, the problem of designing the grounding and shielding practices for the electronic equipment becomes a straightforward procedure.

Grounding and shielding practices. To establish an acceptably high signal-to-noise ratio demands that normal-mode noise be maintained at levels below a specified tolerable value.

Choice of conductor. The construction geometry of the signal conductor (**Fig. 4**) plays a prominent part. Of the many possible varieties, the four most commonly used are random-selected wires in a multiconductor cable, parallel-wire twin-lead, twisted-pair, and coaxial lines. Because of the superiority of the twisted-pair relative to random wires or twin-lead at modest cost increase, its use is universal when noise

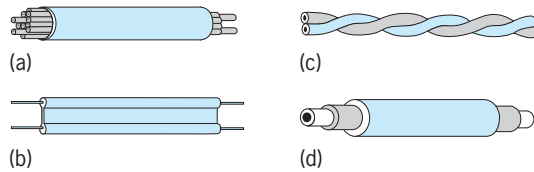


Fig. 4. Common signal-circuit geometries. (a) Two random wires. (b) Straight twin-lead. (c) Twisted pair. (d) Coaxial line.

reduction is a factor. A half twist every inch is both typical and effective.

An encircling metal raceway accomplishes substantial additional attenuation, particularly if it consists of steel conduit. In securing high-quality rejection of capacitively coupled noise voltages, it is common practice to avoid any grounding connection on the signal conductor (coaxial lines excepted) and to incorporate a balanced impedance connection to the central-equipment ground reference terminal in the input circuits of a differential amplifier.

The coaxial line (Fig. 4d) is a theoretically superior construction but may be faced with operating conditions which deteriorate these qualities. As it is inherently an unbalanced line, the exterior shell is commonly connected to the reference ground terminal at the central equipment. The existence of any connection to adjacent ground along the coaxial line (not easy to avoid) creates the ground loop situation described in Fig. 2c with a likely inevitable insertion of normal-mode error signal (noise). Even in the absence of a second ground connection, large-magnitude varying electric fields (Fig. 2b) can create measurable values of charging current flowing along the outer shell to the grounding connection at the central equipment, and in so doing can create an IR voltage drop along that conductor which will appear as a normal-mode noise voltage. An enclosing metal raceway grounded to the central-equipment ground reference terminal can be used to effectively suppress this noise source.

Conductor shielding. A multiplicity of signal circuits may be run within a common metal raceway, providing that the signal magnitude on these conductors is not great enough to create cross-coupled interference (cross talk). Of the various conductor geometries considered in Fig. 4, the random wire case in *a* would be most susceptible and would offer little hope for correction. With the conductor patterns in *b* and *c* the application of wrapped-on shielding tape (or equivalent concentric conducting shell) to each signal circuit, grounded at the central-equipment grounding terminal, would only accomplish high-quality cross-talk suppression. The coaxial line is inherently free of cross-talk interference if the outer conductor is grounded only at the central equipment.

It is possible (and not unlikely) that the outer terminal sensor may itself respond to ambient magnetic field noise (Fig. 2a) or electrostatic field noise (Fig. 2b) or both, and create normal-mode, signal-circuit noise. The magnetic field problem can be resolved by the application of appropriate magnetic

shielding at the sensor, the enclosure of the complete sensor within a closed shell of conductive metal, or a combination of the two. The electrostatic problem can be resolved by enclosing the sensor winding and its circuit leads within a thin-wall shielding shell maintained at zero reference ground potential. The enclosing shell can be the same one used for magnetic induction suppression. The required connection to a zero potential reference point may require a unique grounding circuit extending from the central equipment (Fig. 5).

Control of grounds. The remaining perplexing noise abatement problem relates to the creation of reference grounding terminals at peripheral equipment locations which appear to be at the same potential as the central-equipment grounding terminal; an example is the electrostatic suppression problem discussed in the previous paragraph. An earlier discussion (Fig. 2c) develops the fact that between a building ground-reference point G_2 and a widely separated central-equipment grounding terminal G_1 there will, more often than not, appear a difference in electrical potential. The difference voltage can create normal-mode, error-signal voltages in the several ways described.

If the remote terminal G_2 is not in the same building as is terminal G_1 , or if it is located in an area where large common-mode noise voltages are expected, the use of insulating or equalizing transformers in the signal lines or the adoption of other common-mode suppression techniques should be considered.

If the remote equipment is in the same metal-frame building or in a location where common-mode voltage difference between reference ground terminals is known not to exceed perhaps 10 V, acceptable attenuation should be possible without the introduction of isolation techniques on the signal wires themselves.

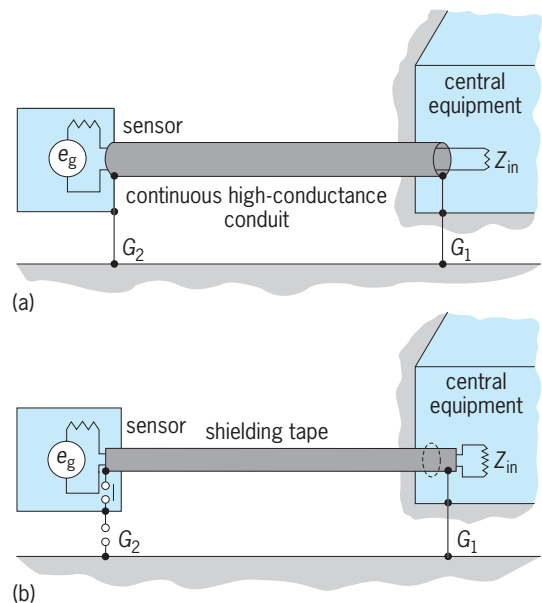


Fig. 5. Practical grounding techniques. (a) High-conductance raceway. (b) Electrostatic shielding.

The installation of a grounding conductor, run with the signal circuits, which interconnects grounding terminal G_2 with G_1 will be helpful but not as effective as might be expected. A substantial fraction of the original noise voltage will remain in the form of an IZ drop along the grounding conductor.

High-conductance raceway. By modifying the form of this grounding conductor, a tremendous gain in attenuation can be achieved. By forming the conductive metal of the grounding conductor into a hollow tube through which the signal conductors are passed, its effectiveness is increased enormously. The resulting tubular grounding conductor, called a high-conductance raceway (Fig. 5a), interconnects grounding terminals G_1 and G_2 as did the cable grounding conductor. Voltage difference $e_{G_1} - e_{G_2}$ appears as an IZ drop along the raceway. The conductance (tube diameter, wall thickness, and type of metal) is selected so that only a small portion of the total IZ drop appears as a resistance drop along the metal tube. Almost the entire voltage difference between grounding terminals G_1 and G_2 is accounted for by a reactive-voltage drop IX created by an induced magnetic field encircling the tubular raceway. This magnetic field also links each signal conductor contained within the raceway and thus induces in each one identically the same IX voltage drop as in the raceway. With the IX voltage component equalized, or canceled out, only the IR drop along the high-conductance raceway appears as a voltage difference between signal conductors and adjacent ground at G_2 .

It is important to note that this cancellation of the IX drop occurs only on those signal conductors run within the high-conductance raceway. If some critical signal circuits must take a different route, an independent high-conductance raceway to contain them will be needed to accomplish the desired high-magnitude attenuation.

The high-conductance raceway technique is applied extensively to interconnect the several independent housings marking up the central electronics equipment to establish a nearly quiet ground reference potential throughout the central equipment. One prominent item of the central equipment is designated as the reference unit and is bonded to the building reference ground, preferably at only one point. All other equipment frames in the group would, by design intent, permit conductive connection only to the central reference structure in the form of high-conductance raceways. Power conductors and other high-noise-level conductors should be run in raceways independent of the sensitive signal circuits.

In extremely severe situations where the IR drop in the high-conductance raceway creates intolerable noise levels, a dual concentric system may be substituted. The exterior high-conductance raceway is identical in design and function with the one just described. Within this is a second metallic tubular raceway, insulated from the former except for an intentional cross bond at the central equipment

grounding terminal G_1 . At the remote equipment locations, the construction consists of an outer enclosure bonded to the building ground terminal G_2 . The electronic equipment chassis containing the active electronics components is insulated from the outer housing and bonded only to the internal concentric raceway. By this modification even the resistance voltage drop along the outer tubular raceway is prevented from appearing as a voltage difference between signal conductors and their enclosing raceway.

Should the remote sensor terminal be small and of simple character, with possible noise being essentially of electrostatic coupling origin, the simple circuit pattern illustrated in Fig. 5b may prove adequate. A lightweight shielding tape enclosing the signal conductors preserves a zero-potential reference plane surrounding the signal conductors. The sensor frame (or enclosing electrostatic shield) is connected to the signal-circuit shielding tape and kept insulated from the local grounding terminal G_2 . The presence of high-intensity electric field noise at the outer terminal may demand a higher-conductance signal-circuit shield to avoid objectionable noise voltages of IR drop origin.

Grounding, electromagnetic shielding (and cancellation), electrostatic shielding, and insulation are intimately intertwined in effecting an acceptable solution.

Data transmission. To an ever-increasing extent the high-density channels of information handling (data transmission) are being accommodated by means other than metallic conductors. The prominent transmission avenues are microwave links and optic fibers. The problems of circuit grounding along these links are simply nonexistent. The interconnecting wire circuits between the distributed sensors and the central processing center, and between the various functional systems at the center, and the transmission to the output terminal at that location will face the same signal-to-noise ratio problem as in the past. See MICROWAVE; OPTICAL COMMUNICATIONS; OPTICAL FIBERS.

Optical isolator. In dealing with the aggravated problems associated with the large magnitudes of common-mode noise voltage encountered with circuits run to an outdoor electric substation, the optical isolator is a valuable tool. This device introduces a short optical transmission path into the electric signal circuit with a capability of withstanding several thousand volts of common-mode voltage without ill effects. See OPTICAL ISOLATOR. R. H. Kaufmann

Bibliography. D. G. Fink and H. W. Beaty (eds.), *Standard Handbook for Electrical Engineers*, 14th ed., 2000; Institute of Electrical and Electronics Engineers, *IEEE Recommended Practice for Powering and Grounding Sensitive Electronic Equipment*, IEEE Std. 1100-1992, 1992; M. Mardigian, *Grounding and Bonding*, 1988; R. Morrison, *Grounding and Shielding Techniques in Instrumentation*, 1986; R. P. O'Riley, *Electrical Grounding: Bringing Grounding Back to Earth*, 5th ed., 1998.

Electronic listening devices

Devices which are used to capture the sound waves of conversation originating in an ostensibly private setting in a form, usually as a magnetic tape recording, which can be used against the target by adverse interests.

There are two kinds of electronic listening devices. One takes advantage of equipment already present on the target's premises, such as a telephone, radio, phonograph, television set, public-address loudspeaker, or tape recorder, to act as a microphone, transmitter, or power supply. The other does not. In the former case, the target's equipment is said to have been compromised.

These practices are unlawful in the United States and Canada except when carried out by law enforcement officers acting under authority of a warrant. In the United Kingdom, electronic eavesdropping by private parties may contravene the laws against trespass and unauthorized use of telephone equipment or radio frequencies, but it is not unlawful in and of itself. The same is generally true in western European countries.

Compromise. Compromise of the target's own equipment takes advantage of the fact that any loudspeaker is capable of functioning just as well as a microphone, that convenient sources of dc power are available within the equipment, or that the equipment is connected to power or signal lines that can transmit the intercepted conversation to some place where recording can conveniently be accomplished.

The equipment most frequently compromised is the telephone handset. The act of compromise may be as simple as bypassing the switch hook with a Zener diode so that the instrument can transmit conversations to a wiretap, or high-impedance parallel connection off the premises, while not signaling a busy tone to someone attempting to dial in. *See* TELEPHONE; ZENER DIODE.

It can be as complex as the infinity transmitter, which is activated by the eavesdropper's dialing in and transmitting, before the telephone rings, an audio tone to a tuned relay concealed in the instrument.

Bugs. Eavesdropping devices that can stand alone are known commonly as bugs. They take advantage of many developments of modern technology, such as microcircuits, miniature ceramic microphones, and miniature batteries. *See* INTEGRATED CIRCUITS; MICROPHONE.

The art of designing bugs has achieved its most advanced state of development in the national intelligence services of the Great Powers. One bug was discovered to have been inserted in the heel of a diplomat's shoe.

The smallest bug available to the private citizen is probably the Hong Kong "spider." It is no larger than a common postage stamp and less than a quarter-inch thick.

Electronically a bug is often just a two-stage frequency-modulated transmitter: an audio amplifier

and a variable-frequency radio-frequency (rf) oscillator. *See* AMPLIFIER; OSCILLATOR; RADIO TRANSMITTER.

Bugs may operate on any frequency from 20 to 1000 MHz, but usually they are designed to operate at a frequency close to that of a powerful local frequency modulation (FM) or very high-frequency (VHF) television station.

A popular hybrid between a compromise device and a bug is the telephone drop-in. In this design, an FM transmitter is made in the form of a telephone microphone. The eavesdropper can casually unscrew the mouthpiece of the target's telephone handset and substitute the drop-in for the original microphone. The range of this device is about 250 ft (75 m). It has the added advantage of drawing its dc power from the telephone company central battery.

Defenses. The telephone line analyzer is used to defeat compromise devices. Under the control of a microcomputer, the line analyzer examines in turn each of the 50 or more telephone lines used by a typical commercial firm. It can detect the tiny amount of current flow into a Zener or other compromise device. A monotonically rising audio signal is then applied to trigger any infinity transmitter. If a compromise device is discovered, a pulse of 800 V is impressed upon the line to burn it out.

Technical surveillance sweeps, as they are known to the trade, are performed alternately with external signal lines connected and disconnected and with telephones on-hook and off-hook.

Often a sniffer is used to sense the presence of a bug. A sniffer may be a simple VHF diode detector with an output meter that reveals the presence of a local carrier. More sophisticated instruments sweep the frequency bands of interest and compare the audio output with locally generated signals, usually a tape recording of simulated business activity. Another type of sniffer affords discrimination between a genuine clandestine device and a strong local broadcasting station.

If the presence of a clandestine device is indicated, the sweep team may employ a panoramic intercept receiver. This instrument displays the amplitude of received signals on a cathode-ray-tube (CRT) trace calibrated in units of frequency. The frequency of the bug can be determined by observing which of the pips on the CRT trace appears to diminish and grow in direct response to the locally generated audio signal. Once its frequency is determined, the bug's physical location can be found by use of a tuned rf field-strength meter having a loop antenna. *See* ANTENNA (ELECTROMAGNETISM); CATHODE-RAY TUBE.

The objective of a sweep such as the one described is to locate the device and discover its technical characteristics with a view toward feeding the device false information that may discredit or compromise the individuals responsible for planting it.

National security forces encounter clandestine devices which are switched on and off by an agent using a control transmitter. Such a device cannot

be detected by the methods described above because they presuppose that the bug is continually on the air or can be activated by locally generated audio.

To detect silent bugs, advantage is taken of the fact that all subminiature transmitters contain one or more semiconductor junctions that are connected to an electromagnetic radiator. Inasmuch as these junctions are nonlinear impedances, a locally generated low-power ultra-high-frequency carrier will be reradiated by the bug, and the reradiated signal will contain strong harmonic components. *See* HARMONIC (PERIODIC PHENOMENA).

Natural devices such as rusty bedsprings will reradiate second-harmonic components, but only a real semiconductor junction will reradiate third (and higher) harmonics.

An ultra-high-frequency (UHF) sweeper is equipped with a meter that produces a negative deflection in the presence of reradiated second-harmonic energy and a positive deflection in the presence of reradiated third-harmonic energy. Such a deflection discloses the presence of a bug, albeit a silent one not detectable by sniffers or panoramic receivers. The UHF sweeper is portable; it bears a resemblance to a household appliance for vacuuming drapes. Increasing positive deflection indicates to the sweep team member that the bug is being approached. *See* HARMONIC ANALYZER.

Advanced devices. Human capacity for designing electronic instruments to spy on one's neighbors appears to be boundless. Some of the more exotic include laser radar which responds to infinitesimal motion of windowpanes caused by conversation in a room, bugs which include miniature television cameras, and minicomputers which are programmed to intercept messages to a resource-sharing computer and to return answers which are spurious.

The so-called pinhole camera, a television camera that uses charged-coupled devices, is less than 2 in. (5 cm) long. Its fast ($f/1.8$) lens provides the capability for excellent infrared vision and its wide-angle (11-mm focal length) characteristics permit viewing a whole room through a pinhole, or a hole no wider than 0.5 in. (13 mm) if a fiber-optic extender must be used to look through a thick wall. *See* CHARGE-COUPLED DEVICES; LENS (OPTICS); TELEVISION CAMERA.

Typical hiding places are in lighting fixtures, radios, television sets, and illuminated signs, and behind switch plates or electrical receptacles. Miniaturized video distribution systems are available that will transmit signals up to 4000 ft (1.22 km) over 24-gauge twisted pair wires. In most jurisdictions these devices have not yet been declared illegal. John M. Carroll

Bibliography. Act IV Security Services Staff, *Surveillance Countermeasures*, 1994; B. Bruno, *Serious Surveillance for the Private Investigator*, 1992; S. French, *The Big Brother Game*, 1976; S. Gowrinathan and J. Shanley (ed.), *Surveillance Technologies II*, 1992; L. Lapin, *How to Get Anything on Anybody*, 1987, bk 2, 1991.

Electronic mail

The asynchronous transmission of messages by using computers and data-communication networks. Historically, electronic mail (or e-mail) referred to any of a number of technologies that allowed people to send documents to one another through electronic means. It was frequently used to describe both wirephoto [the precursor of the facsimile (fax) machine] and telegraphy. Subsequently, usage of the term focused upon the narrower sense given above. *See* FACSIMILE; TELEGRAPHY.

Early systems. The term e-mail began to disseminate in the early 1970s, when time-sharing on large mainframe computers became common. Many systems had the ability to send messages from one time-sharing user to another on the same system by directing it to the receiver's sign-on name. *See* MULTIACCESS COMPUTER.

With the creation of modems capable of being dialed under computer control, the relatively inexpensive creation of e-mail networks became possible. This development frequently was fairly haphazard, one computer dialing others a few times a day and transferring all the messages that had accumulated since the last phone call. As this process was repeated by all the cooperating computers, messages would get to their destinations. *See* MODEM.

Sending messages among a number of machines over dial-up connections required that the addressing scheme be expanded to give a computer identifier as well as the receiver's sign-on name. Since there was no central routing database, users were forced to route messages by hand within the address. For example, a!b!c!smith would send the message from system a to b, from b to c, and there deliver it to user smith on system c. In fact, there could be multiple computers named c, but only one to which b knew how to route mail.

Impact of personal computers. The availability of the personal computer in the 1980s further advanced the technology in three ways. First, since computers became more accessible, e-mail service bureaus became commercially viable. These e-mail systems used the mainframe messaging technology centrally and allowed subscribers to call in via the long-distance telephone network to send and receive their messages.

Second, personal computers coupled with local-area networks produced a number of proprietary e-mail networking systems for work groups. These systems offered better user interfaces and more functionality than the more general electronic-mail solutions at the expense of demanding uniform systems throughout the e-mail network. That is, work groups could communicate with each other only if they were all running identical software. *See* LOCAL-AREA NETWORKS.

Finally, there have always been two distinct actions in the delivery of e-mail: the delivery across the network to the recipient's mailbox, and the recipient's removing a mail item from the box to read it. Originally, both these actions would occur on the

same service bureau computer, where the recipient would connect online to the mail server and read the mail while connected. This forced the e-mail user to learn a completely different computing environment to read or send e-mail and to remain connected to the server while reading or writing messages. With the advent of personal computers with window-based point-and-click user interfaces, the operations of being a mail repository and a mail interface were split. New "mail clients" such as Eudora, Exchange, and Netscape Mail were created. These used special access methods such as POP (Post Office Protocol) or the newer IMAP (Internet Message Access Protocol) to manipulate mail remotely, including moving mail en masse to a recipient's workstation. Once there, the recipient could read it at leisure without being connected to the network. Outgoing messages would be kept until the next time the computer was connected to the Internet, when they would be sent on their way. This allowed the messages to be processed in a variety of new ways such as threading. Threading presents messages in chronological order within a topic rather than just chronologically. This allows a complete topic to be read (or discarded) as a unit. *See* HUMAN-COMPUTER INTERACTION; MICRO-COMPUTER.

Message identification and routing. The use of electronic mail grew continuously until the late 1980s but never achieved widespread use outside of work groups or corporations. The limiting factor was the complicated addressing that had to be worked out before a message could be successfully transmitted.

There were two proposed methods to solve the problem of mail-system identification and routing. The Organization for International Standardization (ISO) formulated the X.400 standard, and the Internet community developed an extended use of the domain name system (DNS). Many impediments to the spread of X.400, such as high software costs and delays in standardization, caused the freely available DNS solution to become the de facto standard.

The DNS describes a worldwide distributed database in which each site maintains its own information about how to route messages to a computer within its administrative domain. A computer wishing to send a message to another asks the DNS for the routing information and uses the information returned to make the connection. This allows a person on virtually any online networking service to send mail to another person by giving only the personal identification and the e-mail system name of the recipient. *See* DISTRIBUTED SYSTEMS (COMPUTERS).

Multimedia transmission. From the time the usage of the term narrowed to exclude facsimile until the early 1990s, generally only coded textual information could be transferred via e-mail. The transmission of nontextual data required special preprocessing, postprocessing, and prior arrangements between the sending and receiving parties. It was very difficult to make these kinds of transfers if the sending and receiving computers were different types.

This restriction was lifted with the adoption of the MIME (Multimedia Internet Mail Enhancements) standard. It described a way of encoding an arbitrary list of media types within a normal textual message in an operating-system-independent manner. Finally, different types of systems could send executable, sound, picture, movie, and other kinds of files to each other via e-mail. *See* MULTIMEDIA TECHNOLOGY.

Attachments. The sender's mail package can attach a file to a standard text message. When the message is sent, the attached files are encoded into a single message. The receiver's e-mail package will decode the message into a base message plus an attachment for each file. The recipient's computer still must have the software required to view or process the file, but the transmission is transparent.

Security. The spread of electronic mail was also hampered by its lack of security. As mail was passed from one site to another closer to its destination, system administrators at each intermediate site could read messages. Also, the source of an e-mail message may be fairly easily forged to make it either untraceable or appear to come from another person. This limited the use of e-mail to so-called friendly applications. Public-key cryptography has been applied to e-mail messaging, notably in PEM (Privacy Enhanced Mail), in response to these security concerns. *See* COMPUTER SECURITY; CRYPTOGRAPHY.

Integration of messaging technologies. Many e-mail-to-fax gateways are in use. These allow messages that originate as e-mail to be converted to fax and then delivered via office fax machines to recipients that do not normally participate in the e-mail community.

Since the communications speeds required for e-mail are quite modest, messages are sometimes transmitted by radio. In the simplest applications of this technology, some pocket paging system providers allow short e-mail messages to be sent to their pagers. In more advanced applications, providers lease notebook-sized systems which can send and receive e-mail through Earth-satellite relay. *See* RADIO PAGING SYSTEMS.

Some voice-mail systems accept e-mail for their clients. The text of the e-mail is passed through a speech synthesizer to artificially read the e-mail text into voice storage. Clients can then hear their e-mail by accessing the voice mail system as they would for voice messages.

Infringement by other technologies. From the time that e-mail was developed until the late 1990s, it was the predominant way of staging a universally available discussion. This was done with a variety of mail forwarding utilities, such as LISTSERV, which would forward a message to a large group of discussion subscribers. Although this method is still in use and has some advantages in circumstances where the asynchronous nature of e-mail is useful, some of these discussions have moved to various chat facilities on the World Wide Web and other technologies. These chat facilities provide a model of a room for a particular topic, and anything typed by someone in the room will be seen in close to real time by

others connected to the room. See COMPUTER; DATA COMMUNICATIONS; DIGITAL COMPUTER; ELECTRICAL COMMUNICATIONS; INTERNET; WORLD WIDE WEB.

Edward Krol

Bibliography. B. Shimmin, *Effective E-mail Clearly Explained: File Transfer, Security, and Interoperability*, AP Professional, 1997; D. Strom and M. T. Rose, *Internet Messaging: From the Desktop to Enterprise*, Harcourt Brace, 1998; D. Wood, *Programming Internet Email*, O'Reilly & Associates, 1999.

Electronic navigation systems

Systems that locate suitably equipped users by means of measurements made by electronic sensors. Typical measurements include distance or bearing from a vehicle to a known point, or direct determination of the present position of the vehicle in a particular coordinate system. From the knowledge of present position, the course and distance to a destination at a known location can be calculated. Most systems are based on the use of electromagnetic (radio) or acoustic waves or inertial sensing. The systems described here are typically used by aircraft, ships, and land and space vehicles. In addition, the use of satellites for land navigation by individuals has become common.

The use of radio waves has been found to be attractive because of their known and nearly constant velocity of propagation, namely, the speed of light (c), which is about 3×10^8 m/s (1.86×10^5 mi/s). Thus, if the time (t) of travel of the radio signal between two points is accurately measured, the distance or range (d) between the points can be accurately determined from $d = ct$. Ships also use systems based on underwater sound waves, with due allowance for the much slower speed of sound in water. Also, electromagnetic waves in the visible (optical) or near-visible spectrum can be used in a similar manner for distance measurement. Antennas designed for reception of radio signals of sufficiently high carrier frequency can also be used to produce tightly focused radiation patterns that can be used in measuring the bearing (angle) between two points. See ELECTROMAGNETIC RADIATION; LIGHT; SONAR; UNDERWATER NAVIGATION; UNDERWATER SOUND.

Parameters. The parameters that characterize the utility of an electronic radionavigation system include accuracy, coverage, fix dimensions, system capacity, signal characteristics, and spectrum.

Accuracy. In navigation, the accuracy of a measured position at a given time is the difference between the indicated position output of the system and the true position at that time. For horizontal accuracy, the circular error probable (CEP) is often used to statistically characterize the random system errors. The CEP represents the radius of a circle that contains 50% of the position fixes. For three-dimensional systems, the spherical error probable (SEP) is often used to statistically represent the random system errors. The SEP represents the radius of a sphere cen-

tered at the origin that contains 50% of errors.

Position accuracy can further be characterized as follows:

1. **Repeatable accuracy:** The accuracy with which a navigator can return to a location, the coordinates of which have been previously measured using the same onboard system.

2. **Absolute accuracy:** The accuracy with which a navigator can determine position in terms of an agreed-upon coordinate system distinct from the navigation system. A common coordinate system is the Earth Center Earth Fixed coordinate system, which is a Cartesian coordinate system whose center is nominally the Earth's center of mass, whose X axis goes through the intersection of the equatorial plane and the prime meridian (which passes through Greenwich, England), and whose Z axis goes through the North Pole. In conjunction with a model of the Earth's shape such as that defined by the WGS-84 ellipsoid (defined in the world Geodetic System 1984, in which positions of satellites in the Global Positioning System are specified), geographic latitude, longitude, and height with respect to the ellipsoidal surface may be determined. See GEODESY.

3. **Relational accuracy:** The accuracy with which a navigator can determine position with respect to another user of the same system. For example, if two ships both having the same navigation system were to rendezvous, the relational accuracy would determine their ability to meet at a common point.

Coverage. The coverage of a navigation system is that surface area or space volume in which the system is adequate to permit the navigator to determine position to a specified level of accuracy. Most navigation systems have been applicable to aircraft, ships, or land vehicles, and coverage has classically been defined in terms of these applications. More generally, coverage should include applicability to space vehicles, submersible craft, and indoor applications.

Fix dimensions. This characteristic defines whether the navigation system provides a linear, one-dimensional line-of-position, or a two- or three-dimensional position fix. The ability of the system to derive other quantities such as velocity, time, or attitude is also included.

System capacity. System capacity is the number of users that a system can accommodate simultaneously.

Signal characteristics. This parameter refers to signal power levels, frequencies, signal formats, coding sequences, data rates, and other characteristics necessary to define the means by which a user derives navigation information.

Spectrum. Spectrum refers to the frequency allocation of the signal. Spectrum allocation issues are addressed by the Federal Communications Commission (FCC) and National Telecommunications and Information Administration (NTIA) within the United States. International spectrum issues are addressed in conjunction with analogous agencies from other countries during the biannual World Radio Conferences. The spectrum allocation of current

Electronic navigation systems					
System	Type ^a	Frequency, MHz	Accuracy (absolute) ^b	Number of stations ^c	Coverage, km (nmi)
Loran C	H	0.100	460 m (1500 ft), 2 drms	60	2200 (1200)
Beacons	R	0.200–1.6	±3° to ±10°, 2 sigma	5000	370 (200)
VOR	R	108–118	±1.4°, 2 sigma ^d	2000	185 (100)
DME	C	960–1215	185 m (600 ft), 2 sigma	1000	185 (100)
TACAN	R, C	960–1215	±1.0°, 180 m (600 ft)	800	185 (100)
ILS (azimuth)	R	108–112	±7.6 m (±25 ft), 2 sigma ^e	2000	33 (33)
ILS (elevation)	R	329–335	±2.1 m (±7.0 ft), 2 sigma ^e	2000	33 (18)
MLS (azimuth)	R	5000–5250	±4.0 m (±13 ft), 2 sigma ^e	44	37 (20)
MLS (elevation)	R	5000–5250	±0.6 m (±2 ft), 2 sigma ^e	44	37 (20)
MLS (DME)	C	960–1220	30 m (100 ft), 2 sigma	44	37 (20)
JTIDS-RelNav	C, PR	960–1215	Variable	Variable	Variable
PLRS	C	420–450	15–25 m (50–80 ft)	Variable	Variable
SSR	R, C	1030, 1090	0.16°, 38 m (125 ft)	800	110–370 (60–200)
SSR/Mode S	R, C	1030, 1090	0.04°, 7 m (24 ft)	150	110–370 (60–200)
TCAS	R, C	1030, 1090	Variable	3500	Variable
GPS	PR	1575, 1227	13 m (43 ft) horizontal; 22 m (72 ft) vertical	24 ^f	Global
GLONASS	PR	1605–1250	30 m (100 ft)	14 ^f	Global
Altimeter	S	4200–4400	2% of altitude		
Mapping radar	S	9375, 9310 ^g	Variable		450 (250)
Map matching	S	Various	Variable		
Doppler radar	S	13,325	0.2% or 1.0 m/s (2 knots)		Global
Inertial	S		0.25–0.5 m/s (0.5–1 knot)		Global

^aH = hyperbolic; R = radial; C = circular, lines of position; PR = pseudorange; S = self-contained; D/H = Doppler history.
^b2 drms (twice distance root mean square) is the 95–98% probability, horizontal position error; 2 sigma is the 95% probability one-dimensional error.
^cWorldwide implementation in 2005.
^dSignal (flight check) accuracy.
^eAt Runway Threshold Approach Reference Datum.
^fSatellites.
^gBeacon frequency.

navigation systems spans the range from very low frequency to C band. See RADIO SPECTRUM ALLOCATIONS.

System classification. From an implementation viewpoint, all navigation systems can be classified as either cooperative or self-contained where cooperative systems exploit external radio or acoustic location signals. Examples of cooperative systems include the very high frequency (VHF) omnidirectional range (VOR), distance-measuring equipment (DME), TACAN, the instrument landing system (ILS), and the Global Positioning System (GPS). Examples of self-contained navigation systems include inertial navigation, Doppler speed-heading navigation, and data-base reference navigation systems (DBRNS). Alternatively, from a navigational viewpoint, systems are frequently classified as either positioning or dead-reckoning systems. Positioning systems determine an instantaneous location fix, whereas dead-reckoning systems accumulate position changes from a known initial location. Most positioning systems are cooperative systems, whereas most dead-reckoning systems are self-contained. The self-contained dead-reckoning systems may use electromagnetic or electromechanical technologies (Doppler radar, gyroscopes, accelerometers, or pressure transducers) to make observations from which the distance traveled and direction from the initial location are determined by continuous mathematical integration of velocity or acceleration. Many modern systems combine the data from cooperative and self-contained sensors to obtain a more accurate solution. Such systems have been called

multisensor, integrated, or hybrid systems. See ACCELEROMETER; DEAD RECKONING; GYROSCOPE; PRESSURE TRANSDUCER.

The major electronic navigation systems (see **table**) will be discussed.

Self-Contained Systems

These systems can be classified as radiating or non-radiating. Radiating systems may be subject to jamming and to homing by radiation-seeking missiles, although the effects of jamming on the systems discussed below are typically small. See ELECTRONIC WARFARE; HOMING; JAMMING.

Radiating systems. These include the radar altimeter, airborne mapping radar, map matching, and Doppler radar.

Radar altimeter. This system (also known as radio altimeter) is a small radar with an antenna mounted on the bottom of an aircraft generating a beam toward the Earth's surface. The signal backscattered by the surface is received and processed by the radar, which generates a reading of altitude above the surface. See ALTIMETER; GROUND PROXIMITY WARNING SYSTEM.

Doppler radar. Radio waves that are transmitted from a moving aircraft toward the ground and backscattered to the aircraft experience a change of frequency, or Doppler shift, directly proportional to the ground speed of the aircraft. A Doppler radar consists of a transmitter and receiver and a frequency tracker that measures and tracks the Doppler shift of the signals in each of three or four antenna beams directed toward the Earth's surface. The time

integral of the Doppler frequency is proportional to the change in position during the time interval. The Doppler velocity measurement is input to a dead-reckoning navigation computation.

Doppler radar navigation systems are widely used on military helicopters. They have also been used in soft lunar and planetary landings. The Doppler velocity measurement concept has been incorporated into airborne search and mapping radars and is also used in sonar systems to measure the ship's velocity. *See* DOPPLER EFFECT; DOPPLER RADAR.

Airborne mapping radars. These radars scan the ground by using specially shaped beams that effectively map the terrain. They can recognize certain features, for example rivers and bridges, and use them for manually or semiautomatic position-fixing updates to a data-base reference navigation system. The output of an airborne mapping radar or a radar altimeter can be used to generate a map or terrain profile, which is then compared with onboard stored maps or terrain profiles to allow an aircraft or missile to automatically fly a prescribed track. *See* AIRBORNE RADAR; GUIDED MISSILE.

Data-base reference navigation systems (DBRNS). A DBRNS utilize a priori maps of a quantity such as bathymetry, topography, gravity, magnetic intensity, radar, or visual images to determine position by correlating or "map matching" to the corresponding sensor measurement aboard the vehicle. Most often they utilize a dead-reckoning navigation system such as an inertial navigator to provide relative navigation between the measurements and thereby correlate a profile of the measurements to the map. The position fix, used to reset an inertial navigation or Doppler dead-reckoning system, is often registered to the midpoint time of the profile. Examples of a DBRNS include bathymetric navigation, digital scene matching area correlation (DSMAC), terrain contour matching (TERCOM), and precision terrain aided navigation (PTAN).

In bathymetric navigation, a fathometer is used to determine a ship's depth below keel. A profile of these measurements is compared with the dead-reckoning system's indicated profile as referenced to the stored a priori map. Accuracy is better than 200 m (650 ft) when accurate maps based on surveys, for example using GPS position as a reference, are available. Coverage depends on the availability of accurate maps and the bathymetry having adequate variations and asymmetry. Flat ocean areas are not amenable to bathymetric navigation. Sonar frequencies of 10–20 kHz are employed.

DSMAC is a target-area missile guidance system. DSMAC utilizes maps made from photographs of the target area that are digitized. Once in the vicinity of the target, DSMAC will match the digitized photograph with the its real-time camera images and correct the missile guidance inertial navigation system. Accuracy is better than 100 m (330 ft).

TERCOM uses radar and barometric altimetry to determine a three-dimensional position by correlating measurement terrain profiles with prestored terrain map profiles. It is analogous to bathymet-

ric navigation above, with the principle application being cruise missiles. Accuracy is better than 200 m (650 ft). Coverage depends on the availability of digitized land areas and the existence of adequate terrain signature.

PTAN is a planned system similar to TERCOM but with expanded coverage. Expanded coverage stems from terrain maps being available from the space shuttle or other space vehicles. PTAN is an alternative guidance technology to GPS for long-range cruise missiles.

Nonradiating systems. Nonradiating systems, such as inertial sensors, offer essentially complete protection against jamming. Inertial navigation equipment based on Newton's second law can determine aircraft acceleration and direction of travel. Acceleration may be determined by measuring the deflection of a spring attached to a known mass. The orientation of the accelerometers is determined by gyroscopes. The acceleration is then doubly integrated in the coordinate frame defined by the gyroscopes to determine the distance traveled since the last position fix. Thus, inertial navigation systems, consisting of accelerometers, gyroscopes, and computers, continuously determine position, velocity, acceleration, direction, and vehicle attitude (pitch and roll). *See* INERTIAL NAVIGATION SYSTEM.

Cooperative Systems

The two general categories of cooperative radio systems are point-source systems and multiple-source systems.

Point-source radio systems. Point-source systems typically determine position in a relative coordinate system by measuring the distance and bearing to a transmitting source at a known location. They may determine distance only or bearing only, if these are the only desired parameters. *See* RHO-THETA SYSTEM.

Bearing measurement. Perhaps the earliest form of a point-source radio system is the direction finder, in which the signal from a single transmitter source is received at two known points or elements. The direction from the vehicle to the source is determined by measuring the differential distance (as indicated by an electronic phase difference) of the signals received at the two points or elements. For operational convenience (size), it is desirable to have the two receiving points close together and to use common circuitry at both measuring points. A loop antenna fulfills both of these requirements. The loop is rotated until the currents in the two vertical arms of the loop are equal in amplitude and phase so that the output of the receiver is zero. The transmitter source is then located 90° from the plane of the loop. In simple loops, there is a 180° ambiguity, but this is typically resolved by temporarily connecting an omnidirectional antenna to the receiver input. As sources for shipboard and airborne direction finders, beacons are the oldest and most numerous navigation aids. Since the frequency bands of beacons are adjacent to the amplitude-modulation (AM) broadcasting band, receivers are easily designed to serve a dual purpose, and they are consequently widely used on

small boats and aircraft of all types. *See* DIRECTION-FINDING EQUIPMENT.

The angular measurement errors, particularly on aircraft and ships, can be quite large (3° or more). Furthermore, lateral navigational error decreases as the transmitter source is approached, a property common to all point-source angle-measuring systems. However, direction finders are still used, particularly for backup navigation and emergency homing, since only a single transmitter source (or beacon) on the ground and a simple rotating antenna and receiver on the vehicle are needed for operation.

Another class of point-source angle-measuring systems is based on the use of rotating or scanning antenna beams at the transmitter source and reception (or reflection) of the transmitted signal by the user vehicle receiver. For example, if a ground transmitter generates a rotating cardioid amplitude pattern at a fixed rate plus an omnidirectional reference signal, the user receiver can measure the relative phase difference between these two signals and thereby determine the bearing to the transmitter source. The very high frequency (VHF) omnidirectional range (VOR), used worldwide for short-distance aircraft navigation, is based on this principle. The bearing function of the military TACAN system is also based on this principle, except that in that system the source transmitter also generates a multilobe pattern for higher accuracy and the resulting ambiguity is resolved by an auxiliary technique.

VOR is the internationally standardized en route navigation system. It meets the Federal Aviation Administration (FAA) en route integrity requirements of a warning within 10 s if there is a failure of ground equipment. Many ground stations have converted to the Doppler mode, which reduces site error due to multipath. *See* DOPPLER VOR; VOR (VHF OMNIDIRECTIONAL RANGE).

Another point-source technique for bearing measurement is based on using a radiated dual-lobe structure, with the carrier in each lobe being modulated with different frequency tones and phases. The user receiver can then detect when the vehicle is operating at the intersection of the two beams. The instrument landing system (ILS), which is used worldwide for aircraft approach and landing, is based on that principle.

ILS is an internationally standardized system that provides a fixed electronic path at major runways throughout the world. It can be affected by close-by buildings and hills and does not allow curved approaches. *See* INSTRUMENT LANDING SYSTEM (ILS).

The known shortcomings of ILS with respect to obscuration and restriction to straight glide paths have led to the development of the microwave landing system (MLS). The use of high microwave frequencies, the nature of the ground-based antenna patterns, and the inclusion of a high-precision version of DME provide higher-accuracy performance and a capability for curved approaches. In view of the promising results with Local Area Differential GPS, also known as Local Area Augmentation Sys-

tem (LAAS), the development of MLS has been terminated in the United States. MLS is installed in about two dozen major airports in the United States and is scheduled to begin decommissioning in 2010. However, because of the unique requirements in the European airspace, acquisition and installation of MLS stations may proceed in European countries. Also, modular combined ILS/MLS/GPS receivers for approach and landing of aircraft are available. The U.S. military services have developed a mobile MLS (MMLS) that can be transported easily.

Secondary surveillance radar (SSR) is an outgrowth of the military Identification of Friend and Foe (IFF) system, using the same frequencies. It is the principal means by which air-traffic controllers identify and track civil and military aircraft worldwide. Beacon ground stations interrogate airborne transponders at 1030 MHz and receive replies at 1090 MHz in order to measure distance using two-way (round-trip) ranging and aircraft bearing using the SSR's narrow antenna beam. The replies are pulse-coded with identity (Mode A) and altitude (Mode C). Mode S, which has been added to the system, provides higher angular accuracy through monopulse techniques, discrete addressing of each aircraft, more efficient code modulation, and a much-higher-capacity data-link capability. *See* MONOPULSE RADAR; RADAR; SURVEILLANCE RADAR.

The Traffic Alert and Collision Avoidance System (TCAS) airborne equipment allows users to interrogate existing unmodified SSR 76 units and tracks the replies in range, bearing (in TCAS ID), and altitude. It generates crew alerts if both range and relative altitude are predicted to be so small that a collision threat exists. Two types of alerts are displayed: traffic advisories (TAs) and resolution advisories (RAs). The resolution advisory shows a vertical escape maneuver (for example "climb") for avoiding a collision. TCAS II uses Mode S, air-to-air, data-link communication for exchanging escape maneuver intentions between aircraft. The simpler TCAS I version, intended for smaller aircraft, does not display resolution advisories. It tracks proximate aircraft in range and altitude only, and displays traffic advisories for intruders that represent a collision threat based on relative range and range rate. *See* AIRCRAFT COLLISION AVOIDANCE SYSTEM.

Distance measurement. Theoretically, there are two possible types of point-source systems for distance determination. One is based on the two-way (round-trip) ranging principle (**Fig. 1**). The interrogator, which may be located on the navigating vehicle, transmits a signal, typically a pulse or pair of pulses, at a known time, which is stored in the equipment. The signal is received at a transponder and, after a fixed known delay, is retransmitted toward the interrogator. When received by the interrogator's receiver, it becomes an input to the ranging circuit. This circuit measures the time difference between the original transmission time and the time of reception (less the known fixed delay), which (when multiplied by the speed of light) is a direct measure of the two-way distance.

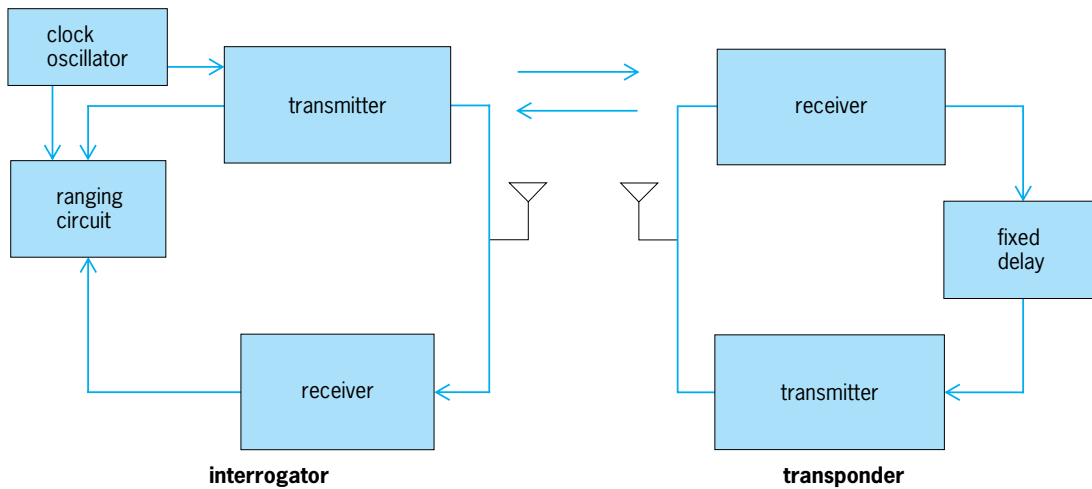


Fig. 1. Two-way (round-trip) ranging system.

An important advantage of this technique is that the signal returns to the point of initial generation for the time-difference measurement process. Therefore, the interrogator's clock oscillator need not be highly stable, since the error due to clock instability depends only on the round-trip time multiplied by the clock drift, and the round-trip time is very short since the signal travels at the speed of light.

Such a system is the basis of the distance-measuring equipment (DME), used for short-range navigation worldwide, and also of the Air-Traffic Control Radar Beacon System (ATCRBS), used for air-traffic control surveillance on a global basis. If the transponder is replaced by a passive reflector, for example an aircraft, the principle of operation is that of all primary surveillance and tracking radars, such as those used for air-traffic control, as well as those used for military applications in ground-based and airborne radars. See AIR-TRAFFIC CONTROL.

DME is an international standard which, together with VOR, has been the basis for the most widely used line-of-sight distance-bearing (rho-theta) aircraft navigation system. Airlines often use airborne equip-

ment in DME/DME mode, obtaining distance to two or more DME transponders rather than in VOR/DME, thereby achieving improved accuracy. See DISTANCE-MEASURING EQUIPMENT.

TACAN is a military modification of DME, using the same channels and adding a bearing capability. The result is a ground antenna system which is more portable than that of the VOR, a property that is particularly useful for systems on aircraft carriers that guide aircraft home. The Vortac system provides DME service to civil aircraft and rho-theta service to military aircraft. See TACAN.

A second possible technique for point-source distance determination is one-way signal-delay (time-of-arrival) measurement between a transmitter source at a known location and a user receiver. In this case, a successful measurement is possible only if the transmitter oscillator (clock) and the user receiver oscillator (clock) are precisely time-synchronized (Fig. 2). Otherwise, true range (delay) cannot be determined, since the time of transmission is not known with respect to the user's clock. Since independent precise time synchronization between two equipments

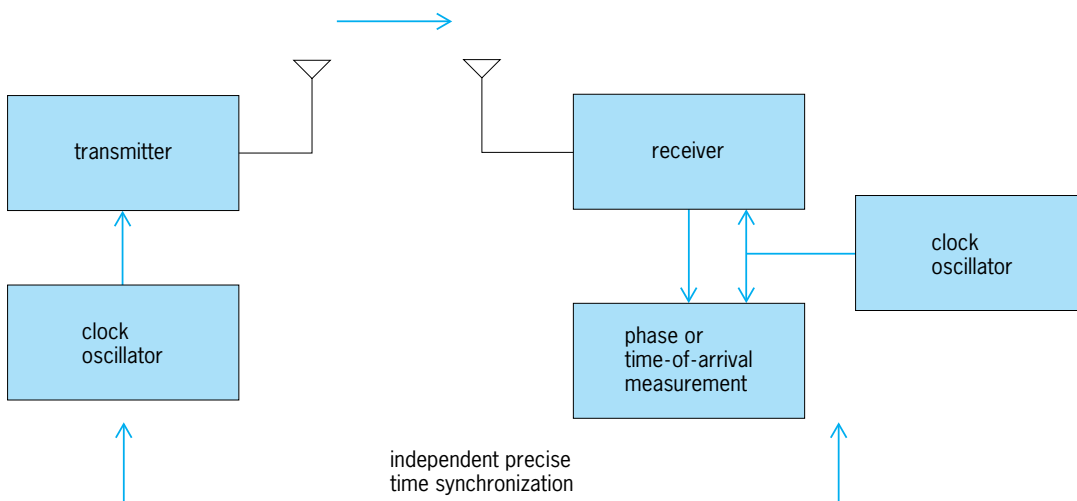


Fig. 2. One-way synchronous ranging system.

is frequently impossible at reasonable equipment cost, no practical point-source distance measurement system based on such synchronization has been developed, but several modern multisensor systems (discussed below) have modes that employ this principle.

Multiple-source radio systems. Multiple-source radio systems determine present position (and vehicle velocity) in an absolute or relative coordinate system by means of measurements to several signal sources at known locations. These systems therefore employ multiple transmitter sources, and user vehicle equipment typically consisting of a receiver or a receiver-transmitter. The four major categories of such systems (with some implementations using combinations of these) are hyperbolic systems, pseudorange systems, one-way synchronous ranging (direct-ranging) systems, and two-way (round-trip) ranging systems.

Hyperbolic systems. These were the first to be developed, but they are still in widespread use. Three or more transmitter sources transmit time-synchronized continuous-wave or pulsed signals. The minimum-size chain (configuration of transmitters), a triad, usually consists of one so-called master and two secondary (slave) stations (Fig. 3). The user receiver measures time differences of arrival of signals from pairs of stations. Loci of constant time differences, or (equivalently) constant differences in distance, from two stations form hyperbolic lines of position (LOPs). The point where two lines of position cross is the position of the user vehicle.

One major advantage of this technique is that the user needs only a receiver and the receiver does not need an expensive high-quality clock oscillator, since only time differences are used. Theoretically, three pairs of sources are needed for a completely unique horizontal position fix, but in practice two pairs usually suffice because the user's altitude is known to acceptable accuracy. The differences in distance can be measured either in terms of differences in times of arrival (for example, of pulses) or in terms of differences in electrical carrier phase, or both. Achievable accuracy is very much a function of the relative geometric location of the sources and the user, and largely depends on the crossing

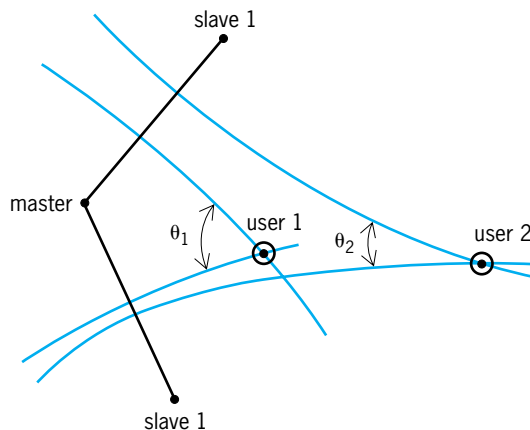


Fig. 3. Effect of relative geometry on accuracy of hyperbolic navigation systems. Accuracy largely depends on θ_1 and θ_2 , the crossing angles at users 1 and 2.

angles of the lines of position (Fig. 3). For a given level of measurement noise, the smaller the crossing angles, the larger the position error. This accuracy degradation is called geometric dilution of precision (GDOP), and is essentially a multiplier on the basic time-difference or line-of-position root-mean-square measurement error. The GDOP is approximately $1/\sin \theta$, where θ is the crossing angle. The Loran C, Omega, and Decca systems are based on this hyperbolic principle, as were the now-terminated Omega and Decca systems. See HYPERBOLIC NAVIGATION SYSTEM.

The Loran C hyperbolic system was originally intended for shipping, both military and civil. The combination of low frequency and pulses provides long range and virtually eliminates multipath propagation interference from sky waves. The use of a form of phase measurement of the carrier cycles provides additional accuracy, and it inherently provides a fine-coarse readout of low inherent ambiguity. Since all Loran-C stations include an atomic clock for timing, precise time and time-interval information can be obtained from the system. For example, telephone systems use this capability for digital switching, signaling, and timing. Loran - C and GPS have been used in hybrid configurations to provide backup in case of outages, obtain precise time, increase overall system integrity, and provide a reference for differential GPS operation. These uses are likely to continue for some years, until GPS is widely accepted as the sole means of navigation. The use of digital computation in Loran C receivers has resulted in a tremendous growth in general aviation use. See LORAN.

Pseudorange. This concept has some similarity in configuration to hyperbolic systems but does not use time differences. Several transmitter sources (for example, space vehicles or terrestrial stations), whose positions are made known to the user, transmit highly time-synchronized signals based on an established system time. With the transmission epochs of the signals from the sources provided or known to the user, the user measures the time of arrival (TOA) of each signal with respect to the user's own clock time, which normally has some offset from system time. The measured time of arrival minus the signal transmission time defines a range measurement that is called pseudorange, since it differs from the true range as a result of the offset of the user's clock relative to the system (transmitter) time. From successive or simultaneous time-of-arrival (pseudorange) measurements from four (or more) sources, the user receiver then calculates the three-dimensional position coordinates and its own clock offset (from system time). This is accomplished by solving four simultaneous quadratic equations (usually the computations are linearized), involving the three known position coordinates of the sources and the four unknowns, namely, the three user position coordinates and the user's clock offset.

Thus, such a system accurately determines not only the user's three-dimensional position but also the user clock offset relative to system time, which is easily related to Universal Time Coordinated (UTC). In addition, by properly combining the known

velocity of the sources and the measured Doppler shift of the signals received from the sources, the three user velocity coordinates and the user clock frequency offset can also be determined. In order to achieve high accuracy in this type of system, the transmitter sources must be highly synchronized through use of clock oscillators of very high stability, for example, atomic clocks. However, the user receiver clock oscillator does not need to be highly stable. *See* ATOMIC CLOCK; DOPPLER EFFECT; TIME.

The pseudorange concept is the basis of operation of two major satellite navigation systems for worldwide use, the U.S. Global Positioning System (GPS) and the Russian Global Navigation Satellite System (GLONASS). In these systems, the satellite-borne clocks are rubidium or cesium atomic clocks, with long-term stabilities of the order of 1 part in 10^{13} . In order to provide high-accuracy time-of-arrival measurement capability, these satellite navigation systems use wide-bandwidth spread-spectrum modulation methods. The pseudorange concept has also been applied in terrestrial navigation systems, notably in one mode of the Joint Tactical Information Distribution System-Relative Navigation (JTIDS-RelNav) system. *See* SATELLITE NAVIGATION SYSTEMS.

Global Positioning System (GPS). This satellite-based radio system was originally developed by the U.S. Air Force to provide worldwide coverage; high-accuracy, three-dimensional position, velocity, and time; and completely passive (receiver-only) operation by all types of dynamic users. It has found wide acceptance by both military and civil users. Two services are available, the Precise Position Service (PPS) for authorized users and the Standard Positioning Service (SPS) for all other users. The root-mean-square value of three-dimensional positioning errors for users of the Precise Positioning Service is specified to be less than 15 m (50 ft), although the achieved accuracy is of the order of 5 m (15 ft) or better in most circumstances. The accuracy of the Standard Positioning Service (SPS) was improved on May 1, 2000, when the intentional degradation of service accuracy known as Selective Availability was discontinued. The specification value of horizontal positioning error that is obtained 95% of the time was reduced by this action from 100 m (330 ft) to 20–30 m (70–100 ft). The accuracy achieved since then is of the order of 5–10 m (15–30 ft), and further reduction is anticipated when additional civil signals become fully operational as soon as 2012.

The specified GPS consists of 24 satellites oriented in six orbital planes including four satellites each. The operational system has typically included 28 or more satellites because of the longer-than-anticipated lifetime of the satellites. Satellites are in 12-h orbits, with four satellites in six orbit planes inclined at 55° , all at an orbital altitude of 20,200 km (10,900 nmi). The satellites transmit highly synchronized, pseudonoise-coded (wide-bandwidth) signals, including data on their ephemerides and clock errors. The user receiver determines at least four pseudoranges by time-of-arrival measurements with respect to its own clock time, and four pseudor-

ange rates (delta ranges) via Doppler measurements with respect to its own clock frequency. From these measurements, the user receiver computes its own three-dimensional position coordinates and its clock bias error, as well as (in some receivers) its three-dimensional velocity coordinates and clock frequency offset. The satellites transmit two codes, the 1.023-megabit-per-second (Mbps) C/A code, which can be tracked by civil users, and the 10.23-Mbps P (or Y) code, which is encrypted for exclusive use by military users. Each satellite has a unique code that distinguishes its signals from the other satellite transmissions that share the GPS frequency allocation by means of code-division multiplexing. System data are modulo-2 added to both codes at 50 bits per second. *See* MULTIPLEXING AND MULTIPLE ACCESS; SPREAD SPECTRUM COMMUNICATION.

Some applications, for example, aircraft landing and airport surveillance, require even higher accuracies than those available from the basic GPS service. For these, the differential GPS (DGPS) concept has been implemented, employing a reference station, whose position is precisely known. Relative positioning accuracy of 3 m (10 ft) is easily achievable with this technique. The Local Area Augmentation System (LAAS) is under development by the FAA to provide DGPS services at airports. Other augmentations to the civil use of GPS include the Wide Area Augmentation System (WAAS), in which two INMARSAT satellites transmit ionospheric delay corrections to users in the contiguous United States (CONUS) that have been determined by combining GPS measurements made at dozens of sites throughout the North America; Maritime and Nationwide DGPS (MDGPS and NDGPS), which transmit pseudorange and (sometimes carrier phase) corrections on the 285–325-kHz maritime radio beacon band; and Continuously Operating Reference Stations (CORS), which are operated by the National Geodetic Survey to support DGPS operation by surveyors, Geographic Information Systems providers, engineers, and scientists, by transmitting non-real-time GPS carrier phase and pseudorange measurements made at precisely surveyed locations.

Global Navigation Satellite System (GLONASS). This system, developed by the Russian government, is very similar in concept to the GPS. It uses 24 satellites in 25,500-km (13,800-nmi) orbits, with eight satellites in three orbital planes at an inclination of 64.8° . This inclination is somewhat higher than that of GPS, thus providing better polar coverage. GLONASS uses frequency-division multiplexing to distinguish signals broadcast from individual satellites. The future of GLONASS will likely be tied to international unification of Global Navigation Satellite System (GNSS) services as they relate to GPS and Galileo.

Galileo. Galileo is an initiative of the European Union, in collaboration with the European Space Agency and European Industry, to launch a European-financed global satellite navigation system under civilian control. It will have fundamental similarities to GPS with additional services for search and rescue and a guaranteed precise-accuracy service. Galileo will consist of a constellation of 30 satellites

in three orbital planes inclined at 54° and at an altitude of around 23,000 km (12,400 nmi). The compatibility, particularly with respect to military frequency allocations, between GPS and Galileo has been a major diplomatic issue. Although the two systems were distinct as of 2006, it is likely that future satellite navigation users will design receivers to utilize both simultaneously, which would further strengthen the dominance of satellite navigation for most terrestrial applications.

JTIDS. The Joint Tactical Information Distribution System (JTIDS) is a decentralized, military, spread-spectrum, data-communication and navigation system, using wide-bandwidth phase coding, frequency hopping, and time-division multiple access. The RelNav function permits all members of a JTIDS network, such as aircraft and ships, to accurately determine their position in both absolute and relative grid coordinates by means of highly precise time-of-arrival measurement of signals received by cooperating units when any member is located in an absolute reference system. Otherwise, all members of a net are positioned relative to one another. The system includes a means for independent precise time synchronization of each unit, and operates in both one-way synchronous ranging and pseudorange modes.

JTIDS with the RelNav function is implemented in a large number of military aircraft and ships of the United States and a number of other countries. An international effort called the Multifunction Information Distribution System (MIDS) is under way toward further miniaturization and modernization of JTIDS terminals. In most aircraft applications, JTIDS-RelNav and an inertial navigation system (INS) are used in a fully integrated multisensor configuration. Combined JTIDS-RelNav/INS/GPS is planned in order to make use of the optimum combination of measurements from these three types of navigation sensors.

One-way synchronous ranging. The implementation of this concept between one of the transmitter sources and the user receiver was discussed above (Fig. 2). In order for a true one-way range measurement to be made, the clock of the user receiver must first be synchronized to that of the transmitter source. In some systems, this is accomplished by an independent means, for example, in the JTIDS-RelNav system and in the Position Location Reporting System (PLRS). The concept is also used in one mode of certain hyperbolic systems, notably Loran C and the decommissioned Omega. When applied to those systems, the concept has been called direct ranging. Two implementations are possible, the rho-rho method or the multiple-rho method.

The rho-rho method requires only two source transmitters, but also requires a highly stable user receiver oscillator (clock) and precise knowledge of the time of transmission from the source stations. A direct range is then developed from each station. The lines of position are circles rather than hyperbolas, with the intersection of the circles representing the user position, thereby leading to more favorable geometric conditions.

The multiple-rho method is an extension of the rho-rho scheme, requiring at least three stations. Using three lines of position permits clock oscillator self-calibration, much as the previously discussed pseudorange concept, and therefore leads to a much less stringent clock oscillator requirement. Again, circular lines of position lead to higher position accuracy because of better geometric behavior.

Two-way (round-trip) ranging. This technique involves multiple two-way distance measurements of the type discussed above (Fig. 1). To obtain a completely unambiguous horizontal position fix, three such two-way ranges are required; however, two are usually sufficient. Since the lines of position are again circular, the geometric accuracy behavior is generally more favorable than for hyperbolic systems. The major disadvantage of this concept is that the user requires a transmitter as well as a receiver (Fig. 1). A typical application of this concept is DME-DME or multiple-DME operation for civil aviation navigation. Other applications are in the military PLRS and for the tracking of aircraft in military training ranges.

PLRS. This is a centralized military, spread-spectrum, position-location and navigation system for aircraft, land vehicles, and personnel. It uses wide-bandwidth phase coding, frequency hopping, and time-division multiple access, and also incorporates data communications capability. It provides military commanders with information on the location of all their elements, as well as their own positions and relative guidance information to each unit. Its operation is based on multilateration using highly precise time-of-arrival measurements of signals exchanged between units in both two-way (round-trip) ranging and one-way synchronous ranging modes.

PLRS and EPLRS, an enhanced version which provides more direct user-to-user communication, are in wide use by the U.S. Marine Corps and Army, respectively. A trend toward distributing the position- and range-bearing computations from the central processors to the individual user units is supported by the decreasing cost of digital processors and by the worldwide availability of GPS. A GPS receiver module can be embedded into a PLRS user unit, and GPS time-of-arrival-based positions can be combined with PLRS time-of-arrival-based measurements, resulting in the increased accuracy and availability of user positioning and navigation information.

Walter R. Fried; Richard L. Greenspan; Marvin May
Bibliography. P. Enge and P. Misra, *Global Positioning System, Signals, Measurements and Performance*, 2d ed., 2006; J. Farrell, *Integrated Navigation Systems*, 5th ed., 2003; E. D. Kaplan and C. J. Hegarty (eds.), *Understanding GPS: Principles and Applications*, 2d ed., 2006; M. Kayton and W. R. Fried, *Avionics Navigation Systems*, 2d ed., 1997; B. W. Parkinson and J. J. Spilker (eds.), *Global Positioning System: Theory and Applications*, 2 vols., 1996; P. G. Savage, *Strapdown Analytics*, 2000; U.S. Department of Defense, *CJCS Master Positioning, Navigation, and Timing Plan*, biannually; U.S. Departments of Transportation and Defense, *Federal Radionavigation Plan*, biannually.

Electronic packaging

The technology relating to the establishment of electrical interconnections and appropriate housing for electrical circuitry. In particular, packaging for microelectronics focuses on interconnections between integrated circuits to achieve higher levels of assembly and electronic function, such as the storage and processing of information.

Electronic packages provide four major functions: interconnection of electrical signals, mechanical protection of circuits, distribution of electrical energy (that is, power) for circuit function, and dissipation of heat generated by circuit function. Thus, electronic packaging is a truly multidisciplinary science, involving such areas as physics, materials science, mechanical and electrical engineering, and computer-aided design.

Printed circuitry. As solid-state transistors started to replace vacuum-tube technology, it became possible for electronic components, such as resistors, capacitors, and diodes, to be mounted directly by their leads into printed circuit boards or cards, thus establishing a fundamental building block or level of packaging that is still in use.

Printed circuitry consists of a relatively stiff core of electrically insulating material in sheet form, with a patterned conducting layer on one or both sides. Holes drilled through the sheet at appropriate locations provide for mounting component leads that are subsequently soldered to ensure both mechanical and electrical connection to the package. Edges of printed circuit cards often provide metal lands (pads) for interconnection to other components. Consumer electronics, such as home video games, make extensive use of very basic printed circuit cards.

Although several different choices are available, the most common material set for printed circuit cards is a fiberglass core with copper conductor. Early patterns were formed by simple silk-screening and etching techniques, but these methods were later replaced by sophisticated photolithographic techniques to achieve very fine dimensions. For advanced applications, cards are now built in multiple layers that consist of conductor layers with intervening insulating layers, with interconnections known as through-vias (electrically conductive, orthogonal pathways between two parallel but isolated conducting layers). These vias are formed by drilling fine holes through the fabricated multilayer card and plating the inside surface of the holes, thus electrically connecting those conductor lines intersected by the drilled holes. The need for increasingly finer features (closer spacing of printed lines and smaller-diameter vias) has brought about significant advances in printed circuit technology, including use of laser-drilled vias and multilayered thin-film patterns on top and bottom of the printed circuit cards. The so-called microvia technology is the most advanced in the printed circuit industry. *See* PRINTED CIRCUIT.

Packaging hierarchy. Complex electronic functions often require more individual components than

can be interconnected on a single printed circuit card. Multilayer card capability was accompanied by development of three-dimensional packaging of so-called daughter cards onto multilayer mother boards, providing interconnections for both power and card signals. Developed for the large mainframe computers during the period 1960–1985, this type of packaging hierarchy is now used for many personal computer applications.

Integrated circuitry allows many of the discrete circuit elements such as resistors and diodes to be embedded into individual, relatively small components known as integrated circuit chips or die. In spite of incredible circuit integration, however, more than one packaging level is typically required, driven in part because of the technology of integrated circuits itself.

Integrated circuit chips are quite fragile, with extremely small terminals. First-level packaging achieves the major functions of mechanically protecting, cooling, and providing capability for electrical connections to the delicate integrated circuit chip. At least one additional packaging level, such as a printed circuit card, is utilized, as some components (high-power resistors, mechanical switches, capacitors) are not readily integrated onto a chip. For complex applications, such as mainframe computers or network servers, a hierarchy of multiple packaging levels is required. A common packaging hierarchy consists of single-chip modules (first-level) soldered to multilayer printed circuit cards (second-level), which in turn are mechanically interconnected to system-level printed circuit boards (third-level) [Fig. 1]. *See* INTEGRATED CIRCUITS.

First-level packaging. In microelectronics, the first-level package is that which provides interconnection directly to the integrated circuit chip. The carrier upon which the chip rests (the substrate) must also provide for interconnection to the next (that is, second-level) package. Major emphasis has been placed on first-level packaging since the early 1970s.

Chip-to-package interconnection is typically achieved by one of three techniques: wire bond, tape-automated bond, and solder-ball flip chip (Fig. 2).

Wire bonding. The wire bond is the most widely used first-level interconnection; it employs ultrasonic energy to weld very fine wires mechanically from metallized terminal pads along the periphery of the integrated circuit chip to corresponding bonding pads on the surface of the substrate. Typical wires are made from aluminum or gold, with small alloying additions to achieve the desired strength required for handling. Wire diameters of only 25–50 micrometers (0.001–0.002 in.) are used because of the small size of the terminals of the integrated circuit chips. Metallization of the substrate and the pad of the integrated circuit chip is typically gold plating and evaporated aluminum, respectively; both are often modified with small alloying additions to customize properties. Bond pads are only a few micrometers thick and are on the order of 50–200 μm (0.002–0.008 in.) in diameter. Wire-bonded integrated circuit chips are

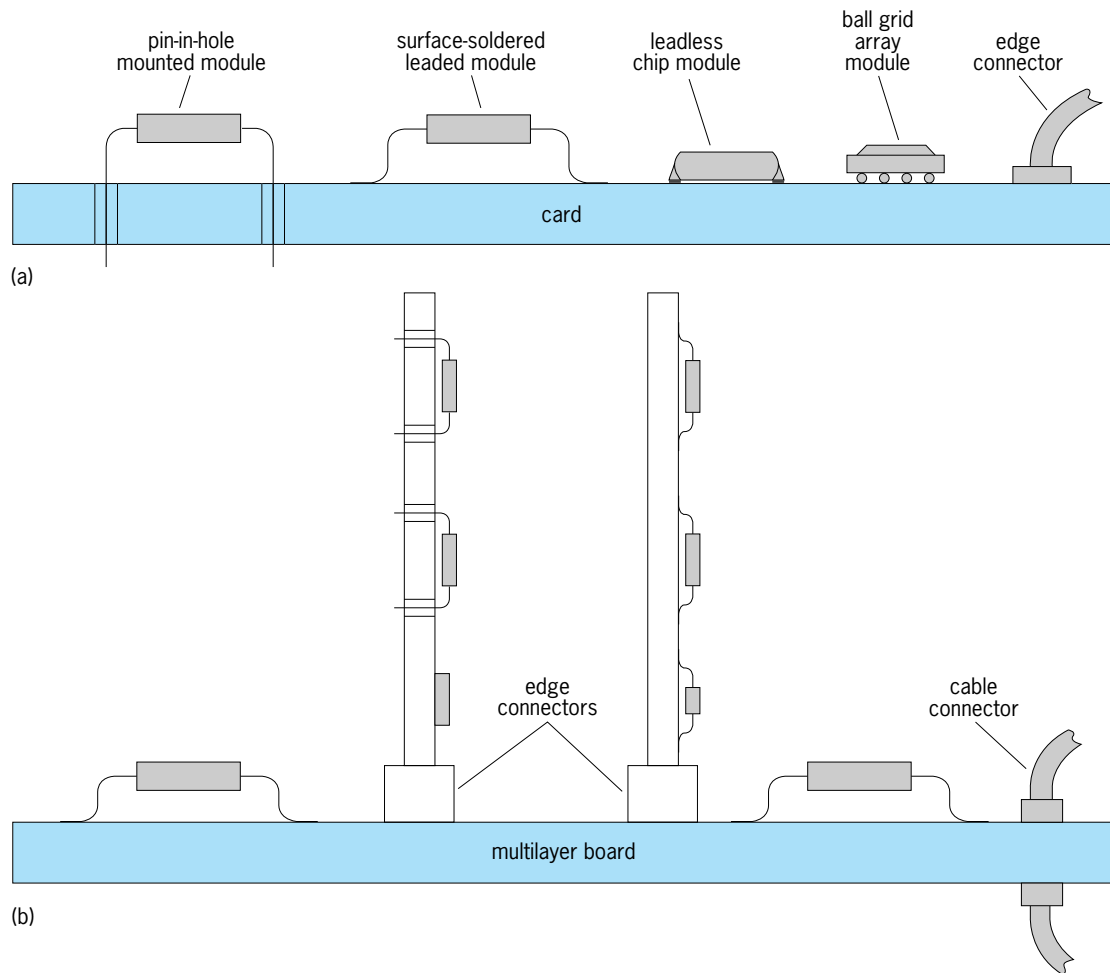


Fig. 1. Packaging hierarchy. (a) A card and single-chip modules. (b) Multilayer board.

attached to the substrate with active circuitry facing out, such that the back of each chip provides a surface for adhesive bonding. Heat is dissipated through the back-bonded interface.

Tape-automated bond. In the tape-automated bond, photolithographically defined gold-plated copper leads are formed on a polyimide carrier that is usually handled like 35-mm photographic roll film, with perforated edges to reel the film or tape. The inner leads are connected to the integrated circuit chip through simultaneous thermocompression bonding of all leads to corresponding gold bump contacts on the terminals of the integrated circuit chips (gang bonding). Outer leads are soldered, thermocompression-bonded, or bonded ultrasonically to the substrate pads, which are similar to wire bond pads in structure.

Solder-ball flip chip. This technique involves the formation of solder bump contacts on the terminals of the integrated chips and reflowing the solder with the chip flipped in such a way that the bump contacts touch and wet to matching pads on the substrate. Lead-tin solder rich in lead is usually used, with flux to ensure adequate wetting to the surface of the substrate. In contrast to wire bonding, an area array configuration rather than peripheral leads is possi-

ble, leading to much greater density of interconnecting terminals per unit area of substrate. These chips are soldered with the circuitry facing toward the substrate (hence the name flip chip) in contrast to more common wire bonding.

Several alternate approaches to achieving flip chip, other than solder ball, are currently practiced. These typically include some form of metallic bump on the integrated circuit die connected to matching pads on the substrate with electrically conductive polymer adhesives.

Substrates. Substrates for first-level packages are quite varied; selection for a given application depends on several factors, including the required number of interconnections, electrical characteristics, reliability objectives, and cost. A major classification is whether the package supports a single integrated circuit chip (single-chip module) or more than one chip (multichip module). The former is by far the most common. Substrate insulator materials for multichip and single-chip modules are selected from one of two broad groups of materials, organics (including epoxy and polyimide) and ceramics (predominantly alumina, but also including silicon carbide and aluminum nitride). Conductors are chosen from a wide range of metals. Specific materials are selected on the

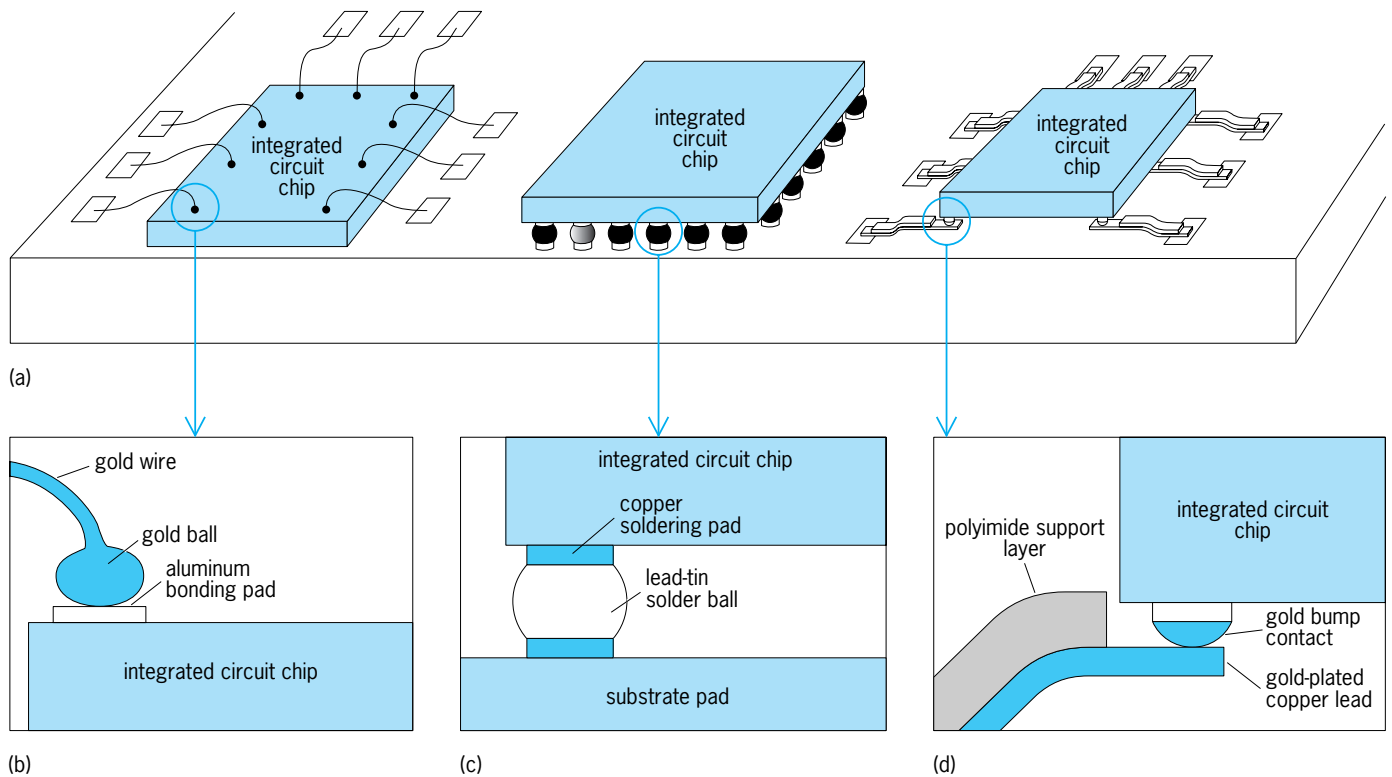


Fig. 2. First-level packaging. (a) Substrate. (b) Wire bond. (c) Solder-ball flip chip. (d) Tape-automated bond.

basis of electrical characteristics (for instance, low-dielectric-strength insulators and low-resistivity conductors) as well as properties required for processability and reliability, such as thermal coefficients of expansion, thermal conductivity, processing temperatures, and thermal stability. *See* CERAMICS; HETERO-CYCLIC POLYMER; POLYETHER RESINS.

Analogous to the advances in printed circuitry, early first-level package substrates were fabricated with a single layer of conductor paste, screened, and then fired onto small pressed ceramic tiles prefabricated with holes through which pins were swaged to provide for second-level interconnections. Ceramic packages today can be made with upward of 60–70 alternating conductor and insulator layers, connected by conductor-filled vias. These complex multilayer ceramics, typically for applications involving multichip modules, are fabricated as green (that is, pressed but unfired particulate) laminates and then sintered to near-theoretical density. Surface metallization, required for first- and second-level interconnections, has advanced so that it can include photolithographically defined patterns.

Significant advances in organic-based substrate technology, notably those utilizing ball grid array interconnect, have led to multilayer plastic packaging for single-chip modules that is competitive in some aspects with traditional ceramic packaging. These substrates utilize many of the same fabrication techniques applied in advanced printed circuit technology discussed earlier.

Interconnection. Means for interconnecting to the second-level package often dictate the general form

of the first-level package (Fig. 3). For instance, if the connection is to be pin-in-hole, the package must have pin-shaped leads that are either peripheral, such as a dual in-line package, or in an area array, such as the pin or ball grid array. If the package is to be soldered on the surface, it can be a leaded carrier, in which leads of various geometries are soldered to pads on the card, or a leadless carrier, in which solder from the card contacts metal pads on the edge of the package. Temperature hierarchy in processing must follow the packaging hierarchy, that is, assembly of first- to second-level packages should occur at

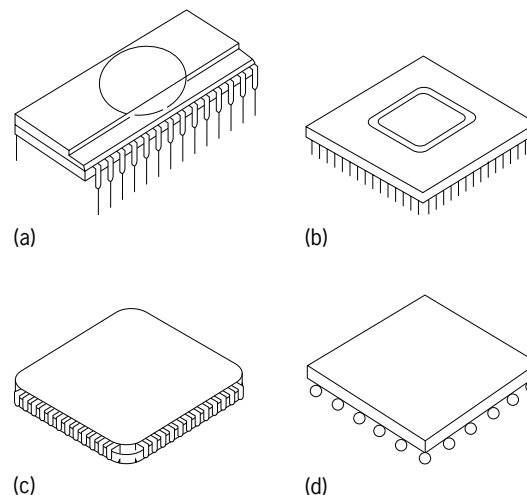


Fig. 3. Typical first-level packages. (a) Dual in-line package. (b) Pin grid array. (c) Leadless chip carrier. (d) Ball grid array.

process temperatures below that at which the chip-to first-level package interconnection was achieved. When more than two packaging levels are required, this guideline can become difficult to follow.

Field data have demonstrated that the most vulnerable portion of the package is the interconnection, either chip-to-substrate or module-to-board. Mismatch between the integrated circuit chip and the substrate materials and substrate-to-card materials caused by thermal expansion can generate cyclic fatigue stresses during machine cycles (that is, each time an electronic system is turned on and off, the integrated circuit chip and its package heat up and cool down). Reliability of both the integrated circuit chip and the interconnections can be enhanced by sealing out environmental moisture. Sealing (encapsulation) techniques range from inexpensive epoxies to high-performance hermetic sealing using soldered or glass-sealed caps to protect the integrated circuit chip and leads. Finally, heat sinks are often attached to the packages to assist in the dissipation of heat generated by operation of the integrated circuit chip. See ELECTRONICS.

Peter J. Brofman

Bibliography. C. A. Harper, *Electronic Packaging and Interconnection Handbook*, 1991; J. Hayward, The problem of interconnect and the impact of printed circuit board technology on integrated circuit package development, *Future Circ. Int.* (London), no. 4, pp. 19-22, 1998; K. Matsuda et al., Simple method for flip-chip bonding on organic substrates, *Int. J. Microcirc. Electr. Packag.*, vol. 20, no. 3, 1997; D. P. Seraphim, R. C. Lasky and C. Li, *Principles of Electronic Packaging: Design and Materials Science*, 1989; R. Tummala and E. J. Rymaszewski, *Microelectronics Packaging Handbook*, 2d ed., 1995.

Electronic power supply

A source of electric power (voltage and current) to operate electronic circuits. Active electronic circuits contain such devices as transistors or vacuum tubes and require external power to amplify, filter, modify, or create electrical signals. The most common source of energy for electronic circuits is obtained by converting the electrical energy available in the conventional alternating-current (ac) electric power mains to an appropriate voltage or current. These converters, or electronic power supplies, can be implemented with a wide variety of circuits. Other power sources include batteries, mechanically driven generators, photovoltaic (solar) cells, and fuel cells. See ALTERNATING CURRENT; CONVERTER.

Most electronic circuits require a direct-current (dc) or constant voltage. If ac power is required, an oscillator or a simple transformer is used. Although some dc-to-dc converters are used, most dc power supplies convert the alternating power from the ac main to dc power. These ac-to-dc power supplies are classified according to the type of circuits used to realize the conversion: rectification, filtering, and regulation. Simple battery chargers are examples of

power supplies that do not require filtering or regulation.

Rectification. An essential step in the conversion of ac to dc is a process called rectification. Rectification converts ac voltage to a waveform with average or dc value by passing only one polarity (half-wave) or by generating the magnitude or absolute value (full-wave). Three types of rectifier (diode) circuits are commonly used. Only one diode is required to obtain half-wave rectification. The diode passes current in one direction and blocks current flow in the other direction. Full-wave rectification can be obtained with four diodes connected in a bridge configuration or with two diodes and a center-tapped transformer (**Fig. 1**). The diodes permit current to flow from each half of the transformer secondary, but only in one direction through the load. The resulting waveform is filtered with a capacitor, and an optional regulator may be added for additional refinement of the dc output voltage. Transformers are normally used at the input of the rectifiers to increase or decrease the voltage and isolate the dc output from the ac input for safety purposes. See DIODE; RECTIFIER; TRANSFORMER.

Filtering. For most applications the ac or alternating portion of the rectified output is unwanted and may cause undesirable results, such as an annoying hum in audio systems. A capacitor can be used to reduce or filter the ac portion of the rectified waveform (**Fig. 1**). The capacitor is charged through the diodes to the peak ac voltage minus the diode forward voltage. Some of the charge stored on the capacitor is delivered to the load each cycle, but the next voltage peak recharges the capacitor. For a capacitance (C), the basic relationship between current (i) and voltage (v), Eq. (1), can be used to estimate the remaining peak-to-peak ripple, according to Eq. (2).

$$i = \frac{Cdv}{dt} \quad (1)$$

$$\text{Peak-to-peak ripple} = V_{pp} \cong \frac{I_{dc}}{fC} = \frac{E_{dc}}{fCR_L} \quad (2)$$

Here, the frequency, f , is the input frequency to the filter, I_{dc} and E_{dc} are the dc current and voltage, and R_L is the load resistance. The value for f can range from 50 Hz for a half-wave rectified European ac line, to several hundred kilohertz for an electronically switched power supply. For a full-wave rectified ac line in the United States, $f = 120$ Hz. Electrolytic capacitors are normally used for ac power-line applications but are not adequate at high frequencies. Other low-pass filter circuits using combinations of capacitors and inductors can be used to pass the dc and attenuate the ac part of the waveform. See ELECTRIC FILTER.

Regulation. Regulators are often used to make the power supply output insensitive to input voltage amplitude variations and further reduce the ripple voltage. The regulator may also be used to adjust or change the dc output voltage and limit the amount of current delivered by the power supply. Regulators are a form of dc-to-dc converter.

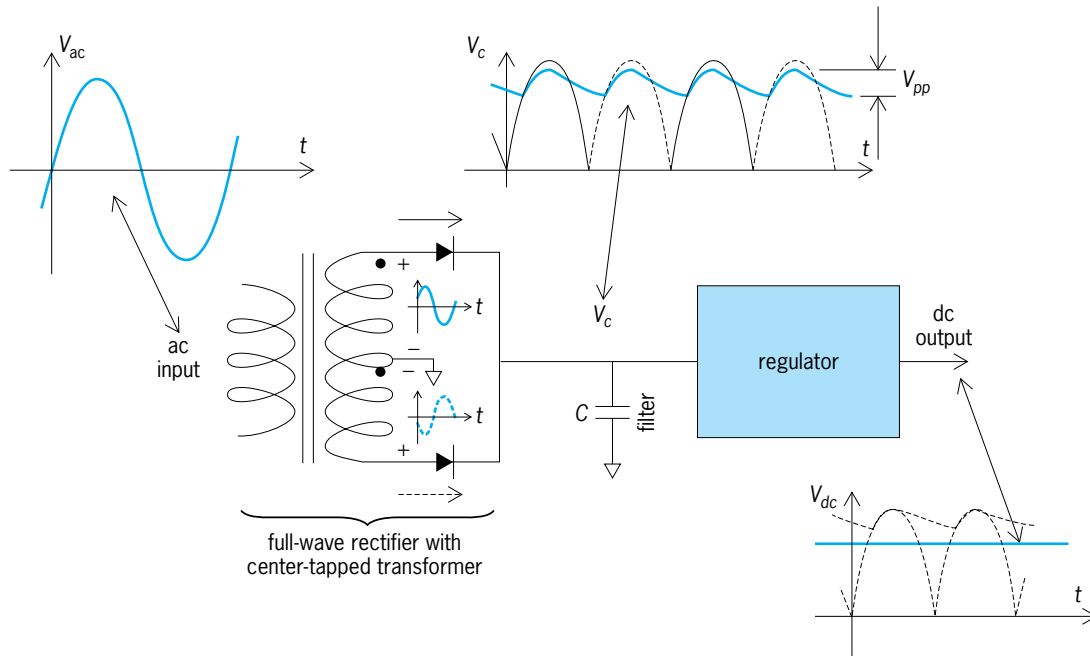


Fig. 1. Typical electronic power supply.

Linear regulators. The oldest and simplest type of regulator is the linear regulator. A simple linear regulator is the shunt type (Fig. 2a). It consists of a Zener diode ($D1$) and a current-limiting resistor (R_s). The Zener diode establishes a fixed, or reference, voltage if it is properly reverse biased. Another common type of linear regulator is the series-pass regulator (Fig. 2b). A negative feedback control system automatically controls the input to the series-pass device (such as an *npn* bipolar transistor, $Q1$) to maintain the output equal to the reference or desired level, V_R . The power dissipation in the R_s , $D1$, or $Q1$ significantly reduces the efficiency of a power supply. See TRANSISTOR; ZENER DIODE.

Switching regulators. More modern power supplies have switching regulators, and are informally called switchers. There are more than a dozen different topologies (basic block diagrams) for switching regulators, but they all use one or more transistors acting as switches; either ON or OFF. In addition to the solid-state (transistor) switches, a typical switching regulator uses capacitors and inductors to store energy and diodes to direct the current. A negative-feedback control system is used to set and maintain the output level, similar to the linear series regulator. For

a switching regulator, the feedback system controls the frequency, or duty cycle, of the switch control voltage. For a buck-or-forward switching regulator (Fig. 3), the output voltage must be smaller than the input dc voltage. When the switch ($Q1$) is on, the current increases through an inductor ($L1$), providing more current to the load, and charges the capacitor ($C1$). When the output voltage rises above the desired value, the control system turns the transistor ($Q1$) off. The current through the inductor continues to flow in the same direction but now through a diode ($D1$). The amplitude of the current decreases. The voltage drops slightly until the control system turns the switch on again. This cycle is repeated automatically to maintain a nearly constant dc output voltage. Other common types of switching regulators are called buck-boost (or flyback) and push-pull or (buck-derived). Switching regulators can be used to generate multiple outputs at different voltage levels. The improved efficiency of switching regulators is due to the fact that energy is stored very efficiently in inductors and capacitors. The remaining losses in the control circuitry, switches, and the diodes are small compared to linear regulators. See INDUCTOR; VOLTAGE REGULATOR.

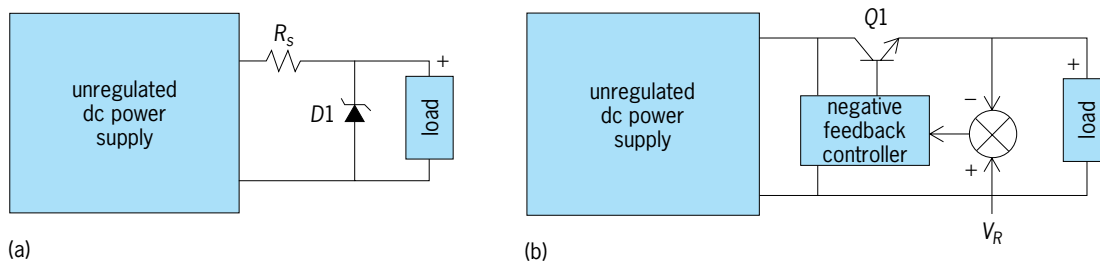


Fig. 2. Linear voltage regulators. (a) Zener-diode shunt type. (b) Transistor series-pass type.

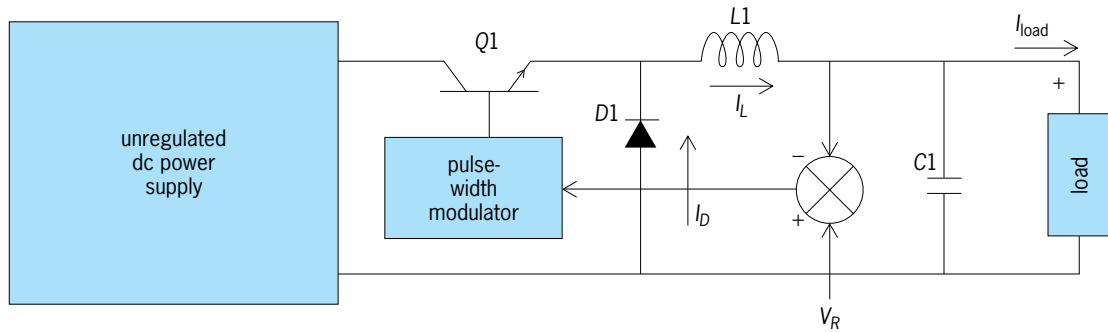


Fig. 3. Simplified circuit diagram of buck-or-forward switching regulator.

Linear versus switching supplies. The major drawback of linear power supplies is the power dissipation of the control device. For example, all of the load current must pass through the pass transistor in a series-pass regulator (Fig. 2b). The resulting power dissipation reduces the efficiency of the supply. Efficiencies of linear dc power supplies range 30–60% for typical output voltages and currents, while switching power supply efficiencies are typically 70–80%. Small linear regulators are relatively inexpensive. Completely integrated linear regulators that can carry up to 1.5 A are available in inexpensive plastic packages. Metal packages may be used for currents in the 5-A range, but power dissipation is a problem for higher-current applications. Size and weight is a problem with larger linear power supplies because large heat sinks and 60-Hz (or 50-Hz) input ac line transformers are required. The typical output power for a 10 cm × 10 cm × 10 cm (4 in. × 4 in. × 4 in.) linear supply is 15 W, while a switching power supply of this size can easily have an output power of more than 150 W. Most switching power supplies use small high-frequency transformers after the switch for isolation, so a large power-line-frequency transformer is not required. Switching regulators operate at a typical frequency of 40 kHz, so much smaller filter capacitors can be used. The increased efficiency of the switching regulator is due to the fact that inductors and capacitors store energy with very little power dissipation or energy loss. Also, the switching devices dissipate very little energy because they are either on or off. When they are off there is no current, and when they carry current in the on state the voltage is very small; thus the average power dissipation is small. Most of the power loss occurs in the transition between states. Therefore, device switching speed is very important. Using faster devices such as field-effect transistors (FETs), insulated-gate bipolar transistors (IGBTs), and metal-oxide-semiconductor (MOS) controlled thyristors results in improved performance over bipolar junction transistors (BJTs). Faster devices allow a higher switching frequency, which results in smaller filter capacitors and transformers.

Other power supplies. Ferroresonant transformer-based power supplies have some advantages for high-current applications such as battery chargers. Ferroresonant power supplies use nonlinear magnetic

properties and a resonant circuit to regulate the output voltage, and they have efficiencies similar to switching power supplies. Power supplies driven by three-, six-, or twelve-phase ac are easier to filter, and they generate much lower harmonic distortion of the current in the ac power system.

Alternating-current power system quality is degraded by most electronic power supplies. This effect may become a serious issue as more electric vehicles are charged in the future or in buildings where electronic systems are the dominant ac load. Flyback power supplies and voltage-multiplier circuits are often used for high-voltage dc power supplies that require less current. National and international regulations for safety, power-line harmonics, and electromagnetic and radio-frequency interference (EMI-RFI) are a very important consideration for modern electronic power supplies of all types. See ELECTRICAL INTERFERENCE; VOLTAGE-MULTIPLIER CIRCUIT.

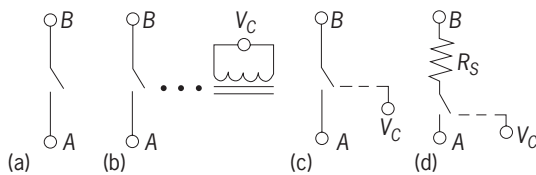
N. G. Dillman

Bibliography. S. G. Burns and P. R. Bond, *Principles of Electronic Circuits*, 2d ed., 1994; G. C. Chryssis, *High-Frequency Switching Power Supplies*, 2d ed., 1989; R. C. Dorf (ed.), *The Electrical Engineering Handbook*, 2d ed., 1997; M. J. Fisher, *Power Electronics*, 1991; D. A. Grant and J. Gowar, *Power MOSFETS, Theory and Applications*, 1989; A. I. Pressman, *Switching Power Supply Design*, 2d ed., 1997; R. S. Ramshaw, *Power Electronics Semiconductor Switches*, 2d ed., 1993.

Electronic switch

An electronic device in which one or more input signals can be routed to one or more outputs by the application of the appropriate electrical control signals. The term is most often applied when analog signals are involved, but the terminology is occasionally used when digital signals are involved.

Operation. Conceptually, an electronic switch can be visualized as a group of one or more mechanical electrical switches (such as light switches used in commercial wiring or toggle switches used in many electronic control panels) in which, instead of mechanically opening or closing the contacts, the physical opening and closing is achieved by applying appropriate electrical control signals to separate



Switching devices. (a) Symbol for mechanical switch. (b) Symbol for electromechanical relay. (c) Symbol for electronic switch. (d) Simple model of electronic switch.

terminals on the switch in much the same way that a relay performs. Unlike the mechanical switch (illus. *a*) or the electromechanical relay (illus. *b*), however, the electronic switch does not contain mechanical contacts but semiconductor devices such as bipolar junction transistors or field-effect transistors. The basic electronic switch is depicted in illus. *c*. *A* and *B* are the terminals of the switch. When a control signal is applied to V_C , the switch closes. When the electronic switch is closed, a small residual resistance R_s remains between the terminals as depicted in the simple model of illus. *d*. The value of this resistance is termed the on resistance. This resistance varies considerably from one electronic switch to another and depends on the manufacturer of the switch. In most applications the nonzero on resistance does not prove problematic, but the user needs to be aware of this limitation. The electronic switch is typically bidirectional in the sense that the terminals *A* and *B* are interchangeable. See ELECTRIC SWITCH; RELAY; TRANSISTOR.

Advantages and disadvantages. Electronic switches offer attractive alternatives to mechanical switches and mechanical relays in several respects. They can be very small, allowing a large number of these devices to be placed in a small area; they are often significantly less expensive than their mechanical counterparts, particularly when small currents are being switched; they can be very fast, with on and off response times which are orders of magnitude faster than can be achieved with mechanical counterparts; they often require only minuscule amounts of energy at the input V_C to control the switch; the control inputs are often directly compatible with the voltage and current levels inherent in most electronic equipment; and they are generally considerably more reliable over a large number of cycles than their mechanical counterparts. On the other hand, they often have more stringent limitations on single voltage or current levels than do mechanical switches, and they do become more costly as the current- or voltage-handling capability increases.

Implementation. Electronic switches are available in stand-alone packages, often termed solid-state switches or analog switches. The term multiplexer, analog multiplexer, or data selector refers to a special class of multiple-input electronic switches in which either digital or analog signals are routed in a predetermined way to one or more outputs, often in a unilateral way (inputs and outputs not interchangeable). The electronic switching or data selection function can also be achieved with relatively

simple circuits involving transistors or operational amplifiers along with a few passive components. See MULTIPLEXING AND MULTIPLE ACCESS; OPERATIONAL AMPLIFIER.
 Randall L. Geiger

Electronic test equipment

Apparatus used for the evaluation of electronic components, subassemblies, and systems.

Early electrical instruments. The galvanometer was an early electromechanical instrument that converted current into angular rotation of a needle pointer across the face of a scale. The instrument consisted of a fine wire coil mechanically connected to a pivoting element. Also connected to the element were the needle pointer and a spring, which held the pointer to an initial position on the scale. The coil was located in the field of a permanent magnet so that when current flowed through the coil the induced magnetic field rotated the element, causing the needle to move across the scale. Greater current caused a larger deflection. This instrument could also be used as a voltmeter, since a higher voltage induces a larger current flow through the coil. See AMMETER; GALVANOMETER.

The advent of the cathode-ray tube (CRT) enabled the development of another widely used instrument, the oscilloscope. The CRT translates voltage into the deflection of an electron beam, which visibly activates a luminescent phosphor inside the face of the tube. Vertical and horizontal deflection plates correspond to *x* and *y* movement of the beam. Usually a sawtooth-shape waveform is connected to the horizontal plates, so that the *x* direction represents time. Therefore, when a voltage is connected to the vertical plates, the visual representation on the CRT face has the amplitude and polarity of the voltage displayed as a function of time. See CATHODE-RAY TUBE; OSCILLOSCOPE.

An important advance in electronic test equipment was the incorporation of circuits that directly converted the analog signal to be measured into a digital reading. Initially, these improvements provided easier and more accurate operator interpretation of the instrument's measurement. The most significant impact of these converters was to enable a computer interface to the equipment to be set up. Many new instruments were developed with computer bus capabilities, allowing direct computer monitoring and control of the instrument. Other equipments were designed with embedded computers, which provided very sophisticated analysis of the data within the instrument itself. See ANALOG-TO-DIGITAL CONVERTER; DIGITAL-TO-ANALOG CONVERTER; EMBEDDED SYSTEMS.

The best example of this enhancement is with spectrum analyzers. Digital signal processing (DSP) electronics have been included in these instruments which provide a frequency spectral analysis of the input signal. These instruments are widely used in the evaluation of radio communication signals, acoustic signal processing, and the analysis of

mechanical forces and vibration. *See* ACOUSTIC SIGNAL PROCESSING; SPECTRUM ANALYZER.

Automatic test equipment. The complexity and sophistication of electronic components also brought about the evolution of the test equipment. As components and subassemblies became more complicated, individual pieces of test equipment such as voltmeters, oscilloscopes, and signal generators could no longer efficiently evaluate the performance of the unit under test. The individual test instruments were integrated to form rack-and-stack automatic test equipment, where the individual instruments were stacked on top of each other and connected to a control computer by a common bus. After the initial test programming effort, these computer-controlled testers were very efficient at providing the input stimulus to the unit under test and monitoring the output response. Improved microelectronic components led to increasing capabilities, speed, and accuracy of stimulus and monitoring instruments, and component cost reductions eventually led to the complete integration of equipments into one complex automatic-test-equipment system (see *illus.*). *See* ELECTRONICS; INTEGRATED CIRCUITS; LOGIC CIRCUITS; MICROPROCESSOR; SIGNAL GENERATOR.

Characterization versus production testing. Two fundamentally different types of testing can require different types of electronic test equipment. Design verification or characterization tests are usually done on a small sample of units to validate the design or the process and to determine how much performance margin exists. Extensive data are taken, but the testing is a nonrecurring task, so test time and data acquisition efficiency are not too important. Product-specific test fixtures are designed, and general-purpose instruments are usually employed in an engineering laboratory environment, where the equipment can be shared or redeployed from one program to another.



Automatic test equipment for testing microelectronic devices. (Rockwell International)

Production testing, however, is done on all the units, so speed and efficiency are important. There is usually little data gathering, since the results of the tests are simply "pass" or "fail." Although there may be many different tests at various parametric conditions within one production test program, they are done in quick succession, and the tests typically require only a few seconds to carry out. Thus, only a few expensive production test systems can test thousands of units per day. The importance of short unit test times is reflected in the design of automatic test equipment. Stimulus- and response-measurement subsystems have been developed with ever briefer settling times (the time required to stabilize after a change in mode), and faster test-control computers have been included in automatic test equipment.

Test diagnosis. The ability of automatic test equipment to sort out good from bad units is only the first requirement of modern production test equipment. Since the repair of defective units can be costly, the design and program development of test equipment frequently must include special provisions to provide failure-mode analysis. Information such as the probable defective component on a printed circuit board or the probable defective board in a system is of great value in efficient repair. Some test systems include fault dictionaries, supplemental tests after first failure, or even artificial-intelligence features to assist in repair. *See* ARTIFICIAL INTELLIGENCE; FAULT ANALYSIS.

For microelectronic devices, where repair is usually impractical, failure-mode analysis is used to improve the manufacturing process. Bad devices can be tested under relaxed conditions to determine if they have a gross defect such as a short or an open, or if they have a parametric failure where they do not operate at rated speed or voltage. This information is used as feedback to the manufacturing process, where adjustments can be made to improve the yield, or number of good units, per manufacturing batch.

Quality level. The economic aspects of electronic system manufacturing have dictated the test thoroughness or quality level that must be achieved by the test equipment. In general, the cost to diagnose and repair a unit under test increases an order of magnitude for each successive level in manufacturing. If it costs \$1 to detect a bad device, it would cost \$10 to detect and replace that device in the next-higher assembly, the printed circuit board; it would cost \$100 to detect and replace the bad board at the system level, and \$1000 to repair that system once installed in the field. Therefore, the earliest detection of a defect is the most economical. Component customers who formerly asked for defect levels slightly less than a 1% now expect defect rates on the order of a few parts per million. The test equipment and the overall test strategy have been affected by this expectation of high quality level. The biggest impact is that testers with higher speed and improved accuracy are required. Worst-case conditions such as high temperature or low voltage must be determined for the units to be tested, and then all units are subject to

these conditions in a single test or even a multiple-pass test. Test equipment must be integrated with artificial environment chambers for the component, printed-circuit-board, and system tests. *See ENVIRONMENTAL TEST; PRINTED CIRCUIT.*

The importance of electronic test equipment far exceeds that of the simple measurements or even the complex components for which they have been designed. The common electronic automotive instrument cluster, for example, combines measurements of speed, angular speed, fuel level, and engine performance in one easily read panel, but behind that panel are the fundamental elements of test instruments.

William R. Mann

Bibliography. K. Brindley, *Automatic Test Equipment*, 1991; M. Collons, *Computer Controlled Testing and Instrumentation: An Introduction to the IEC-625: IEEE 488 Bus*, 1983; Hewlett Packard, *Test and Measurement Catalog*, 1995; R. J. Traister, *Meters and Scopes: How to Use Test Equipment*, 1988; D. Wobschall, *Circuit Design for Electronic Instrumentation*, 2d ed., 1987.

Electronic warfare

The art and science of preserving the use of the electromagnetic spectrum for friendly forces while denying its use to the enemy. This is accomplished via use of techniques, devices, and equipment by an adversary to deny or counteract an enemy's use of radar, communications, guidance, or other radiowave devices. The principal techniques that are employed are called electronic countermeasures (ECM). Because of the growing use of optical and infrared techniques for communications, guidance, detection, and control, they are sometimes called electromagnetic, rather than electronic, countermeasures to convey more adequately the idea that countermeasures are not confined to the portion of the spectrum where electronic techniques alone are applicable but may be used throughout the electromagnetic spectrum. Other techniques include the use of anti-radiation weapons designed specifically to home on and destroy radars, and the attempts to reduce the electromagnetic signatures of aircraft, missiles, and ground vehicles.

In a broader sense, electromagnetic warfare refers to the worldwide struggle by major military powers for supremacy in the electromagnetic spectrum. This largely silent, secret struggle occurs continually, independent of whether its participants are at war with one another, as was the case during the Cold War. Participants seek control of the spectrum because of the increasing dependence of military forces on its use for surveillance, communications, detection, measurement, guidance, and control. With mastery of the spectrum, one adversary could achieve an indispensable ingredient for conquering an enemy or discouraging a potential aggressor.

Since the late 1990s, there has been a close association of electronic warfare with information warfare (IW). Information warfare can be described as

taking action to protect one's information systems from exploitation and corruption while simultaneously exploiting, corrupting, or destroying the adversary's information systems. Electronic warfare is considered the implementation tool or active component of information warfare.

Categories of ECM. Traditionally, ECM equipment and techniques are categorized as active and passive, depending on whether or not they radiate their own energy. The passive category includes reconnaissance or surveillance equipment that detects and analyzes electromagnetic radiation from radar and communications transmitters in a potential enemy's aircraft, missiles, ships, satellites, ground installations, and communications systems. The reconnaissance devices may be used to identify and map the location of emitters without in any way altering the nature of the signals they receive. Other types of passive ECM do seek to enhance or change the nature of the energy reflected back to enemy radars, but they do not generate their own energy.

Active ECM equipment generates energy, either in the form of noise to confuse an enemy's electromagnetic sensors or by radiating false or time-delayed signals to deceive radio or radar equipment and their operators. *See ELECTRICAL NOISE.*

Surveillance systems. Reconnaissance or surveillance ECM systems are carried by Earth-orbiting satellites, aircraft, ships, crewless (drone) aircraft (commonly known as unmanned aerial vehicles or UAVs), and automotive vehicles. Some are located on the ground (**Fig. 1**). A few systems are small enough to be carried by a person. Reconnaissance systems are

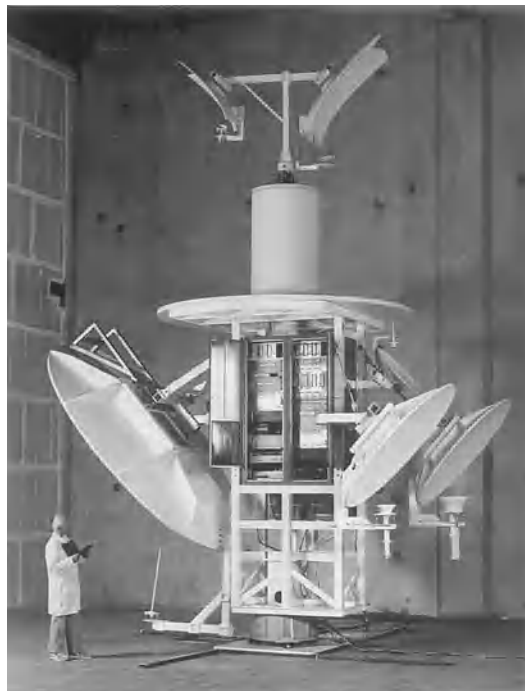


Fig. 1. Fully automatic land-based ELINT system, which monitors and documents the nature and location of hostile electromagnetic emitters and detects any changes in their location. (*Litton Applied Technology*)

interchangeably called ferret or electronic intelligence (ELINT) systems. *See* DRONE; MILITARY SATELLITES.

They consist of sensitive receivers electromechanically or electronically tuned over desired portions of the spectrum in search of transmissions of interest. Bearing to an intercepted signal can be determined by direction-finding techniques. Once secured, the signals can be displayed for analysis by an operator or stored on tape or other storage media for subsequent analysis or both. *See* DIRECTION-FINDING EQUIPMENT.

The type of information obtained by modern surveillance systems is illustrated by a standard system in use by the United States Air Force. It has a superheterodyne receiver and five separate radio-frequency tuning heads that span the spectrum from 1 to 18 GHz. (More advanced systems extend the range up to 40 GHz to cope with higher-frequency threats.) The receiver sequentially scans through those five heads at sweep rates up to 20 Hz, selectable by an operator. The operator has a panoramic display of the five separate amplitude-versus-frequency traces projected by a multigun cathode-ray tube. When a signal of interest is observed, a cursor can be superimposed over it, causing the receiver to discontinue the sweep mode and lock onto the intercepted signal. At this point, the display automatically switches into an analytical mode, offering an enlarged time display with repetitive pulsed signals on five separate logarithmically spaced time scales. Then the pulse shape, its repetition rate, and pulse width can be visually identified. Separately, an automatic digital indication of the signal's frequency in 100-kHz increments, pulse duration to 0.1 microsecond, and pulse repetition period to 1 microsecond can be obtained. *See* RADIO; RADIO RECEIVER.

Not all reconnaissance equipment is intended for spectral observation. Another form of reconnaissance equipment, carried by some U.S. Air Force aircraft, aids photo interpretation by detecting and partially identifying electromagnetic signals. Intelligence information gathered by the set is directly recorded on the film of a reconnaissance camera. The system, which covers six military frequency bands, locates an emitter in any one of eight sectors of the photograph. It helps in identifying the frequency, pulse-repetition frequency, and pulse width of the unknown intercepted radar signal.

Radar warning receivers. Radar warning-receiver (RWR) systems, which came into widespread use on United States tactical and transport aircraft during the Vietnamese war, are a more limited form of the ELINT system, intended as a means of self-protection. Unlike the latter, which is intended to search for signals over a broad range of the spectrum, the warning receiver is programmed to alert a pilot whose aircraft is being illuminated by a specific radar signal above predetermined power thresholds anticipated by ELINT systems. A typical airborne threat-warning system consists of a processor, receivers, power supply, antennas, display, and controls, and generates both visual and aural warnings of radar threats. When

the pilot has been alerted, the aircraft can be maneuvered to evade the threat or initiate counteraction with onboard ECM capability. The principal radar threats, which appeared in the Vietnamese war and in the Middle East wars of 1971 and 1973, and again during Operation Desert Storm in the Persian Gulf in 1991, were Soviet-made surface-to-air missile (SAM) radar and anti-aircraft gunfire direction radars. Although both are ground-based threats, the warning-receiver concept is equally applicable to radar threats from opposing aircraft. Similarly, warning receivers are carried by ships or other vehicles.

Airborne warning systems usually are crystal video or superheterodyne receivers that indicate the presence of a radar threat and offer a coarse relative bearing to it. An associated receiver and extra antennas augment this capability, automatically determining direction to the threat and giving the pilot a means of homing on the target for weapons delivery by nulling an indicator that centers the aircraft on the signal.

The warning systems also exploit a particular weakness in SAM weapons, which requires that a command radar be turned on prior to missile launch so that the missile can be guided to its airborne target after launch. The receiver can detect characteristic changes in power level of the radar, thereby warning the pilot that a missile is about to be fired. SAM commanders often counter this technique by simply turning on the command radar frequently even when there is no impending target intercept. This can prompt aircraft to prematurely jettison fuel tanks or weapons in an effort to evade the supposed imminent launch of a SAM missile. *See* GUIDANCE SYSTEMS; GUIDED MISSILE; MISSILE.

These procedures illustrate how each advance in radar technology has a countermeasure and how each countermeasure is followed by a counter-countermeasure. Usually the countermeasure and the counter-countermeasure involve a technological advance, but frequently simple ingenuity or changes in tactics supply the new measure.

Intruding aircraft have a time advantage over defense radars because they can detect enemy radar at considerably greater distances than the hostile radar can spot them. The reason for this is that energy emitted by radar decays as a function of the square of the range from the radar to the target. *See* RADAR.

As technology progresses, the difference between the warning system and limited-capability ELINT receivers is narrowing with the advent of small digital processors, particularly for tactical military aircraft. The tactical aircraft needs to be able to react quickly to changes in power levels of specific radar threats, movement of transportable SAM radars, and frequency shifts by frequency agility or frequency hop radars.

Reflectors. One of the oldest passive ECM techniques is the use of chaff. These are metallic strips cut to lengths resonant at the defense radar frequency so that they return spurious radar echoes to enemy radar (**Fig. 2**). Chaff can confuse an enemy by generating false targets, or noise, forcing the enemy to take time to analyze the returns and sort real from false



Fig. 2. Metallized glass chaff, which returns spurious radar echoes. (Lundy Electronics and Systems)

targets. Chaff can screen or mask aircraft or higher-speed ships so that the enemy is unable to determine their presence, or can help an aircraft break track once it is alerted by its warning receiver that it is being tracked by radar.

Chaff has evolved from the aluminum strips hand thrown from bombers during World War II to 0.001-in.-diameter (25-micrometer) metallized glass fibers that are automatically ejected in great quantities from aircraft by electromechanical, pyrotechnic, or pneumatic dispensers. Typically, dispensing mechanisms aboard modern military aircraft forcibly eject cardboard packages, roughly $3 \times 5 \times \frac{1}{2}$ in. ($8 \times 13 \times 1.3$ cm) either manually as initiated by the pilot or automatically if the RWR is interfaced directly with the chaff/flare system. Each package, containing thousands of strands of aluminum-coated glass fibers, is torn open at release. The contents then blossom into clouds within the aircraft's slipstream.

In Vietnam, naval pathfinding aircraft would seed air corridors with chaff through which attacking aircraft would subsequently fly. The chaff masked their specific presence from enemy radars. In addition, strategic bombers are outfitted with chaff-dispensing rockets fired ahead of the bombers to assure that an adequate screen blooms before the aircraft arrives.

By the mid-1970s, U.S. Air Force and Navy tactical fighters were being outfitted to carry large chaff pods capable of dispensing 300 lb (136 kg) of lightweight chaff contained on six separate drums and blown out through separate exit pipes at the rear of the pod (Fig. 3). The dispensing rate and the duration and frequency of chaff bursts can be programmed into the pod on the flight line.

To assure availability aboard an aircraft of chaff cut to proper lengths, the military services have inves-

tigated dispensing systems that cut chaff to lengths comparable to the frequencies of radars that warning receivers indicate are illuminating the aircraft. A technique known as delayed-opening chaff (DOC), by which chaff packages are deployed from aircraft or ballistic missiles approaching their targets, also was developed. A DOC package is descended by parachute and then blown open by a time fuse. This might deceive an enemy into believing that there are targets far removed from the dispensing aircraft.

Other passive ECM techniques include the use of special radar-absorbing materials, such as ceramics or ferrites, which reduce reflection coefficients so that the amount of radar energy returned to the illuminating radar is reduced. A standard technique is the special shaping of bodies, specifically in missile reentry systems and airframes, that reduces the vehicle's radar cross section (the area of the vehicle that effectively returns echoes to the transmitter). The Air Force has successfully applied these and other signature-reducing techniques, called stealth, to diminish the radar observability of such aircraft as the F-117 fighter and the B-2 bomber. During the Gulf war, the F-117 achieved great success in Desert Storm, and was credited with covering 40% of the strategic targets in Iraq while flying only 2% of the combat sorties by coalition aircraft. The aircraft evaded radar and infrared tracking threats and was credited with being the only aircraft to strike in heavily defended downtown Baghdad. See MILITARY AIRCRAFT; RADAR-ABSORBING MATERIALS.

The use of corner reflectors, or Luneberg lenses, which concentrate the energy they reflect back to the radar, can provide a useful ECM aid by confusing an enemy. These reflectors, fixed to drone aircraft

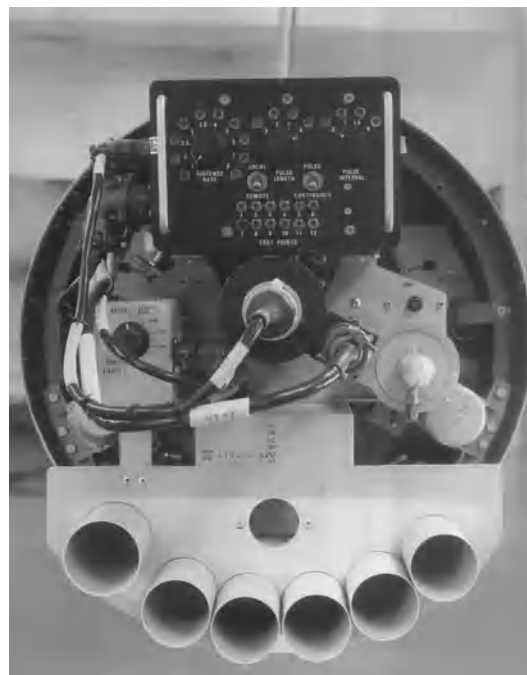


Fig. 3. Bulk-chaff dispensing pod which ejects chaff through six exit pipes.



Fig. 4. Countermeasures-dispensing system aboard an air transport used in humanitarian relief efforts ejecting flares to foil attacks by ground-launched infrared guided missiles. (Lockheed Aircraft Service Co.)

released from bombers, produce large returns on enemy radar screens, causing the drones to appear on radar like the larger bombers that released them. See ANTENNA (ELECTROMAGNETISM).

Infrared protection. United States aircraft carry sensitive infrared ECM receivers to protect them from heat-seeking infrared missiles. Because such missiles normally are launched toward the rear of the target so that they can home on the heat from the aircraft's jet exhaust, the receivers are located on the aircraft's tail. When the pilot is warned of impending danger from the infrared missile, infrared flares can be ejected through the chaff dispenser either automatically or manually in order to divert the heat-seeking missile.

Following the loss of an Italian transport to infrared guided missiles during a humanitarian relief mission into Sarajevo, Bosnia, in 1993, various United Nations countries began equipping their air-relief transports with missile-warning and flare-dispensing systems (Fig. 4). The flares are intended to lure away infrared guided missiles fired from shoulder-held launchers as the aircraft are descending on approach to their destinations.

U.S. Navy and Army aircraft are equipped with active infrared transmission devices designed to induce infrared missiles tracking an aircraft's engine exhaust to break target lock. These transmitters generate high-data-rate infrared pulses to jam signals derived within the missile seeker from the rotating action of the seeker's reticle in chopping incoming infrared energy.

U.S. Army and Marine rotary-wing aircraft have been outfitted with infrared countermeasures systems for protection from heat-seeking missiles, especially at low altitudes where they are vulnerable to attack by ground-fired, infrared-guided missiles (Fig. 5). The energy source in the transmitter of one such system radiates heat to foil heat-seeking missiles. The heat pattern generated by the 13-in.-

high (33-cm), 28-lb (13-kg) transmitter confuses the missile's guidance and tracking system, throwing the threatening weapon off course.

Active ECM. The many active ECM techniques can be classified broadly either as noise or deception jamming. The former is the oldest, simplest, and most straightforward, but requires higher average power levels and is more expensive. Deception jamming is the more artful and sophisticated technique, operating on the characteristics of the pulse train generated by threat radars.

Noise jamming has a disadvantage in that it can alert enemies to the fact that they are being jammed. Yet in tactical situations where the enemy has little reaction time and noise contributes to confusion, this distinction has little importance.

Communications-jamming techniques include use of white-noise jamming, in which artificially produced white noise modulates the output of a transmitter. This is one of the simplest and easiest approaches, despite its excessive use of power. White-noise jamming is one of the more prevalent forms of noise jamming. In sweep-through jamming the transmitted carrier frequency is swept through a part of the spectrum at a high rate, producing pulses masking the incoming signal. This makes more economical use of transmitter power. See ELECTRICAL NOISE.

Deception-jamming techniques are predicated on the idea of operating on pulses received from the enemy so that the signal reradiated from the target deceives the enemy radar or its operators (Fig. 6). For instance, the ECM deception set may receive an enemy radar pulse, circulate it through a delay line, amplify it, and reradiate it back toward the enemy. Because the enemy determines the position of the target by the round-trip transit time of the radar energy, its radar decision circuitry will conclude that the target is at a greater distance than it actually is because of the deceptive pulse delay inserted in that round-trip period by the active ECM set. This technique is known as range gate walk-off (RGWO). Velocity gate walk-off (VGWO) is accomplished by applying some form of frequency translation technique such as serrodyning to the returned pulse. Serrodyning is an operation employed by the velocity deception electronic countermeasures (ECM) jammer to counter those radars utilizing a target's



Fig. 5. Transmitter of infrared countermeasures system, providing 360° protection, mounted atop fuselage of UH-1H utility helicopter. (Sanders Associates)

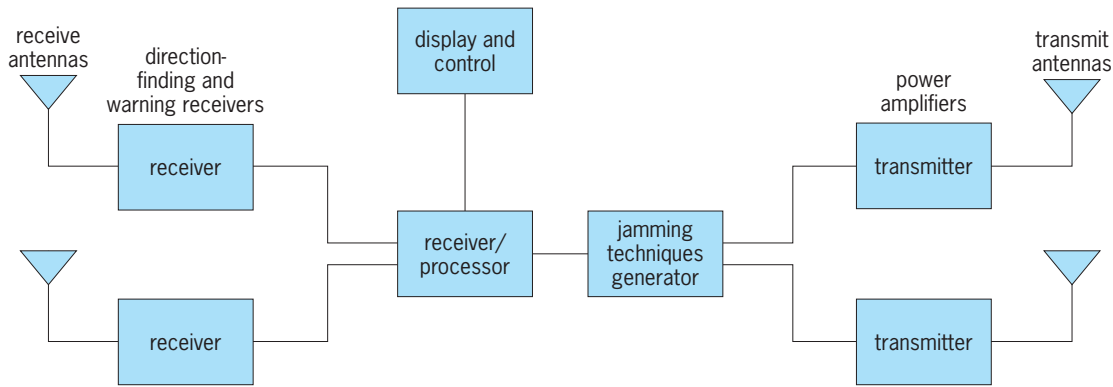


Fig. 6. Typical pulse repeater jamming system.

Doppler shift information. The jammer will frequently translate an incoming signal by a desired amount and at a specific rate during retransmission. If the radar follows the ECM signal, it is left without a target in its velocity gate when the ECM signal is turned off. The operation can be accomplished by either a traveling wave tube (TWT) or a digital phase shifter (DPS).

Similarly, the ECM deception set may operate on the radar pulse train, returning many pulses instead of one, in an effort to deceive the enemy into believing there are many targets spaced at different positions. The genuine return may be blanked out, or false signals may be generated at higher levels so that the enemy may think the genuine return is simply noise. By changing the character of the pulses in the train, the ECM set can deceive the antagonist in a number of ways. Newer ECM systems combine both noise and deception features.

The level of sophistication achieved in modern deceptive electronic countermeasures systems is indicated by a system deployed on some U.S. Navy fighter and attack aircraft. This microprocessor-controlled system employs multiple jamming techniques in countering air, ground, or sea-based pulse radar threats. When threat signals are received and identified by matching them against an emitter library in an electronic data file, the system chooses the optimum active technique to defeat the threat. In a high-threat-density environment, the 190-lb (86-kg)

system automatically selects the appropriate countermeasures technique and power level for responding to the threat in accordance with threat priorities. Identification and prioritization of the threat and the choice of jamming response are all handled by the reprogrammable data file. The reprogramming capability enables the user to add or change threat characteristics as they become known.

ECM system. Because of the complex electromagnetic environment in warfare, no single technique is satisfactory for all occasions. A sophisticated warplane may have a penetration-aids subsystem with a number of separate systems operating interdependently to give aural and visual warning of radar threats. A cryogenically cooled infrared receiving set with its own decision circuitry located in the top of the aircraft's vertical stabilizer warns the pilot of threats from heat-seeking missiles. An ECM deception set automatically generates false target returns to tracking-radar threats in several frequency bands. A pneumatically driven chaff dispenser, working with other subsystems, selects and automatically ejects disposable packages including chaff and infrared flares (Fig. 7). A wideband radar homing and warning set detects enemy ground and airborne radar threats, identifies and determines the detected threat location, and calculates azimuth and elevation of the threat. It monitors the effectiveness of the ECM deception set and automatically supplies signals to deploy disposables. It also provides bearing



Fig. 7. Programmable countermeasures-dispensing system, which can dispense chaff, flares, and active expendables to protect aircraft from threats. (Tracor)

inputs to the aircraft's inertial navigation computer and drives the lead computing gunsight in an attack mode so that the aircraft can attack a threatening emitter detected by the warning set.

ECM aboard missiles. Many ECM techniques developed for aircraft, ships, and ground operations were investigated and in some instances adapted by the United States to ensure penetration capability of its ballistic missiles. The techniques are similar to what was previously described, although active ECM techniques are less practical because of the severe environments and power limitations. Several ways of using active ECM have been investigated, however, including the possibility of installing jammers within small rocket-fired precursor decoys which would be fired ahead of a warhead. The jammers, each of which is tuned within different frequency bands, spread confusion among enemy radars.

The difficulties in securing extremely rugged microwave transmitting tubes to withstand missile environments and in generating a maximum amount of broadband power with a minimum of weight are formidable. A major hurdle was overcome during the mid-1960s by developing sources of kilowatt dc power having a high specific weight and high specific volume. See MICROWAVE TUBE.

To use active ECM properly, an adversary must have prior knowledge about the defensive radars so that ECM jammers can be tuned. In an aircraft or on a ship, there is more likely to be space and power for devices to sense hostile signals and tune transmitters. In a missile, where space and power are at a premium, this ability is far more costly.

In general, however, the military relies for missile penetration more heavily on passive ECM techniques, such as chaff and inflated metallized balloons, which reinforce enemy-radar returns. One concept is to use tethered chaff, which would be ejected from a missile system on atmospheric reentry and carried along near the vehicle to create large masking radar echoes. Regenerative chaff, a technique by which a missile or powered decoy would sequentially eject clouds of chaff, is another.

Trends in ECM. Newer ECM systems display greater integration among various subsystems to tighten up defense capability. They incorporate infrared warning receivers and in some cases new active infrared jamming systems, as well as attack warning from a data link.

The Navy is integrating its standard modular electronic warfare systems carried by a wide variety of ships. In one version, a passive ECM system would be coupled to rocket-propelled decoy and chaff packages to lure away radar-guided antiship missiles detected by the ECM system. A small hand-held active radar decoy can be used against monopulse radars that are particularly difficult to jam by conventional angle jamming.

As part of a continuing practice of updating the ECM capability of military aircraft to meet newly recognized threats, the U.S. Air Force has modified repeatedly its strategic bomber force. Its B-52 bombers have been retrofitted with improved jammers to



Fig. 8. Antenna for shipborne deception ECM system mounted on vessel's superstructure. (Hughes Aircraft Co.)

counter higher-frequency radars carried by interceptors or located on the ground for control of surface-to-air missiles.

A number of allied aircraft were shot down or damaged by ground-launched infrared missiles during the Desert Storm operation. This experience reinforced efforts to find better sensors, either infrared or pulse-Doppler, to provide rear-sector warning of impending missile threats to aircraft, and to develop suitable missile alert warning systems.

The sinking of an Israeli destroyer by Soviet-made Egyptian cruise missiles in 1967 revealed a glaring inadequacy in the defensive systems of Western naval vessels that prompted a resurgence in shipboard ECM activities. This concern was heightened when the American destroyer *Stark* was taken unawares and was accidentally struck by two cruise missiles fired from an Iraqi aircraft during the Iran-Iraq war. A primary element of defense against a cruise missile is the brief interval of time, as little as 90 s, that a ship has from the instant a threatening missile rises above the horizon until it is upon the ship. This briefness requires that a ship quickly detect and identify the threat, determine its direction of flight, and take appropriate countermeasures. This may best be accomplished by integrating ship sensors, including ECM signal sensing, and identification and direction-finding capabilities, with antiaircraft weapons and improved command-and-control capability.

The Navy has developed systems to protect aircraft carriers and other large warships from radar-guided cruise missiles. One of these is a deception system that senses the threatening radar signal from the missile seeker and radiates a false electronic image of the

target ship, offset by a safe distance from the actual ship. Then the missile will home on this ghost image of the target and fall harmlessly into the sea (Fig. 8).

Military events have established and reaffirmed the crucial role played in modern warfare and in strategic and tactical defense by electronic warfare. This was apparent in the Persian Gulf conflict in 1991. Coalition electronic warfare aircraft were a significant factor in suppressing enemy air defenses with active jamming, passive location systems, and anti-radiation missile delivery ability. See MILITARY AIRCRAFT.

Barry Miller; Paul J. De Lia

Bibliography. D. L. Adamy, *EW 101: A First Course in Electronic Warfare*, 2001; D. L. Adamy, *EW 102: A Second Course in Electronic Warfare*, 2004; J. P. R. Browne and M. T. Thurbon, *Electronic Air Warfare*, 1997; E. J. Chrzanowski, *Active Radar Electronic Countermeasures*, 1990; F. Neri, *Introduction to Electronic Defense Systems*, 2d ed., 2001; R. Poisel, *Introduction to Communication Electronic Warfare Systems*, 2002; A. Price, *The History of U.S. Electronic Warfare*, vol. 1, 1984, vol. 2, 1989; D. C. Schleher, *Introduction to Electronic Warfare*, 1986; R. J. Wiegand, *Radar Electronic Countermeasures System Design*, 1991; R. G. Wiley, *Electronic Intelligence: The Analysis of Radar Signals*, 2d ed., 1993; R. G. Wiley, *Electronic Intelligence: The Interception of Radar Signals*, 1985.

Electronics

The technological area involving the manipulation of voltages and electric currents through the use of various devices for the purpose of performing some useful action with the currents and voltages. This large field is generally divided into two primary areas, analog electronics and digital electronics.

Analog electronics. In analog electronics, the signals to be manipulated take the form of continuous currents or voltages. The information in the signal is carried by the value of the current or voltage at a particular time t . Some examples of analog electronic signals are amplitude-modulated (AM) and frequency-modulated (FM) radio broadcast signals, thermocouple temperature data signals, and standard audio cassette recording signals. In each of these cases, analog electronic devices and circuits can be used to render the signals intelligible. Processing signals in an analog fashion has advantages and disadvantages, discussed below.

Commonly required manipulations include amplification, rectification, and conversion to a nonelectronic signal. Amplification is required when the strength of a signal of interest is not sufficient to perform the task that the signal is required to do. For example, the signal obtained by the piezoelectric transducer on an old-fashioned record player is very weak and does not contain sufficient power to drive a set of speakers hard enough that humans can understand the information that is contained in the signal. The transducer signal must be amplified. The amplification process involves making a replica of

the original signal containing sufficient power to perform the required task. However, the amplification process suffers from the two primary disadvantages of analog electronics: (1) susceptibility to replication errors due to nonlinearities in the amplification process and (2) susceptibility to signal degradation due to the addition, during the amplification process, of noise originating from the analog devices composing the amplifier. These two disadvantages compete with the primary advantage of analog electronics, the ease of implementing any desired electronic signal manipulation. The decision concerning whether to implement an electronic operation in analog form will depend upon the relative weights given to this advantage and the disadvantages cited above. See AMPLIFIER; DISTORTION (ELECTRONIC CIRCUITS); ELECTRICAL NOISE.

Historically, analog electronics was used in large part because of the ease with which circuits could be implemented with analog devices. This facility was primarily a result of the large variety of devices available, including vacuum tubes, transistors, and some special-purpose integrated circuits. However, as signals have become more complex, and the ability to fabricate extremely complex digital circuits has increased, the disadvantages of analog electronics have increased in importance, while the importance of simplicity has declined. Thus, modern-day circuit designers have opted more frequently to implement their designs digitally. This tendency is also the result of the error-correction possibilities inherent in digital electronics. Digital electronics is the most important class of electronics.

Digital electronics. The advent of the transistor in the 1940s made it possible to design simple, inexpensive digital electronic circuits and initiated the explosive growth of digital electronics. However, computers and other digital circuits had been implemented with electron tubes for at least two decades. Digital signals are represented by a finite set of states rather than a continuum, as is the case for the analog signal. Typically, a digital signal takes on the value 0 or 1; such a signal is called a binary signal. Because digital signals have only a finite set of states, they are amenable to error-correction techniques; this feature gives digital electronics its principal advantage over analog electronics. See ELECTRON TUBE; TRANSISTOR.

In common two-level digital electronics, signals are manipulated mathematically. These mathematical operations are known as boolean algebra. The operations permissible in boolean algebra are NOT, AND, OR, and XOR, plus various combinations of these elemental operations. The NOT operation performs an inversion of the digital signal; for example, a 1 signal is taken to 0, and a 0 signal is taken to a 1. The AND operation in boolean algebra is similar to the well-known multiplication operation. The digital signals are multiplied together, and the result is a 0 unless all of these signals are 1; in that case, the result is a 1 as well. The OR operation is closely related to addition. If any single input signal, or combination of input signals, is 1, then the result is a 1;

otherwise, if the inputs are all 0, the result is also 0. The XOR (exclusive or) operation resembles common addition, but the carry information is thrown away. In this case, if just one of the inputs is a 1, the output is a 1; otherwise the output is 0. Digital circuits exist to implement all of these operations, as well as more complex operations comprising a combination of these simple ones. As a result, signals can be manipulated in a large variety of ways; for example, signals can be added to each other or multiplied by each other, or transformations such as the Fourier transform can be made, allowing one type of signal to be converted into another. *See* BOOLEAN ALGEBRA; FOURIER SERIES AND TRANSFORMS; LOGIC CIRCUITS.

Electronic devices. Several common electronic devices serve useful purposes in circuit design. The devices most often used include resistors, capacitors, inductors, diodes, and transistors.

Resistor. A resistor is a two-terminal device that impedes the flow of electrons from one of its terminals to the other. The ratio of the voltage applied across the terminals of the resistor to the current which flows from one terminal to the other is the value of the resistance and is measured in units of ohms. One ohm is 1 volt per ampere. *See* ELECTRICAL RESISTANCE; RESISTOR.

Capacitor. A capacitor is a two-terminal device consisting of a set of parallel plates of conductive material on which electric charge can be stored. The ratio of the voltage applied across the terminals of the device to the steady-state charge stored on the plates is the capacitance of the device. Capacitance is measured in units of farads. One farad is 1 volt per coulomb. *See* CAPACITANCE.

Inductor. An inductor is a two-terminal device consisting of a series of conductor loops which generate a magnetic field. As current flows through the inductor, a magnetic field is created in which energy can be stored. The ratio of the voltage across the terminals of the device to the time-rate-of-change of the current flowing between the terminals is the inductance of the device. Inductance is measured in units of henrys. One henry is 1 volt-second per ampere. *See* INDUCTANCE; INDUCTOR.

Diode. A diode is a two-terminal device that permits electric current to flow in only one direction. The most common modern implementation of a diode takes the form of a semiconductor *pn* junction, that is, a junction between a semiconductor doped *p*-type and one doped *n*-type. *See* DIODE; JUNCTION DIODE; SEMICONDUCTOR DIODE.

Transistor. A transistor is a three-terminal device that uses the voltage on one terminal to control the current flowing between the two remaining terminals. There are several major classes of transistors, including field-effect transistors (FETs), bipolar junction transistors (BJTs), and minor devices such as point-contact transistors. Each type of device has characteristics which make it particularly suitable for certain applications. For example, the FET is relatively easy to fabricate in an integrated-circuit format, and so finds wide use in those types of applications. Generally, transistors are used to control one voltage or

current in a circuit with another, in switching or amplifying circuits.

Logic families. A logic family consists of circuits fabricated to operate according to a particular specification. When integrated circuits of the same logic family are wired together on a printed circuit board, typically no buffering of the outputs or inputs is required for proper functionality. Thus, a multichip circuit can be built up almost exclusively from integrated circuits, circumventing the need for most discrete components. Some of the logic families are resistor-transistor logic (RTL), diode-transistor logic (DTL), transistor-transistor logic (TTL), emitter-coupled logic (ECL), and complementary metal-oxide-semiconductor (CMOS) logic.

The most successful logic family is the CMOS family. The widespread use of CMOS arises from its low power consumption and its acceptance of a wide range of voltages, giving it versatility that is not found in the other logic families. The low power consumption characteristic of CMOS is a result of the fact that it comprises FET devices which have extremely high input impedances and low operating bias currents. CMOS logic's ability to accept a wide range of supply and operating voltages is also due to the FET technology used, and serves to simplify interfacing the logic input and output points to the data collection and use points.

Electronic circuit design. Electronic circuits are composed of various electronic devices. In circuits built from discrete components, the components are typically soldered together on a fiberglass board known as a printed circuit board. On one or more surfaces of the printed circuit board are layers of conductive material which has been patterned to form the interconnections between the different components in the circuit. In some cases, the circuits necessary for a particular application are far too complex to build from individual discrete components, and integrated-circuit technology must be employed. Integrated circuits are fabricated entirely from a single piece of semiconductor substrate. It is possible in some cases to put several million electronic devices (transistors, resistors, capacitors) inside the same integrated circuit. Modern computer microprocessors are an example of this type of very large scale integration. Special techniques must be used to design such circuits because of the limited characteristics of some of the devices that can be realized. For example, in an integrated circuit it is very difficult to implement an inductor, and so priority is placed on circuit designs that function without inductors. Many integrated circuits can be fabricated on a single wafer of silicon at one time, and at the end of the fabrication process the wafer is sawed into individual integrated circuits. These small pieces, or chips as they are popularly known, are then packaged appropriately for their intended application. *See* INTEGRATED CIRCUITS; PRINTED CIRCUIT.

The microprocessor is the most important integrated circuit to arise from the field of electronics. This circuit consists of a set of subcircuits that can perform the tasks necessary for computation and are

the heart of modern computers. The microprocessors in current use can copy digital data elements from one transistor array (memory) location to another, load an instruction from memory, execute that instruction (typically to add, or otherwise manipulate, the data stored in memory), and copy the results of various instructions into memory, as well as carry out other processing necessary for computation. Microprocessors are programmed, or told what to do, via a series of instructions which together are called a program. Thus, the microprocessor is an example of a circuit that can be reconfigured on the basis of a set of instructions it receives from the programmer and executes. Microprocessors that understand large numbers of instructions are called complete instruction set computers (CISCs), and microprocessors that have only a very limited instruction set are called reduced instruction set computers (RISCs). Microprocessors have found wide acceptance in a large variety of applications and are used to control everything from the temperature of a private dwelling to the softness (calcium-ion concentration) of the water used in the fabrication foundries of the integrated-circuit manufacturers. See DIGITAL COMPUTER; MICROPROCESSOR.

Other circuit designs have been standardized and reduced to integrated-circuit form as well. An example of this process is seen in the telephone modem, used for performing digital communications over telephone or radio-frequency-optical links. Modulation techniques have been standardized to permit the largest possible data-transfer rates in a given amount of bandwidth, and standardized modem chips are available for use in circuit design. See MODEM.

The memory chip is another important integrated electronic circuit. This circuit consists of a large array of memory cells composed of a transistor and some other circuitry. As the storage capacity of the dynamic random-access memory (DRAM) chip has increased, significant miniaturization has taken place, to the point where modern DRAM integrated circuits have device feature sizes of approximately 0.2 micrometer. See CIRCUIT (ELECTRONICS); SEMICONDUCTOR MEMORIES.

David R. Andersen

Bibliography. N. R. Malik, *Electronic Circuits: Analysis, Simulation, and Design*, 1995; A. S. Sedra and K. C. Smith, *Microelectronic Circuits*, 3d ed., 1992.

Electronvolt

A practical unit of energy often used in atomic, nuclear, and particle physics. It is the energy acquired by an electron (or any singly charged particle) when it is accelerated by a potential difference of 1 V. The unit of energy in the International System (SI) is the joule: 1 eV is equal to 1.6022×10^{-19} J. Although the electronvolt is a non-SI unit, it is often expressed using SI prefixes, for example, meV (10^{-3} eV), keV (10^3 eV), MeV (10^6 eV), GeV (10^9 eV). See ELECTRICAL UNITS AND STANDARDS; ELECTRON; METRIC SYSTEM; PHYSICAL MEASUREMENT; UNITS OF MEASUREMENT.

The electronvolt is also commonly used for photon energies. Wavelength and photon energy E are related by equation below.

$$E(\text{eV}) = \frac{1239.84}{\text{wavelength (nm)}}$$

See PHOTON; QUANTUM MECHANICS. Fred Schlachter

Electrooptics

The branch of physics that deals with the influence of an electric field on the optical properties of matter, especially in crystalline form. These properties include transmission, emission, and absorption of light.

An electric field applied to a transparent crystal can change its refractive indexes and, therefore, alter the state of polarization of light propagating through it. When the refractive-index changes are directly proportional to the applied field, the phenomenon is termed the Pockels effect. When they are proportional to the square of the applied field, it is called the Kerr effect. See KERR EFFECT; POLARIZED LIGHT; REFRACTION OF WAVES.

Pockels cell. The Pockels effect is used in a light modulator called the Pockels cell. This device (Fig. 1) consists of a crystal C placed between two polarizers P_1 and P_2 whose transmission axes are crossed. A crystal often used for this application is potassium dihydrogen phosphate (KH_2PO_4). Ring electrodes bonded to two crystal faces allow an electric field to be applied parallel to the axis OZ , along which a light beam (for example, a laser beam) is made to propagate. The crystal has been cut and oriented in such a way that, in the absence of an electric field, the polarization of light propagating along OZ does not change, and therefore no light is transmitted past P_2 . However, when an electronic driver V applies an electric field of the proper magnitude

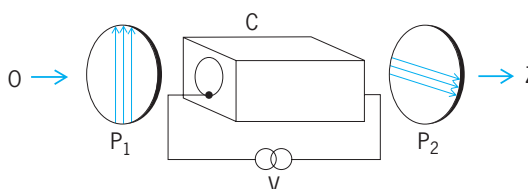


Fig. 1. Typical Pockels cell light modulator. The arrows (OZ) represent a light beam.

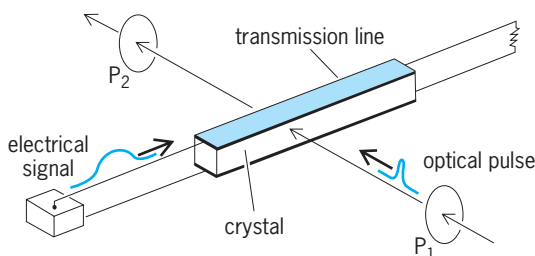


Fig. 2. Electrooptic sampling arrangement using a transmission-line Pockels cell. Examples of optical and electrical waveforms are shown.

(2000 V/cm or 5000 V/in. is a typical value) across the crystal, the crystal becomes birefringent, and light that has been vertically polarized by P_1 undergoes a 90° change in polarization in the crystal, and is now well transmitted past P_2 . Pockels cells of this type can be switched on and off in well under 1 nanosecond. See BIREFRINGENCE; CRYSTAL OPTICS.

Electrooptic sampling. The linearity and high-speed response of the Pockels effect within an electrooptic crystal make possible a unique optical technique for measuring the amplitude of repetitive high-frequency (greater than 1 GHz) electric signals that cannot be measured by conventional means. The technique is known as electrooptic sampling and employs a special traveling-wave Pockels cell between crossed polarizers, P_1 and P_2 (Fig. 2). The electrooptic crystal (commonly lithium tantalate, LiTaO_3) is placed between electrodes that are part of a high-speed transmission line along which the unknown electric signal travels. Ultrashort optical pulses (on the order of 1 picosecond) from a mode-locked laser, synchronized with the electric signals and directed through the crystal, are used to repetitively sample a small portion of the electrical waveform. After passing through the second polarizer P_2 , the optical pulses have an intensity proportional to the amplitude of that portion of the electrical waveform experienced as the pulse passed through the crystal. By varying the relative delay between the arrival of the optical pulse and electric signal in the crystal and graphing the transmitted intensity as a function of delay, an accurate replica of the electrical waveform is obtained. The availability of subpicosecond optical pulses makes it possible to measure electric signals containing frequencies up to 1 THz with similar temporal resolution. Electrooptic sampling is an important technique for the analysis of ultrafast electric signals such as those generated by high-speed transistors and optical detectors. See LASER; MICROWAVE SOLID-STATE DEVICES; OPTICAL DETECTORS; OPTICAL MODULATORS; OPTICAL PULSES.

Michel A. Duguay; Janis A. Valdmanis

Self-electrooptic-effect device. A self-electrooptic-effect device (SEED) is a combination of a quantum-well electrooptic modulator with a photodetector which, when light shines on it, changes the voltage on the modulator. Although the device relies internally on an electrooptic effect, the output from the modulator is controlled by the light shining on the photodetector, giving an optically controlled device with an optical output. Most of these devices rely on the quantum-confined Stark effect in semiconductor quantum-well heterostructures as the electrooptic mechanism and utilize the changes in optical absorption resulting from this mechanism.

The quantum-confined Stark effect is conveniently observed by passing a light beam through a diode containing quantum wells in an undoped region (Fig. 3) between p - and n -doped regions. In the simplest self-electrooptic-effect device, this structure also can function simultaneously as the photodetector. The wavelength of the incident light can be so chosen that, as the voltage across the modulator de-

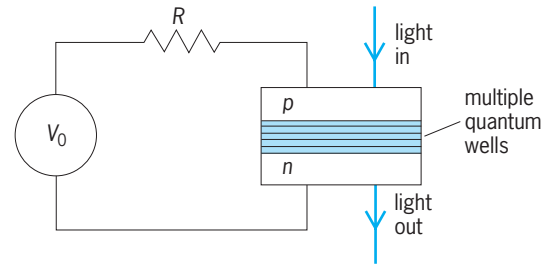


Fig. 3. Sample circuit for a self-electrooptic-effect devices.

creases, the absorption of light by the modulator increases. When light is shone on the diode in the circuit of Fig. 3, photocurrent is generated, resulting in a voltage drop across resistor R . Consequently, the voltage across the modulator falls, increasing its optical absorption. This absorption increase itself results in a further increase in photocurrent, which in turn leads to a further reduction in voltage across the diode, and so on. This positive feedback can be so strong that it leads to a discontinuous switching into a highly absorbing state, giving the optically bistable input-output characteristic shown in Fig. 4. Thus this simple device can function as an optically controlled memory or switch with an optical output.

Many different configurations of self-electrooptic-effect devices have been proposed for different functions. Substituting an inductor for the resistor gives an oscillator. Substituting a conventional photodiode for the resistor can give a bistable device whose switching thresholds are controlled by the light shining on the photodiode (a diode-based SEED or D-SEED). Substituting another quantum-well diode for the resistor can give a device that is bistable in the power ratio of the two light beams incident on the two diodes (a symmetric SEED or S-SEED). Transistors and more complex circuitry can also be incorporated to give yet more complex relations between optical inputs and outputs. Such devices can be used to perform all the elementary logical functions on light beams.

Self-electrooptic-effect devices can offer operating energy densities comparable to electronic devices when the modulators, photodetectors, and

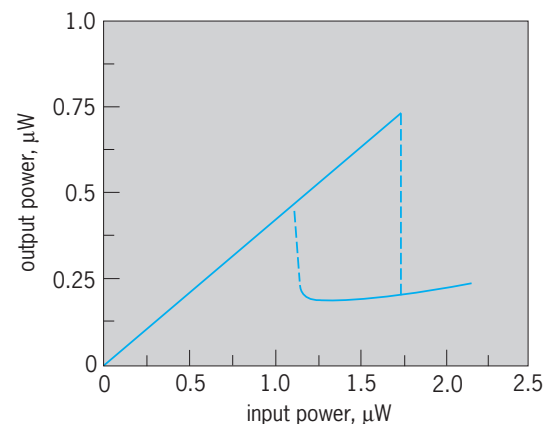


Fig. 4. Set of optical input-output characteristics for an integrated device.

any other circuitry are so integrated as to minimize any undesired capacitance. They can also be fabricated in two-dimensional arrays compatible with optical processing schemes. See ARTIFICIALLY LAYERED STRUCTURES; EXCITON; SEMICONDUCTOR HETEROSTRUCTURES; STARK EFFECT. David A. B. Miller

Bibliography. F. Argullo-Lopez, *Electro-Optics: Phenomena, Materials, and Applications*, 1994; H. Haug (ed.), *Optical Nonlinearities and Instabilities in Semiconductors*, 1988; M. A. Karim, *Electro-Optical Devices and Systems*, 1990; M. A. Karim (ed.), *Electro-Optical Displays*, 1990; R. Waynant and M. N. Ediger (eds.), *Electro-Optics Handbook*, 1994; A. Yariv, *Optical Electronics*, 5th ed., 1997.

Electrophilic and nucleophilic reagents

Electrophilic reagents are chemical species which, in the course of chemical reactions, acquire electrons, or a share in electrons, from other molecules or ions. Although this definition embraces all oxidizing agents and all Lewis acids, electrophilic reagents are ordinarily thought of as cationic species, such as H^+ , NO_2^+ , Br^+ , or SO_3 (or carriers of these species such as HCl , CH_3COONO_2 , or Br_2), which can form stable covalent bonds with carbon atoms. Electrophilic reagents frequently are positively charged ions (cations). See ACID AND BASE.

Nucleophilic reagents are the opposite of electrophilic reagents. Nucleophilic reagents give up electrons, or a share in electrons, to other molecules or ions in the course of chemical reactions. Nucleophilic reagents frequently are negatively charged ions (anions). Typical nucleophilic reagents are hydroxide ion (OH^-), halide ions (F^- , Cl^- , Br^- , and I^-), cyanide ion (CN^-), ammonia (NH_3), amines, alkoxide ions (such as CH_3O^-), and mercaptide ions (such as $C_6H_5S^-$). See SUBSTITUTION REACTION.

Joseph F. Bunnett

Electrophoresis

The migration of electrically charged particles in solution or suspension in the presence of an applied electric field. Each particle moves toward the electrode of opposite electrical polarity. For a given set of solution conditions, the velocity with which a particle moves divided by the magnitude of the electric field is a characteristic number called the electrophoretic mobility. The electrophoretic mobility is directly proportional to the magnitude of the charge on the particle, and is inversely proportional to the size of the particle. An electrophoresis experiment may be either analytical, in which case the objective is to measure the magnitude of the electrophoretic mobility, or preparative, in which case the objective is to separate various species which differ in their electrophoretic mobilities under the experimental solution conditions.

Tiselius cell. The phenomenon of electrophoresis was first observed in 1807 by the Russian physi-

cist F. F. Reuss, but electrophoresis was not employed as an experimental technique until the introduction of a new electrophoresis apparatus by Arne Tiselius in 1937. The apparatus of Tiselius detected electrophoretic motion by the moving-boundary method, in which a boundary is created between the solution of particles to be examined and a sample of pure solvent. As the particles migrate in an electric field, the boundary between solution and solvent can be observed to move, and if there are a number of species in the solution with different electrophoretic mobilities, a series of boundaries of various shapes and magnitudes can be detected. Using his apparatus, Tiselius demonstrated the heterogeneity of human blood plasma, and showed for the first time that the globulin molecules could be separated into different classes, which were designated alpha, beta, and gamma globulin. The moving-boundary method was used for three decades to separate complex mixtures of charged macromolecules in solution and to study the physical characteristics of solutions of proteins and other macromolecules of biological and industrial importance.

Gel techniques. The resolving power of electrophoresis was greatly improved by the introduction of the use of gel supporting media. The gel matrix prevents thermal convection caused by the heat which results from the passage of electric current through the sample. The absence of convection reduces greatly the mixing of the various parts of the sample, and therefore allows for more stable separation. The dimensions of the cross-links of the gel may also provide a molecular sieving effect, which increases the resolving power of the electrophoretic separation of molecules of different size. In addition, the gel media may support a gradient of a separate reagent, which assists in the separation of macromolecules. Gradients of pH and of reagents of various types may be combined in two-dimensional arrays for even greater resolving power. A very successful derivative of the gel technique is the determination of the molecular weights of protein molecules by electrophoresis of the molecules in a gel medium which contains substantial amounts of detergent. The detergent denatures the protein molecules, changing them from globular, compact structures to long, flexible polymers which are coated with detergent molecules. These polymers move in the electric field through the gel medium with a velocity which is determined by the length of the polymer, and therefore by the molecular weight of the protein unit. This method is the most common technique for the determination of molecular weights of proteins in biochemical studies. See GEL; PH; PROTEIN.

Isoelectric focusing. An important variation of the electrophoresis technique is isoelectric focusing. In this technique the medium supports a pH gradient which includes the isoelectric pH of the species being studied. Many charged macromolecules have both positive and negative charges on their surfaces, and the electrophoretic mobility is related to the net excess of charge of one type or the other. As

the pH becomes more acidic, the number of positive charges increases, and as the pH becomes more basic, the number of negative charges increases. For each molecule of this type, there is one pH at which the net charge on the surface is zero, so that the molecule does not move when an electric field is applied and thus has an electrophoretic mobility of zero. This pH is called the isoelectric pH. If the molecule is introduced into a pH gradient which includes its isoelectric pH, it will migrate to the position of the isoelectric pH and then become stationary. In this way, all molecules of a given isoelectric pH will migrate to the same region—hence the term isoelectric focusing. The method of isoelectric focusing is particularly good for the analysis of microheterogeneity of protein species and other species which may differ slightly in their chemical content. *See* ISO-ELECTRIC POINT.

Isotachopheresis. One important variant of electrophoresis is the phenomenon of isotachopheresis, in which ionic species move with equal velocity in the presence of an electric field. The isotachopheretic condition is maintained when particles of different mobilities form boundaries within the solution so that the most mobile ions form the leading edge of the moving sample, followed by the less mobile ions in the order of their mobility. The reason that the ions with different mobility can move with the same velocity is that the electric field in each of the regions is inversely proportional to the mobility of the ionic species in order to maintain a constant current throughout the sample. Isotachopheresis has a number of important analytical and preparative applications which are featured by its advantages of high resolving power, sensitivity, speed, and the ability to concentrate rather than to dilute the components which are being analyzed.

Particle electrophoresis. The methodology of electrophoresis may be modified considerably when the particles undergoing analysis are of sufficient size to be viewed either with the naked eye or with the assistance of an optical microscope. This general area of particle electrophoresis has its most important applications in the analysis of the surface charge of living cells and in the study of various types of particles used in industrial coating processes. The most straightforward technique, called optical cytopherometry or microelectrophoresis, is that in which a human experimenter views the particles in an electric field under an optical microscope and determines manually the amount of time necessary to traverse a given distance. Although time-consuming and tedious, this technique has had many important applications. Attempts to modernize this method have included the introduction of high-speed photography, television technology, and the study of particle electrophoresis by the laser Doppler effect. Particle electrophoresis can be conducted under many of the conditions which were originally developed for smaller macromolecules. The use of gradients of density and pH and the methods of isoelectric focusing and isotachopheresis are commonly applied to

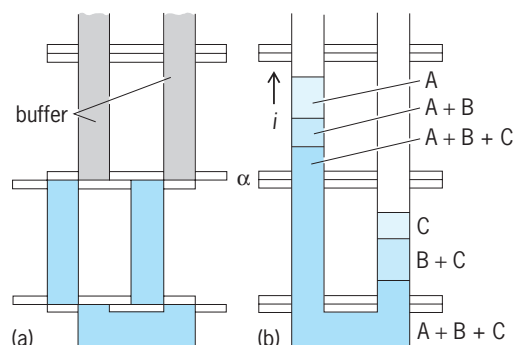


Fig. 1. Diagram of Tiselius electrophoresis cell. (a) Formation; buffer plus proteins A, B, and C. (b) Motion of electrophoretic boundaries.

particles which are even large enough to be viewed with the naked eye (**Fig. 1**).

A common interference effect in the performance of particle electrophoresis is the sedimentation of the particles in the field of the Earth's gravity. This effect can be minimized in importance by performing the experiments in a medium of equal density, by performing the electrophoresis in a vertical direction and accounting for the effect of sedimentation, or by performing the experiment in such a short period of time that the degree of sedimentation in the Earth's gravity is not a significant interfering effect.

Laser applications. Application of the optical laser to electrophoretic detection resulted in the development of a technique which can be used for analytical electrophoresis experiments on particles of all sizes. The basic principle is that the highly monochromatic (single-frequency) laser light impinges upon the particles, and is scattered from the particles in all directions. When observing the laser light which has been scattered from a moving particle, one can detect that there is a slight shift in the frequency of the light as a result of the motion of the particle. This is the Doppler effect, which causes the change in the apparent tone of passing trains or cars and which is the operating principle of other familiar techniques such as radar. The application of the laser Doppler principle to electrophoresis experiments, often called electrophoretic light scattering (ELS), is an important method for the rapid determination of electrophoretic velocities. The complete electrophoretic mobility distribution of a sample of many particles can be determined in a time as short as 1 s with a precision heretofore unobtainable by standard technology. *See* DOPPLER EFFECT; LASER.

In an electrophoretic light-scattering apparatus with a sample Doppler spectrum (**Fig. 2**), the very slight Doppler shifts caused by the electrophoretic motions are detected by the electronic "beats" which result when the scattered light is incident simultaneously with an unshifted reference beam, or "local oscillator," on a photodetector. An important application of electrophoretic light scattering is the characterization of leukemic cells (**Fig. 3**). The normal lymphocytes and leukemic cells isolated in the same way have distinctly different electrophoretic mobilities and mobility distributions, as shown by

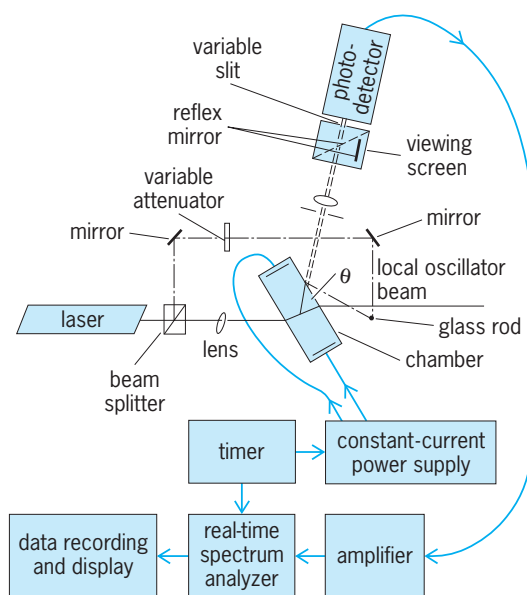


Fig. 2. Diagram of an electrophoretic light-scattering apparatus. Doppler-shifted light from the particles moving in the chamber is mixed with an unshifted reference beam (the local oscillator). The beat frequencies from the photodetector are spectrum-analyzed to produce the electrophoretic mobility histogram. (After G. M. Hieftje, ed., *New Applications of Lasers to Chemistry*, ACS Symposium Ser. 85, 1978)

the two spectra in Fig. 3a. The spectrum in Fig. 3b shows the simultaneous detection of leukemic and normal cells in a mixture, based solely on their Doppler-detected electrophoretic mobilities. Electrophoretic light scattering has been used for the study of many types of living cells, cell organelles, viruses, proteins, nucleic acids, and synthetic polymers. See ELECTROLYTIC CONDUCTANCE; SCATTERING OF ELECTROMAGNETIC RADIATION. B. R. Ware

Capillary electrophoresis. Electrophoresis can be performed in a capillary format. A typical system consists of two reservoirs and a capillary filled with a buffer solution (Fig. 4a). A high voltage is applied across the capillary by using a high-voltage power supply. The very small diameter capillaries (typically 5–100 micrometers) employed in this technique allow for efficient heat dissipation. Therefore, much higher voltages can be employed than those used in slab gel electrophoresis, leading to faster, more efficient separations. Compounds are separated on the basis of their net electrophoretic mobilities.

Most often, the detector is placed on line and analytes are detected as they flow past the detector (Fig. 4a). Spectroscopic detection (ultraviolet and laser-based fluorescence) is usually performed in this manner by using the capillary itself as the optical cell. Alternatively, detectors can be placed off line (after the column; Fig. 4b). In this case, the detector is isolated from the applied electric field through the use of a grounding joint. Electrochemical detection and mass spectroscopic detection are generally accomplished in this manner, since the electric field can interfere with the performance of these detectors. The several different modes of capillary

electrophoresis include capillary zone electrophoresis, micellar electrokinetic chromatography, capillary gel electrophoresis, isoelectric focusing, and isotachopheresis.

Capillary zone electrophoresis. This is the simplest and most widely used form of capillary electrophoresis. The capillary is filled with a homogeneous buffer, and compounds are separated on the basis of their relative charge and size. Most often, fused silica capillaries are employed. In this case, an electrical double layer is produced at the capillary surface due to the attraction of positively charged cations in the buffer to the ionized silanol groups on the capillary wall. In the presence of an electric field, the cations in the diffuse portion of this double layer move toward the cathode and drag the solvent with them, producing an electroosmotic flow. The resulting flow profile is flat rather than the parabolic shape characteristic of liquid chromatography. This flat flow profile causes analytes to migrate in very narrow bands and leads to highly efficient separations. The electroosmotic flow is also pH dependent, and it is highest at alkaline pH values.

In most cases, the electroosmotic flow is the strongest driving force in the separation, and all

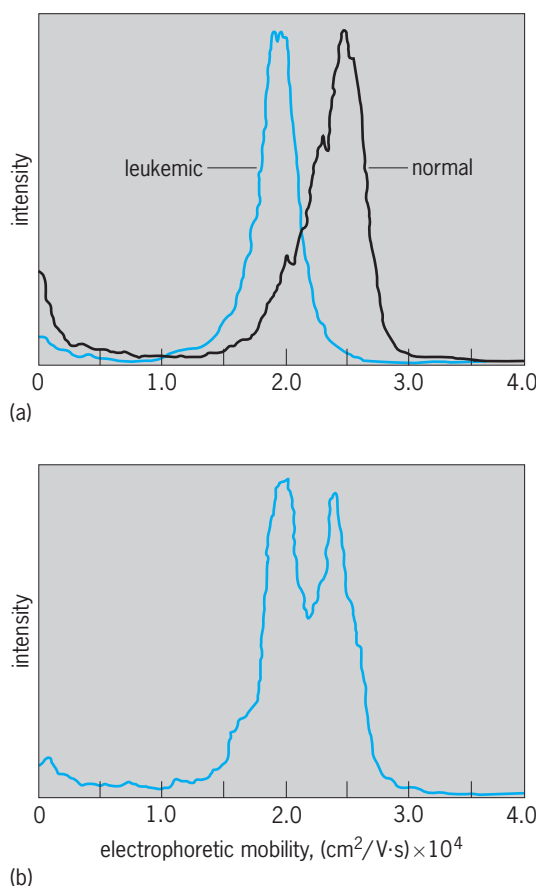


Fig. 3. Characterization of leukemic cells with electrophoretic light scattering. (a) Spectra of leukemic cells and normal cells. (b) Spectrum for mixture of the two cell types. (After B. A. Smith, B. R. Ware, and R. S. Weiner, *Electrophoretic distributions of human peripheral blood molecular white cells from normal subjects and from patients with acute lymphocytic leukemia*, *Proc. Nat. Acad. Sci. USA*, 73:2388, 1976)

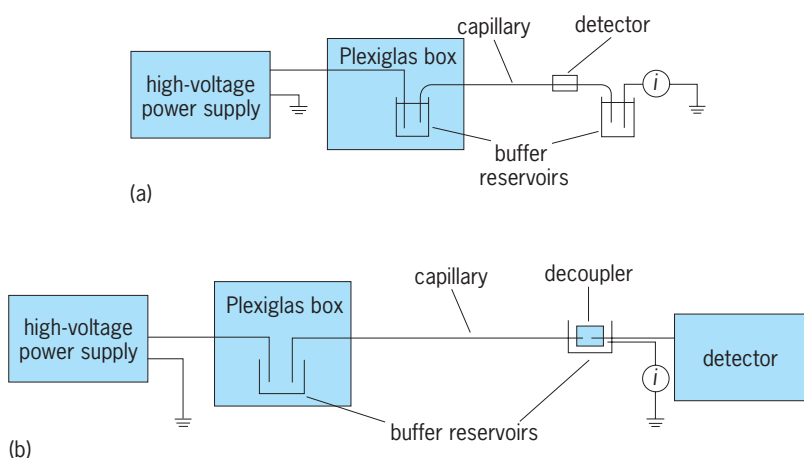


Fig. 4. Capillary electrophoresis. (a) Arrangement of a basic system with on-line detection. (b) System with off-line detection.

analytes, regardless of charge, migrate toward the cathode. Therefore, it is possible to separate and detect positive, negative, and neutral molecules in the same electrophoretic run, if the detector is placed at the cathodic end. Negatively charged compounds are attracted to the anode but are swept up by the electroosmotic flow and elute last. Neutral molecules, which are not separated from each other in capillary zone electrophoresis, elute as a single band with the same velocity as the electroosmotic flow. Positive compounds have positive electrophoretic mobilities in the same direction as the electroosmotic flow, and they elute first (Fig. 5). Capillary zone electrophoresis is generally employed for the separation of small molecules, including amino acids, peptides, and small ions, and for the separation of drugs, their metabolites, and degradation products.

Although the electroosmotic flow that is characteristic of the fused silica capillaries makes it possible to separate and detect both positive and negative ions in a single run, the presence of the negatively charged

silanol groups can lead to irreversible adsorption to the capillary surface of proteins containing a large number of positive functional groups (high pI, the isoelectric point of the protein). Therefore, for many of the applications geared toward protein separations, the inside capillary surface is modified to mask the free silanol groups. This modification minimizes protein adsorption, but it also dramatically reduces the electroosmotic flow, leading to much longer separation times.

Additives to the buffer in the reservoirs (the run buffer) can also be used to change the mechanism of separation. For example, cyclodextrins have been employed for the separation of a number of optical isomers. Other additives, such as organic solvents and metal ions, have been utilized to improve selectivity or to change the electroosmotic flow. See ELECTROKINETIC PHENOMENA.

Micellar electrokinetic chromatography. In micellar electrokinetic chromatography, an ionic surfactant is added to the run buffer above the critical micelle concentration. Analytes partition in and out of the micelle based on their hydrophobicity (and charge). The migration time is determined not only by site and charge, as in capillary zone electrophoresis, but also by the amount of time that the analyte spends in the micelle relative to that spent in the run buffer. The most common surfactant employed is dodecyl sodium sulfate, which produces negatively charged micelles. In this case, hydrophobic compounds are retained longer on the column because of the high negative electrophoretic mobility of the anionic micelles. The resulting separation is very similar to that obtained by reversed-phase chromatography. This technique is especially useful for the separation of neutral compounds. However, charged compounds can exhibit additional ionic interactions with the micelles, which can affect their migration time. See MICELLE.

Capillary gel electrophoresis. This is the version of capillary electrophoresis most similar to conventional

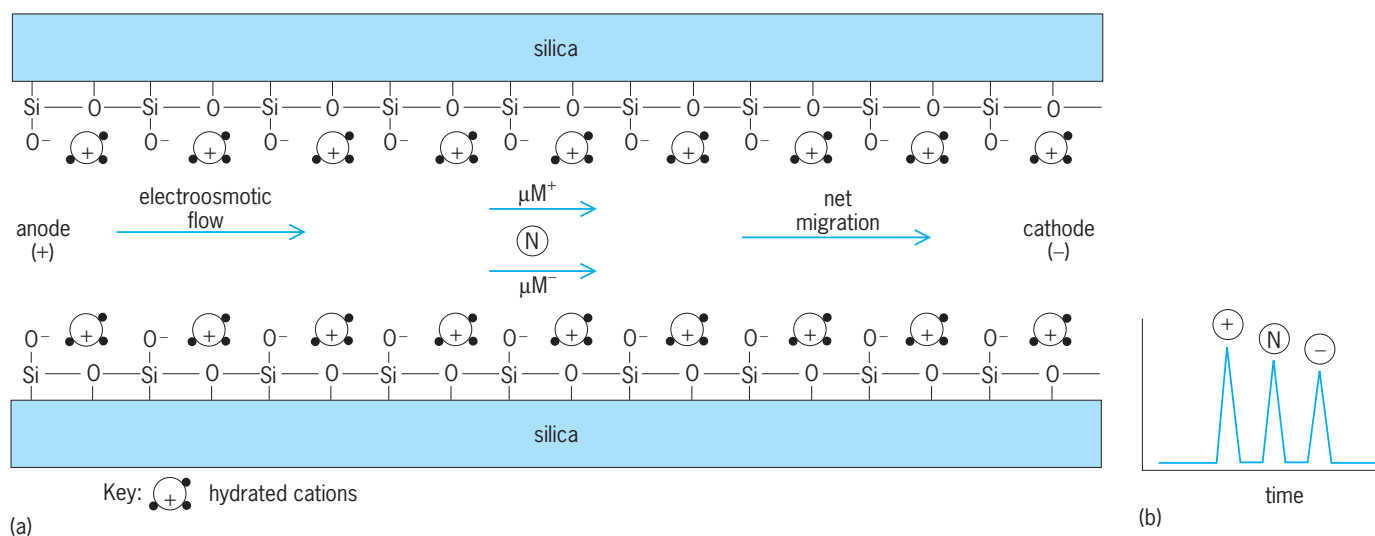


Fig. 5. Capillary zone electrophoresis. (a) Separation mechanism showing electrophoretic mobility of the positive ion (μM^+) and negative ion (μM^-); N is a neutral molecule. (b) Migration order of the ions.

slab gel electrophoresis. It has been applied extensively to the separation of oligonucleotides and proteins. In this form of capillary electrophoresis, a polyacrylamide gel (or other polymeric material) is placed inside the fused silica capillary, which has been modified so that the electroosmotic flow is essentially zero. The gel can be cross-linked, but separations with non-cross-linked gels have also been accomplished. Separation is based on size and charge. Since oligonucleotides exhibit high negative electrophoretic mobilities, their detection takes place at the anodic end of the capillary. Separations of oligonucleotides differing by only one base can easily be obtained with this technique. *See* OLIGONUCLEOTIDE.

Capillary isoelectric focusing. Isoelectric focusing can also be accomplished in the capillary electrophoresis format. As with capillary gel electrophoresis, a capillary with reduced electroosmotic flow must be employed. A pH gradient is created in the capillary by filling it with a mixture of the solute and a solution of an ampholyte, an electrolyte having both acidic and basic characteristics. For the focusing to take place, one buffer reservoir is filled with a basic solution and the other with an acidic solution. The high voltage is then applied. Under these conditions, both the solute and the ampholytes migrate until they reach a region where the pH is equal to their pI. When the isoelectric focusing process is completed, all charge is neutralized and the separation current goes to zero. The zones can then be mobilized for detection by pressure or through the addition of salt to one of the reservoirs. Proteins differing by less than 0.005 pI unit have been separated by using this technique.

Capillary isotachopheresis. This method, also known as displacement electrophoresis, is a completely different mode of capillary electrophoresis in which all the analytes move at the same velocity. Only one class of ions (anions or cations) can be separated in capillary isotachopheresis. In the case of an anion separation, the capillary is filled with an electrolyte containing an anion that has a higher mobility than any of the analytes. The sample solution is then loaded on the capillary, and the cathodic end of the capillary is submerged in an electrolyte solution containing anions of lower mobility than any of the analytes in the mixture. When the high voltage is applied, the anions in the lead electrolyte, which have the highest electrophoretic mobilities, move quickly toward the anode, producing a gap of low conductivity and high electric field strength. The electric field causes analyte ions to move toward the gap, with those exhibiting the largest electrophoretic mobilities moving first. Eventually, a steady state is reached in which each solute is isolated in its own band of differing field strength and all bands are moving at the same velocity. This method has also been employed as a preconcentration step for other modes of capillary electrophoresis.

All the capillary electrophoresis methods have the advantage of the ability to analyze small volumes (typical injection volumes are 1–50 nanoliters). This makes it possible to analyze very small samples or to

use the same sample for several different analyses. One unique application of this technique is the determination of amino acids and neurotransmitters in single cells. Susan Lunte

Alternating-field electrophoresis. Alternating-field agarose gel electrophoresis is a technique for separating very large molecules of deoxyribonucleic acid (DNA); fragments of DNA ranging in size from 30 to 10,000 kilobasepairs (kb) can be resolved. For the molecular biologist, this is a considerable advance over conventional agarose gel electrophoresis, which is limited to the resolution of less than 50 kb.

Conventional gel electrophoresis employs a single pair of electrodes to generate an electric field that is constant in both time and direction and that is uniform across the gel. DNA molecules are negatively charged with a uniform charge-to-mass ratio and thus migrate steadily toward the positive electrode. Although DNA is a linear molecule, in solution it tends to collapse into a random coil configuration. Agarose is a porous material that acts like a sieve, retarding the movement of the DNA; the larger the molecule (and therefore the larger the random coil), the more the retardation, and thus the molecules separate on the basis of size. However, above approximately 50 kb, the dimensions of the random coil are larger than the pore size of the agarose. The DNA can no longer be sieved through the gel, and resolution is lost.

The DNA can be visualized as a long, snake-like molecule. It tends to align itself parallel to an external electric field because it is uniformly charged along its length. This is the principle behind alternating-field electrophoresis. Contrary to conventional electrophoresis alternating-field electrophoresis does not use a constant electric field but one which regularly alternates in direction. With each change in field direction, the DNA molecules attempt to reorient themselves. When this happens an end or a small loop of the molecule, which has dimensions much smaller than those of the random coil of the entire molecule, may be positioned by a pore in the agarose. The electric field can then pull the DNA by the end through the hole. When the field regularly alternates from one direction to another, the DNA regularly reorients and is pulled, in a snakelike fashion, through an adjacent hole.

Not all molecules in the system will make equal progress under these conditions, because not all molecules can reorient themselves with equal speed. The larger the molecule, the more time it requires in a given field strength (determined by the applied voltage) to change directions and the less time it has left to move. This is the basis for the separation seen in alternating-field electrophoresis. Small molecules will reorient themselves quickly and spend most of the time migrating. Larger molecules will reorient slowly and be unable to move very far. Very large molecules may require such long times for reorientation that there is no time remaining for migration. Kathleen Gardiner

Bibliography. P. Bocek et al. (eds.), *Analytical Isotachopheresis*, 1988; P. Camilleri (ed.), *Capillary*

Electrophoresis: Theory and Practice, 2d ed., 1997; N. Guzman (ed.), *Capillary Electrophoresis Technology*, 1993; B. D. Hames (ed.), *Gel Electrophoresis of Proteins: A Practical Approach*, 3d ed., 1998; R. Weinberger, *Practical Capillary Electrophoresis*, 2d ed., 2000; R. Westermeier, *Electrophoresis in Practice: A Guide to Methods and Applications of DNA and Protein Separations*, 4th ed., 2005.

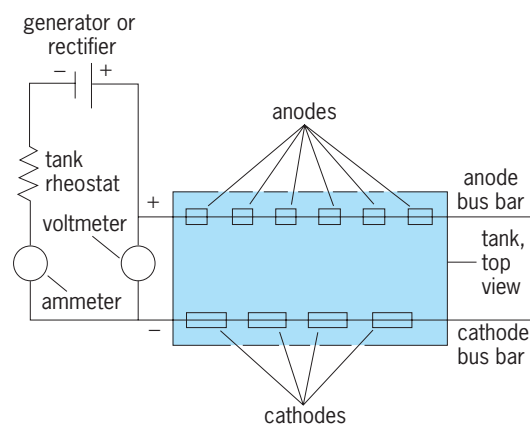
Electroplating of metals

A process for applying a coating to a metal surface by electrochemical means. Electroplating is extensively used to produce printed circuit boards. Its main advantage is that the circuit can be produced directly rather than having to be etched out of a piece of copper sheet. Electroplating is also widely used to impart corrosion resistance. Most parts of automobile bodies are zinc plated for corrosion resistance. Since zinc is more readily attacked by most corrosive agents that automobiles encounter than steel, it provides galvanic or sacrificial protection. An electrolytic cell is formed in which zinc, the less noble metal, is the anode and the steel, the more noble one, is the cathode. The anode corrodes and the cathode is protected. Zinc also provides a good base for paint. If a metal is more noble than the one on which it is electroplated, it provides protection against corrosion only if it is completely continuous. If a small area of the substrate is exposed such as under a pin-hole, corrosion occurs there, rapidly forming a pit. See ELECTROCHEMISTRY.

An example of an electroplated coating applied primarily for wear resistance is hard chromium on a rotating shaft. Electroplating is also used to build up worn or undersized parts. Gold-plated jewelry is an example of a decorative application. Gold and also palladium are electrodeposited on electrical contacts. Here the absence of an oxide film avoids the rise in the electrical resistance of the contact. Nickel and chromium are plated for some decorative applications; however, their wide use in the automotive industry has diminished considerably, primarily because of the associated environmental problems. Magnetic components made of such alloys as permalloy can be manufactured by electroplating.

Process

The electroplating process consists essentially of connecting the parts to be plated to the negative terminal of a direct-current source and another piece of metal to the positive pole, and immersing both in a solution containing ions of the metal to be deposited (see *illus.*). The part connected to the negative terminal becomes the cathode, and the other piece the anode. In general, the anode is a piece of the same metal that is to be plated. Metal dissolves at the anode and is plated at the cathode. If the current is used only to dissolve and deposit the metal to be plated, the process is 100% efficient. Often, fractions of the applied current are diverted to other reactions



Typical connections for a simple electroplating process.

such as the evolution of hydrogen at the cathode; this usage results in lower efficiencies as well as changes in the acidity (pH) of the plating solution. In some processes, such as chromium plating, a piece of metal that is essentially insoluble in the plating solution is the anode. When such insoluble anodes are used, metal ions in the form of soluble compounds must also be added periodically to the plating solution. The anode area is generally about the same as that of the cathode; in some applications it is larger. See ELECTRODE; ELECTROLYSIS.

Most plating solutions are of the aqueous type. There is a limited use of fused salts or organic liquids as solvents. Nonaqueous solutions are employed for the deposition of metals such as aluminum that have overvoltages lower than hydrogen. Such metals cannot be plated in the presence of water, as hydrogen would be preferentially reduced. See OVERVOLTAGE.

In addition to metal ions, plating solutions contain relatively large quantities of various substances used to increase the electrical conductivity, to buffer, and in some instances to form complexes with the metal ions. Relatively small amounts of other substances, which are called addition agents, are also present in plating solutions to level and brighten the deposit, to reduce internal stress, to improve the mechanical properties, and to reduce the size of the metal crystals or grains or to change their orientation.

The quantity of metal deposited, that is, the thickness, depends on the current density (A/m^2), the plating time, and the cathode efficiency. The current is determined by the applied voltage, the electrical conductivity of the plating solution, the distance between anode and cathode, and polarization. Polarization potentials develop because of the various reactions and processes that occur at the anode and cathode, and depend on the rates of these reactions, that is, the current density. If the distance between anode and cathode varies because the part to be plated is irregular in shape, the thickness of the deposit may vary. A quantity called the throwing power represents the degree to which a uniform deposit thickness is attained on areas of the cathode

at varying distances from the anode. Good throwing power results if the plating efficiency is low because of polarization where the current density is high. A plating solution such as an alkaline cyanide bath, in which the polarization of the cathode increases strongly with increasing current density, has good throwing power. A plating solution such as an acid-sulfate copper bath, which is almost 100% efficient, has poor throwing power.

Two other deposition processes are closely related to electroplating: electroless plating and displacement plating. Both processes require no applied current. Electroless plating is a process which, once initiated, continues; that is, it is autocatalytic. Displacement plating occurs when the metal deposited is more noble than the substrate and the substrate dissolves. The reaction ceases when the substrate is completely covered by a pore-free deposit. *See* ELECTROLESS PLATING.

Preparation

In order for adherent coatings to be deposited, the surface to be plated must be clean, that is, free from all foreign substances such as oils and greases, as well as oxides or sulfides.

Cleaning. Three principal methods are employed to remove grease and attached solids. (1) In solvent cleaning, the articles undergo vapor degreasing, in which a solvent is boiled in a closed system and its vapors are condensed on the metal surfaces. Many solvents used in the past have been found to be toxic. (2) In emulsion cleaning, the metal parts are immersed in a warm mixture of kerosine, a wetting agent, and an alkaline solution. (3) In electrolytic cleaning, the articles are immersed in an alkaline solution, and a direct current is passed between them and the other electrode, which is usually steel. Heavily soiled articles are cleaned in solvent or alkaline spray machines. Cleaning solutions may contain sodium hydroxide, carbonate, phosphate, and metasilicate, plus wetting agents and chelating agents. More highly alkaline solutions are used for steel than for other metals. Most of the cleaning is accomplished by the scrubbing action of the evolved gases and the detergency of the components of the solution. The articles may be connected as anodes, as is usual for steel, or as cathodes, as is usual for other metals. Electrolytic cleaning is usually the last cleaning step.

Ultrasonic cleaning is also used extensively, especially for blind holes or gears packed with soils. Ultrasonic waves introduced into a cleaning solution facilitate and accelerate the detachment of solid particles embedded in crevices and small holes. Frequencies from about 18,000 to 24,000 Hz are usually employed. They produce cavitation, which causes rapid local circulation. *See* CAVITATION; ULTRASONICS.

Pickling. In this process, oxides are removed from the surface of the basis metal. For steel, warm, dilute sulfuric acid is used in large-scale operations because it is inexpensive; but room-temperature, dilute hy-

drochloric acid is also used for pickling because it is fast acting. In cathodic pickling of steel, attack of the metal is retarded while the oxide is being dissolved. In addition to rough pickling, acid treatments to activate the surface just prior to plating are often used.

Hydrogen embrittlement may be caused by the diffusion of hydrogen into steel during pickling and also in the certain plating operations. Especially with high-carbon steels, hydrogen causes cracking (a reduction in the fatigue strength and ductility). Hydrogen is gradually evolved on standing, and more rapidly evolved by heating to about 390°F (200°C). *See* DIFFUSION; EMBRITTLEMENT.

Electropolishing. Electropolishing is the preferential, electrochemical dissolution of asperities on surfaces. It results in a smoother and brighter surface. Thin electrodeposits tend to copy the morphology of the surface on which they are plated. Therefore, to attain a bright deposit, the surface on which it is plated must also be bright, a property which can be attained by electropolishing. The electropolished surface is smooth, clean, and free of deformed material, making it an excellent surface for plating. *See* ELECTROPOLISHING.

Equipment

The two classes of electroplating equipment are electrical and mechanical.

Electrical equipment. Most plating operations are conducted with direct current at potentials of 3–12 V. Since electricity is delivered as alternating current at 220 V, it is necessary to reduce the voltage and to rectify the current. Both motor generators and rectifiers are used. Selenium, germanium, and silicon rectifiers have been used extensively in plating. Special generators are used for pulse plating; the pulse variables are often computer controlled. *See* DIRECT-CURRENT GENERATOR; SEMICONDUCTOR RECTIFIER.

Mechanical equipment. The plating tanks are either poly(vinyl chloride) or steel, which requires no lining for alkaline solutions. For neutral and acid solutions, steel tanks are lined with rubber or plastic. For chromic acid baths, lead or plastic linings are used. Most large plating operations are conducted in conveyor tanks. In semiautomatic conveyors, the cathode racks are carried only through the plating tank. In fully automatic conveyors, the cathodes are carried successively through tanks that contain cleaning, pickling, and plating solutions, with intermediate rinses.

Plating of continuous strip is another important operation. Rolls of copper foil are produced by this method. A rotating hard-chrome-plated drum is partially immersed in the plating solution, where it is the cathode. Copper is deposited on the immersed part of the drum. As the copper foil does not adhere to chromium, it peels off as it emerges from the solution. After rinsing and drying, the foil can be rolled up. Tin is also plated on steel by the continuous-strip method.

Small objects are plated in “barrels,” usually of hexagonal shape, with perforated plastic sides. The

barrels rotate on a horizontal axis in a tank containing the plating solution and anodes. The articles contact cathode connections as they tumble during the rotation of the barrels.

Plating of Specific Metals

Most metals can be electroplated from either aqueous or fused-salt solutions. The more important metals plated from aqueous baths are chromium, copper, gold, nickel, silver, tin, and zinc. Alloys can also be electroplated. Electrodeposited copper-zinc and lead-tin alloys are used extensively.

Chromium. Electroplated chromium is used primarily to produce wear- and corrosion-resistant coatings. Chromium is not deposited for decorative purposes as extensively as in the past because of the associated pollution problem from the discharge of hexavalent chromium. Chromium plating solutions consist primarily of chromic oxide, sulfuric acid, and water. The solutions in which a hard chrome finish is deposited contain a lower ratio of chromic oxide to sulfuric acid than the ones used for the decorative type, and are maintained at a somewhat higher temperature. Decorative chromium is plated in a high state of internal stress, generally leading to cracking except in very thin deposits. If there are only a few cracks, the undercoat, which is generally nickel, may be subjected to very high local corrosion currents. To avoid this problem, the chromium deposits are intentionally microcracked or made porous in order to increase the exposed area of the undercoat and thus reduce the current density. The throwing power and current efficiency of chromium plating solutions are rather poor. *See* CHROMIUM.

Copper. Electroplated copper is used extensively in the manufacture of printed circuit boards. As copper cannot be directly plated on the insulating-material substrates, they must be rendered electrically conductive first. The main advantage of using electrolytically deposited copper to produce printed circuits is that the areas of the board that are made conductive can be controlled. The actual circuit can be produced by selective etching or selective plating involving techniques such as photosensitizing, photoresist, and etch resist. Areas exposed through a mask can become conductive by coating them with electroless copper after activation with a solution of stannous and palladium chloride. Suitable organic materials are also used to render the board conductive. The areas made conductive serve as the substrate for copper electrodeposition. Some circuits are produced only with electroless plating of copper. Only the through holes are plated, generally with electroplated copper over electroless plated copper when the boards are made by laminating copper sheet to the plastic substrate. The copper segments of the printed circuit may be coated with an electroplated tin-lead alloy to facilitate subsequent soldering and also to protect them from oxidation. Gold deposits are also used for this purpose. There are other uses of electroplated copper in the electronic industry, such as in the production of microchips. Electroplated copper is also the undercoat for decorative

nickel-chromium deposits. *See* COPPER; PRINTED CIRCUIT.

Gold. Gold plating is used for electrical contacts which must remain free of oxides, connections for microcircuits, certain information-storage devices, solid-state components, and jewelry. The use of gold for electrical and electronic application exceeds that for decorative purposes. Electroplated palladium is being substituted for gold in some applications. *See* GOLD; PALLADIUM.

Nickel. Nickel coatings covered with chromium provide corrosion-resistant and decorative finishes for steel, brass, and zinc-base die castings. The most widely used plating solution is the Watts bath, which contains nickel sulfate, nickel chloride, and boric acid. All-chloride, sulfamate, and fluoroborate plating solutions are also used. The sulfamate solution is used for low-stress applications. There are a number of compounds, mostly organic, which can be added to nickel-plating baths. When specularly bright nickel is plated, at least one from each of two groups of compounds called class I and class II brighteners must be added to the plating solution. Class I compounds, for example, benzene and naphthalene disulfonic acids, contain sulfur. Class II brighteners, for example, butyne diol, do not contain sulfur. Class II brighteners can also cause leveling, which is producing a deposit that is smoother than the surface on which it is plated. When class II brighteners are added to plating solutions without sulfur-containing agents, they result in semibright nickel deposits. For improved corrosion protection, a duplex nickel deposit, that is, a relatively thick, semibright, sulfur-free plate covered with a thinner bright coating, is frequently used. A thin chromium layer is plated over the bright nickel. If corrosion begins through a pore or crack in the chromium, it will penetrate to the more noble semibright nickel and then spread laterally. Thus, rust spots that result from penetration to the steel are less likely to develop. Exposure of the skin to nickel-plating solutions can cause a rash and severe itching. The sensitivity of persons to nickel salts varies considerably.

Electroless nickel deposits are used for decorative purposes on nonmetallic substrates. Adhesion occurs mainly mechanically by the penetration of the deposit into crevices caused by a prior selective-etch treatment of the substrate. The reducing agents in electroless nickel-plating solutions are either hypophosphites or borohydrides. The deposits are highly supersaturated solid solutions of phosphorus or boron in nickel. The grain size of electroless nickel containing less than 7% by weight of phosphorus is extremely small. Deposits of higher phosphorus content can be considered to be amorphous and are also nonmagnetic. Electroless nickel can also be precipitation-hardened by heat treatment. Electroless nickel plated on metallic substrates produces a uniform thickness over irregularly shaped parts. *See* HEAT TREATMENT (METALLURGY); NICKEL; SOLID SOLUTION.

Silver. The principal use of silver electrodeposits is for tableware because of their corrosion resistance

(except to sulfur-containing foods) and pleasing appearance. Other important uses are for bearings and electrical circuits, waveguides, and hot-gas seals. The plating solutions are of the cyanide type and generally contain additives that produce bright deposits. See SILVER.

Tin. Tin is used in electrodeposition as a component of solders. Solder is electroplated over copper to protect it from oxidation and to facilitate subsequent joining operations. The advantage of electroplating the solder is that it can be applied only where it is needed. Because of the desirability of eliminating lead in solders, those with higher tin contents are used.

Tin-plated steel for cans in which food is preserved has limited use, because lacquers are preferred to prevent contact between steel and food. Tin plating is employed for refrigerator coils and bearing surfaces. Bright pore-free surfaces can be produced by melting the electroplated tin and allowing it to "reflow."

Stannous chloride and stannous sulfate are the main components of acid tin-plating solutions. The alkaline solution contains sodium stannate or potassium stannate. Both types are used in the production of electrolytic tin plate. The main components of solutions for the electrodeposition of solders are lead fluoroborate and tin fluoroborate. Sulfonates are also used. See TIN.

Zinc. The sacrificial protection of steel against corrosion is the main reason for plating zinc. It is used extensively in the automobile industry for this purpose. Screws, bolts, and washers are barrel-plated with zinc also for corrosion protection. Continuous zinc electroplating of wire and strip is another application. The advantage of electrodepositing zinc over hot dipping is the ability to apply thinner coatings and higher purity. A conversion film formed by dipping in a chromate solution inhibits the formation of white corrosion products. Acid plating solutions are used preferentially over alkaline cyanide ones because of the pollution associated with the latter. See ZINC.

Properties of Electrodeposits

In the various applications of electrodeposits, certain properties must be controlled and, therefore, must be measured. Specifications have been supplied by the American Society for Testing and Materials, and by other technical societies and some governmental agencies, establishing the various properties and the methods of determining them. The properties of electroplated metal which should be considered, depending on the use of the deposit, are thickness, adhesion to the substrate, brightness, corrosion resistance, wear resistance, the mechanical properties of yield strength, tensile strength, ductility, and hardness, internal stress, solderability, density, electrical conductivity, and the magnetic characteristics.

Thickness. Several properties depend on the thickness of an electrodeposit and how it varies over a work piece. Several types of gages and also micro-

scopic, gravimetric, and chemical methods are used to measure thickness.

X-ray fluorescence is most often used to measure deposit thicknesses. When the electrical conductivity of the substrate and the deposit differ considerably and neither is ferromagnetic, eddy-current gages can be used to measure deposit thickness. Beta backscatter gages can measure the thickness of a deposit if its atomic number is considerably different from that of the substrate, as is the case with gold and nickel. See NONDESTRUCTIVE EVALUATION; X-RAY FLUORESCENCE ANALYSIS.

Microscopic methods of thickness determination require that the cross section be metallographically prepared and be precisely perpendicular at the interface between substrate and deposit or at a known angle of inclination. The thickness of thin deposits should be determined by scanning electron microscopy, while optical microscopes have adequate resolution for thick ones.

Chemical thickness measurements involve a calibration of the time needed for a reagent to penetrate to a certain depth. Then the time for the same reagent to reach the substrate is measured and the deposit thickness calculated. If the substrate can be dissolved without attacking the deposit, the latter can be weighed. If the area and the density are accurately known, the average thickness can be calculated.

Adhesion. Good bonding between substrate and deposit is very important in almost all plating applications. Poor adhesion results in peeling or blistering of the deposit. Generally, the adhesion between two metals is strong. Poor bonding results if the substrate surface was not clean prior to plating, if foreign substances were absorbed, or if brittle phases form as a result of interdiffusion at elevated temperatures. The bond between a nonmetallic substrate and the electroless plate is primarily mechanical and can therefore vary considerably. Adhesive coatings on the substrate can improve the bonding to electroless deposits.

Adhesion tests are not well standardized. In the peel test, which is widely used, part of the deposit is loosened from the substrate, and the force required to pull off the rest is measured. A tensile test, in which the bonded interface is placed perpendicular to the applied force, is also used. Scraping tests and attempts to interpose a knife edge between substrate and deposit are examples of qualitative tests.

Brightness. For an electrodeposit to be specularly bright, the hills and valleys of the surface morphology should not vary in height by more than the wavelength of light. This requirement is generally fulfilled when the crystal size is very small. Addition agents in the plating solutions are also absorbed on faster-growing sites, permitting the microscopic recessed region to catch up. Quantitative measurements of brightness are rarely performed; comparisons by eye are generally made.

Corrosion. With sacrificial coatings, corrosion protection occurs as long as enough of the coating remains; thus protection depends mainly on the thickness of the plate. In order for a substrate to be

protected against corrosion by a deposit that is more noble, the coating must be completely free of pores, cracks, or discontinuities. The substrate must not be exposed at any point to the corrosive medium. When the less noble substrate is completely covered, the corrosion resistance of the coating itself is enhanced if the structure and chemical compositions are uniform.

Corrosion testing is performed either under actual or simulated service conditions or under accelerated exposure. The conditions under which accelerated exposure tests are performed are more severe than those encountered in service, so that the results can be seen more quickly. Salt-spray and acetic acid-modified salt-spray accelerated corrosion tests are most frequently used. The copper-accelerated acetic acid-salt-spray test and a test in which parts are exposed to a slurry of clay in a salt solution are examples of the second type. These tests are used to duplicate under accelerated conditions the corrosion of plated automotive parts in winter conditions. Electrochemical measurements of the corrosion current density and potential can also give an indication of corrosion resistance. *See* CORROSION.

Wear resistance. Although a high hardness value is an indication of good wear resistance, there are factors relating to wear resistance that are not well understood. There are also several types of wear, the two most common types being adhesive and abrasive wear. Corrosion and lubrication can greatly affect wear. Wear resistance is generally determined by moving two contacting surfaces relative to one another and noting the weight loss after some period of time. *See* WEAR.

Mechanical properties. In many respects, the mechanical properties of electrodeposits are similar to those of metals that have been heavily plastically deformed, that is, cold-worked, in that they exhibit higher strength and hardness and lower ductility than the same materials in the as-cast or annealed state. The number of crystal defects such as dislocations and lattice vacancies in electrodeposits is of the same order of magnitude as that found in severely cold-worked metals. Electrodeposited metals also recrystallize on heating. The small crystal size which electrodeposited metals frequently possess, and the inclusion in the crystal lattice of foreign materials from the plating-bath additives and side reactions, contributes to the strength. Many electrodeposited metals also have a large number of growth twins, which can increase strength. *See* CRYSTAL DEFECTS.

In general, the same factors that increase strength reduce ductility. Annealing electrodeposits can reduce strength and improve ductility. However, both strength and ductility are reduced upon annealing if a brittle grain-boundary phase is formed, as in the case of nickel deposits from plating solutions to which certain sulfur-containing agents have been added. Pores due to included gas also lower both strength and ductility. *See* GRAIN BOUNDARIES; METAL, MECHANICAL PROPERTIES OF.

The mechanical properties of electrodeposits can be determined by a tensile test. However, it is not

advisable to use the tensile specimen used to test metal sheet, because the ratio of the width to the thickness is generally too large. It is preferable to scale the dimensions. Care must be taken so that the specimen is accurately aligned and the edges are not serrated from cutting. The bulge test is also used. Hardness is usually measured by an indentation test on a metallographically prepared cross section of the deposit. The Knoop indenter is particularly useful for evaluating relatively thin deposits. *See* HARDNESS SCALES.

Internal stress. Most electrodeposits are in a state of internal tensile stress. The stress can be reduced or made compressive by adding certain materials to the plating bath. The addition of saccharin to nickel-plating solutions has this effect. The causes of internal stress are not fully understood. There are indications that the coalescence of crystals or crystallites is one important cause of tensile stress. The diffusion of hydrogen so as to cause the surface layers to shrink can be another cause. Filling of voids with hydrogen gas so as to expand them can result in compressive stress.

Internal stress can be measured from the bending of a strip plated on only one side or from the change in length of a specimen plated on both sides. A long strip wound into a spiral, with one end fixed and the other free, is the most frequently used stress-measuring device. If this spiral is plated on one side, stress causes the free end to rotate. The rotation is calibrated to yield a stress value.

Solderability. Solderability is an important property, especially for plated segments of printed circuits. Solderability depends on the cleanliness of the surface, especially the degree of oxidation after prolonged storage. Coating copper deposits by plating them with solder, tin, or gold can prevent oxidation and thus improve solderability. Solderability is generally determined by how well molten solder pools spread. The spreading is also affected by the temperature and the flux. Because of the different coefficients of thermal expansion of the printed-circuit substrate and the copper deposit, high stresses can develop on heating for soldering. These stresses can result in the fracture of thin and brittle plated segments, especially at through holes. *See* SOLDERING.

Density, and electrical and magnetic properties. There are no special methods of measuring the density of electrical and magnetic properties of electrodeposits. The density of electrodeposits can be reduced because of the presence of gas or other foreign materials. Lattice vacancies can result in higher electrical resistivities of plated metals. The very fine grain size of electrodeposited or electroless plated ferromagnetic metals results in high coercive forces.

Special Applications

Special processes associated with electrodeposition include electroforming, anodizing, pulse plating, and laser plating.

Electroforming. Electroforming is a special type of plating in which thick deposits are subsequently removed from the substrate, which acts as a mold. The

process is particularly suitable for forming parts that require intricate designs on inside surfaces, for example, waveguides.

Intricate machining operations can be performed much more easily on outside surfaces. First, the outside surface of the substrate mold is machined; then, the contours of the design can be transferred to the inside surface of the deposit; and, finally, the deposit is separated from the mold. The mold or matrix can be either metallic or nonmetallic. Nonmetallic molds must be rendered electrically conductive by the application of a powder or by chemical reduction, electroless plating, or vapor deposition. For nonadherent deposits, substrate removal is easy; otherwise, the substrate must be dissolved or melted away. Important applications of electroforming are in the production of phonograph record masters, printing plates, and some musical instruments as well as in waveguides. See WAVEGUIDE.

Anodizing. In anodizing, a process related to plating, an oxide is deposited on a metal that is the anode in a suitable solution. The process is primarily used with aluminum, but it can be applied to beryllium, magnesium, tantalum, and titanium. Relatively thick oxide deposits can be produced, even though they are electrically insulating. The presence of small pores through which the solution can reach the metal surface permits the continuation of the reaction. The electrical current is carried by the electrolyte in the pores.

Solutions for anodizing aluminum generally contain either chromic or sulfuric acid. After anodizing, the pores should be sealed, to improve protection of the substrate, by a hot-water or steam treatment, which causes hydration and a resulting volume expansion of the oxide. Colored coatings can be produced by the incorporation of dyes.

Pulse plating. The maximum current density in electrode position is limited by the replenishment by diffusion or convection of the metal ions in the solution layer adjacent to the work piece. If the current is pulsed so that it is off for a certain time interval, higher current densities can be used because ions are replenished during the off time. The resulting higher deposition potentials can result in grain refinement, which, in turn, can lead to brighter and stronger deposits without the use of addition agents. Internal tensile stresses can also be reduced by pulse plating.

Laser plating. In many electronic applications it is necessary to produce very fine lines of deposits. When laser radiation is directed onto a substrate, the adjacent solution is heated and the plating rate so greatly increased that for all practical purposes deposition is limited to the illuminated area. Thus a laser can be used to draw the desired deposit design on the substrate. Directing a jet of plating solution onto the substrate can accomplish essentially the same result. Laser and jet plating can also be combined to give even better control over the dimensions and shape of the deposit.

Rolf Weil

Bibliography. J. W. Dini, *Electrodeposition: The Materials Science of Coatings and Substrates*, 1993;

L. J. Durney (ed.), *Graham's Electroplating Handbook*, 4th ed., 1984; F. Lowenheim (ed.), *Guide to the Selection and Use of Electroplated and Related Finishes*, 1982; L. T. Romankiew and D. R. Turner (eds.), *Electrodeposition Technology: Theory and Practice*, 1987; W. H. Safranek, *Properties of Electrodeposited Metals and Alloys*, 2d ed., 1986; P. W. Wild, *Modern Analysis for Electroplating*, 2d ed., 1991.

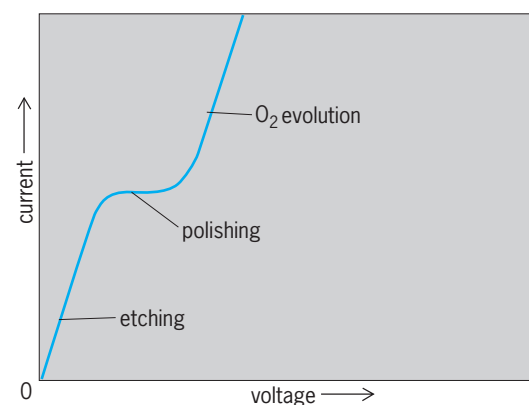
Electropolishing

A method of polishing metal surfaces by applying an electric current through an electrolytic bath in a process that is the reverse of plating. The metal to be polished is made the anode in an electric circuit. Anodic dissolution of protuberant burrs and sharp edges occurs at a faster rate than over the flat surfaces and crevices, possibly because of locally higher current densities. The result produces an exceedingly flat, smooth, brilliant surface.

Electropolishing is used for many purposes. The brilliance of the polished surface makes an attractive finish. Because the polished surface has the same structural properties as the base metal, it serves as an excellent surface for plating. Electropolishing avoids causing differential surface stresses, one of the requirements for the formation of galvanic cells which cause corrosion. Because no mechanical rubbing is involved, work hardening is avoided. Contaminants, which often are associated with the use of abrasives and polishing compounds, are also avoided. The surface is left clean and may require little or no preparation for subsequent treatment or use. Electropolishing also minimizes loss of high-temperature creep-rupture strength.

In electropolishing, the work is submerged in an electrolyte and connected to the positive terminal (anode) of a source of direct-current (dc) power. The negative terminal is connected to an electrode that will resist chemical interaction with the electrolyte. Carbon is often used. The electrolyte is usually a concentrated acid, although alkaline solutions, dilute acids, and salts have been successfully employed.

Current density on the work surface is crucial. The relationship of anode current to anode voltage



Anode current versus anode voltage in electropolishing.

Electropolishing processes				
Metal to be polished	Electrolyte	Temperature, °F (°C)	Current density, A/ft ² (A/m ²)	Polishing time, min
Aluminum	Phosphoric acid and ethylene glycol	180 (82)	140 (1500)	2–5
Brass and copper	Copper cyanide, zinc cyanide, and potassium cyanide	80 (27)	50 (500)	0.1–0.5
Molybdenum and niobium	Sulfuric acid and hydrofluoric acid in methyl alcohol	80–130 (27–54)	4500 (50,000)	0.05–0.4
Silver	Silver (powder), potassium cyanide, and potassium carbonate	80–125 (27–52)	15–20 (160–200)	0.1–0.2
Steel	Sulfuric acid and phosphoric acid, chromic or humic acids may be added	160 (71)	15–25 (160–270)	2–5
Tantalum and tungsten	Sulfuric acid and hydrofluoric acid in methyl alcohol	80–130 (27–54)	4500 (50,000)	0.05–0.4
Zinc*	Chromic oxide	65–90 (18–32)	1500 (16,000)	1–1.5

*After treatment, zinc surfaces must be successively dipped in dilute solutions of potassium chromate (2 min) and hydrochloric acid (5 s).

is shown in the **illustration**. Below a certain level of voltage, etching occurs. Above that level a constant-current region is reached in which polishing occurs. One mechanism proposed to explain this effect suggests that a film, or a systematic structuring of solution contents and ions, forms over the work surface and that this film, being thinner over burrs and sharp protuberances, allows selective dissolution of the surface to cause flatness. Whether or not this is correct, raising the voltage to a level sufficient to cause oxygen evolution results in solution agitation at the work surface and thereby interferes with polishing.

Acids commonly used as electrolytes, either alone or in combination, include acetic, chromic, citric, hydrofluoric, nitric, phosphoric, and sulfuric, with phosphoric playing a major role in many processes. In most processes the temperature of the electrolyte is usually maintained in the range 90–250°F (30–120°C). Current densities vary widely, but 100–500 A/ft² (1100–5400 A/m²) of surface will remove moderate amounts of metal. Current densities up to 5000 A/ft² (54,000 A/m²) may be employed if large amounts of metal are to be removed. Polishing times vary from a few seconds to 15 min or longer (see **table**).

In general, higher electrolyte temperatures, higher current densities, and longer polishing times all tend to produce brighter finishes. The range of finishes includes all gradations from satin to mirror, control being obtained through selection of electrolyte, solution concentration, operating temperature, and polishing time.

Numerous metals and alloys electropolish well. Aluminum, both wrought and cast, can be electropolished and anodized to yield a finish that rivals chromium plate. Brass, copper, chromium plate, gold, nickel, nickel alloys, carbon steel, the 300 and 400 stainless steels, and zinc can all be electropolished. See ELECTROPLATING OF METALS.

William W. Snow

Bibliography. R. M. Burns and W. W. Bradley, *Protective Coatings for Metals*, 3d ed., 1975; D. R. Gabe, *Principles of Metal Surface Treatment and Protec-*

tion, 2d ed., 1978; V. Sedlacek, *Metallic Surface, Films, and Coatings*, 1993.

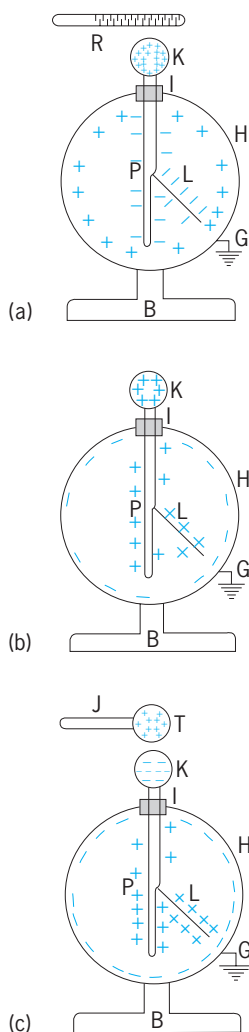
Electroscope

An instrument for detecting the presence and sign of an electric charge. It is the simplest type of ionization chamber. See IONIZATION CHAMBER.

The **illustration** shows a common type of simple gold-leaf electroscope. Gold leaf, shown as L, is used because it is an extremely thin conducting foil which has low mass per unit area and is very flexible. Hence, it responds quickly and vigorously to small electrostatic forces. Aluminum foil is almost as satisfactory as gold foil. In the illustration P is a metal support post for L; I is an insulator of high quality through which P passes and terminates in a metal knob K; and H is a cylindrical metal housing with flat ends and windows so located that the motion and final position of L are visible. H serves as an electrostatic shield that is connected to the environment, as well as a shield against air currents. The base B supports the electroscope.

The hard-rubber rod R (illus. *a*), which has been given a net charge by rubbing, has set up the charge distribution shown, by the process of electrostatic induction, since the isolated structure consisting of K, P, and L has to remain in total electrically neutral. P and L have the same sign of charge and so L is repelled from P. With R still present, K is connected to H, and then K, P, L, and H have a total induced positive charge. If K is disconnected from H and finally R is taken away, the remaining positive charge on K, P, and L induces a negative charge on H (illus. *b*). At this stage the electroscope has a positive charge on its leaf system. This charge and the leaf deflection decay in minutes through leakage currents across the surface of I.

If an electroscope has a charge of known sign, as in illus. *b*, it can be used to test the sign of an unknown charge, as in illus. *c*, where the metal test ball T, with its insulating handle J, has the unknown charge. In the situation pictured, L moves farther away from P



Electroscope. (a) Being charged by induction by negative charge on hard-rubber rod R. (b) Positive charge left on its leaf after induction process is complete. (c) Testing the sign of an unknown charge on test ball T.

as T is brought slowly up toward K, showing that T has a positive charge. If T had a negative charge, L would move toward P, as T slowly approaches K. The converse situation, if the leaf system in illus. c had a negative charge initially, can be readily visualized.

Although electroscopes have been built with a wide variety of geometries, the principle of operation is essentially the same for all. If an electrostatic voltmeter has a scale, permitting quantitative measurements, it is called an electrometer or electrostatic voltmeter. For information on electrometers See ELECTROMETER; VOLTMETER.

Ralph P. Winch; Bryan P. Kibble

Bibliography. A. F. Kip, *Fundamentals of Electricity and Magnetism*, 2d ed., McGraw-Hill, 1969; Physical Science Study Committee, *Physics*, 7th ed., Kendall Hunt, 1991.

Electrostatic discharge testing

A stress testing procedure for determining the electrostatic discharge (ESD) susceptibility level of electrical and electronic devices and systems.

Beginning in the mid-1970s, active electronic devices such as semiconductor components and integrated circuits became susceptible to fast transient events called electrostatic discharge events. The cores of these more sensitive electronic devices began failing after human handling because there was no built-in protection against the fast transients from charged humans. Over the next 25 years, ESD design engineers were able to design-in ESD protection circuits to match the advanced technology changes. However, at the end of 1990s the dramatic increase in the combination of high speed, low power, high performance, and high pin count (over 1000 pins) resulted in extremely dense circuits in devices and system-in-package (SIP) assemblies. These, when combined in one device or assembly, led to increased sensitivity of the devices and therefore increased susceptibility to the fast ESD transients. See ELECTRIC TRANSIENT; ELECTROSTATICS; INTEGRATED CIRCUITS; STATIC ELECTRICITY.

There are four such ESD events: (1) the transient from a charged human, (2) the transient from a charged machine or piece of equipment, (3) the transient from the charged package of the device, and (4) the transient from a charged human holding a piece of conducting material. Each ESD event is unique, and therefore the associated pulses have very different parameters.

Since the mid-1970s, ESD stress testing methods have measured the differences in the failure voltage thresholds or the failure current thresholds of these ESD-sensitive devices to determine their susceptibility. ESD stress testing simulates (using a tester) what happens when a real-world ESD event "strikes" the pin of an ESD-sensitive device. The core of such a device is susceptible to the fast ESD transients, so ESD protective circuits are placed at every pin or pad of the device.

There are four ESD testing instruments in use to stress the electronic devices and assemblies, one for each of the different ESD events. The testing techniques are detailed in four corresponding ESD standards, which describe methods or procedures for stressing the devices and determining their susceptibility. With the advent of high-performance and high-pin-count devices came the need for modern ESD testing instruments and techniques, the major issue being one of added parasitics (added capacitance and inductance) in the ESD tester. These older tester parasitics modified the shape and other parameters (such as pulse rise time and pulse peak current) associated with the pulse from the ESD event.

Basic principles of stress testing. The basic principles associated with ESD stress testing are founded on the four different ESD events. The first three events each cause a single device or component to fail, and the testers that simulate these events are therefore used to qualify the devices as to their level of susceptibility before they are shipped. The fourth event (human + metal) has a different transient that causes assemblies to temporarily stop functioning (and requires rebooting) but does not destroy the devices, and the tester that simulates this event is

Data measured at 500 V for the four ESD model events					
ESD transients	RC (Ω/pF)	Peak current, A	Rise time, ns	Decay time, ns	Scope bandwidth
Human body (fingers)	1500/100	0.33	2–10	150	500 MHz
Machine	0/200	8.75	8–15	66–90	350 MHz
Charged device	25/4	7.50	<0.2	<0.4	1–6 GHz
Human + metal	330/150	2.375	0.7–1.0	30	1.0 GHz

therefore used to qualify assemblies. Each device type must be stressed by each of the first three ESD transient types, and at several voltage levels, beginning with the lowest (hundreds of volts) and ending with the highest (thousand of volts). At each voltage level, several devices are stressed to ensure that the results are repeatable and reproducible and that process variance is properly taken into account.

The **table** details the differences between the pulse parameters of the four transient ESD events. The values of the main circuit elements (resistance and capacitance) in the tester that simulate each ESD

event are shown in the second column. The next three columns show the peak current, rise time, and decay time for each event. The scope displays the transient pulse from which the parameters can be recorded, and the different bandwidths confirm that each of the four ESD events is very different from the others.

Traditional ESD testers. Figures 1 and 2 show the basic simulators with the resistance–capacitance (RC) network set up to test the devices. The actual testers have more parasitics (stray capacitances and inductances) than shown here, and although the ESD standards specify limits, the parasitic effects are large enough to modify the waveforms. Figure 1 can be used to represent three ESD simulators, depending on the RC network specified. For example, substituting the RC value of 1500 $\Omega/100$ picofarads will produce the human-body circuit for simulation. Figure 2 represents the charged-device simulator–tester, where charging the device occurs through the V_{ss} (ground) bus pins or the V_{cc} (high-voltage) bus pins of the device. Here, R_D is the device resistance, C_D is the device capacitance, and L_D is the inductance. (The value $C_D = 15$ pF is an example from actual stressing of a real device, while the value of 4 pF in the table is used as a calibrator.)

Need for new stress testers. In addition to the inadequacies of additional parasitics in traditional ESD testers, the actual data provided only pass/fail results. Using this go/no-go procedure, the device fabrication processes are changed after a failure, and then the devices are redesigned and retested. This iteration is continued until the desired level for device operation is achieved. However, the increased number of I/O pins coupled with the need for faster and denser structures resulted in increased ESD sensitivity. To design a good ESD protection structure, the designer must understand the basic device failure mechanisms as well as the mechanisms of the various integrated circuit protection structures. By the mid-1980s, the sensitivity of the cores of electronic devices increased to such a high level that a new ESD test method, the transmission line pulse (TLP) ESD test method, had to be developed. This method used the cleaner and more reliable rectangular or square pulse to stress the device, and was employed in these early years to simulate and correlate with the human-body discharge event. The TLP method provided additional information beyond the human-body go/no-go test method.

The TLP ESD test method was used as an in-house laboratory diagnostic tool for engineering evaluation, characterization, and development of the

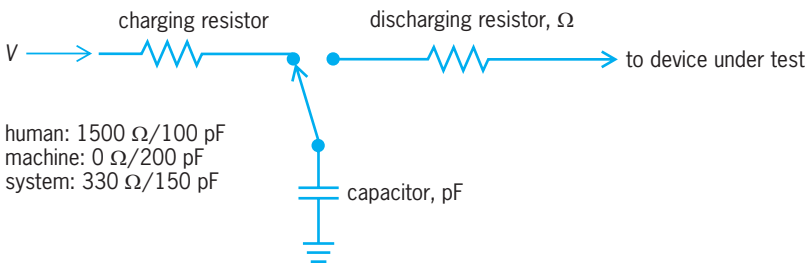


Fig. 1. RC (resistance–capacitance) network for three ESD events: human body, machine, and human + metal (system). Resistance and capacitance are shown for each network.

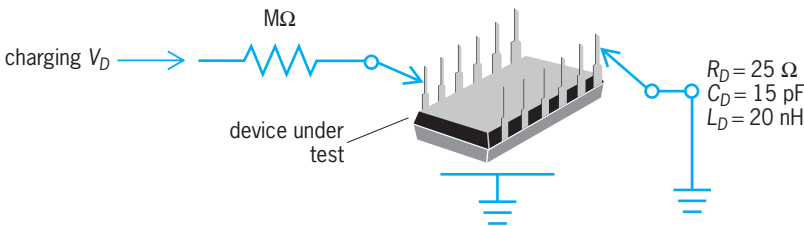


Fig. 2. RCL (resistance–capacitance–inductance) network for charged-device ESD event.

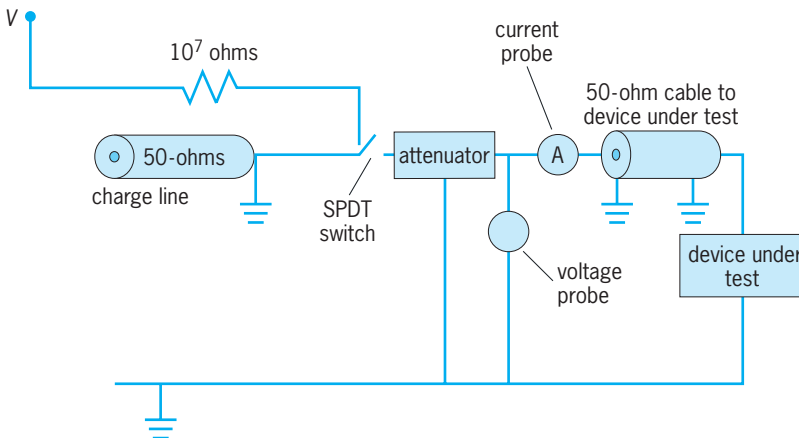


Fig. 3. TLP (transmission line pulse) tester in constant-impedance (50-ohm load) mode.

protection structures, with the goal of improving the robustness of the ESD protective structures in integrated circuits. It did so by allowing the engineer to determine precise current paths in the device and hence to develop improved and new circuits. The failure criteria in TLP testing were linked to the failure criteria of the traditional human-body ESD testing. The TLP pulse parameters and testing results were compared with those from the traditional human-body simulator in order to show that correlation exists between the two methods of ESD testing.

However, the need for higher pin count devices combined with the increased sensitivity of high-performance semiconductor devices led to the need for more improvement in the TLP test method. The first commercially developed TLP tester-simulator was completed in 1999, and the first TLP standard test procedure was developed and published in 2003. It has since evolved from a laboratory diagnostic tool to a standard practice used for engineering evaluation and development, including failure analysis and subsequent redesign for increased robustness of the protective structures.

TLP tester. There are two basic TLP ESD testing configurations: the constant-current (500 → 1000-ohm load) mode and the constant-impedance (50-ohm load) mode (Fig. 3). The 50-ohm constant-impedance mode has inherent higher current capability and is charged by the 50-ohm charge cable line. The attenuator is used to reduce the unwanted reflections. The single-pole double-throw (SPDT) switch is used to go between charging and discharging from the high-voltage charge line. The TLP tester stresses the device under test in voltage steps starting from less than 1.0 V to as high as 500 V. Since a short circuit in this configuration produces 10 amperes, the impedance of the device will determine the device response current and corresponding response voltage.

The ESD protection structure of the device under test responds to the stresses by producing voltage and current values, so a plot of the device data results in a special current-voltage (I - V) curve, much like that seen on a curve tracer plot. Here, however, additional points and parameters [trigger voltage (V_{t1}) and trigger current (I_{t1}), holding voltage (V_{sb}), and device impedance slope] are obtained (Fig. 4). The TLP stress pulses are increased until the device produces a secondary breakdown point called V_{t2} and I_{t2} , where failure begins.

The current-voltage-leakage (I - V - L) plot tells about the complete performance of the ESD protection structure. The leakage current $I_{L, device}$ is plotted against the device response current I_{device} . This dynamic I - V - L curve shows one "leakage current" point for each stress voltage. Figure 4 shows that there is no change in the leakage current (top x -axis scale) up to a specific current value I_{t2} (y -axis scale). Then there is a catastrophic change in the leakage current to coincide with the device second breakdown voltage, V_{t2} , also referred to as the total failure point. The gradual device leakage current evolution (from picoamperes

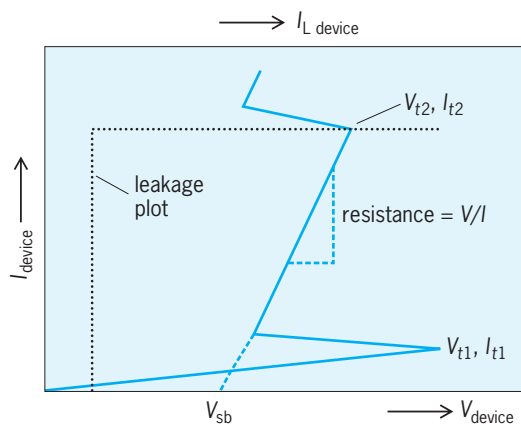


Fig. 4. Current-voltage-leakage (I - V - L) characteristic curve from a transistor device.

to nanoamperes) can change drastically well before the usually defined total failure point. The leakage is taken after each stress pulse because it can determine if a failure has occurred or if an early failure is to be expected. Leakage data are collected at the dc operating voltage of the device. These data are critical because I - V plots alone cannot give any information about possibly early failures; only the leakage plot can.

General ESD stress testing will continue to be the simplest way to check and ensure the voltage or current level at which devices will continue to operate after being stressed by an ESD (transient) pulse. The human-body tester has been the main simulator for the human-body ESD event, but the new TLP test method provides information that is not available using the traditional methods. See ELECTRONIC TEST EQUIPMENT.

Leo G. Henry

Bibliography. J. E. Barth et al., TLP calibration, correlation, standards, and new techniques, *IEEE Trans. Electr. Pack. Manuf.*, 24(2):99-108, April 2001; ElectroStatic Discharge Association, *Electrostatic Discharge (ESD) Technology Roadmap*, 2005; H. Geske, DVI complaint ESD protection to IEC 61000-4-2 level standard, *Conformity*, pp. 12, 14-17, September 2004; L. G. Henry, All ESD testing standards are not created equal: A Rosetta Stone analysis, *Compliance Eng.*, July-August 2003; L. G. Henry et al., Charged device model metrology: Limitations and problems, *Microelectr. Reliab.*, 42:919-927, 2002; L. G. Henry and M. Chaine, Component level ESD testing, chap. 6 in ElectroStatic Discharge Association, *ESD Phenomena and the Reliability for Microelectronics* (white paper), 2002; T. J. Maloney and N. Khurana, Transmission line pulsing techniques for circuit modeling of ESD phenomena, *EOS-7: EOS/ESD Symposium*, p. 49, 1985.

Electrostatic lens

An electrostatic field with axial or plane symmetry which acts upon beams of charged particles of uniform velocity as glass lenses act on light beams. The action of electrostatic fields with axial symmetry is

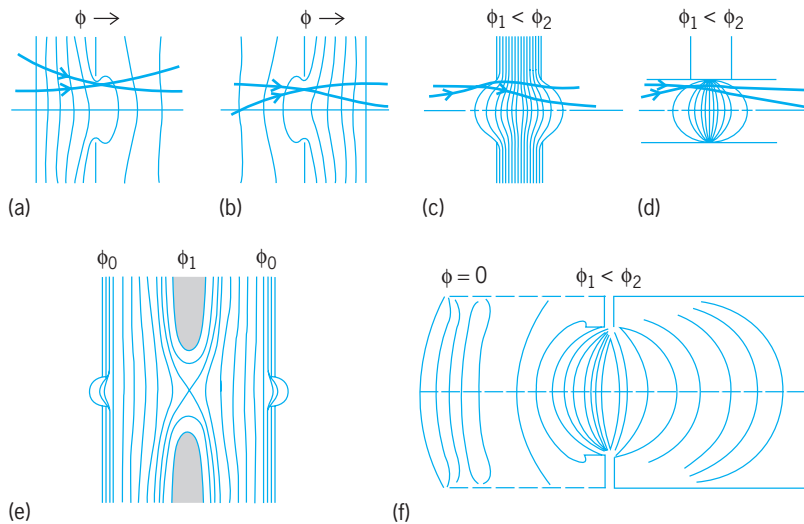


Fig. 1. Axially symmetric electrostatic lenses. (a) Single-aperture lens (decreasing field). (b) Single-aperture lens (increasing field). (c) Two-aperture lens. (d) Two-cylinder lens. (e) Unipotential lens. (f) Cathode lens (image tube). (After E. G. Ramberg and G. A. Morton, *J. Appl. Phys.*, vol. 10, 1939, and V. K. Zworykin et al., *Electron Optics and the Electron Microscope*, Wiley, 1945)

analogous to that of spherical glass lenses, whereas the action of electrostatic fields with plane symmetry is analogous to that of cylindrical glass lenses. Plane symmetry as used here signifies that the electrostatic potential is constant along any normal to a family of parallel planes.

The action of an electrostatic lens on the paths of charged particles passing through it is most readily visualized with the aid of an equipotential plot of the fields in a plane of symmetry of the lens. The equipotential lines in the plot indicate the intersection with the plane of the drawing of surfaces on which the electrostatic potential is a constant. The paths of charged particles in the electrostatic field are bent toward the normals of the equipotentials as the particles are accelerated, and away from the normals as the particles are decelerated. See ELECTRON MOTION IN VACUUM.

Axially symmetric lenses. These lenses are generally formed at or between circular apertures and cylinders maintained at suitable potentials. A number of such lenses are shown with characteristic path plots in **Fig. 1**. For any of these it is possible to define focal points, principal planes, and focal lengths in the same manner as for light lenses and to determine with their aid image magnification for any object position (**Fig. 2**). For a thin electrostatic lens in particular, that is, a lens for which the extent of the variation in potential is small compared to its focal length, the object side focal length f_o and the image side focal length f are given by Eq. (1). Here

$$\frac{\phi_o^{1/2}}{f_o} = \frac{\phi_i^{1/2}}{f_i} = \frac{3}{16}(\phi_o\phi_i)^{1/4} \int \left(\frac{\phi'}{\phi}\right)^2 dz \quad (1)$$

$\phi(z)$ is the potential along the axis of the lens, ϕ' its derivative with respect to z (that is, the electric field along the axis), and ϕ_o and ϕ_i are the potentials in object and image space, respectively. The integra-

tion is extended over the lens field. The quantity ϕ is here normalized so that it is equal to the accelerating potential of the particle in question.

Axially symmetric lenses are commonly divided into the four classes that follow.

Simple aperture lenses. These are the lens fields formed about circular apertures in a plane metallic electrode at potential ϕ with different electrostatic fields $-\phi'_o$ and $-\phi'_i$ on the two sides. In most cases the focal length f of such a lens is given to a sufficient degree of accuracy by the Davisson-Calbick formula for an aperture, Eq. (2). Simple aperture lenses are

$$\frac{1}{f} = \frac{\phi'_i - \phi'_o}{4\phi} \quad (2)$$

encountered as parts of more complex electrostatic lens systems, as well as at the mesh openings of metal screens employed as electrostatic shields in vacuum tubes.

Bipotential, or immersion, lenses. In these lenses image space and object space are field-free, but at different potentials. Typical examples are the lenses formed between apertures or cylinders at different potentials (**Fig. 1c** and **d**). If the separation d of the two apertures is large compared to their diameters and if each component aperture lens satisfies the conditions for a validity of the Davisson-Calbick formula, the focal lengths of the bipotential aperture lens are given by Eq. (3). The distances of the principal planes

$$\begin{aligned} \frac{1}{f_o} &= \left(\frac{\phi_i}{\phi_o}\right)^{1/2} \frac{1}{f_i} \\ &= \frac{3}{8d} \left[1 - \left(\frac{\phi_o}{\phi_i}\right)^{1/2}\right] \left(\frac{\phi_i}{\phi_o} - 1\right) \end{aligned} \quad (3)$$

from the plane of symmetry are given by Eqs. (4) and (5). Generally, the principal planes are displaced

$$b_o = \frac{-d}{2} - \frac{4d\phi_o}{3(\phi_i - \phi_o)} \quad (4)$$

$$b_i = \frac{d}{2} - \frac{4d\phi_i}{3(\phi_i - \phi_o)} \quad (5)$$

from the plane of symmetry toward the low-potential side, with the image-side principal plane closer to object space than the object-side principal plane.

For two cylinders of equal diameter D , whose difference in potential is small compared to their mean

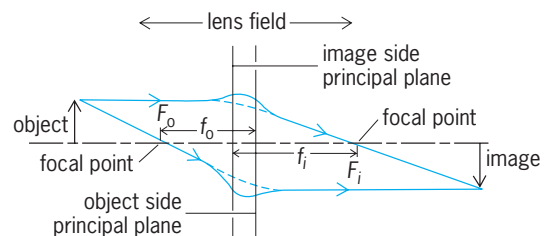


Fig. 2. Definition of principal planes, focal points, and focal lengths for axially symmetric electrostatic lenses.

potential, Eq. (6) gives the focal lengths.

$$\begin{aligned} \frac{1}{f_o} &= \left(\frac{\phi_i}{\phi_o} \right)^{1/2} \frac{1}{f_i} \\ &= \left(\frac{\phi_i}{\phi_o} \right)^{1/4} 0.66 \left(\frac{\phi_i - \phi_o}{\phi_i + \phi_o} \right)^2 \frac{1}{D} \end{aligned} \quad (6)$$

Bipotential lenses, in particular lenses formed between two cylinders at different potentials, find wide application in beam-focusing devices such as electron guns. Like unipotential lenses, they invariably act as converging lenses.

Unipotential lenses. For this type the potentials are equal in object and image space. In their simplest form these lenses consist of three apertures of which the outer two are at a common potential ϕ_o and the central aperture is at a different, generally lower, potential ϕ_1 . For such lenses with a central aperture of diameter D and the two outer apertures, of smaller diameter, separated a distance D from the plane of symmetry, the weak-lens focal length is given by Eq. (7).

$$\frac{1}{f} = \frac{0.2}{D} \left(\frac{\phi_o - \phi_1}{\phi_o} \right)^2 \quad (7)$$

As ϕ_1 approaches zero, the quantity $1/f$ increases more rapidly than indicated by this formula; it attains a value of $0.7/D$ for $\phi_1 = 0$. Unipotential lenses operated at high potentials (for example, $\phi_o = 50$ kilovolts, $\phi_1 = 0$) are employed as objectives and projection lenses in electrostatic electron microscopes. The electrodes are commonly made out of stainless steel and given a high polish. See ELECTRON MICROSCOPE.

Cathode lenses or immersion objectives. Here the lens field extends from the emitter surface up to field-free image space. Examples are the cathode region of an electron gun, the electron-optical system of an electrostatic image tube or image converter, and the objective of an emission electron microscope. In the electron gun the cathode lens converges the electrons emitted by the cathode to a small spot, the crossover, which is imaged by a second electron lens as the scanning spot on the cathode-ray tube screen. See CATHODE-RAY TUBE.

In the image tube the cathode itself—a transparent photoemissive surface on which a light picture is projected—is imaged on a fluorescent screen beyond the cathode lens. Frequently a cathode lens consists essentially of a uniform accelerating field followed by a short lens. The image magnification m is then given by Eq. (8). Here u is the distance between cathode

$$m = \frac{v}{2u} \quad (8)$$

and short lens and v is the distance between short lens and image. The quantity v is given by Eq. (9),

$$1v = \frac{1}{f} - \frac{1}{2u} \quad (9)$$

where f is the focal length of the short lens.

Lenses of plane symmetry. These lenses, analogous to cylindrical glass lenses, are formed between paral-

lel planes and at slits, replacing the circular cylinders and apertures of lenses with axial symmetry. For the simple slit in an electrode at potential ϕ separating two regions of field $-\phi'_o$ and $-\phi'_i$, the focal length is given by the Davisson-Calbick formula for a slit as shown by Eq. (10).

$$\frac{1}{f} = \frac{\phi'_i - \phi'_o}{2\phi} \quad (10)$$

Edward G. Ramberg

Electrostatic precipitator

A device used to remove liquid droplets or solid particles from a gas in which they are suspended. The process depends on two steps. In the first step the suspension passes through an electric discharge (corona discharge) area where ionization of the gas occurs. The ions produced collide with the suspended particles and confer on them an electric charge. The charged particles drift toward an electrode of opposite sign and are deposited on the electrode where their electric charge is neutralized. The phenomenon would be more correctly designated as electrodeposition from the gas phase. The practical aspects of the electrostatic precipitator were demonstrated in 1906 by F. G. Cottrell.

Construction. In its simplest form the experimental arrangement may consist of a vertical tube containing an insulated concentric wire (Fig. 1). A direct current potential of 10–100 kV may be applied to the center wire. A corona discharge occurs in a small area surrounding the wire. The suspended particles are ionized in the corona discharge and migrate to the wall of the vertical tube. If the suspended particles

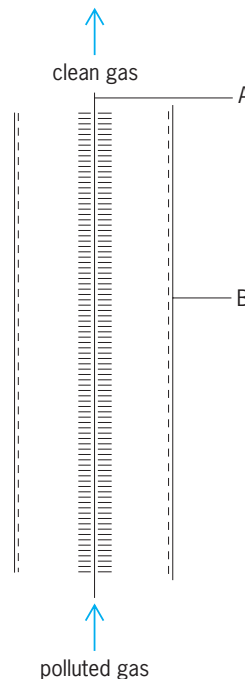


Fig. 1. Diagram of simple precipitator. A, corona wire; B, grounded tube.

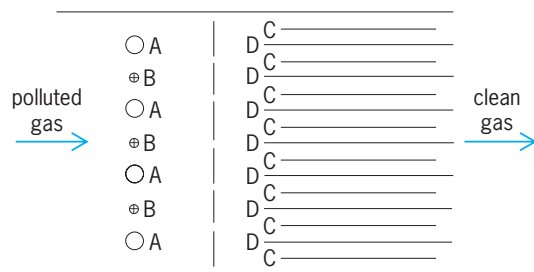


Fig. 2. Diagram of two-stage precipitator. A, grounded cylinders; B, corona wires; C, grounded collector plates; D, charged plates of polarity similar to B.

are liquid, they will accumulate on the wall and coalesce into droplets which can be drained away from the bottom of the tube. Solid particles may be displaced by mechanical vibration or scrapers and discharged into a conical collector at the bottom of the tube.

In more elaborate arrangements the ionization may occur in one vessel, whereas the deposition occurs in a second stage. A simplified two-stage apparatus is illustrated in Fig. 2. In the first chamber the particles become charged but are prevented from depositing on the grounded cylinders by suitable adjustment of the rate of flow of the gas. In the second chamber (consisting of alternately charged, loosely packed, parallel plates) precipitation can be achieved satisfactorily by applying a lower potential than is necessary in the charging chamber, since a corona discharge is not a requirement.

The corona discharge is usually produced by making the center wire negative, because precipitation efficiency is higher under these conditions. However, less ozone is produced by reversing the polarity, and a positive wire is commonly employed in the cleaning of air when the presence of ozone may be objectionable. The high-voltage direct current is commonly produced by mercury vapor or vacuum tube rectifiers. The power requirements vary from 2 to 5 kWh per 1,000,000 ft³ (28,300 m³) of gas being treated, depending upon the amount and nature of the particles being removed. Considerable difficulty is encountered in electrical leakage across the insulators because drops of liquid or solid particles are deposited on them.

H. J. White in 1951 stated that suspended particles become charged by two different mechanisms. The bombardment mechanism results from direct collisions between the suspended particles and ions produced in the corona discharge. The diffusion mechanism consists of the attachment of ions to the suspended particles by ion diffusion. Although both mechanisms operate simultaneously, the former is more important. However, the latter effect may become predominant for smaller particles, perhaps in the 0.1–0.2-micrometer range.

Application. The use of electrostatic precipitators has become common in numerous industrial applications. Each installation, however, is a separate design problem. There is no comprehensive theory applicable in every case for the separation of particles or droplets from a moving gas stream which is usually in

turbulent flow. Consequently, empirical methods are generally used for designing precipitators. In fact, it is sometimes customary to make an accurate model of a proposed precipitator installation and to adjust and correct the model performance in the laboratory to meet specifications prior to the manufacture of a scaled-up unit.

In view of the lack of a firm theoretical basis the efficiency of a given unit must be stated carefully in reference to the particular unit, with respect to operating voltage, design geometry, flow rate, and particle removal. The efficiency is usually expressed as the weight percentage of the material removed from the input stream. Such an expression inadequately represents the effect of particle size. It is usually desirable to remove the smallest particles from effluent gases, since their higher degree of opacity makes their discharge more visible and therefore more objectionable to the public eye; units are designed accordingly. The large particles are therefore those which escape; however, the reverse condition may also be achieved by suitable design. The efficiency is thus greatly influenced by the particle or droplet size distribution in the stream, and an efficiency percentage based on weight must be referred to the size distribution, most correctly, before and after precipitation. The efficiency depends exponentially on the stream residence time in the precipitator, so that slower flow rates or longer vessels give better precipitation efficiencies.

Among the advantages of the electrostatic precipitator are its ability to handle large volumes of gas, at elevated temperatures if necessary, with a reasonably small pressure drop, and the removal of particles in the micrometer range. Some of the usual applications are (1) removal of dirt from flue gases in steam plants; (2) cleaning of air to remove fungi and bacteria in establishments producing antibiotics and other drugs, and in operating rooms; (3) cleaning of air in ventilation and air conditioning systems; (4) removal of oil mists in machine shops and acid mists in chemical process plants; (5) cleaning of blast furnace gases; (6) recovery of valuable materials such as oxides of copper, lead, and tin; and (7) separation of rutile from zirconium sand. *See* DUST AND MIST COLLECTION. George S. Mill; W. O. Milligan

Electrostatics

The class of phenomena recognized by the presence of electrical charges, either stationary or moving, and the interaction of these charges, this interaction being solely by reason of the charges and their positions and not by reason of their motion. *See* ELECTRIC CHARGE.

At least 90% of the topics that are normally classified as electrostatics are concerned with the manipulation of charged particles by electric fields. When a particle becomes charged by rubbing or other means, it has either a surplus or a deficit of electrons. A body with a surplus of electrons is said to be negatively charged; a body with a deficiency,

positively charged. The amount or quantity of charge on a body is expressed in coulombs (positive or negative). A coulomb is an enormous amount of charge, and in most electrostatic situations charge levels of a small fraction of a coulomb give rise to significant effects. Electrostatic forces always exist between charged bodies. Bodies with like charge experience repulsive forces, while oppositely charged bodies experience attraction.

Coulomb's law. If two bodies are charged to Q_1 and Q_2 coulombs and are separated in vacuum by a distance of r meters, the force F in newtons between them is given by Eq. (1).

$$F = \frac{Q_1 Q_2}{4\pi\epsilon_0 r^2} \quad (1)$$

During the eighteenth century, Henry Cavendish carried out an ingenious experiment with two concentric metallized spheres which were electrically connected. He showed that after charging the outer sphere no charge was detectable on the inner sphere. He concluded that his observations proved that the force between charged bodies depended inversely on the square of the distance separating them. See ELECTRICAL SHIELDING.

Experimental work on the direct measurement of force between charged bodies was carried out by C. A. Coulomb several years later, and the law which expresses mathematically the force between charged bodies is known as Coulomb's law. In electrical science, ϵ_0 is an important constant known as the permittivity or dielectric constant of free space, and is also sometimes called the primary electric constant. It has the value $\epsilon_0 = 8.85416 \times 10^{-12}$ farad per meter. See COULOMB'S LAW; PERMITTIVITY.

Electric fields. Coulomb's law shows that a body charged to Q_1 experiences a force due to the presence of another body charged to Q_2 . Q_2 may be considered to influence the whole of space surrounding it, because if Q_1 were to be positioned anywhere it would experience a force due to the presence of Q_2 . The property of a charge to influence the whole of space can be modeled by a three-dimensional force field permeating the whole of the space surrounding the charge Q_2 . This field is called the electric field, and is often symbolized E . It is a vector quantity having both magnitude and direction. From Coulomb's law, the force acting on Q_1 can be written as in Eq. (2).

$$F = Q_1 E \quad (2)$$

From Eq. (2) it follows that the electric field E may be defined as the coulombic force in newtons per coulomb of charge. The more familiar units of electric field are volts per meter. It again follows that the field E due to charge Q_2 is given by Eq. (3). When

$$E = \frac{Q_2}{4\pi\epsilon_0 r^2} \quad (3)$$

there are many charged bodies present in an environment, the force that would be exerted on a charged particle at any location can be found by calculating the field at the location due to the presence of each

charged body separately, and the net field is obtained by adding up the individual components. An entire region may be considered in this way so that the electrostatic force acting on a particle at any position is readily determined. Hence the influence of charged bodies upon the trajectories of charged particles can be calculated. In a cathode-ray tube, for example, various fixed bodies in the form of deflection and focusing electrodes are appropriately charged to set up an electrostatic field which controls the trajectories of electrons. See CATHODE-RAY TUBE; ELECTRIC FIELD.

Potential energy. A system of charged particles or bodies is unstable unless the particles are prevented from moving, since the like-charged particles will repel each other until they are infinitely far apart, and oppositely charged bodies will attract one another and come together. The system has potential energy. The potential energy of two charged particles separated by a distance r can be shown to be given by Eq. (4). See ENERGY; POTENTIALS.

$$\text{Potential energy} = \frac{Q_1 Q_2}{4\pi\epsilon_0 r} \quad \text{joules} \quad (4)$$

Charging methods. The three principal methods of applying electric charge to objects are corona charging, induction charging, and tribocharging. See ENERGY.

Corona charging. The corona-charging method, which relies upon the impact of charged atoms or molecules (ions) on charged bodies, is particularly important in a number of technological applications. Copious quantities of ions may be generated by a corona discharge, which is a region in which an intense electric field acts upon air molecules and ionizes them so that free ions are produced. A sharply pointed electrode maintained at a high positive or negative potential induces a stream of positive or negative ions which may be used for charging surfaces. The stream of ions from a corona point is usually so intense that neutral air molecules become entrained in the flow to produce a corona wind which can be felt on the back of a hand or can deflect a candle flame.

Ions from a corona discharge may be used to charge isolated bodies, insulating surfaces or particles by simply directing a corona wind onto the surface to be charged. In the case of particles, it is normally sufficient for them to pass through a corona discharge region to receive a significant charge from ion-particle collisions. See CORONA DISCHARGE.

Induction charging. Surfaces may be charged by exposure to an electrostatic field. If the surface is a liquid and it is disrupted into droplets, they will be electrically charged. An example of induction charging of a liquid is illustrated in Fig. 1. With no voltage applied to the capillary tube, the liquid drips from the end of the tube. The drops are not significantly charged. When a voltage of a few hundred volts is applied to the capillary tube, an electric field is set up in the drop formation region. As each drop forms, the electric field induces surface charge, and an

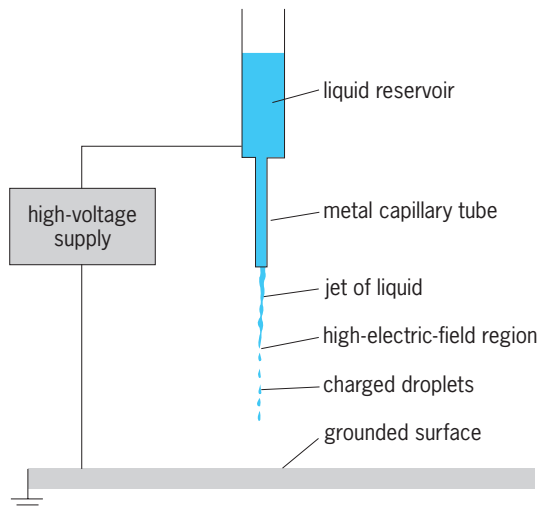


Fig. 1. Induction charging of a liquid.

electrostatic downward force acts on each drop in addition to the gravitational force. As a result, drops break off from the capillary tube before they are fully developed and are therefore smaller than normal. Thus the presence of the field increases the drip rate and, most importantly, causes the drops to be electrically charged. If the capillary-tube voltage is raised to a few thousand volts, the dripping changes to a jet or many jets of liquid which may subsequently break up into small charged droplets constituting a spray. See ATOMIZATION.

Induction charging of equipment and personnel may occur when they are exposed to an electric field. Personnel charged in this way may generate electrostatic discharges when approaching grounded surfaces. Sensitive microelectronic devices can be damaged and computer data can be corrupted by such discharges.

Applications. Electrostatics is put to good use in a wide variety of applications, many involving solid or liquid particles.

Air cleaners. The electrostatic precipitator enables smoke emissions from power-station chimneys, smelting plants, and other industrial plants to be reduced to relatively low, acceptable levels. On a smaller scale, efficient filters exist for removing dust from the air in offices, public places, and the home. In some filters, dust particles undergo corona charging as they are sucked by a fan through a duct, and are then collected on grounded electrodes; in others, permanently electrified filter material is used, made from thin plastic sheets which have been treated by surface bombardment from a corona ion source. Ions may attach so tenaciously to the surface of the sheets when conditions are carefully controlled that the treated material acquires a permanent electric charge. Such electrified material is known as an electret. An electret is the electrical equivalent of a permanent magnet, and will attract dust particles, fibers, and so forth, in the same way as a magnet attracts iron filings. See AIR FILTER; ELECTRET; ELECTROSTATIC PRECIPITATOR.

Spraying. In several applications which utilize electrostatics, solid or liquid particles are charged and sprayed onto grounded objects. Dry powder coating is used in many industries in preference to the wet-paint-spraying process. Crop spraying is another important application in which electrostatic forces help to efficiently apply herbicides or insecticides (Fig. 2). For the dry-powder-coating process, pigmented powder particles, which may incorporate resin, modifiers, and possibly a curing agent, are dispersed pneumatically and charged using a corona point at the end of a spray gun. Liquids are dispersed hydraulically, pneumatically, or by electrostatic forces alone, and droplets are charged by induction. As the sprayer nozzle is operated at several thousand volts, an intense electric field is set up between the nozzle and the grounded object that is to be sprayed. Charged particles tend to follow the field lines and are even guided to deposit on the far side of targets. This wrap-around effect minimizes

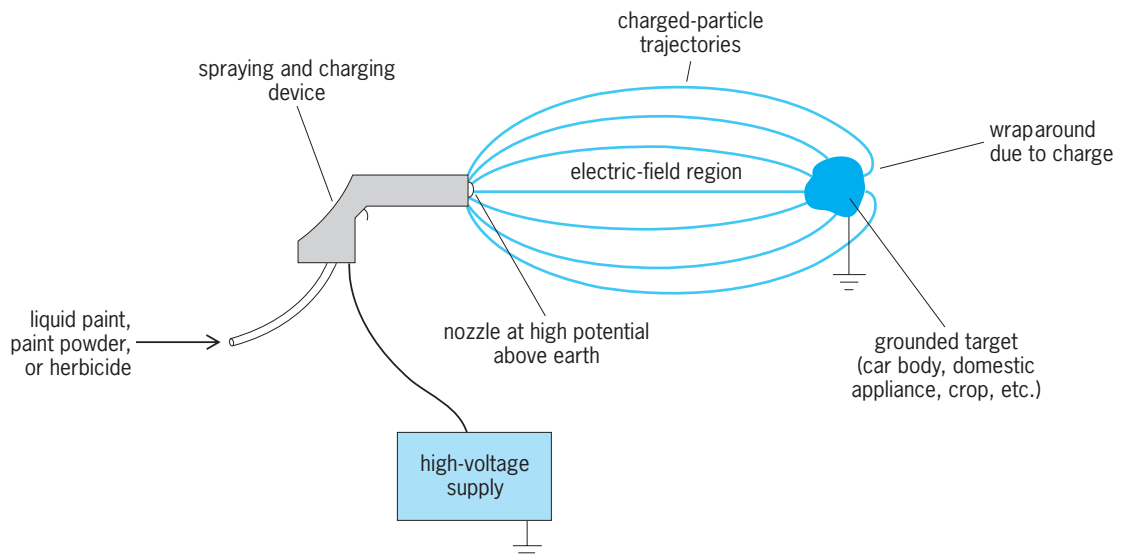


Fig. 2. Electrostatic spraying of powders or droplets.

target overspray, and is a particularly important feature of electrostatic spraying. In the powder-coating process, the powder adheres to the grounded workpiece by electrostatic forces known as image forces. A charged particle close to a conducting surface is attracted by an equal but opposite charge induced as a mirror image on the surface. Coated objects are baked or treated with infrared radiation, causing thermoplastics to melt and run to form a continuous film. In thermosetting powders, so-called cross-linking occurs, and once the coating is cured, it will not remelt. In the liquid-paint- and crop-spraying processes, centrifugal forces are sometimes used in conjunction with electrostatic fields, producing sprays consisting of droplets that are fairly uniform in size. This enables deposition to be more accurately controlled.

Research into electrostatic spraying, sometimes called electrohydrodynamic atomization, is leading to new applications for the deposition of ceramic, glass, and polymer films and for powder particle production of special materials. The electrostatic spraying of materials is also used for analysis by means of mass spectrometry, as the electrostatic spray process is gentle and does not disrupt delicate complex molecules.

By using induction charging, it is possible to disperse fuel oils into charged sprays or mists. Charged electrodes allow some measure of control over the combustion of such mists; that is, it is possible to produce electrically controlled flames.

Electrophotography. In electrophotography an optical system is used to project the image to be copied onto a light-sensitive semiconducting surface precharged by a corona source. Exposure of the surface to light reduces the electrical conductivity of the material and allows surface charge to leak away to a back plate in proportion to the intensity of the light, so that bright parts of the image are regions that have lost most of the original charge while dark zones remain fully charged. Gray shades produce intermediate levels of charge. A mixture of very fine black toner particles (approximately 2 micrometers in diameter) and coarser carrier particles (approximately 100 μm in diameter) is then brought into contact with the charged surface. Transfer of only charged toner particles onto the latent charged surface occurs. The toner charges as it separates from the carrier particles. The amount of toner that attaches by electrostatic forces to each zone of the surface depends directly on the surface charge density. Dark regions of the image are highly charged and thus attract a relatively thick layer of toner. A sheet of paper is then laid over the toner-covered surface, and transfer of toner to paper occurs so that an image remains on the paper when it is peeled off the surface. The paper is then treated to fix the image. Electrophotography is technologically so advanced that color copying with excellent quality is possible. See PHOTOCOPYING PROCESSES.

Ink-jet printers. Modern ink-jet printers are of various types, not all of which utilize electrostatic principles. One type uses hydraulic pressure via arrays of acoustically excited nozzles which produce fine jets of ink, usually by hydraulic pressure. The jets

travel a few millimeters and then break up into drops. By ensuring that drops are formed in the presence of an electrostatic field, they become charged and may be deflected electrostatically to a printing surface. Ink-jet printers are highly developed machines that print rapidly and silently. Colored inks may be used not only for alphanumeric printing but for high-resolution printing of pictures which approach photographic quality. See INK; PRINTING.

Metallic ion and droplet beams. Another development being commercially exploited is the production of metallic ion or droplet beams using electrostatic forces acting upon a liquid-metal surface. Considerable success has been achieved with many molten metals including gold and silver. Either ion or charged droplet beams may be formed depending on the operating conditions of the source. The beams so formed may be very well defined and directed with great accuracy onto targets where they can be used for ion implantation or for the formation of conducting tracks in the fabrication of microelectronic circuits. See INTEGRATED CIRCUITS; ION IMPLANTATION.

Electrostatic coalescers. Crude oil from wells, especially offshore, is invariably mixed with water or brine, with water content sometimes in excess of 80% of the volume. The extraction process subjects the mixture to considerable shearing, and stable oil-water emulsions form, which do not readily separate out. An oil continuous phase containing water drops will support an intense electric field which polarizes the individual drops. Each drop essentially becomes an electrostatic magnet with positively and negatively charged ends, aligned in the electric field. By virtue of this polarization, drops in proximity are attracted toward each other and coalesce to form larger drops. In this way a population of fine drops within an oil can be induced to coalesce rapidly to form much larger drops, which are then able to settle out under gravity from the oil. Electrostatic treaters using electric fields not only to coalesce drops but also to move and deposit inorganic particles of sand, mud, and clay and organic particles have been developed.

Space exploration. Ion engines which produce thrust by electrostatically accelerating mercury or cesium ions have been successfully operated in space. Colloid thrusters, operating on exactly the same principles as electrostatic paint or crop sprayers, have also been developed. In these a propellant such as glycerol is atomized and accelerated from a nozzle by an electrostatic field. See ION PROPULSION. A. G. Bailey

Bibliography. A. G. Bailey, *Electrostatic Spraying of Liquids*, 1988; A. G. Bailey, The science and technology of electrostatic powder spraying, *J. Electrostat.*, 45:85-120, 1998; I. Berta, *Industrial Electrostatics*, 1994; J. S. Chang, A. J. Kelley, and J. M. Crowley (eds.), *Handbook of Electrostatic Processes*, 1995; J. A. Cross, *Electrostatics: Principles, Problems and Applications*, 1987; J. M. Crowley, *Fundamentals of Applied Electrostatics*, 1986, reprint 1991; Special issue on electrohydrodynamic atomization, *J. Aerosol Sci.*, 30(7):823-978, 1999.

Electrostriction

A form of elastic deformation of a dielectric induced by an electric field; specifically, the term applies to those components of strain which are independent of reversal of the field direction. Electrostriction is a property of all dielectrics and is thus distinguished from the converse piezoelectric effect, a field-induced strain which changes sign upon field reversal and which occurs only in piezoelectric materials. *See* DIELECTRIC MATERIALS; PIEZOELECTRICITY.

Electrostrictive strain is approximately proportional to the electric susceptibility, elastic compliance, and the square of the field strength, and is extremely small for most materials.

The electrostrictive effect in certain ceramics is employed for commercial purposes in electromechanical transducers for sonic and ultrasonic applications. *See* MICROPHONE. Robert D. Waldron

Electrotherapy

The use of electric current to treat disease. Electrotherapy is based on principles developed during the nineteenth and twentieth centuries following the first demonstration of "animal electricity" by Luigi Galvani in the eighteenth century. This article covers low-frequency electrotherapy, diathermy and hyperthermia, and electroconvulsive therapy.

Low-Frequency Electrotherapy

Low-frequency electrotherapy uses relatively weak alternating electric currents. They are delivered by electrodes that are placed under or on the surface of the body and are connected to pulse-shaping circuits that are located inside or outside the body.

Excitable tissues. Electrodes that stimulate electrically excitable cells, such as those in muscle and nerve tissues, are usually placed directly in or on the tissue by surgery or are inserted through a vein by catheterization. There are many applications for electrode stimulation: irregular heart rhythms can be controlled by pacemakers; muscles, such as those of the diaphragm and urinary bladder, that become paralyzed can be made to contract electrically; and long-term pain can be relieved by implanting electrodes in the spinal canal. When coupled to appropriate braces, sensors, and programmed computers, electrodes in contact with the muscle groups of a lower extremity have been shown to help persons with spinal cord injury become ambulatory. Surface electrodes are widely used for temporary relief of pain, a technique known as transcutaneous electrical nerve stimulation; for preventing muscle atrophy after injury or immobilization; and for treatment of curvature of the spine, or scoliosis. *See* CARDIAC ELECTROPHYSIOLOGY; MUSCLE; PAIN.

Nonexcitable tissues. A wide variety—if not all—of the body's nonexcitable cells alter their function in specific ways and at specific times in response to appropriate, very weak electrical stimuli. (That observation reemphasizes the central roles of physics

in all living processes and the interaction of electric charges as the basis of the chemical aspects of those processes.) Much of the progress in electrotherapy has evolved from the observation that both hard and soft tissues, such as bone and arteries, become electrically charged when they are cyclically deformed by mechanical or hydrodynamic forces. Weak voltages, in the range of microvolts to millivolts per centimeter, occur because both cells and tissues contain piezoelectric molecules that respond to deforming forces by becoming electrically polarized. Furthermore, electrically charged entities such as ions, cells, and molecules are transported by hydrodynamic forces past sites of structurally fixed electric charge and, in the process, create a voltage. That phenomenon is an example of an electrokinetic event. By establishing the patterns and values of those stress-generated voltages in bone, researchers have been able to develop three methods of influencing the behavior of nonexcitable cells that are involved in the repair of nonunited fractures. The oldest method was an invasive procedure, and was followed by two noninvasive procedures. *See* BONE.

Electrode implantation. The first clinical method for treating nonunited fractures employed surgically implanted electrodes. Once placed at the fracture site, they delivered constant, direct currents similar in amplitude to those that occur naturally in bone after fracture, known as injury potentials, or as a response to mechanical deformation. Unfortunately, surgical methods carry a risk of infection, and direct currents can result in adverse electrochemical reactions around the electrodes.

Pulsed electromagnetic fields (PEMFs). Two noninvasive electrical methods have proved effective in altering cellular behavior. The first involves the placement of dynamically charged, insulated plates outside tissue-culture vessels or the fractured limbs of animals. Broad-scale application has not yet proven practical in humans, however. The second method uses one or more coils of wire coupled to a pulse generator to create a weak time-varying magnetic field that penetrates the body, much as radio waves enter a closed building. The field characteristics are designed to induce pulsing electric currents in the tissue, with waveforms, frequencies, and amplitudes similar to those produced normally by skeletal tissues during high-impact exercise. The waveforms of pulsed electromagnetic fields are quite asymmetrical and contain a broad range of frequencies, which are characteristics that distinguish them from power-line, radio-frequency, and microwave fields.

Depending upon the energy patterns generated by the magnetic fields in tissue, the function of cells involved in abnormal processes can be changed without producing heat. In nonunited fractures, the normal repair has been interrupted at an intermediate stage of healing, and soft, rather than hard, tissue forms a bridge between the bone fragments. The final stages of bony repair are not achieved until the soft tissues undergo calcification. Certain types of pulsed electromagnetic fields can initiate calcification, aid in the ingrowth of new blood vessels, and

increase bone formation, each of which is important in restarting the healing process and ultimately achieving successful union. Pulsed electromagnetic fields have been widely used to treat nonunited fractures, many of which had failed to heal after one or more operations and the person faced amputation of the affected part. In contrast to other electromagnetic fields, the pulse types do not carry any known risks, and hospitalization and a surgical procedure are usually unnecessary, since pulsed electromagnetic fields are applied in the doctor's office. The method, therefore, is substantially less costly than most other treatments available.

As the understanding of their mechanisms of action at the cellular and subcellular levels has increased, pulsed electromagnetic fields have been used successfully to treat other problems of bone and its surrounding soft tissue. For example, when avascular necrosis of the hip in young adults results in bone tissue death, revitalization has been achieved and hip function restored following pulsed electromagnetic field therapy. In older persons, shoulder pain from chronic inflammation of tendons that has proved resistant to classical forms of treatment has responded to pulsed electromagnetic fields. Therapeutic requirements for selective cellular effects appear to parallel those involved in the treatment of disease states with drugs. Future success will lie in encoding the appropriate physical (as opposed to chemical or drug) "information" in the pulsed electromagnetic field waveform to produce changes in cell function and thereby control or correct specific pathologic processes. Other medical applications for pulsed electromagnetic field, such as in cardiology, remain to be explored.

C. Andrew L. Bassett

Diathermy and Hyperthermia

The therapeutic benefits of heat have been known for centuries, and modern medicine has used technology to provide controlled heat to diseased tissues.

Diathermy. Therapy for afflicted deep tissues that do not respond to conventional methods, such as infrared heating or heating pads, can often be achieved with diathermy. Heating results from the electrical resistance of tissue to high-frequency or microwave currents. Increasing the tissue temperature to a range of 106–113°F (41–45°C) increases the physiologic response and therapeutic benefit, which includes increased extensibility of collagen tissues in joint contractures, decreased joint stiffness in rheumatoid arthritis, and pain relief and reduction of muscle spasms through the local action of heat on nerve endings and receptors. Warming can also resolve inflammatory infiltrates, edema, and exudates and increase blood flow in diseased or damaged tissue. Heating has been used in cancer therapy under proper temperature control.

Various instruments have been developed to accommodate placement of the heating element on the body. They include (1) coupling of 13.56- or 27.12-MHz shortwave currents from electric fields by way of insulated or noncontacting capacitor plates or specially shaped electrodes for placement in natural

body cavities; (2) induction of 13.56- or 27.12-MHz electric currents by magnetic fields from solenoidal coils that enclose the afflicted body member or from flat coils that are placed near the surface of the afflicted tissues; and (3) transmission of 433-, 915-, or 2450-MHz microwave energy into the tissues by radiation from dipole or cavity antennae. The shortwave magnetic and the 433- or 915-MHz microwave applicators are superior to the shortwave capacitive and 2450-MHz microwave applicators for treating tissues beneath subcutaneous fat layers, because they minimize the undesirable selective heating of the fat and the variation in heating levels between patients. Diathermy is contraindicated for areas that are anesthetized or noninnervated, for areas with inadequate vascular supply, in the presence of acute inflammation accompanied by the formation of pus and fever, whenever there is a possibility of hemorrhage, and with malignancies, when there is inadequate monitoring of tissue temperature.

Hyperthermia. Hyperthermia is an experimental method of treating malignant tumors that uses heat alone, heat in combination with ionizing (x-ray) radiation, or heat with chemotherapy. One form of heating involves the application of radio-frequency energy by using methods similar to—but more sophisticated and more carefully controlled than—those in diathermy treatment. The effective temperature range of hyperthermia is very narrow, 108–111°F (42.5–44°C); the benefits are minimal at lower temperatures, and damage to normal cells is probable at higher temperatures. Several mechanisms are thought to account for the selective destruction of tumor cells: (1) selective heating caused by the lower rate of blood flow in tumor tissues; (2) greater sensitivity of tumor cells to heat due to their hypoxic, acidic, and poor nutritional state; (3) synergism of ionizing radiation and hyperthermia due to thermal killing of cells that are hypoxic and are at those critical stages of growth when they are most resistant to radiation; and (4) increase in cell membrane permeability and sensitivity to chemotherapeutic drugs.

Exposures at frequencies other than those authorized for diathermy are performed in shielded rooms to prevent interference with radio communication facilities. The selective heating of tumor tissues uses the same basic techniques as those for diathermy: shortwave capacitive, shortwave inductive, and microwave applicators. However, because of the potentially life-threatening nature of malignant disease, more drastic measures can be justified, including the application of low-frequency currents by surface-contact and invasive electrodes and the implantation of magnetic seeds that can be heated by low-frequency magnetic fields. These approaches enhance the localization of tumor heating and minimize the destruction of normal tissues.

The frequencies used exceed 100 kHz to prevent excitation of nerves and muscles. Heating may be further localized through use of additional frequencies and more sophisticated applicators, such as phased-antenna arrays and cylindrical antennae designed for positioning in natural body cavities. Extensive

multipoint temperature-measuring equipment with automatic feedback power-control systems is necessary to maintain a sufficiently high temperature in the tumor without destructive temperature changes in normal tissues. See THERMOTHERAPY. Arthur W. Guy

Electroconvulsive Therapy

Electroconvulsive therapy is a procedure for treating severe psychiatric disorders. It is usually given as a series of treatments, typically numbering between 6 and 12, over a period of a few weeks. Each 10-min treatment is preceded by the administration of an anesthetic agent, to render the patient unconscious throughout the procedure, and a drug that blocks muscle movement in the body, to confine the provoked seizure to the brain and prevent convulsive movements in the body. Electrodes are then placed on the scalp and a small jolt of electricity is applied. The electric current provokes a generalized brain seizure, much like that experienced spontaneously by patients with grand mal epilepsy. The patient's brain waves are usually monitored with electroencephalography during the seizures, which typically last about 1 min.

In use since about 1940, electroconvulsive therapy has been shown to be effective in treating a specific set of psychiatric disorders, the most common of which is depressive illness. Depressed patients often undergo electroconvulsive therapy after failing to respond to other forms of therapy, particularly antidepressant medications. For others, electroconvulsive therapy is preferred because the medical risks associated with drug treatment are too great: in some elderly patients, the frequency of medical injury and death from antidepressants is believed to exceed that with electroconvulsive therapy. The fact that electroconvulsive therapy produces rapid improvement and has a high probability of success is another reason for its use. In patients at risk for death from suicide or medical conditions associated with depressive illness, electroconvulsive therapy may be preferred over other therapeutic methods. It has also been shown to be an effective treatment for patients with acute mania and schizophrenia. Less common uses are the treatment of a small set of medical and neurological disorders that are resistant to conventional therapies, such as Parkinson's disease, intractable epilepsy, and psychosis associated with toxic states. See AFFECTIVE DISORDERS; PARKINSON'S DISEASE; SCHIZOPHRENIA.

Electroconvulsive therapy has a characteristic set of side effects. Following each treatment, patients experience a period of confusion that may last from several minutes to several hours, and immediately following the complete course of treatments, they frequently have memory difficulties. Memory problems include difficulty in remembering newly learned information and in recalling events from the recent past. At the same time, other aspects of intellectual functioning, such as the ability to perform a task and solve a problem, are unchanged or show improvement over that displayed before electroconvulsive therapy. The improvements can be attributed to the

beneficial effects of electroconvulsive therapy on the psychiatric conditions that caused the diminished intellectual performance. Objective evidence of persistent difficulties in learning or recalling new information is difficult to obtain several weeks after treatment, but some patients may experience permanent gaps in their memory for events that occurred in the weeks prior to, during, and following the treatment course. The beneficial and adverse effects of electroconvulsive therapy can likely be traced to changes in brain physiology and biochemistry. Studies in humans and animals indicate it is unlikely that electroconvulsive therapy causes structural damage to the brain, but the precise ways in which it produces beneficial and adverse effects are unknown. See ELECTROCONVULSIVE THERAPY. Harold A. Sackheim

Bibliography. R. Abrams, *Electroconvulsive Therapy*, 3d ed., 1997; A. Chiabrera, C. Nicolini, and H. P. Schwan (eds.), *Interaction Between Electromagnetic Fields and Cells*, 1985; T. Fried, R. Johnson, and W. McCracken, *Archives of Physical Medicine and Rehabilitation*, 65:228-231, 1984; P. L. Ludmer and N. Goldschlager, *N. Engl. J. Med.*, 311:1671, 1984; C. Polk and E. Postow, *Handbook of Biological Effects of Electromagnetic Fields*, 2d ed., 1996.

Electrothermal propulsion

Vehicular propulsion that involves electrical heating to raise the energy level of the propellant. In contrast, chemical rockets use the chemical energy of one or more propellants to heat and accelerate the decomposition products (monopropellants) or combustion products (bipropellants) for thrusting purposes. In both instances, the high-energy propellant gases are exhausted through a nozzle where they are accelerated to a high velocity (U_e), and thrust is produced by reaction. By decoupling the heating or energy addition process from the restraints of propellant chemistry considerations, electrothermal devices can be operated on a wide variety of materials, many of which would not otherwise be considered to be propellants. Water and space station liquid-waste streams are two examples of such propellants being considered for electrothermal propulsion purposes. See ROCKET PROPULSION.

Practical electrothermal thrusters come in two forms, resistojets and arcjets. In resistojets the electrical energy is first deposited in a heater or resistive element and then transferred to the propellant. The need to first heat a material limits the maximum operating temperature and the maximum enthalpy of the propellant. As a consequence, the essential simplicity of the device is balanced by well-defined limitations on exhaust velocity or specific impulse I_s . (Specific impulse is a measure of propellant utilization efficiency. It is defined as the thrust per unit weight of propellant expended to produce the thrust. Numerically, $I_s = U_e/g$, where g is the gravitational acceleration at the Earth's surface, 9.8 m/s² or 32.2 ft/s².) Arcjets circumvent this limitation by using the propellant as the heater element. High temperatures and

specific-impulse values can be achieved but only at the price of design complexity. See SPECIFIC IMPULSE.

One further and fundamental distinction is that electrical thrusters are power-limited whereas chemical thrusters are energy-limited. This distinction profoundly influences the form and applications of these two classes of propulsion systems. By definition, electrical propulsion systems must have an associated power supply for operation. Solar panels or nuclear power supplies can supply well-defined power to the thruster for essentially unlimited time, which is limited by the availability of propellant for propulsion. Consequently, although the power level is well constrained, the total energy available is virtually boundless. Chemical systems are exactly opposite. In this case, the total energy available for propulsion is well defined by the propellant volume. The propellant tanks (liquid) or casing (solid) can hold only a certain amount of propellant. However, the rate at which this propellant is used, the rate of energy usage per unit time, or the power can be exceedingly high. All of the space shuttle's propellant is exhausted in about 10 min, whereas electrical thrusters operate for months or years. In other words, electrical propulsion devices are characterized by low thrust (power), high exhaust velocity, and parsimonious use of propellant (their primary feature of interest), and chemical rockets are high-thrust systems that are profligate in the use of propellant. Chemical thrusters are ideal for escaping the Earth's gravity well; electrical thrusters are ideal for moving payloads in low-acceleration conditions removed from gravity wells, that is, for the more ambitious missions far removed from low Earth orbit.

Resistojet thrusters. These have been under development for many years and are now flying in a stationkeeping role on many communications satellites. An electric heating coil is used to add energy to the propellant (Fig. 1). The maximum temperature, and therefore the thruster performance limit, is set by the maximum temperature of the coil material and correspondingly structural requirements. The heat of the coil is coupled either radiatively or directly to the propellant gas. In the radiatively coupled resistojet (Fig. 1a), the coil is placed in an evacuated cavity and

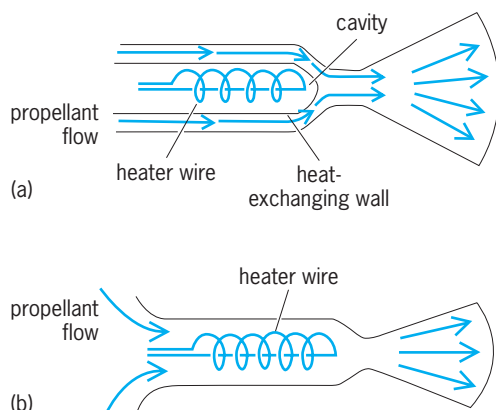


Fig. 1. Resistojet configurations. (a) Radiatively coupled. (b) Directly coupled.

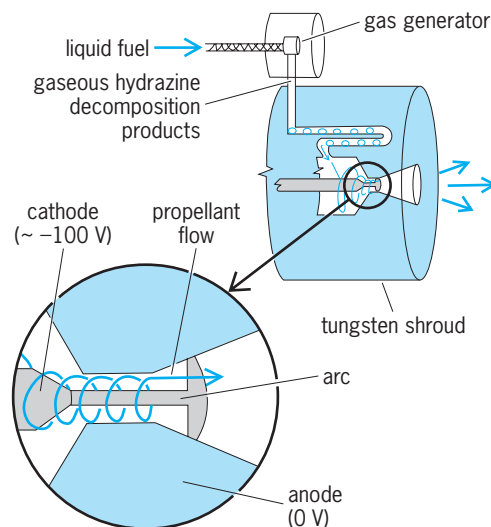


Fig. 2. Hydrazine arcjet thruster, showing enlargement of heat transfer region.

transfers energy by radiation to the inner cavity wall. The outer cavity wall acts as a heat exchanger which, in turn, is in direct contact with the propellant. Evaporation of the heater element is a life-limiting issue in this type of thruster. In a directly coupled resistojet (Fig. 1b) the heater element is fully immersed in the propellant and transfers energy principally by conduction and convection. Direct heating is more efficient but exposes the coil to the propellant, which can lead to erosion or corrosion. As in the radiatively coupled thruster, such processes limit the coil lifetime.

In principle, resistojets can be designed to operate at any power level. Existing devices typically require 300 to 800 W of electrical power and produce specific impulses in the range of 250 to 330 s (exhaust velocities from 2.5 to 3.3 km/s or 1.5 to 2.1 mi/s) with thrust levels that are typically on the order of 0.5 newton (0.1 lbf). Many different types of propellants, including hydrazine (N_2H_4), ammonia (NH_3), molecular hydrogen (H_2), and the organic waste gases such as carbon dioxide (CO_2) and methane (CH_4), have been successfully used. The omnivorous nature of the resistojet has led to its consideration as reboost propulsion for the NASA crewed space station using waste gases as propellants; and resistojets using water (H_2O) as a propellant have been selected for the reboost propulsion of the industrial space facility.

Arcjet thrusters. These are similar in concept to resistojets, but rather than transferring energy to the gas through a resistive coil, the arcjet utilizes an electric arc discharge which passes directly through the propellant. While this avoids the materials limitations imposed by the heater coils and the heat exchanger walls of the resistojet, a small portion of the gas must be ionized to carry the current. The interaction of this plasma with the electrodes is an important design consideration that has to be properly addressed to assure adequate longevity. See ARC DISCHARGE.

A schematic of an arcjet thruster is shown in Fig. 2. Low-speed propellant enters the arcjet plenum in a

tangential fashion to produce a vortical motion to assist plasma formation on the centerline. The gas accelerates around the cathode and through the constrictor and finally expands supersonically to ambient space. A direct-current arc is struck between the cathode and the anode (that is, the nozzle body). The high-temperature, high-velocity plasma which forms the laminar arc column is stabilized as it passes through the constrictor region by both energy (radiation and conduction) and mass transfer (diffusion and convection) to and from the propellant gas which surrounds it. The thruster walls are thus shielded from the hot arc plasma by the propellant. Electrons within the arc core are readily accelerated by the applied electric field and transfer energy via collision to their much more massive neutral and ionic neighbors. This process tends to cool the arc plasma and heat the surrounding propellant flow. Although the temperatures in the low-density central region may be in excess of 20,000 K, the bulk or mass-averaged temperature of the propellant is typically between 5000 and 8000 K, a level that substantially exceeds that of resistojets. As the arc plasma enters the diverging section of the nozzle, it is rapidly cooled and eventually expands to the anode wall. With proper design, high wall-current densities can be avoided so that the process of arc attachment remains is nondestructive.

Electrons produced at the cathode by a combination of thermionic and field emission provide most of the current in the arcjet. An equal number of ions, however, is generated in the anode fall region because of collisions of electrons with neutral gas molecules as they accelerate toward the wall. Thus, except in the cathode and anode fall regions (that is, within a few mean free paths of these surfaces), the arcjet plasma is spatially neutral. Langmuir probe measurements in arcjet plumes indicate that electron temperatures are on the order of 1 eV and number densities are on the order of 10^{13} cm⁻³. For the typical arcjet flow rates, this represents an electron/ion mole fraction of less than 1%. See FIELD EMISSION; PLASMA DIAGNOSTICS; THERMIONIC EMISSION.

Arcjet thrusters can operate over a wide range of power by scaling the physical dimensions of the device. They may be broadly classified by their power requirements. If a nominal cathode voltage of -100 V with respect to the anode is assumed, the definitions given in the **table** apply. Typical thrust levels range from 0.25 N (0.05 lbf) for low-power arcjets to 2.5 N (0.5 lbf) for high-power arcjets. Arcjets are therefore considered as candidates for missions ranging from stationkeeping to orbit transfer. Like the resistojet, an arcjet can operate with many different types of pro-

pellant. Because the arcjet is an electrothermal device, operation is relatively insensitive to propellant selection; operation has been demonstrated by using argon (Ar), helium (He), molecular hydrogen (H₂), ammonia (NH₃), and hydrazine (N₂H₄). Depending on the choice of propellant, exhaust velocities can range from 4 to 9 km/s (2.5 to 5.5 mi/s) with corresponding specific impulses between 400 and 900 s. Development of prototype engines has been undertaken, with flights projected in the early 1990s. See SPACECRAFT PROPULSION.

G. W. Butler

Bibliography. M. N. Hirsh and H. J. Oskam, *Gaseous Electronics*, 1978; R. G. Jahn, *Physics of Electric Propulsion*, 1968; *Proceedings of the AIAA/ASME/SAE/ASEE 24th Joint Propulsion Conference*, Boston, 1988; *Proceedings of the 1986 JANNAF Propulsion Meeting*, New Orleans, 1986; *Proceedings of the 19th International Electric Propulsion Conference*, Colorado Springs, 1987; J. A. Stone, D. C. Byers, and P. Q. King, *The NASA Electric Propulsion Program*, NASA Tech. Mem. 101324, 1988.

Electroweak interaction

One of the three basic forces of nature, along with the strong nuclear interaction and the gravitational interaction. The terms "force" and "interaction between particles" are used interchangeably in this context. All of the known forces, such as atomic, nuclear, chemical, or mechanical forces, are manifestations of one of the three basic interactions.

Until the early 1970s, it was believed that there were four fundamental forces: strong nuclear, electromagnetic, weak nuclear, and gravity. It was by the work of S. Glashow, S. Weinberg, and A. Salam that the electromagnetic and the weak nuclear forces were unified and understood as a single interaction, called the electroweak interaction. This unification was a major step in understanding nature, similar to the achievement of J. C. Maxwell and others a century earlier in unifying the electric forces and magnetic forces into the electromagnetic interactions. A goal of theoretical physics is to achieve a further simplification in understanding nature and describe the presently known three basic interactions in a unified way, usually referred to as the grand unified theory (GUT). Whether this is possible remains to be seen. See ELECTROMAGNETISM; FUNDAMENTAL INTERACTIONS; GRAND UNIFICATION THEORIES; GRAVITATION; MAXWELL'S EQUATIONS; STRONG NUCLEAR INTERACTIONS; WEAK NUCLEAR INTERACTIONS.

Properties of basic interactions. Some of the properties of the basic interactions are summarized in **Table 1**. The strong nuclear forces are the strongest, electroweak is intermediate, and gravity the most feeble by a huge factor. The ranges, that is, the distances over which the forces act, also differ greatly. The strong nuclear and the weak interactions have a very short range, while electromagnetism and gravity act over very large distances. Thus, at very short subatomic distances the strong nuclear force, which holds the atomic nucleus together and

Classification of arcjet thrusters			
Range	Power, kW	Current, A	Mass flow rate, g/s
Low	1-2	10-20	0.04
Medium	5-10	50-100	0.12
High	20-30	200-300	0.24

TABLE 1. Basic forces in nature

Interaction	Relative strength	Property acted on	Force carrier	Range
Strong nuclear	1	Color charge (r, g, b)	Gluon (g)	10^{-13} cm
Electroweak	Electromagnetic	Electric charge (q)	Photon (γ)	∞
	Weak nuclear	Weak charges (t_3, y)	Bosons (W^\pm, Z^0)	10^{-16} cm
Gravity	10^{-40}	Mass (m)	Graviton (G)	∞

governs many interactions of the subnuclear particles, dominates. At larger distances the electromagnetic forces dominate, and hold the atom together and govern chemical and most mechanical forces in everyday life. At even larger scales, objects such as planets, stars, and galaxies are electrically neutral (have an exact balance of positive and negative electric charges) so that the electromagnetic forces between them are negligible, and thus the gravitational forces dominate in astronomical and cosmological situations.

Each of the basic forces acts on, or depends on, different properties of matter. Gravity acts on mass, and electromagnetic forces act on electric charges that come in two kinds, positive and negative. The strong nuclear forces act on a much less well-known property, called color charge, which come in three kinds, r , b , and g (often called red, blue, and green). The weak nuclear forces act on equally esoteric properties called weak isospin t and hypercharge y . While the mass and the electric charge are properties that are recognized in everyday situations, the color charge and the weak isospin and hypercharge have no correspondence in the large-scale everyday world. See COLOR (QUANTUM MECHANICS); ELECTRIC CHARGE; HYPERCHARGE; ISPIN; MASS.

Fundamental constituents of matter. All known forms of matter are made of molecules and atoms, which are made up of the nucleus (protons and neutrons) and orbital electrons. These in turn can be understood to be made up of the fundamental constituents, the quarks and the leptons. Each of these comes in six kinds (Table 2). The six leptons are

usually organized into three families:

$$\begin{pmatrix} \nu_e \\ e^- \end{pmatrix} \begin{pmatrix} \nu_\mu \\ \mu^- \end{pmatrix} \begin{pmatrix} \nu_\tau \\ \tau^- \end{pmatrix}$$

Similarly, the six quarks can be organized into three families:

$$\begin{pmatrix} d \\ u \end{pmatrix} \begin{pmatrix} s \\ c \end{pmatrix} \begin{pmatrix} b \\ t \end{pmatrix}$$

The masses of the leptons and quarks are frequently specified in units of electronvolts (1 MeV = 10^6 eV, and 1 GeV = 10^9 eV). In more usual units, the electron mass of 0.5 MeV is equal to 9×10^{-31} kg. For many years the masses of the three neutrinos (ν_e, ν_μ , and ν_τ) were believed to be zero. However, an effect called neutrino oscillations indicates that they have small but nonzero masses. The values of these masses are not yet known; only experimentally measured upper limits on them are available. See NEUTRINO.

All of the quarks and leptons have gravitational and weak interactions since they have nonzero values of mass and weak isospin and hypercharge. The particles with zero electric charge have no electromagnetic interactions, and the leptons have no strong nuclear interactions since they carry no color charge. See LEPTON; QUARKS.

Exchange forces and gauge bosons. The present understanding is that the basic forces are not contact forces but act over distances larger than the sizes of the particles (action at a distance). In this picture,

TABLE 2. Leptons and quarks, the fundamental constituents of matter

Particle	Symbol	Mass	Electric charge (q)	Color charge	Weak charges	
					Isospin (t_3)	Hypercharge (y)
e neutrino	ν_e	≤ 12 eV	0	0	+1/2	-1
Electron	e^-	0.51 MeV	-1	0	-1/2	-1
Mu neutrino	ν_μ	≤ 0.2 MeV	0	0	+1/2	-1
Muon	μ^-	106 MeV	-1	0	-1/2	-1
Tau neutrino	ν_τ	≤ 18 MeV	0	0	+1/2	-1
Tau	τ^-	1777 MeV	-1	0	-1/2	-1
Down quark	d	3-7 MeV	-1/3	r, b, g	-1/2	1/3
Up quark	u	1.5-3 MeV	+2/3	r, b, g	+1/2	1/3
Strange quark	s	100 MeV	-1/3	r, b, g	-1/2	1/3
Charm quark	c	1.2 GeV	+2/3	r, b, g	+1/2	1/3
Bottom quark	b	4.3 GeV	-1/3	r, b, g	-1/2	1/3
Top quark	t	174 GeV	+2/3	r, b, g	+1/2	1/3

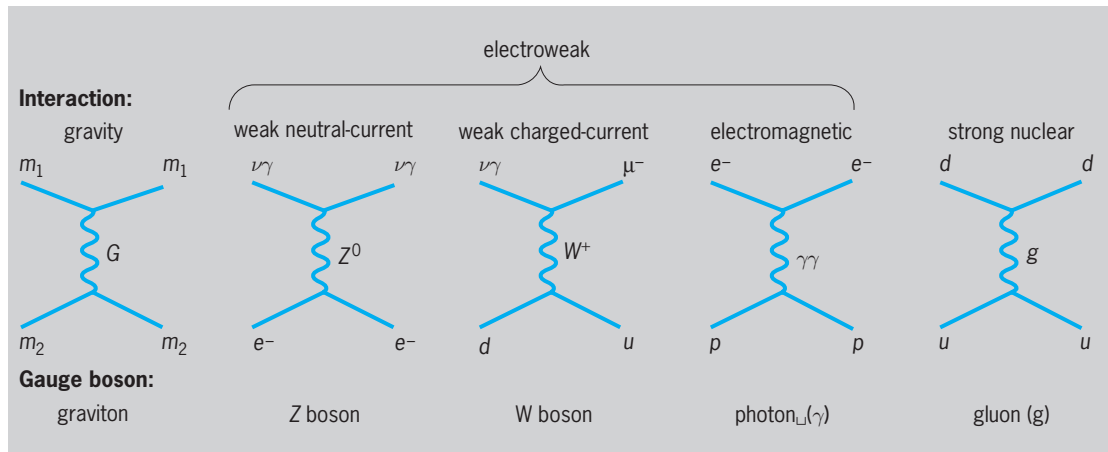


Fig. 1. Basic interactions mediated by gauge bosons.

based on field theory, the forces are carried or mediated by intermediate particles that are called gauge bosons. For example, the electromagnetic force between an electron and a proton is carried by the quantum of the electromagnetic field called the photon (γ). The strong nuclear force is carried by the gluon (g), and the gravitational force is carried by the graviton (G) [Fig. 1]. The weak nuclear force comes in two categories: the charge-changing (charged-current, for short) mediated by the W^\pm bosons, and the neutral-current weak interactions mediated by the Z^0 boson (Fig. 1). See GLUONS; GRAVITON; INTERMEDIATE VECTOR BOSON; PHOTON.

Helicity and parity violation. All of the fundamental constituents, the quarks and the leptons, carry one-half unit of angular momentum (spin = 1/2) as if they were spinning around their own axis. (Such particles are called fermions.) By the rules of quantum mechanics, the direction of this spin is quantized to be either parallel or antiparallel to the direction of motion of the particle. Particles with spin direction parallel to their direction of motion have helicity +1 and are called right-handed, and particles with antiparallel spin have helicity -1 and are called left-handed. See ANGULAR MOMENTUM; HELICITY (QUANTUM MECHANICS); SPIN (QUANTUM MECHANICS).

One of the symmetries of nature is called parity, which is a symmetry between right-handed and left-handed coordinate systems. If parity symmetry holds, left-handed and right-handed particles must have the same interactions. In 1956 T. D. Lee and C. N. Yang proposed that parity symmetry is violated in the weak interactions, and this proposal was soon verified experimentally. It was found that the left-handed and the right-handed particles have different weak interactions. The weak isospin and hypercharge assignments given in Table 2 are for left-handed quarks and leptons. The right-handed particles have somewhat different assignments. In particular, the right-handed particles have no weak isospin, and thus only the left-handed particles participate in the charged-current weak interactions. See PARITY (QUANTUM MECHANICS).

Electroweak unification. Until the early 1970s, the electromagnetic and the weak interactions were believed to be separate basic interactions. At that time the Weinberg-Salam-Glashow model was proposed to understand these two interactions in a unified way. The model was based on an $SU(2) \times U(1)$ gauge symmetry in which the $SU(2)$ part corresponds to a weak isospin triplet of gauge bosons, the W^+ , W^0 , and W^- , and the $U(1)$ corresponds to a singlet, the B^0 . The W 's couple to the property called weak isospin, and the B^0 couples to weak hypercharge. See GAUGE THEORY.

In its original form, this model, based on an unbroken gauge symmetry, led to some physically unacceptable features such as zero masses for all the constituent particles and predictions of infinities for some measurable quantities. Through the pioneering work of G. 'tHooft, M. Veltman, and others, it was shown that the theory can be made renormalizable, removing the infinities and providing masses to the particles, by spontaneous breaking of the gauge symmetry and the introduction of one new particle, the Higgs boson. See HIGGS BOSON; RENORMALIZATION; SYMMETRY BREAKING.

The neutral gauge bosons, the W^0 and B^0 , form a quantum-mechanical mixture, which produces the two physically observable gauge bosons, the γ and the Z^0 , as given by Eqs. (1). The γ is the well-known

$$\begin{aligned}\gamma &= \sin\theta W^0 + \cos\theta B^0 \\ Z^0 &= \cos\theta W^0 - \sin\theta B^0\end{aligned}\quad (1)$$

photon that mediates the electromagnetic interactions. The Z^0 mediates the neutral-current weak interactions, and the W^\pm mediate the charged-current weak interactions (Fig. 1). In this way, all of these interactions are described by a common unified theory. The mixing angle θ in Eqs. (1), forming the γ and the Z^0 , is called the weak mixing angle and is the fundamental parameter of the theory. The strength and nature of the interactions of the particles are determined by the vector and axial vector coupling constants g_v and g_A . In the electroweak model all of

these couplings are given in terms of the single parameter of the theory, the weak mixing angle, and the properties of the leptons and quarks given in Table 2. The model also gives a relationship between the electric charge q and the weak charges t_3 and y , Eq. (2).

$$q = t_3 + 1/2 y \quad (2)$$

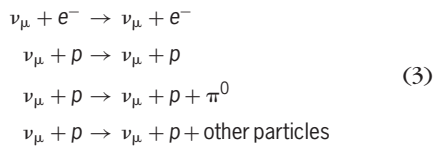
See NONRELATIVISTIC QUANTUM THEORY; QUANTUM MECHANICS.

The coupling constants that govern the electro-weak interactions of all of the particles can be summarized as:

1. Electromagnetic interactions: $g_v = q, g_A = 0$
2. Charged current weak interactions:
 $g_v = -g_A = t$
3. Neutral current weak interactions:
 $g_v = t_3 - 2q \sin^2 \theta, g_A = -t_3$

In the above expressions, t stands for the magnitude of the weak isospin, and t_3 is its projection along an axis of quantization.

The electroweak theory has great predictive power. Its first and most striking prediction was the existence of neutral-current weak interactions mediated by the Z^0 boson. Until the time of this prediction, the weak interactions were believed to be of charged-current nature only, with no neutral-current component. Some examples of neutral-current processes predicted by the theory were reactions (3).



See NEUTRAL CURRENTS.

The experimental discovery of such neutral-current weak interaction processes was considered very strong support for the validity of the theory. From the experimental measurements of the cross sections (that is, interaction probabilities) of such processes, an early estimate of the weak mixing angle was derived to be $\sin^2 \theta \approx 0.23$. With this value of $\sin^2 \theta$, the theory was able to predict the masses of the hypothetical W^\pm and Z^0 gauge bosons to be approximately $M_W \approx 81$ GeV and $M_Z \approx 92$ GeV.

A second major triumph for the electroweak theory was the discovery of the W and Z bosons in 1983 at the proton-antiproton collider at the CERN Laboratory in Geneva, Switzerland, with masses very close to the values predicted by the theory. At this time the validity of the theory was considered to be firmly established. See PARTICLE ACCELERATOR.

Precision tests. A great deal of experimentation followed the discovery of the W and Z bosons to carry out detailed precision tests of the electroweak theory. These were done in a wide variety of contexts such as parity-violating effects in atomic physics, neutrino interactions, polarized electron scattering,

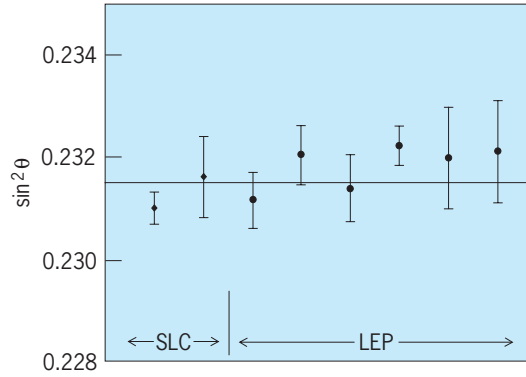
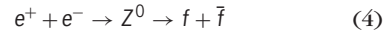


Fig. 2. Results of measurements of the weak mixing angle, θ , at the SLC (Stanford Linear Collider) and the LEP (Large Electron-Positron) Collider. The horizontal line indicates the average of all the measurements: $\sin^2 \theta = 0.23155 \pm 0.0018$.

and precision measurements of the W and Z boson masses. All of the experimental results were in excellent agreement with the predictions of the theory. During the 1990s the most precise tests of the theory were carried out at two positron-electron (e^+e^-) colliders, the LEP (Large Electron-Positron) Collider at CERN and the SLC (Stanford Linear Collider) at Stanford Linear Accelerator Center (SLAC) in Stanford, CA, using reaction (4), where f can be any



one of the leptons or quarks, and \bar{f} stands for their antiparticles. In these experiments, the value of the $\sin^2 \theta$ parameter was measured in a large variety of ways. The agreement among the results is quite good (Fig. 2), providing the most sensitive test of the theory, and the most precise value of the weak mixing angle: $\sin^2 \theta = 0.23155 \pm 0.00018$.

Prospects. The successful electroweak theory, combined with quantum chromodynamics (QCD), the theory describing the strong nuclear interactions, forms the so-called standard model of particle physics. This standard model has been brilliantly successful in accurately predicting and describing all experimental results over a huge energy range, varying from the electronvolt energies of atomic physics to the 100-GeV energy range of the largest existing particle colliders. As such, it represents a landmark achievement of both experimental and theoretical physics. See QUANTUM CHROMODYNAMICS; STANDARD MODEL.

However, in spite of these great successes, the story is not complete, and two major problems remain to be solved in this field. The first one has to do with the realization that the standard model cannot be complete in its present form since it cannot explain the masses of the fundamental constituents, the quarks and leptons. These masses vary over a large range (Table 2), from a few electronvolts to 174 GeV. The basic gauge symmetry on which the standard model is based would indicate that these masses should all be the same. There must therefore be an additional piece of the puzzle, which is usually referred to as the source of the electroweak

symmetry breaking, that remains to be found. Hypothetical ideas about this missing element of the model range from the prediction of a single additional particle, the Higgs boson, to complicated models such as supersymmetry that predict dozens of new elementary particles. The search for this new physics that will lead to a more complete theory motivates research in this field, including the construction of high-energy particle accelerators and colliders. See HIGGS BOSON; SUPERSYMMETRY; SYMMETRY BREAKING.

The second outstanding problem in this field is the search for a theory that not only describes the strong nuclear and the electroweak interactions but includes gravity as well. The standard model is based on the principles of quantum mechanics, while the current understanding of the gravitational forces is based on Einstein's theory of general relativity. No one so far has been able to combine these two theories; that is, a quantum theory of gravity does not, as yet, exist. The search for such a grand unified theory is a major focus of activity in theoretical physics. See ELEMENTARY PARTICLE; QUANTUM GRAVITATION; RELATIVITY.

Charles Baltay

Bibliography. S. L. Glashow, Towards a unified theory: Threads in a tapestry, *Rev. Mod. Phys.*, 52:539–543, 1980; P. Renton, *Electroweak Interactions*, Cambridge University Press, 1990; A. Salam, Gauge unification of fundamental forces, *Rev. Mod. Phys.*, 52:525–538, 1980; S. Weinberg, Conceptual foundations of the unified theory of weak and electromagnetic interactions, *Rev. Mod. Phys.*, 52:515–523, 1980.

Element (chemistry)

An element is a substance made up of atoms with the same atomic number. Some common elements are oxygen, hydrogen, iron, copper, gold, silver, nitrogen, chlorine, and uranium. Approximately 75% of the elements are metals and the others are nonmetals. Most of the elements are solids at room temperature, two of them (mercury and bromine) are liquids, and the rest are gases.

Occurrence and classification. A few of the elements are found in nature in the free (uncombined) state. Some of these are oxygen, nitrogen, the noble gases (helium, neon, argon, krypton, xenon, and radon), sulfur, copper, silver, and gold. Most of the elements in nature are combined with other elements in the form of compounds. The most abundant element on the Earth is oxygen; the next most abundant is silicon. The most abundant element in the universe is hydrogen and the next most abundant is helium. See ELEMENTS, GEOCHEMICAL DISTRIBUTION OF.

The elements are classified in families or groups in the periodic table. Elements are also frequently classified as metals and nonmetals. A metallic element is one whose atoms form positive ions in solution, and a nonmetallic element is one whose atoms form negative ions in solution. See PERIODIC TABLE.

Atoms of a given element have the same atomic

number, but may not all have the same atomic weight. Atoms with identical atomic numbers but different atomic weights are called isotopes. Oxygen, for example, is made up of atoms whose atomic weights are 16, 17, and 18. Hydrogen is made up of isotopes 1, 2, and 3; the isotopes of masses 2 and 3 are called deuterium and tritium, respectively. Carbon is made up of isotopes 11, 12, 13, and 14. Carbon-14 is radioactive and is used as a tracer in many chemical experiments.

All the elements have isotopes, although in certain cases only synthetic isotopes are known. Thus, fluorine exists in nature as ^{19}F , but the artificial radioactive isotope ^{18}F can be prepared. Many of the isotopes of the different elements are unstable, or radioactive, and hence they disintegrate to form stable atoms either of that element or of some other element. See ATOMIC MASS; RADIOACTIVITY.

Origin and uses. The origin of the chemical elements is believed to be the result of the synthesis by fusion processes at very high temperatures (in the order of 100,000,000°C or 180,000,000°F and higher) of the simple nuclear particles (protons and neutrons) first to heavier atomic nuclei such as those of helium and then on to the heavier and more complex nuclei of the light elements (lithium, boron, and so on). The helium atoms bombard the atoms of the light elements and produce neutrons. The neutrons are captured by the nuclei of elements and produce heavier elements. These two processes—fusion of protons and neutron capture—are the main ones forming the chemical elements. Furthermore, the energy of the Sun and the stars is derived primarily from the fusion of hydrogen nuclei and electrons to form helium nuclei. It is believed that this element-producing fusion process is occurring even today in hot stars. See NUCLEOSYNTHESIS.

The elements form the raw materials for the great chemicals industry today. Various metals are used for structural materials, protective coatings, ornamental devices, jewelry, and tableware. Such nonmetals as chlorine, bromine, hydrogen, sulfur, and nitrogen are important for the manufacture of many of the common chemicals of commerce. Neon is used to make neon light signs, and radon is used as a source of radioactive rays for therapy.

A number of elements, found in only very slight traces or not at all in nature, have been synthesized. Those are technetium, promethium, astatine, francium, and all the elements with atomic numbers above 92. These elements have been synthesized by a variety of nuclear reactions that involve transmuting atoms of one element into atoms of another by bombarding that element with neutrons or fast-moving particles (protons, deuterons, and alpha particles) which will change the atomic number to that of the new element. Not only have these elements been synthesized, but isotopes of all the other elements also have been synthesized. See ATOMIC STRUCTURE AND SPECTRA; ISOTOPE; TRANSURANIUM ELEMENTS.

Alfred B. Garrett

Bibliography. F. A. Cotton et al., *Advanced Inorganic Chemistry*, 6th ed., Wiley-Interscience, 1999;

N. N. Greenwood and A. Earnshaw, *Chemistry of the Elements*, 2d ed., 1997; S. Hofman, *On Beyond Uranium Journey to the End of the Periodic Table*, 2002.

Element 112

Element 112 was discovered in 1996. It should be a heavy homolog of the elements mercury, cadmium, and zinc, and is expected to be the last element in the *6d* shell. See CADMIUM; HALF-LIFE; MERCURY (ELEMENT); ZINC.

1																	18
H																	He
2																	2
3	4															10	
Li	Be															Ne	
11	12															18	
Na	Mg	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36
K	Ca	Sc	Ti	V	Cr	Mn	Fe	Co	Ni	Cu	Zn	Ga	Ge	As	Se	Br	Kr
37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52	53	54
Rb	Sr	Y	Zr	Nb	Mo	Tc	Ru	Rh	Pd	Ag	Cd	In	Sn	Sb	Te	I	Xe
55	56	71	72	73	74	75	76	77	78	79	80	81	82	83	84	85	86
Cs	Ba	Lu	Hf	Ta	W	Re	Os	Ir	Pt	Au	Hg	Tl	Pb	Bi	Po	At	Rn
87	88	103	104	105	106	107	108	109	110	111	112	113					
Fr	Ra	Lr	Rf	Db	Sg	Bh	Hs	Mt	Ds	Rg							

lanthanide series	57	58	59	60	61	62	63	64	65	66	67	68	69	70
	La	Ce	Pr	Nd	Pm	Sm	Eu	Gd	Tb	Dy	Ho	Er	Tm	Yb

actinide series	89	90	91	92	93	94	95	96	97	98	99	100	101	102
	Ac	Th	Pa	U	Np	Pu	Am	Cm	Bk	Cf	Es	Fm	Md	No

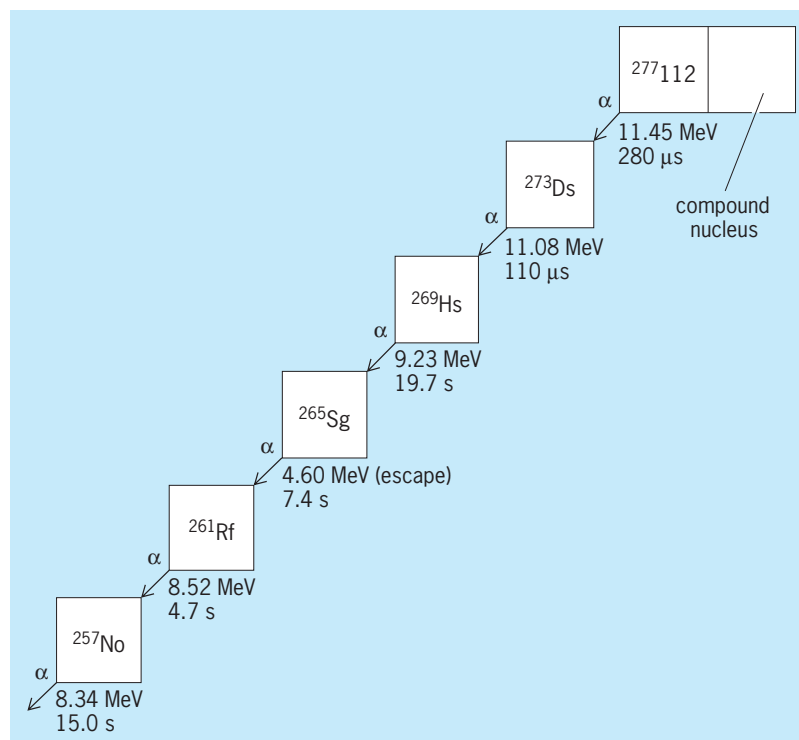
A long-standing claim to having synthesized element 112 is made by a Jerusalem-European Organization for Nuclear Research (CERN) collaboration. A spontaneous fission emitter with a half-life of several weeks was found in 1971 in a tungsten-beam stop at CERN. It followed the chemistry of mercury, and was assigned to element 112. A two-step reaction mechanism was proposed. From the tungsten-beam stop of the 24-GeV particle accelerator strontium-88 is produced by high-energy fission of fissionable spallation products. In a second reaction the ^{88}Sr ions induce a fusion reaction with another tungsten nucleus. The nucleus $^{272}112$ (the isotope of element 112) is synthesized at very low excitation energy by radiative capture and is detected after chemical separation by its fission. See NUCLEAR REACTION; SPALLATION REACTION; STRONTIUM; TUNGSTEN.

Element 112 was discovered on February 9, 1996, at GSI (Gesellschaft für Schwerionenforschung), Darmstadt, Germany, by detection of $^{277}112$, which was produced by fusion of a zinc-70 projectile and a lead-208 target nucleus following the cooling down of the fused system by emission of a single neutron. The fused system was observed at an excitation energy of 12 MeV. Sequential alpha decays to darmstadtium-273, hassium-269, seaborgium-265, rutherfordium-261, and nobelium-257 allowed unambiguous identification by using the known decay properties of the last three members of the chain. In the decay chain (see illustration), the first three members are new isotopes. Isotope $^{277}112$ has a half-life of 0.7 ms, and it is produced with a cross section of $0.5 \times 10^{-37} \text{ cm}^2$. The new isotopes of darmstadtium and hassium are

of special interest. Their half-lives differ by more than four orders of magnitude, and alpha energies are very different, which is characteristic of a closed-shell crossing. At the neutron number $N = 162$, a closed shell was theoretically predicted, and this is verified in the decay chain observed. Isotope ^{269}Hs has a half-life of 9 s, which is long enough to allow studies on the chemistry of this element. The crossing of the neutron shell at $N = 162$ is an important achievement in the field of research on superheavy elements. The stabilization of superheavy elements is based on high fission barriers, which are due to corrections in the binding energies found near closed shells. The shell at $N = 162$ is the first such shell predicted, and is now verified. See ALPHA PARTICLES; DARMSTADTIUM; HASSIUM; LEAD; NEUTRON; NOBELIUM; NUCLEAR STRUCTURE; RADIOISOTOPE; RUTHERFORDIUM; SEABORGIUM.

The methods used to produce element 112 were the same as those used for darmstadtium and roentgenium. The decay chain of the new element was observed in an irradiation time of about 3 weeks. The trend toward smaller cross sections has continued. See ROENTGENIUM.

Element 112 is unnamed, but there is no doubt that it was discovered. A further decay chain of $^{277}112$ was detected in May 2000 by the GSI group, ending by spontaneous fission decay. In 2004, two more decay chains ending by spontaneous fission of ^{261}Rf were found at RIKEN, Japan. The decay chain of $^{277}112$ feeds the isotope ^{269}Hs , which was produced in the reaction $^{248}\text{Cm}(^{26}\text{Mg},5n)$ by the GSI nuclear



Sequence of decay chains that document the discovery of element 112. Numbers below boxes are alpha energies and correlation times. Element 112 is produced in the reaction $^{70}\text{Zn} + ^{208}\text{Pb} \rightarrow ^{277}112 + 1n$.

chemistry group in 2003. The decay chains observed confirm the pattern of isotopes from $^{277}112$.

Theoretical model calculations of the element 112 chemistry indicate that it may exhibit rather unusual properties. On the basis of its projected position in group 12 of the periodic table it should behave like a heavy metal similar to its homolog mercury (atomic number 80). However, it is predicted that an extra stabilization of its binding electrons may result in a relatively high volatility, possibly giving rise to some similarity with the noble gas radon (atomic number 86). See RADON.

Peter Armbruster; M. Schädel
Bibliography. S. Hofmann et al., The new element 112, *Z. Phys. A*, 354:229–230, 1996.

Elementary particle

The idea that everything is made from a few basic elements originated in ancient Greece. In the nineteenth century the elementary pieces of matter were the atoms of the chemical elements, but in the first half of the twentieth century atoms were found to be compounds of electrons, protons, and neutrons. These became known as elementary particles, that is, particles that are not compounds of other particles. See ELECTRON; NEUTRON; PROTON.

Historical overview. The electrical attraction of opposite electrical charges is what grips negatively charged electrons around the positively charged atomic nucleus; it is the protons that give a nucleus its charge. The amount of positive charge on a proton is equal in magnitude but opposite in sign to that of an electron. This is crucial for the fact that matter in bulk is electrically neutral, whereby gravity controls the motion of planets and our attraction to the Earth's surface. However, why these two particles carry such precisely counterbalanced electrical charges is a mystery. See COULOMB'S LAW; ELECTRIC CHARGE; GRAVITATION.

A proton is some 1836 times as massive as an electron, their masses being respectively 939 and $0.511 \text{ MeV}/c^2$, where c is the speed of light. A neutron has a mass of $940 \text{ MeV}/c^2$ and is very similar to a proton. In the immediate aftermath of its discovery in 1932, the neutron was thought to be a version of a proton with no electrical charge, though today their relationship is understood to be more profound.

An electron or a proton is stable, at least on time scales longer than the age of the universe. When neutrons are in the nuclei of atoms such as iron, they too may survive unchanged for billions of years. However, an isolated neutron is unstable, with a mean life of 886 s. Neutrons can also be unstable when large numbers of them are packed with protons in a nucleus, which leads to natural radioactivity of many elements. Such neutrons undergo beta decay, a result of the weak force, which converts a neutron into a proton and emits an electron and a neutrino. See BETA PARTICLES; RADIOACTIVITY.

Neutrinos have no electric charge and masses that are too small to measure with present techniques.

The neutrino emitted in the neutron beta decay has a mass that is less than $3 \text{ eV}/c^2$, that is, less than $1/100,000$ the mass of the electron. See NEUTRINO.

In the 1930s the known elementary particles were the above four together with the photon (γ), the quantum particle of electromagnetic radiation. The electron and neutrino are known as leptons and are still recognized as elementary in the sense that they are not composed of more fundamental constituents whereas the proton and neutron are built from quarks. Two varieties, or flavors, of quark are required to make protons and neutrons. They are the up quark, (denoted u), with electrical charge that is a fraction $+2/3$ of a proton's charge, and the down quark (d), whose fractional charge is $-1/3$. The combination uud is sufficient to make a proton and udd makes a neutron. See LEPTON; PHOTON; QUARKS.

The up and down quarks, together with the electron and the above neutrino, form a basic family of what is known as the standard model of fundamental particles. These particles are fermions, defined by having an intrinsic angular momentum or spin of $1/2$ in units of Planck's quantum. The photon is a boson, with integer spin of 1. The exchange of photons between electrically charged particles transmits the electromagnetic force. See ANGULAR MOMENTUM; ELECTROMAGNETISM; QUANTUM STATISTICS; SPIN (QUANTUM MECHANICS); STANDARD MODEL.

This simple picture began to break down around 1950 with the discoveries of new forms of particle, first in cosmic rays and then in experiments at high-energy particle physics accelerators. With modern hindsight it is possible to classify the discoveries into two classes: leptons and hadrons. See COSMIC RAYS; PARTICLE ACCELERATOR.

The muon is a lepton with the same spin and electric charge as an electron but some 207 times more massive at $106.5 \text{ MeV}/c^2$. Today six members of the lepton family are known: three that are electrically charged, the electron (e), muon (μ) and tau (τ) [mass $1777 \text{ MeV}/c^2$]; and three varieties of neutrino known as the electron-neutrino (ν_e), muon neutrino (ν_μ), and tau neutrino (ν_τ) [Table 1]. The masses of the neutrinos have not yet been directly measured, but they are believed to have small but nonzero masses. All of these are fundamental particles. Leptons are unaffected by the strong (nuclear) force but feel the weak force and, if electrically charged (electron, muon and tau), the electromagnetic force. See STRONG NUCLEAR INTERACTIONS; WEAK NUCLEAR INTERACTIONS.

Particles that feel the strong force are known as hadrons. The cosmic rays revealed the existence of unstable hadrons, some of which became known as strange particles, such as the K meson and lambda hyperon (Λ). Experiments at particle accelerators enabled such particles to be produced when the kinetic energy of particles colliding with atomic nuclei was converted into ephemeral particles, which were revealed by their trails in bubble chambers (or nowadays by electronic detectors). Among these were short-lived heavier versions of the proton and

TABLE 1. The fundamental particles^a

Gauge bosons		$J_C^P = 1^-$	Self-conjugate except $\overline{W^+} = W^-$				
Name	Symbol	Charge ^b	Mass and width, GeV		Couplings		
Photon	γ	0	0		$A \Rightarrow \gamma A$		
Gluon ^c	g	0	0		$A \Rightarrow gA'$		
Weak bosons							
Charged ^d	W^\pm	± 1	80.4, 2.1		$d \Rightarrow W^+u$		
Neutral ^e	Z^0	0	91.2, 2.5		$A \Rightarrow Z^0A$		
Fermions ^f		$J = 1/2$	All have distinct antiparticles, except perhaps the neutrinos				
Name	Charge ^b	Symbol and mass, GeV	Symbol and mass, GeV		Symbol and mass, GeV		
Leptons							
Neutrinos	0	ν_e	$< 2 \times 10^{-9}$	ν_μ	$< .0002$	ν_τ	$< .02$
Charged leptons ^g	-1	e	.00051	μ	.106 ^h	τ	1.777 ^h
Quarks ^c							
Up type	$2/3$	u	.0015-.004	c	1.2	t	175 ⁱ
Down type	$-1/3$	d	.004-.008	s	.10	b	4.5

^aThe graviton, with $J_C^P = 2^+$, has been omitted, since it plays no role in high-energy particle physics.
^bIn units of the proton charge.
^cThe gluon is a color SU₃ octet (8); each quark is a color triplet (3). These colored particles are confined constituents of hadrons; they do not appear as free particles, which is indicated by placing their entries in boxes.
^dThe branching ratios (%) of the decay modes of the W^+ are:
 $u\bar{d}, c\bar{s}$ 34 each
 $\nu_e\bar{e}^+, \nu_\mu\bar{\mu}^+, \nu_\tau\bar{\tau}^+$ 11 each
^eThe branching ratios (%) of the decay modes of the Z^0 are:
 $d\bar{d}, s\bar{s}, b\bar{b}$ 15.6 each
 $u\bar{u}, c\bar{c}$ 11.6 each
 $\nu_e\bar{\nu}_e, \nu_\mu\bar{\nu}_\mu, \nu_\tau\bar{\nu}_\tau$ 6.7 each
 $e^+e^-, \mu^+\mu^-, \tau^+\tau^-$ 3.4 each
^fThe three known families (generations) of fermions are displayed in three columns.
^gAny further charged leptons have mass greater than 100 GeV.
^hThe μ and τ leptons are unstable, with the following mean life and principal decay modes (branching ratios in %):
 μ $\tau_\mu = 2.2 \times 10^{-6}$ s $e\bar{\nu}_e\nu_\mu$ 100
 τ $\tau_\tau = 2.9 \times 10^{-13}$ s $\mu\bar{\nu}_\mu\nu_\tau$ 17, $e\bar{\nu}_e\nu_\tau$ 18, (hadrons) $\bar{\nu}_\tau$ 65
ⁱThe t quark has a width = 2 GeV, with dominant decay to Wb .

neutron, known as resonances, whose lifetimes of the order of 10^{-23} s are similar to the time it takes a beam of light to traverse a proton. All of these hadrons are composed of quarks. See HADRON; PARTICLE DETECTOR.

There are six flavors of quark. The up (u), charm (c), and top (t) have electrical charge $+2/3$; the down (d), strange (s), and bottom (b) have charge $-1/3$ (Table 1). As individual quarks cannot appear in isolation, a direct measure of their mass is not possible, but approximate scales of mass have been determined. The up and down quarks, when trapped inside hadrons, have energies of around 300 MeV; most of this is due to their motion, their masses being at most $8 \text{ MeV}/c^2$. The strange quark is some $100 \text{ MeV}/c^2$ more massive, the charm quark mass is around $1.2 \text{ GeV}/c^2$ ($1200 \text{ MeV}/c^2$), the bottom around $4.5 \text{ GeV}/c^2$, and the top around $175 \text{ GeV}/c^2$. The reason behind these values, or even their qualitative pattern, is not understood.

Strange hadrons contain one or more strange quarks. Hadrons containing charm or bottom quarks are known as charm and bottom hadrons (including a special category known as charmonium and bottonium). Top quarks are so heavy that no hadrons containing them have been identified. It is even possi-

ble that top hadrons cannot form, as the top quark is so unstable that it decays (in a form of beta decay, $t \rightarrow b\nu_c$) before being able to grip to other quarks to make observable hadrons.

Fundamental particles and interactions. The fundamental forces that act on the particles are gravity, the electromagnetic force, and the strong and weak forces. These forces act by the exchange of particles, known as gauge bosons (Table 1). For the electromagnetic, strong, and weak forces, these have been identified, respectively, as the photon, gluon, and weak bosons; the graviton, the quantum of the gravitational force, is firmly predicted by theory, but the prospect of direct observation is exceedingly remote. The gravitational force is so feeble between individual particles that it can be ignored for all practical purposes.

The interactions of particles are responsible for their scattering and transformations (decays and reactions). Because of interactions, an isolated particle may decay into other particles. Two particles passing near each other may transform, perhaps into the same particles but with changed momenta (elastic scattering) or into other particles (inelastic scattering). The rates or cross sections of these transformations, and so also the interactions responsible

for them, fall into three groups: strong (typical decay rates of 10^{21} – 10^{23} s $^{-1}$), electromagnetic (10^{16} – 10^{19} s $^{-1}$), and weak ($<10^{15}$ s $^{-1}$). See FUNDAMENTAL INTERACTIONS; GLUONS; GRAVITON.

The photon, gluon, and weak bosons have spin 1. The photon and gluon have no mass, whereas the weak bosons are massive. The weak bosons may be electrically charged (the W^+ and W^-), with masses of 80.4 GeV/ c^2 , or neutral (the Z^0), with mass of 91.2 GeV/ c^2 . The fundamental couplings are known as gauge couplings (Table 1), in which a gauge boson is absorbed or emitted by a fundamental particle (fermion or gauge boson). See GAUGE THEORY; INTERMEDIATE VECTOR BOSON.

A massless photon with little or no energy can be emitted and absorbed by electrons in atoms or molecules within the normal laws of energy and momentum conservation. However, when a neutron undergoes beta decay, it converts to a proton by emitting a “virtual” W^- , violating energy conservation by some 80 GeV. The Heisenberg uncertainty principle of quantum mechanics states that such a violation can occur but only for a limited time, within which the W^- converts into an electron and neutrino, which are the real end products of beta decay. The energy violation of 80 GeV is so extreme that the ensuing time scale, and in consequence the distance through which the W^- can travel, are exceedingly short. This limits the effects of the forces associated with W exchange whereby it appears “weak.” See ENERGY; UNCERTAINTY PRINCIPLE.

Data show that the fundamental couplings of a photon or of a W or Z^0 to a lepton or quark are similar in magnitude such that the electromagnetic and weak forces are now regarded as two manifestations of a single electroweak force. The apparent difference in their strengths is because the weak bosons are massive whereas the photon is massless. These differences in masses, and consequent differences in the effective forces, are examples of symmetry breaking. Theorists believe that the source of this broken symmetry is a spinless particle known as the Higgs boson. It is the interaction between the Higgs boson and the fundamental fermions and bosons that gives them their masses. The details of this interaction remain to be understood, and the first step is to confirm the existence of the Higgs boson. Theory predicts that its mass exceeds 100 GeV/ c^2 and that it should be manifested at experiments at CERN in Geneva, Switzerland, commencing from 2007. See ELECTROWEAK INTERACTION; HIGGS BOSON; SYMMETRY BREAKING.

Stability. Most particles are unstable and decay into particles with lower masses. Massless particles, such as the photon, are stable. Neutrinos, long thought to be massless and stable, are now known to have the possibility of small masses; the lightest neutrino is the lightest fermion and is stable as its decay would be into bosons, which could not conserve angular momentum. The present view is that the only massive particles that are strictly stable are the electron and the lightest neutrino or neutrinos. The electron is the lightest charged particle; its decay

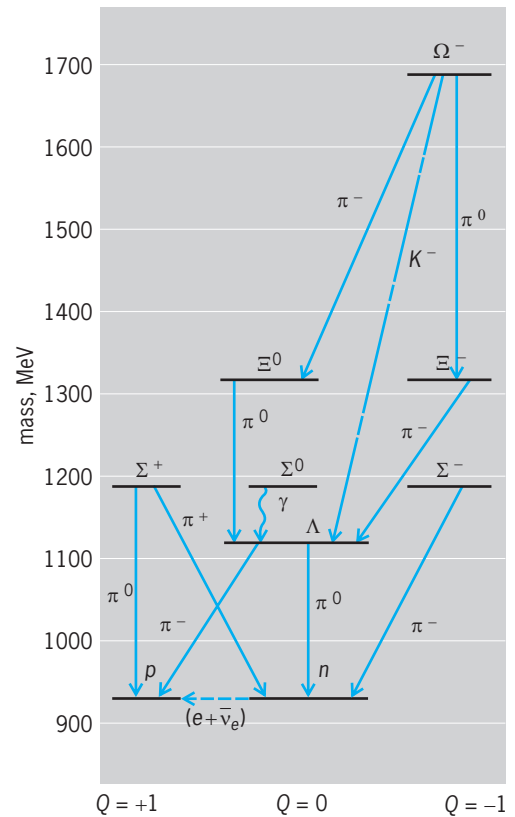


Fig. 1. Decay modes of quasistable baryons made from u , d , or s quarks.

would be into neutral particles and could not conserve charge. The proton may be stable, though there is no fundamental reason why it need be; it is to all practical purposes stable with a half-life exceeding 10^{30} years.

Unstable elementary particles must be studied within a short time of their creation, which occurs in the collision of a fast (high-energy) particle with another particle. Such high-energy particles occur naturally in the cosmic rays, but their flux is small; thus most elementary particle research is based on high-energy particle accelerators.

Hadrons can be divided into the quasistable and the unstable. The quasistable hadrons (Fig. 1) are those that are too light to decay into other hadrons by way of the strong interactions, such decays being restricted by the requirement that flavors be conserved. The quasistable hadrons that decay through weak interactions have long mean lives—more than 10^{10} times the characteristic time of strong interactions, $\hbar/(m_\pi c^2) = 0.5 \times 10^{-23}$ s, where m_π is the mass of the π meson or pion, and \hbar is Planck's constant divided by 2π . Three hadrons, π^0 , η , and Σ^0 , can decay by way of the electromagnetic interaction. These three have mean lives of the order of $10^5 \times \hbar/(m_\pi c^2)$.

The important practical distinctions in the experimental study of interactions are among (1) the stable massive particles (electrons and nuclei), which can be used as target particles as well as in beams; (2) the particles with mean lives greater than 10^{-8} s (γ , ν , μ^\pm , π^\pm , K_L , K^\pm), which can be used only

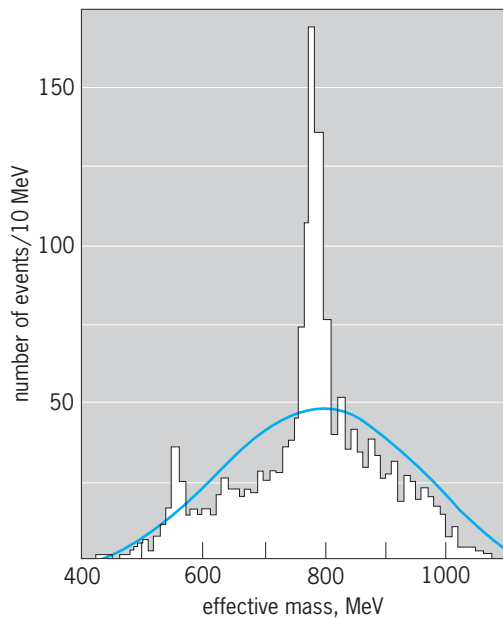


Fig. 2. Observation of the η meson ($m = 550$ MeV) and the ω meson ($m = 783$ MeV) as resonances in the reaction $\pi^+p \rightarrow \pi^+p \pi^+ \pi^- \pi^0$. Effective mass is the total relativistic energy of three of the emerging pions in their center-of-mass coordinate system. Peaks in the data indicate that collisions can create a short-lived particle of corresponding mass, which decays into three pions.

in beams; (3) the quasistable hadrons with mean lives of the order of 10^{-10} s ($\Xi^{0,-}$, Λ , Σ^\pm , K_S , Ω^-), which have only a small but usable chance of interacting in matter before decaying; and (4) the remaining hadrons, which have a vanishingly small chance of reinteracting except when produced within a nucleus.

The unstable hadrons are also called particle resonances. Their lifetimes, of the order of $\hbar/(m_\pi c^2)$, are much too short to be observed directly. Instead they appear, through the uncertainty principle, as spreads in the masses of the particles—that is, in their widths—just as in the case of nuclear resonances (Fig. 2).

The first unstable hadron observed was the $\Delta(1232)$ in pion-nucleon scattering. Flavorless vector mesons can be formed as resonances in collisions between electrons and positrons. This technique has led to the discovery of charmonium and bottomonium spectroscopies: quasiatomic systems made from a charm or bottom quark and its corresponding antiquark.

Antiparticles. To each kind of particle there corresponds an antiparticle, or conjugate particle, which has the same mass and spin but has opposite values of charge, strangeness, charm, or bottomness (quantum numbers that are conserved additively). Antiparticles are denoted by putting a bar over the symbol of the corresponding particle: antiproton \bar{p} , antiquark \bar{q} and so forth. Thus a baryon, such as the proton, made of three quarks implies that an antibaryon is made of three antiquarks, $\bar{q}\bar{q}\bar{q}$. All baryons and antibaryons have half-integer values of spin and are fermions. Mesons are examples of hadrons that have integer

spins and are made of an equal number of quarks and antiquarks. The simplest example of such a meson is thus $q\bar{q}$ and scores of examples are known. It is theoretically possible to form mesons made of $qq\bar{q}\bar{q}$, known as tetraquarks or molecules, or even of no quarks, made entirely from gluons, known as glueballs. There is some evidence that mesons with no spin, f_0 and a_0 , having masses 980 MeV/ c^2 , are examples of $qq\bar{q}\bar{q}$, and that heavier examples of spinless mesons f_0 with masses of 1500 and 1710 MeV/ c^2 may be quantum mixtures of $q\bar{q}$ and a glueball. See ANTIMATTER; BARYON; MESON.

Nomenclature. Hadrons of the same spin, of similar mass, and distinguished by their electrical charges are distinguished by writing charge as a superscript, for example, pi meson (π) and sigma hyperon Σ have π^+ and Σ^0 . (But p and n are usually written instead of N^+ and N^0 .)

With the discovery of large numbers of unstable hadrons, it has become impossible to name each with a different letter; instead letters are used to denote classes of hadrons. First, hadrons with a common value of spin (denoted J), whose quantum wave functions behave the same under space and charge parity, P and C , and whose masses are similar, form families with the same isospin, I . In recognition of these patterns, a naming scheme has developed (Table 2). For instance, π_2 is used for any meson with the properties $I = 1$ and $J^{PC} = 2^{-+}$, and a particular hadron is distinguished by giving its approximate mass, in MeV [for instance, $\pi_2(1680)$]. If the hadron is the lightest one of its class, the mass can be omitted from its name; it is always omitted if the hadron is hadronically stable. The mass can be replaced by a specification of the quark state of the hadron if this is known, for instance $\psi(2S)$ instead of $\psi(3686)$. The scheme is described in full in the *Review of Particle Physics*, which is published biennially and updated more frequently on the World Wide Web. See ISPIN; PARITY (QUANTUM MECHANICS); SYMMETRY LAWS (PHYSICS).

Since hadrons have an extended structure, they can have rotationally excited states, as do molecules and deformed nuclei. Such states form a sequence of hadrons with increasing spins ($J_0, J_0 + 1, \dots$) and masses but with the same values of other quantum numbers (except for space and charge parity, P and C , which alternate in sign). For historical reasons such sequences of hadrons are called Regge trajectories.

TABLE 2. Quantum numbers and names of mesons

J^a	PC	$I = 1$	$I = 0$	$s\bar{s}$	$c\bar{c}$	$b\bar{b}$	$t\bar{t}^b$
0,2,4,...	$++$	π	η	η'	η_c	η_b	η_t
1,3,5,...	$+-$	ρ	ω	ω'	ω_c	ω_b	ω_t
1,2,3,...	$--$	ρ	ω	φ	ψ	Υ	θ
0,1,2,...	$++$	a	f	f'	χ	χ_b	χ_t

^aCertain mesons, termed exotic states, may have other values of J . For example, the π_1 has $J^{PC} = 1^{-+}$.

^bThe t and \bar{t} quarks may decay before the mesons in this column can form.

Quantum chromodynamics. In addition to their electrical charges and flavors, quarks carry a property called color, which is the source of the force between quarks. As photons couple to electric charge and mediate the electromagnetic force, so massless gluons transmit the color force. The quantum theory of the color force is called quantum chromodynamics (QCD). *See* COLOR (QUANTUM MECHANICS).

The gluon field is coupled to color. The coupling of the gluon to a particle is fixed by the color of the particle and just one universal coupling constant g , analogous to the electronic unit of charge e . There are three varieties of color for quarks and three of opposite sign for antiquarks. The rules of attraction and repulsion for electrostatics have an analogy in chromostatics: like colors repel and unlike colors attract. The latter property has two manifestations. One is directly analogous to electrostatics, in that a quark with a given color attracts an antiquark carrying the opposite sign of color charge. The other possibility is attraction between quarks, each carrying a different color, and in a quantum state that is antisymmetric under interchange of the two particles.

As electric charges of opposite sign lead to neutral atoms, so do the color charges of quarks and antiquarks combine to form overall colorless hadrons. Such states are known as color singlets in the language of the mathematical theory underpinning QCD. As hadrons are colorless, the long-range forces observed between them are no different than those between other particles. The two simplest combinations of quarks that can be colorless are $q_1\bar{q}_2$ and $q_1q_2q_3$; these are found in nature as the basic structure of mesons and baryons, respectively. The exchange of gluons between any of the quarks in these colorless combinations gives rise to an attractive force, which binds them together.

Not just quarks and antiquarks but also gluons carry color, and therefore gluons can mutually attract and repel one another. This situation is very different from electromagnetism, where the photon does not carry charge. The consequence of this self-coupling is that the interaction between two colored particles through the gluon field, which at short distances is an inverse-square Coulomb force, proportional to g^2/r^2 (where r is the distance between the particles), becomes stronger than this inverse-square force at larger r . This property is interpreted by saying that the coupling strength g is effectively larger at larger r ; this coupling strength defines the so-called running coupling constant $g(r)$. Computer simulations of the QCD theory imply that this coupling becomes infinitely strong at distances of the order of 10^{-15} m, leading to confinement of the constituent quarks within colorless hadrons, in accord with experiment. The mutual interaction among gluons is predicted to lead to colorless hadrons made from gluons without the need for quarks. As mentioned above, such hadrons are known as glueballs.

This situation is very different from that in quantum electrodynamics (QED). There the long-range force is precisely of the inverse-square form at large

distance, with coefficient equal to the product of the charges, e^2 . In dimensionless form, this is $e^2/\hbar c \approx 1/137$, the dimensionless parameter α of QED. In QED, vacuum polarization shields a charge—just as in any polarizable medium—resulting in an increase in the apparent charge as one gets closer, penetrating the cloud of shielding charge; in QCD, vacuum polarization—the effect of the gluon self-coupling—is antishielding, resulting in an apparent decrease in the color strength as one gets closer. This short-distance property of QCD is known as asymptotic freedom. *See* FINE STRUCTURE (SPECTRAL LINES); FUNDAMENTAL CONSTANTS; QUANTUM ELECTRODYNAMICS.

At small r , or equivalently at large momentum transfer, the running coupling becomes small and consequently perturbation theory becomes reliable. Calculations in perturbative QCD, such as the gluonic radiative corrections to the cross section for $e^+e^- \rightarrow$ hadrons, are in good agreement with experiment. This is an important reason for the present view that QCD is the correct theory of the hadron glue. The ability to make precise computations in QCD for hadron interactions involving large transfers of momentum underpins much of the experimental planning at facilities such as the Large Hadron Collider (LHC). *See* QUANTUM CHROMODYNAMICS.

Quarkonium. According to the so-called naïve quark model, hadrons are bound states of nonrelativistic (slowly moving) quarks, analogous to nuclei as bound states of nucleons. The interactions between the quarks are taken qualitatively from QCD, namely a confining central potential and (exactly analogous to electromagnetic interactions) spin-spin (hyperfine) and spin-orbit potentials; quantitatively, these potentials are adjusted to make the energy levels of the model system fit the observed hadron masses. This model should be valid for hadrons composed of heavy quarks but not for hadrons containing light quarks (u , d , s), but in fact it succeeds in giving a good description of many properties of all hadrons. One reason is that many of these properties follow from so-called angular physics, that is, symmetry-based physical principles that transcend the specific model. *See* BARYON; MESON.

As mentioned above, mesons with the composition $c\bar{c}$ and $b\bar{b}$ are called charmonium and bottomonium. These names are based on the model of positronium, the atomic bound state of a positron and an electron; the generic name for flavorless mesons, $q\bar{q}$, is quarkonium (**Fig. 3**). The mesons are classified by their intrinsic angular momentum or spin, J , and the behavior of their quantum wavefunctions under parity P and charge conjugation C . *See* POSITRONIUM.

Some states of quarkonium are cleanly observable as resonances in e^-e^+ collisions: $e^-e^+ \rightarrow \gamma^* \rightarrow (q\bar{q})_1^{--}$, where γ^* is a virtual photon that materializes into a state of a quark and antiquark with $J^{PC} = 1^{--}$ (the quantum numbers of a photon). This vector quarkonium meson is made at rest in the center-of-mass frame of the colliding e^- and e^+ and is

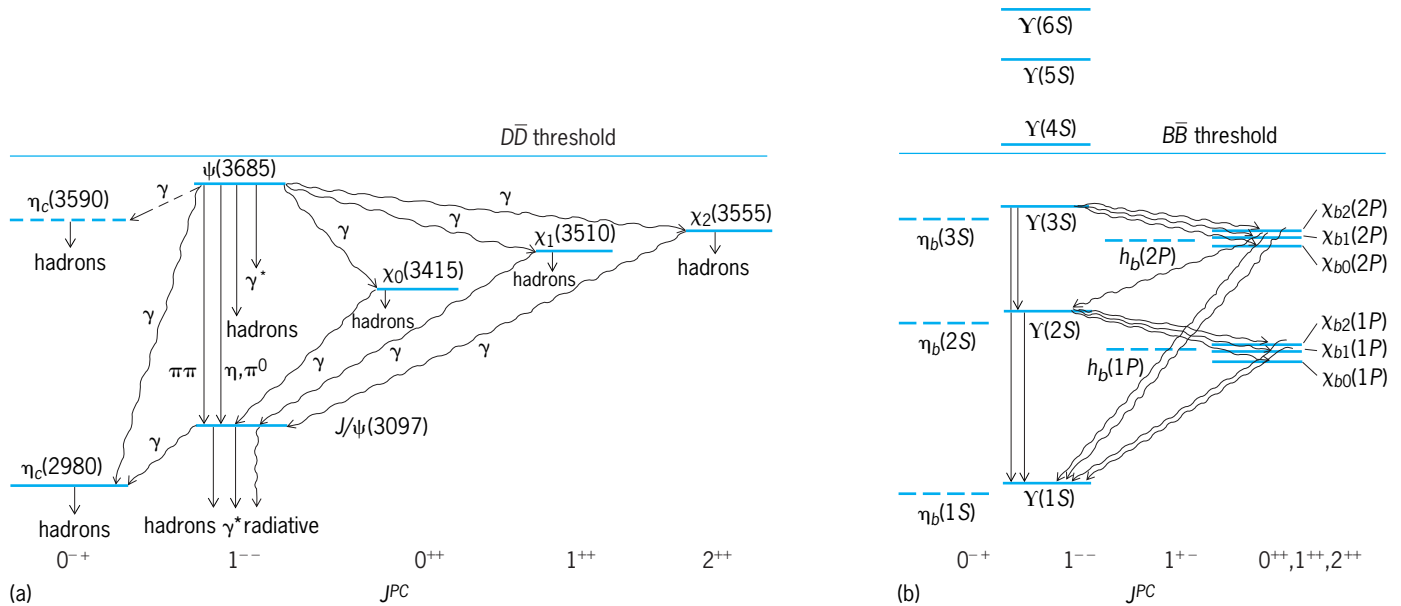


Fig. 3. Level schemes and transitions of (a) charmonium and (b) bottomonium. Observed radiative transitions are indicated by wavy lines, observed states and other observed transitions by solid lines, and uncertain states and transitions by broken lines. The notation γ^* refers to decay processes involving intermediate virtual photons, including decays to e^+e^- and $\mu^+\mu^-$. (After Particle Data Group, *Review of particle properties*, *Phys. Lett.*, 179B:1-350, 1986)

unaccompanied by other hadrons. Quarkonium states of other quantum numbers can be observed as decay products of excited 1^{--} mesons. Charmonium has been produced also by the annihilation of antiprotons on protons, and observed in the decay products of bottom mesons (B), as in the process $B \rightarrow K + J/\psi(3097)$.

Since both heavy quarkonium and positronium are systems of a fermion bound to its antifermion by a central force, they are qualitatively very similar (quantitatively, their masses and excitation energies are in the ratios of $\sim 10^4$ and 10^7 , respectively). However, the potential between the quarks is not $\sim 1/r$ (coulombic), which gives a different ordering in the energy levels for quarkonium and positronium. The pattern observed for quarkonium shows that the quark-antiquark potential is consistent with r (linear potential), which agrees with the prediction of QCD, as calculated by the lattice Monte Carlo method. Information on the spin-spin and spin-orbit parts of the $q\bar{q}$ interaction also comes from the quarkonium energy levels.

Just as positronium decays by annihilation of its constituents into photons, quarkonium can decay by annihilation of its q and \bar{q} into gluons, which then fragment into hadrons. This mechanism results in a strong C dependence of decay rates; for example, the $\eta_c(2981)$ decays ~ 200 times faster than the $J/\psi(3097)$. Vector quarkonium also has the decays $(q\bar{q})_1^{--} \rightarrow \Upsilon^* \rightarrow x\bar{x}$, where x is any charged fundamental fermion, lighter than the quark q ; this decay mode accounts for about 28% of the decays of the J/ψ and 17% of the decays of the upsilon (Υ) meson. See J/PSI PARTICLE; UPSILON PARTICLES.

Weak decays of quarkonium, in which one of the heavy quarks is transformed to a lighter quark, are

negligibly slow; for example, the chance that a J/ψ meson decays this way is $\sim 10^{-8}$.

Excited states of quarkonium have decay modes in which the heavy quarks neither annihilate nor transform. If the heavy quarks remain bound, this is a deexcitation, a decay into a lower state of the quarkonium with emission of one or more hadrons or a Υ (radiative decay). Radiative decays are the main source of the η_c , χ_0 , χ_1 , and χ_2 mesons.

Highly excited quarkonium states can decay into pairs of mesons. An example is the bottomonium decay $\Upsilon(10,575) \rightarrow B^-B^+$. The B meson is the lightest hadron which contains a b quark, and so $2m_B = 10,545$ MeV is the energy threshold for such decays; the $\Upsilon(10,575) = \Upsilon(4S)$ decays to $B\bar{B}$ nearly 100% of the time, with a decay rate 2000 times faster than the total decay rate of the next lower vector bottomonium state, $\Upsilon(10,355) = \Upsilon(3S)$. Generally, autoionization decay modes are the fastest decay modes of quarkonium whenever they are kinematically possible. Electron-positron annihilation at an energy of 10,355 MeV has been used as a source of large amounts of B mesons. Comparison of the properties of B and its antiparticle counterpart \bar{B} is used to investigate possible asymmetries between matter and antimatter.

Heavy-quark hadrons. The known heavy quarks are the c (charm), b (bottom), and t (top) quarks, whose masses are larger than the mass of a proton, ~ 1 GeV. A hadron which contains a single heavy quark resembles an atom where the heavy quark sits nearly at rest at the center, and is a static source of the color field, just as the atomic nucleus is a static source of the electric field. Just as an atom is changed very little (except in mass) if its nucleus is replaced by another of the same charge (an isotope), a

heavy-quark hadron is changed very little (except in mass) if its heavy quark is replaced by another of the same color. This is called heavy-quark symmetry. So, for example, the D and B mesons are similar, except in mass, which plays an important role in the quantitative analysis of their weak decays.

The only hadrons so far observed that contain more than one heavy quark are mesons containing a conjugate pair: $c\bar{c}$ and $b\bar{b}$. The creation of hadrons such as $c\bar{c}$ or ccu from ordinary matter, which contains only u and d quarks, requires the creation of two pairs of heavy quarks at nearly the same place and time; this is highly improbable, hence the difficulty in producing such states.

Bottom mesons B made from $b\bar{u}$ or $b\bar{d}$ and their antiparticle counterparts \bar{B} made of $\bar{b}u$ or $\bar{b}d$ have special interest as their weak decays are not symmetric between particle and antiparticle. It is hoped that the study of this asymmetry can shed light on the nature of the large-scale asymmetry between matter and antimatter in the universe. To this end, over 10^9 B and \bar{B} (“B-bar”) pairs have been made at customized “B-factories”: electron-positron colliders tuned to an energy of around 10.5 GeV where the production rate of these mesons is optimal. Such experiments have taken place at Stanford, California, and the Belle facility at KEK in Tsukuba, Japan, and definitive results have begun to emerge.

Weak interaction. The weak interaction changes the flavor of particles as in its most familiar manifestation, the beta decay of the neutron, which transmutes into a proton, electron, and neutrino. Neutrinos, which have no color or electric charge, and hence do not respond to strong or electromagnetic forces, are ideal probes of the weak interaction.

All experimental results at low energies are consistent with the weak reactions being point four-Fermi interactions, that is, occurring when four fermions are at one point in space-time. This was the essence of Fermi’s original theory of beta decay. As noted above, it is now known that such beta decays are the result of the exchange of a massive charged vector (spin-1) particle, the W^\pm (W^+ and W^- are antiparticles of one another), which is emitted and absorbed by the so-called weak charged current. From the quark point of view, the so-called charged-current weak interactions are flavor-changing; here flavor is meant in a broad sense which includes leptons. This particle is part of a renormalized field theory, which is a spontaneously broken gauge theory (by the Higgs mechanism). This is the electroweak theory of S. Weinberg, S. Glashow, and A. Salam, which now is accepted as the standard model of electroweak interactions. It has four gauge fields, whose quanta are the W^+ and W^- , a massive neutral vector particle Z^0 , and the photon.

The Z^0 is coupled to all particles with strengths controlled by a single parameter, the Weinberg (or weak) angle θ_w , whose value is determined by observation of the neutral-current weak interactions, which result from exchange of the Z^0 . These were

unknown before this theory suggested their existence. The particle that emits a Z^0 remains unchanged, as is the case with the emission of a photon.

Universality. An important property of the weak interaction is universality. This property was first observed in the fact that the interactions responsible for the beta decays $\mu \rightarrow e\nu_e\nu_\mu$ and $n \rightarrow p e\nu_e$ are nearly equal; that is, the couplings $W\nu_\mu\mu$ and Wud of the W^\pm are nearly equal. (Likewise it is found that the couplings $W\mu\nu_\mu$ and $W e\nu_e$ are equal.) This near-equality points to a deep similarity of leptons and quarks.

The universality is also seen in Wcs and Wtb couplings, which are also almost equal. In addition to these couplings within a generation, there are feeble couplings between the generations, at least for the quarks. This was originally found in the existence of two couplings involving the u quark, namely Wud (strangeness-conserving) and Wus (strangeness-changing), the latter about one-fourth as large as the former. This contrasts with the leptons, where the ν_e is involved in only one coupling, namely $W\nu_e e$ (about equal in strength to Wud), and $W\nu_e\mu$ vanishes. Nonetheless, an overall universality still occurs when all the possible transitions are accounted for in the mathematical matrix formalism of N. Cabibbo, M. Kobayashi, and T. Maskawa (involving the so-called CKM matrix). The numbers in the CKM matrix cannot, in general, be made real. As a result, the W coupling leads to an asymmetry between quarks and antiquarks, with potential implications for the large-scale asymmetry observed between matter and antimatter in the universe. However, while the values of the CKM matrix elements have been experimentally measured, so far there is no theory for them.

There is indirect evidence that a similar matrix formalism might apply to the leptons and that transitions such as $W\nu_e\mu$ occur. Such ideas also suggest that transitions between the electron and the muon might occur, albeit very rarely. This is an active area of research.

Weak bosons. The weak bosons W^+ , W^- ($=\overline{W^+}$), and Z^0 were first observed in 1983 in proton-antiproton collisions. The W^\pm , and similarly the Z^0 , are made by the annihilation of a quark pair in a $p\bar{p}$ collision, $q\bar{q} \rightarrow W \rightarrow$ decay products, where the q and \bar{q} are contained in an incident proton p and antiproton \bar{p} , respectively. The branching ratios for the decays into quarks (Table 1) have been predicted from theory and confirmed in experiment. What the weak bosons do not decay into is important information about the nonexistence of additional fundamental particles (quarks, leptons, gauge bosons, or Higgs particles) with masses less than 45 GeV. The decay modes $e^\pm\nu_e$ and $\mu^\pm\nu_\mu$ of the W^+ , and e^+e^- and $\mu^+\mu^-$ of the Z^0 , are easy to observe—despite the rarity of the W or Z production (about 1 per 106 $p\bar{p}$ collisions)—because no other process makes many leptons with momenta as high as 40 GeV, at a large angle to the incident p and \bar{p} beams. Because there are three times as many kinds of quarks as

leptons (each flavor of quark comes in three colors), most decays of the W and Z are into pairs of quarks, which fragment into 40-GeV jets; but these decays are hard to observe because of the background of similar jets made by the scattering process $q\bar{q} \rightarrow q\bar{q}$.

The Z is made cleanly by e^+e^- annihilation, enabling its decay modes to be easily seen. The decays $Z \rightarrow \nu_l \nu_l$, where ν_l is a neutrino, cannot be seen directly, but at least the total branching ratio into these modes can be deduced by observation of the total width of the Z . With the assumption of universality, this determines the number of light (mass < 45 GeV) neutrinos to be 3; that is, there are none beyond those presently known: ν_e , ν_μ , and ν_τ . By implication, if every neutrino is accompanied by a charged lepton, and this lepton pair is also accompanied by a pair of quarks (charges $2/3$ and $-1/3$), the three generations of leptons and quarks are the totality of such fermions.

Grand unified theory. The negative charge of the electron and the positive charge of a proton are equal to better than 1 part in 10^{19} . As the electron appears to be a fundamental fermion, as are the quarks that make up the proton, this equality is remarkable and hints at some deeper relation between leptons, such as the electron, and the quarks. The balance is achieved by a further act of balance: the electrical charges of quarks come in fractions that are $1/3$ of an integer (in units of a proton), and the strong color forces acting on the quarks cluster them in threes. The three-ness of the color charges conveniently conspires therefore with the third-integer nature of the electric charges, suggesting some relation between the electric and color charges.

These hints are further solidified when we look at the quantum field theories of the forces. The carriers of the electromagnetic and color forces are massless spin-1 particles, the photon and gluon. The weak force also is carried by spin-1 particles, the W and Z , though here they are massive. Apart from this mass effect that spoils the symmetry, and the fact that the weak interactions do not respect mirror symmetry (parity), the underlying quantum theories of the three forces show profound similarities.

Their strengths differ at room temperature, due to the different strengths of the interaction between the spin-1 particles and the fermions and also to the fact that the W and Z are massive. However, in experiments at higher energies these differences have been observed to die away. Experiments at the LEP accelerator at CERN have shown that when electrons and positrons annihilate at a total energy of 90 GeV, the Z boson is produced with a strength that is essentially the same as that of the electromagnetic interaction. Studies of the W boson at these extreme energies show that here too the weak interaction has an intrinsic strength that is united with the electromagnetic and has been obscured at low energies by the large masses of the W and Z , in contrast to the massless photon. Thus, as noted above, it

is now accepted that the electromagnetic and weak forces are two manifestations of a single electroweak force.

The most familiar form of the weak interaction is the beta decay, such as that of a neutron into a proton, electron, and neutrino, which is ultimately driven by a down quark transmuting to an up with the emission of a virtual W^- . The coupling of the W to quarks appears to be identical to its interaction with leptons: were the W boson the only probe of leptons and quarks, they would be indistinguishable. This identity further strengthens the belief that leptons and quarks are profoundly related.

In theories that postulate such a relation, the leptons and quarks occur together in multiplets of the large symmetry group called families (or generations). The known fundamental fermions do seem to fall into three families (Table 1). Each family consists of a doublet of leptons (neutrino [charge 0] and charged lepton [charge $-e$]), and a color triplet of doublets of quarks (up-type [charge $2/3 e$] and down-type [charge $-1/3 e$]). (From the total width of the Z^0 , it is deduced that there are no more families whose neutrino has mass < 45 GeV.) However, there is no known connection between the leptons and the quarks of each family; for example, it is only because (ν_e, e) is the lightest lepton doublet and (u, d) is the lightest quark doublet that they are regarded as members of the same family.

The main difference between them is that quarks carry color and respond to the strong force whereas leptons do not. Here, though, there are further hints of unity at high energies. The strength of the color forces in QCD weakens with increasing energy; the strength of the analogous electroweak interaction, according to the quantum theory, increases. Extrapolating to higher energies, the theories imply that these forces all have a similar intrinsic strength at 10^{15} GeV, known as the grand unification energy. In grand unification theories, this energy is the mass m_{lq} of hypothetical superheavy "leptoquark gauge bosons," analogous to the W^\pm and Z^0 bosons of the electroweak subtheory. At momentum transfers larger than this, the mass of the leptoquark bosons is irrelevant and the large symmetry is unbroken; for smaller momentum transfers, as at low energies, the symmetry is broken and the three forces appear different, the more so the lower the momentum transfer.

The couplings of the leptoquark gauge bosons turn leptons into quarks, or vice versa (this is the reason for the name leptoquark), or quarks into antiquarks. There are many theoretical ideas on the mathematical structure of the true unified theory. Different approaches that agree with known phenomena have widely different implications for the properties of new particles that are implied by the theory. Nonetheless, there are certain common features, such as that the exchange of a leptoquark boson can therefore result in the transformation $qqq \rightarrow lqq$, for example, $p \rightarrow e^+ \pi^0$. This baryon- and lepton-number-violating interaction is a

much weaker interaction than the analogous ordinary weak interaction, because leptoquark bosons are much heavier than weak bosons. The implication that protons are unstable, albeit with a half-life that exceeds the age of the universe by many orders of magnitude, is a rather general consequence of attempts to build an empirically successful unified theory. Studies of large samples of matter in the hope of finding evidence for a single proton decay have been made, but no definitive evidence for proton decay has yet been found. *See* GRAND UNIFICATION THEORIES.

Accelerators. For a century, beams of particles have been used to reveal the inner structure of atoms. These beams have progressed from naturally occurring alpha and beta particles, courtesy of natural radioactivity, through cosmic rays, to intense beams of electrons, protons, and other particles at modern accelerators. By smashing the primary beams into a target, some of the energy can be converted into new particles, which can themselves be accumulated and made into secondary beams. Thus, beams of pions, kaons, muons, and neutrinos have been made, as well as antiparticles such as positrons and antiprotons. There are even beams of heavy ions—atoms stripped of their electrons—which enable violent collisions between heavy nuclei to be investigated.

There has also been a renewed interest in cosmic rays, where nature provides particles at energies far beyond anything that we can contemplate achieving on Earth. The problem is that such rays come at random and are much less intense than beams made at accelerators.

The basic principle is that electrically charged particles are accelerated by electric forces. If enough electric force is applied to an electron, say, it will be accelerated along a straight path, as in the linear accelerator at Stanford in California, which can accelerate electrons to energies of 50 GeV in a distance of 3 km (2 mi).

Under the influence of a magnetic field, the path of a charged particle will curve. By using electric fields to speed them, and magnetic fields to bend their trajectory, particles can be guided around circles over and over again. This is the basic idea behind huge rings, such as the 27-km (17-mi) LEP circular accelerator at CERN, which accelerated electrons and positrons to 100 GeV each.

Originally, the beams were directed at static targets. In a linear accelerator aimed at a static target, the debris of the collision is propelled forward, just as a stationary car is shunted forward when another car crashes into its rear. When a beam hits a stationary target, the hard-won energy of the beam particles is transferred largely into energy of motion—into moving particles in the target—and is effectively wasted. This problem is overcome if particles can be brought to collide head-on, so that their energy can be spent on the interaction between them. In such a collision the debris flies off in all directions, and the energy is redistributed with it—none is wasted in setting stationary lumps in motion. Thus, there has been an

increasing strategy of making counterrotating beams of particles and antiparticles, such as electrons and positrons, or protons and antiprotons, and colliding them head-on.

The major application has been to enable collisions between particles and antiparticles, principally protons and antiprotons or electrons and positrons. A limit to the energy attainable at a circular accelerator of electrons or positrons is that these high-energy particles radiate away energy when they travel on a circular path. This “synchrotron radiation” is greater the tighter the radius of the orbit and the higher the energy of the particle. Protons and antiprotons also emit synchrotron radiation; but with a mass of nearly 2000 times that of an electron, they can reach much higher energies before the amount of energy lost becomes significant, and they also pack a bigger punch. Hence they are the prime choice when the aim is to reach out to previously unexplored higher energies, as at the Large Hadron Collider. *See* SYNCHROTRON RADIATION.

At CERN the 27-km ring of magnets that accelerated electrons and positrons to 100 GeV has been replaced by the superconducting magnets of the Large Hadron Collider, which can guide protons up to energies of 8000 GeV (8 TeV). Counterrotating beams will collide head-on at a total energy of 16 TeV. Experiments at this facility were scheduled to begin in late 2007.

In order to carry out fine-detail studies of new particles discovered at the Large Hadron Collider, there are plans to produce large numbers of them under more controlled conditions. To do so, the plan is to generate electron-positron collisions at the optimum energy. As this energy is expected to be several hundred GeV, two linear accelerators will be required—one for the electrons and the other for the positrons—which are aligned so as to produce head-on collisions of the beams.

Supersymmetry. The fundamental particles of matter, the leptons and quarks, are all fermions with spin $\frac{1}{2}$. The forces that act on them are mediated by the photon, the gluon, and the *W* and *Z* bosons, all with spin 1. Grand unification is based on this common feature of bosons as force carriers acting on the basic fermions. Supersymmetry (SUSY) theory implies that there is a further symmetry between the forces and the matter particles, such that the known fermions are partnered by new bosons, and the known bosons by new fermions, with novel forces transmitted by these fermions. It is hypothesized that this partnering may lead to a more complete unification between particles and forces.

In SUSY the families of bosons that twin the known quarks and leptons are known as superquarks and superleptons, more commonly referred to as squarks and sleptons. If SUSY were an exact symmetry, each variety of lepton or quark would have the same mass as its squark or slepton sibling. However, the selectron empirically must have a mass greater than 100 GeV, which implies that it must be hundreds of thousands of times more massive than the electron. Similar remarks follow for all the sleptons and

squarks, implying that SUSY is a very badly broken symmetry.

An analogous statement applies to the superpartners of the known bosons. The naming convention here is to add -ino (pronounced "eeno") to denote the superfermion partner of a standard boson. Thus there are predicted the photino, gluino, wino, and zino. The hypothetical graviton, carrier of gravity, is predicted to have a partner, the gravitino. Here too the SUSY is badly broken, and these "inos" must have masses far greater than their conventional counterparts.

Although SUSY must be badly broken, in the sense that the masses of the particles and sparticles differ greatly, the basic ideas are appealing mathematically and merge as a fundamental symmetry of space and time as encoded in Einstein's theory of relativity and the quantum theory. The resulting pattern of sparticles solves technical problems in the quantum theories of the forces at high energies and the response of the particles to those forces, in particular preventing nonsensical predictions such as that certain processes occur with infinite probability. Quantum effects, due to SUSY particles emerging fleetingly as virtual particles in accord with the uncertainty principle, can affect measurable quantities in present experiments. From such data it is predicted that the lightest sparticles may occur with masses of a few hundred GeV and as such be observable at experiments at the Large Hadron Collider. There is also the tantalizing possibility that the lightest sparticles are electrically neutral, such as the photino or gluino, and are metastable. As such they could form large-scale clusters under their mutual gravitational attraction and form a substantial part of the dark matter of the universe. See DARK MATTER; SUPERSYMMETRY; WEAKLY INTERACTING MASSIVE PARTICLE (WIMP).

Mass. The masses of the elementary particles have been discussed previously. Here, a description will be given of how these masses affect physical phenomena, and some of the enigmas associated with mass will be summarized.

The mass of matter in bulk is almost entirely due to the nucleons, protons, and neutrons, that form the nuclei of atoms. The mass of the proton is caused by the confining effects of QCD, whereby the quarks and gluons that make up the proton are entrapped within a sphere of radius approximately 10^{-15} m. Heisenberg's uncertainty principle implies that when spatially constrained to such a distance, the constituents gain momentum and energy of hundreds of MeV. The equation $E = mc^2$ then implies that this energy acts as an inertia or mass, leading to the 938-MeV/ c^2 mass of the proton.

The individual up and down quarks within the nucleons carry masses of only a few MeV/ c^2 . Being permanently confined within hadrons, a direct measure of their masses is not possible, but they can be determined indirectly to be in the range 1.5–4 MeV for the up and 4–8 MeV for the down. It is the greater mass of the down quark that gives the neutron its slightly larger mass relative to the proton.

The strange quark is approximately 100 MeV/ c^2 more massive than the up and down, which gives the extra mass to strange hadrons relative to their nonstrange counterparts. Charm, bottom, and top quarks have larger masses, approximately 1.2, 4.5, and 175 GeV/ c^2 , respectively. The fundamental leptons also have a spread of masses. The neutrinos have masses that are too small to determine, and they had long been thought to be massless; however, the phenomenon of neutrino oscillations shows that at least two of the three have nonzero masses, of the order of a few eV/ c^2 at most, possibly only a fraction of this. The electron, muon, and tau leptons have masses of about 0.5, 106, and 1777 MeV/ c^2 , respectively. The reason for these patterns in the masses is not understood.

The symmetry properties of unified theories would be realized if all the fundamental particles were massless. In the case of the photon, gluon, and graviton, this is the case empirically, but the W and Z bosons have masses of 80.4 and 91.2 GeV/ c^2 , respectively. As described above, it is their large masses that enfeebls the weak forces transmitted by W and Z relative to the electromagnetic force transmitted by the massless photon.

Relative to the grand unification energy of 10^{15} GeV, all of these masses are negligible, and it is theorized that the true unified symmetries of nature would be exact but for the effect of masses. The nature and origin of these masses is thus a fundamental problem of great interest, and the phenomenon is known as electroweak symmetry breaking. The standard model explains mass by proposing that it is due to a new field called the Higgs field, after Peter Higgs who in 1964 recognized this theoretical possibility. Mass is then the effect of the interaction between the fundamental particles and the Higgs field. Photons do not interact with the Higgs field and thereby remain massless; the W and Z bosons do interact and thereby gain their large masses. The quarks and leptons are also presumed to gain their masses by such interactions.

The Higgs field is manifested in quantum theory by the appearance of particles known as Higgs bosons. In Higgs's original formulation, the boson has spin 0. However, if supersymmetry is realized in nature, there should be a family of such particles including the fermion super-Higgs or "Higgsino." The Higgs particles are predicted to occur with masses between 100 and 2000 GeV/ c^2 . The Large Hadron Collider at CERN has been constructed so as to be able to cover this entire range of possibilities and thereby to identify the origin of electroweak symmetry breaking. Frank E. Close

Bibliography. F. Close, *The New Cosmic Onion*, Taylor and Francis, 2006; F. Close, *Particle Physics: A Very Short Introduction*, Oxford, 2004; B. Greene, *The Elegant Universe*, Jonathan Cape, 1999; S. Weinberg, *Dreams of a Final Theory*, Vintage, 1993; W. S. C. Williams, *Nuclear and Particle Physics*, Oxford, 1991; W.-M. Yao et al. (Particle Data Group), Review of particle physics, *J. Phys. G*, 33:1, 2006.

Elements, cosmic abundance of

The average chemical and isotopic composition of the solar system is appropriately referred to as cosmic, since this elemental abundance distribution is found to be nearly the same for interstellar gas and for young stars associated with gas and dust in the spiral arms of galaxies. The Sun makes up more than 99.9% of the mass of the solar system, so the bulk chemical composition of the solar system is essentially the same as that of the Sun. The cosmic abundances of the nonvolatile elements are determined from chemical analyses of a type of meteorite known as CI chondrites, whereas the relative abundances of the volatile elements are determined from quantitative measurements of the intensities of elemental emission lines from the Sun's photosphere. In most silicate-rich meteorites and the Earth, Moon, Venus, and Mars, the most abundant elements are

oxygen, magnesium, silicon, iron, aluminum, and calcium. Average solar-system composition consists of 70.7 wt % hydrogen, 27.4 wt % helium, and only 1.9 wt % of all remaining elements, lithium to uranium. Cosmic abundances are now widely referred to as standard abundances in the astrophysical literature. See ASTRONOMICAL SPECTROSCOPY; ELEMENT (CHEMISTRY).

The cosmic abundances of the elements are given in numbers of atoms of one element relative to a certain number of atoms of a reference element (see **table**). Silicon is commonly taken as the reference element in the study of the composition of the Earth and meteorites, and the data are given in atoms per 10^6 atoms of silicon. Photospheric abundances in the Sun and other stars are often given in the logarithmic "dex" scale, relative to 10^{12} atoms of hydrogen [$A_{\text{El}} = \log(N_{\text{El}}/N_{\text{H}}) + 12.00$], where A_{El} is the photospheric abundance of element El, N_{El} is the number

Cosmic abundances of the elements (expressed as number of atoms per 10^6 atoms of silicon)*

Atomic number (Z)	Element symbol	CI chondrites [†]	Solar photosphere [‡]	Cosmic abundance [‡]	Atomic number (Z)	Element symbol	CI chondrites	Solar photosphere	Cosmic abundance
1	H	5.29×10^6	2.79×10^{10}	2.79×10^{10}	44	Ru	1.86	1.93	1.86
2	He	6.04×10^{-10}	2.72×10^9	2.72×10^9	45	Rh	0.344	0.367	0.344
3	Li	57.1	0.351	57.1	46	Pd	1.39	1.36	1.39
4	Be	0.73	0.70	0.73	47	Ag	0.486	0.243	0.486
5	B	16.9	9.9	16.9	48	Cd	1.61	1.64	1.61
6	C	7.58×10^5	9.23×10^6	9.23×10^6	49	In	0.184	1.27	0.184
7	N	5.99×10^4	2.32×10^6	2.32×10^6	50	Sn	3.82	2.79	3.82
8	O	7.66×10^6	1.88×10^7	1.88×10^7	51	Sb	0.309	0.279	0.309
9	F	843	1012	843	52	Te	4.81		4.81
10	Ne	2.36×10^{-12}	3.35×10^6	3.44×10^6	53	I	0.90		0.90
11	Na	5.74×10^4	5.96×10^4	5.74×10^4	54	Xe	3.50×10^{-13}		4.7
12	Mg	1.074×10^6	1.06×10^6	1.074×10^6	55	Cs	0.372		0.372
13	Al	8.49×10^4	8.22×10^4	8.49×10^4	56	Ba	4.49	3.76	4.49
14	Si	$\equiv 1.00 \times 10^6$	$\equiv 1.00 \times 10^6$	$\equiv 1.00 \times 10^6$	57	La	0.4460	0.412	0.4460
15	P	1.04×10^4	7.85×10^3	1.04×10^4	58	Ce	1.136	1.06	1.136
16	S	5.15×10^5	5.96×10^5	5.15×10^5	59	Pr	0.1669	0.143	0.1669
17	Cl	5.24×10^3	8.81×10^3	5.24×10^3	60	Nd	0.8279	0.881	0.8279
18	Ar	9.62×10^{-12}	7.00×10^4	1.01×10^5	62	Sm	0.2582	0.285	0.2582
19	K	3.77×10^3	3.67×10^3	3.77×10^3	63	Eu	0.0973	0.0902	0.0973
20	Ca	6.11×10^4	6.38×10^4	6.11×10^4	64	Gd	0.3300	0.367	0.3300
21	Sc	34.2	41.2	34.2	65	Tb	0.0603	0.0221	0.0603
22	Ti	2.40×10^3	2.92×10^3	2.40×10^3	66	Dy	0.3942	0.385	0.3942
23	V	293	279	293	67	Ho	0.0889	0.0507	0.0889
24	Cr	1.35×10^4	1.30×10^4	1.35×10^4	68	Er	0.2508	0.237	0.2508
25	Mn	9.55×10^3	6.84×10^3	9.55×10^3	69	Tm	0.0378	0.0279	0.0378
26	Fe	9.00×10^5	8.81×10^5	9.00×10^5	70	Yb	0.2479	0.335	0.2479
27	Co	2.25×10^3	2.32×10^3	2.25×10^3	71	Lu	0.0367	0.0320	0.0367
28	Ni	4.93×10^4	4.95×10^4	4.93×10^4	72	Hf	0.154	0.211	0.154
29	Cu	522	452	522	73	Ta	0.0207		0.0207
30	Zn	1.26×10^3	1.11×10^3	1.26×10^3	74	W	0.133	0.359	0.133
31	Ga	37.8	21.1	37.8	75	Re	0.0517		0.0517
32	Ge	119	71.6	119	76	Os	0.675	0.785	0.675
33	As	6.56		6.56	77	Ir	0.661	0.624	0.661
34	Se	62.1		62.1	78	Pt	1.34	1.76	1.34
35	Br	11.8		11.8	79	Au	0.187	0.285	0.187
36	Kr	1.64×10^{-13}		45	80	Hg	0.34		0.34
37	Rb	7.09	11.1	7.09	81	Tl	0.184	0.221	0.184
38	Sr	23.5	26.0	23.5	82	Pb	3.15	2.48	3.15
39	Y	4.64	4.84	4.64	83	Bi	0.144		0.144
40	Zr	11.4	11.1	11.4	90	Th	0.0335		0.0335
41	Nb	0.698	0.733	0.698	92	U	0.0090	<0.0094	0.0090
42	Mo	2.55	2.32	2.55					

*CI chondrite data are from Anders and Grevesse (1989), except for B for which data of Zhai and Shaw (1994) were used. Solar photosphere data are from Grevesse and Sauval (1998). Cosmic abundances are from Anders and Grevesse (1989), except for B, for which the newer data of Zhai and Shaw (1994) were used, and C, N, and O, for which the newer photosphere values of Grevesse and Sauval (1998) were used. For elements with long-lived radioactive isotopes, abundances are present-day values; abundances of these elements 4.5×10^9 years ago were higher.

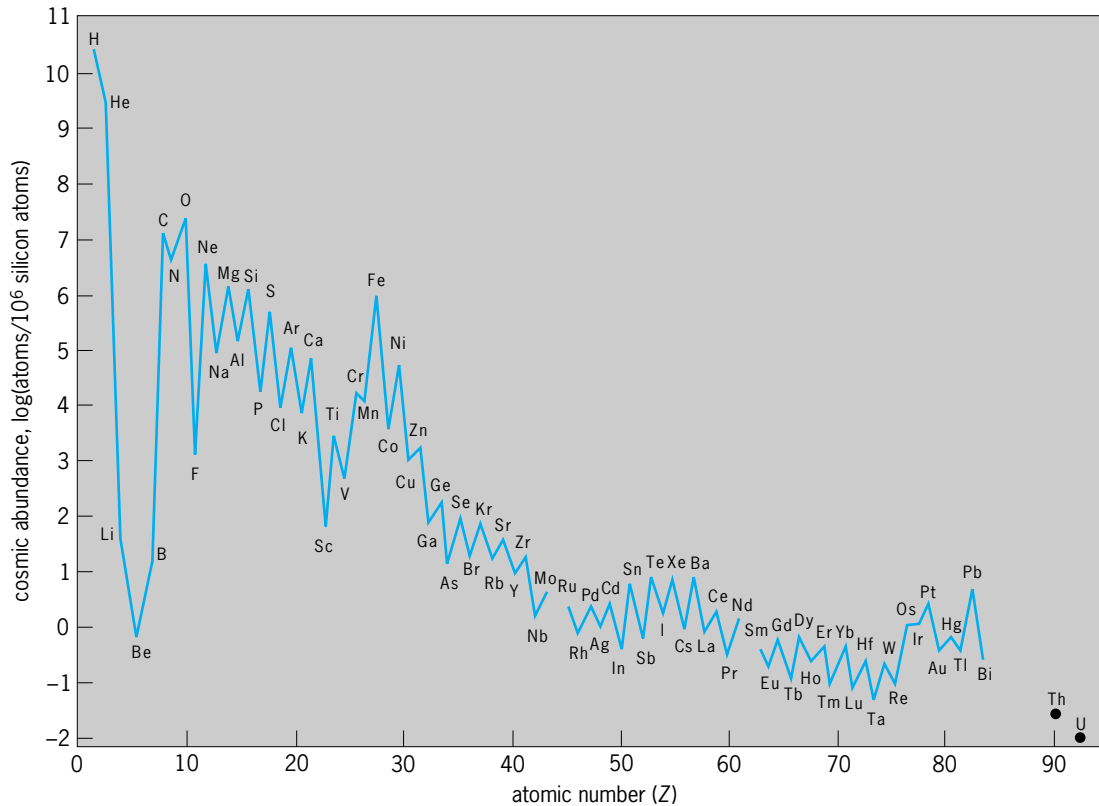


Fig. 1. Cosmic abundances of the elements, plotted as a function of atomic number (Z). Abundances are given in atoms per 10^6 atoms of silicon. For radioactive elements, abundances are those at the origin of the solar system 4.5×10^9 years ago. Elements with an even number of protons are more abundant than those with an odd number. (After E. Anders and N. Grevesse, *Abundances of the elements: Meteoric and solar*, *Geochim. Cosmochim. Acta*, 53:197-214, 1989)

of atoms of element E_l per 10^{12} atoms of hydrogen, and N_H is the number of atoms of hydrogen (10^{12}). See ELEMENTS, GEOCHEMICAL DISTRIBUTION OF.

Uses. Cosmic abundances of elements have several important uses. First, by comparing cosmic abundances to chemical analyses of various types of meteorites, inferences can be made about chemical fractionation processes that occurred in the primitive solar nebula, such as condensation and vaporization. Also, by comparing them to rock compositions, inferences can be made about processes that occurred early in the history of rocky planets, such as separation of a metallic core and differentiation of silicates into mantle and crust. Second, cosmic abundances serve as a standard of comparison for spectroscopic measurements of elemental abundances of the photospheres of other stars and for measurements of elemental and isotopic abundances in cosmic rays. Finally, nucleosynthesis occurs in many different stellar environments. Explanations of nucleosynthesis must account for how elements and isotopes from various astrophysical sources are made and then mixed to form the solar system's average chemical and isotopic composition. See NUCLEOSYNTHESIS; SOLAR SYSTEM.

CI chondrites. While there may be no obvious reason why a meteorite should have the average composition of the solar system, there are strong arguments that a particular group of meteorites, the CI chondrites, do in fact have the average nonvolatile

elemental abundances of the solar system. The CI chondrites are primitive meteorites that do not actually contain chondrules; rather, they are composed primarily of clay minerals. Their textures indicate some hydrous alteration from an unknown precursor, but this alteration apparently did not affect bulk elemental composition. See CLAY MINERALS.

The processes of chemical fractionation in primitive meteorites are well known. The major sources of chemical fractionation are separation of elements by volatility, due to partial evaporation or incomplete condensation of the components within meteorites; separation of siderophile (metal-seeking) from lithophile (rock-seeking) elements, by magnetic or density separation; and separation by the passage of liquid water in parent bodies, whereby elements in water-soluble form, such as sodium (Na), potassium (K), calcium (Ca), and bromine (Br), can be fractionated from elements in water-insoluble forms. All of these processes could potentially have fractionated CI chondrites relative to the true average chemical composition of the solar system. However, when CI chondrite abundances are compared with those of the solar photosphere, the small deviations that are found are independent of volatility, siderophile character, or susceptibility to aqueous alteration. The only significant deviations found are that the non-volatile elements lithium (Li), beryllium (Be), and boron (B) are depleted in the Sun relative to CI chondrites. Nuclear reactions in the Sun are known to

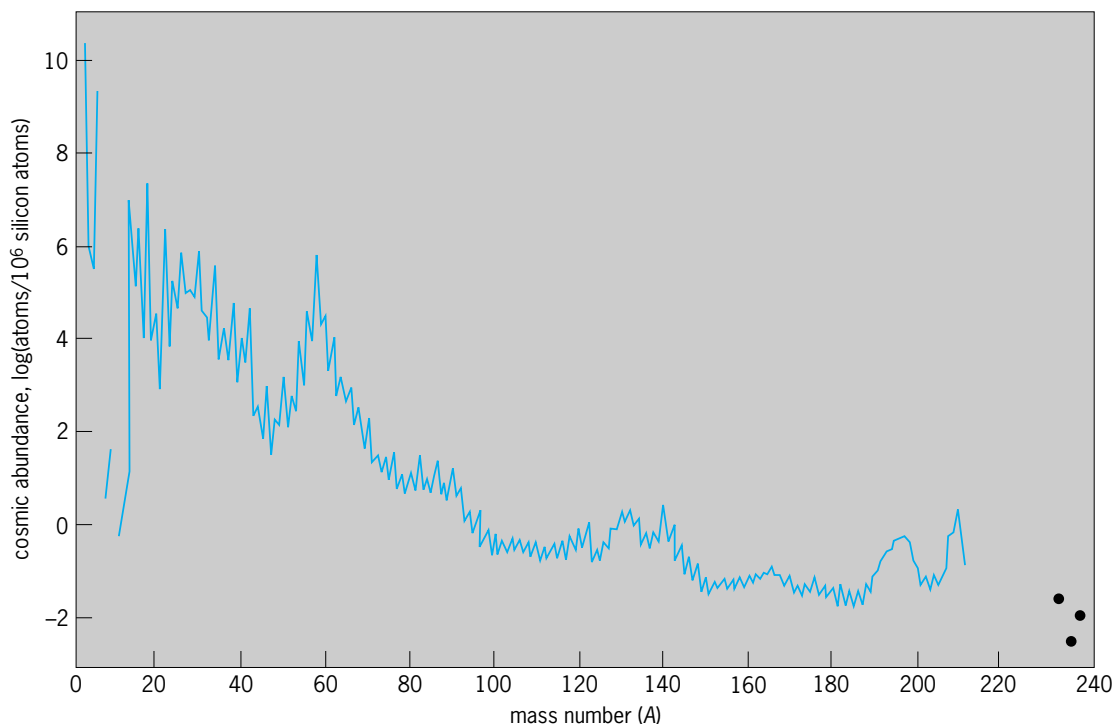


Fig. 2. Cosmic abundances of the isotopes, plotted as a function of mass number (A). Abundances are given in atoms per 10^6 atoms of silicon. For radioactive nuclides, abundances are those at the origin of the solar system 4.5×10^9 years ago. For some mass numbers, there is more than one element. A smooth curve can be drawn through the abundances of the odd-mass isotopes. (After E. Anders and N. Grevesse, *Abundances of the elements: Meteoric and solar*, *Geochim. Cosmochim. Acta*, 53:197–214, 1989)

consume these elements. Similarly, when abundances of odd-mass isotopes are plotted as a function of mass, there are no systematic deviations of groups of elements according to their chemical properties. See METEORITE.

Systematic patterns. There are a number of interesting features of cosmic abundances when they are plotted as a function of atomic number (Z) or mass number (A) (Figs. 1 and 2). When they are plotted as a function of Z , there are maxima at hydrogen (H), oxygen (O), iron (Fe), xenon (Xe), and lead (Pb); and there is a pronounced odd-even effect where elements with an even number of protons are more stable. There are a number of peaks and valleys that are seen more readily when cosmic abundances are plotted as a function of mass. Hydrogen-1 (^1H) and helium-4 (^4He) are the most abundant isotopes; there is a gap between ^4He and carbon-12 (^{12}C) where Li, Be, and B have very low abundances; there is a decrease in abundances from $A = 12$ (carbon) to $A = 45$ (scandium); a rise to a peak at $A = 56$ (iron) followed by a steep, then gradual decrease in abundances among the heavier elements, interrupted by several small peaks. After $A = 45$, it becomes readily apparent that even-mass isotopes are more abundant than odd-mass isotopes.

The idea that there might be some systematics in the abundances of the elements as a function of mass has a long history. It has been determined that elemental abundances in terrestrial rocks have been chemically fractionated by a variety of processes in the 4.5 billion years of Earth history. In 1947, it was found that the abundances of nuclides, especially of

odd mass, are a smooth function of mass number. The most widely used compilation of cosmic abundances of elements indicates a level of precision where deviations from the smoothness of the odd-mass abundance versus mass-number curve (Fig. 2) have become apparent (see table). See ISOTOPE.

Comets and interplanetary dust. Other potentially useful samples of cosmic abundances are comets and interplanetary dust. Abundances in dust ejected from Halley's Comet were measured by mass spectrometers on the *Giotto* spacecraft that flew past Halley. Abundances match the solar photosphere for many elements, but there are two notable discrepancies: Halley is low in iron and high in silicon. Halley's Comet, at least the dust portion of it, does not appear to represent pristine solar nebular matter. The most pristine samples of interplanetary dust have been collected by aircraft flying at high altitude. Small grains (less than 15 micrometers in diameter) show no evidence of heating upon atmospheric entry. These grains have the mineralogy and chemistry of CI and CM chondrites. CM chondrites λ (type 2) have a bulk chemical composition indicating some nebular volatility fractionation when compared with CI chondrites, but are more primitive than other types of meteorites. Interplanetary dust particles are so small that analyses are not of high enough precision to judge whether they are a better sample of cosmic abundances than CI chondrites. See COMET; COSMOCHEMISTRY; HALLEY'S COMET; INTERPLANETARY MATTER; SPACE PROBE.

Andrew M. Davis

Bibliography. E. Anders and N. Grevesse, *Abundances of the elements: Meteoritic and solar*,

Geochim. Cosmochim. Acta, 53:197–214, 1989; N. Grevesse and A. J. Sauval, Standard solar composition, *Space Sci. Rev.*, 85:161–174, 1998; H. Suess and H. C. Urey, The abundance of the elements, *Rev. Mod. Phys.*, 28:53–74, 1956; M. Zhai and D. M. Shaw, Boron cosmochemistry, Part I: Boron in meteorites, *Meteoritics*, 29:607–615, 1994.

Elements, geochemical distribution of

The distribution of the chemical elements within the Earth in space and time. Knowledge of the geochemical distribution of the elements in the Earth, particularly in the Earth's crust, and of the processes that lead to the observed distributions make it possible to locate and use efficiently essential elements and minerals and to predict their dispersal patterns when they reenter the natural environment after use.

Formation of the Earth. To understand the present-day distribution of the elements in the Earth, it is necessary to go back to the time of Earth formation approximately 4.5 billion years ago. It is generally believed that the Earth and the other planets in the solar system formed by agglomeration of smaller fragments of solid material orbiting around the Sun. This material had precipitated from a cooling hot gas cloud (the solar nebula), with the most refractory materials condensing out first, the most volatile last. The distribution of elements in the solar system in this early phase thus had much to do with volatility, and the solid material that aggregated to form the planets was a mix of volatile and nonvolatile materials. In a general sense, this is reflected in the overall compositions of the planets: the volatile-rich material accumulated in the cooler regions away from the Sun and formed the gaseous planets such as Jupiter and Saturn; the less volatile material dominated in the so-called terrestrial planets of the inner solar system, for example, Earth and Mars. The material which accreted to form the Earth was probably similar to that found today in meteorites, a mixture of iron-nickel metal and silicate and oxide minerals. See COSMOCHEMISTRY; ELEMENTS, COSMIC ABUNDANCE OF; SOLAR SYSTEM.

Although the Earth may have been an approximately homogeneous mixture of accreted materials at the time of its formation, it is now made of many chemically distinct parts. At the fundamental level, these are the core, the mantle, and the crust. While chemical fractionation in the solar nebula depended upon volatility, chemical differentiation within the Earth took place by the separation of molten material from unmelted residue under the influence of gravity. Because large amounts of energy were released from accreting fragments, the early Earth was very hot, and during the accretion stage itself, temperatures in some parts exceeded the melting point of iron metal. Pools of dense molten iron, with dissolved nickel and other elements, aggregated and sank through the Earth under gravity to form the core, leaving behind a mantle of silicate and oxide minerals. The present core constitutes about 32.4% of the Earth's mass. The distinct parts of the Earth possess unique overall compositions (Table 1).

Dynamics of mantle. Many geological phenomena on the surface of the Earth are caused by the convection of rocks within the mantle. The dynamic activity of the mantle arises because temperature differences of only a few hundred degrees Celsius reduce the density of large volumes of rock in the mantle, causing them to rise toward the surface in the form of plumes. When the pressure on the heads of such rising plumes of hot rocks decreases, the melting temperature of the minerals is reduced, and melting occurs even though the temperature of the rocks has not increased.

The silicate and oxide minerals of rocks in the mantle melt over a considerable temperature range. The resulting liquid, called magma, is less dense than the unmelted residue, thereby allowing the magma to rise to the surface, where it is erupted to form volcanoes. The Hawaiian Islands and most of the island chains in the Pacific Ocean are examples of magma formation by decompression melting in mantle plumes. See SILICATE MINERALS.

This is the process by which the crust has formed over geologic time, and it still continues as is evident from the existence of active volcanoes. The crust constitutes only 0.4% of the Earth's mass; the mantle,

TABLE 1. Element distribution among some of the major subdivisions of the Earth*

Element	Continental crust	Oceanic crust	Upper mantle	Core [†]
Oxygen	45.3	43.6	44.2	—
Silicon	26.7	23.1	21.0	—
Aluminium	8.39	8.47	1.75	—
Iron	7.04	8.16	6.22	85.5
Calcium	5.27	8.08	1.86	—
Magnesium	3.19	4.64	24.0	—
Sodium	2.29	2.08	0.25	—
Potassium	0.91	0.13	0.02	—
Titanium	0.68	1.12	0.11	—
Nickel	0.011	0.014	0.20	5.5
Sulfur	NA	NA	NA	9.0

*Estimates of element abundances are in percent by weight and are arranged in order of decreasing abundance in the continental crust. Sulfur contents are not well known and are designated "not applicable."

[†]The estimate for the core is just one of several models. Others substitute light elements such as oxygen, carbon, or silicon for some or most of the sulfur shown here.

1 H																	2 He
3 Li	4 Be											5 B	6 C	7 N	8 O	9 F	10 Ne
11 Na	12 Mg											13 Al	14 Si	15 P	16 S	17 Cl	18 Ar
19 K	20 Ca	21 Sc	22 Ti	23 V	24 Cr	25 Mn	26 Fe	27 Co	28 Ni	29 Cu	30 Zn	31 Ga	32 Ge	33 As	34 Se	35 Br	36 Kr
37 Rb	38 Sr	39 Y	40 Zr	41 Nb	42 Mo	43	44 Ru	45 Rh	46 Pd	47 Ag	48 Cd	49 In	50 Sn	51 Sb	52 Te	53 I	54 Xe
55 Cs	56 Ba	57–71 La–Lu	72 Hf	73 Ta	74 W	75 Re	76 Os	77 Ir	78 Pt	79 Au	80 Hg	81 Tl	82 Pb	83 Bi	84	85	86
87	88	89	90 Th	91	92 U												
rare earths	57 La	58 Ce	59 Pr	60 Nd	61	62 Sm	63 Eu	64 Gd	65 Tb	66 Dy	67 Ho	68 Er	69 Tm	70 Yb	71 Lu		

Key:

□ lithophile ■ atmophile □ chalcophile ■ siderophile

Geochemical classification of the elements in relation to the periodic table. The short-lived radioactive elements (43, 61, 84, 85, 86, 87, 88, 89, 91) have been omitted.

from which it and the core have separated, makes up 67.2%. See EARTH CRUST.

Large-scale distribution. The large-scale distribution of the elements in the Earth depends on the affinity of each element for specific compounds or phases. Those elements that alloy easily with iron, for example, are mostly sequestered in the Earth's core; those which form oxides and silicate minerals tend to be concentrated in the Earth's crust and mantle. Although many elements display multiple characteristics depending on the chemical environment, a classification according to geochemical affinity is nevertheless useful (see *illus.*). The categories in this classification include atmophile (elements that are gases and concentrate in the atmosphere), lithophile (elements that form silicates or oxides and are concentrated in the minerals of the Earth's crust), siderophile (elements that alloy easily with iron and are concentrated in the core), and chalcophile (elements such as copper which commonly form sulfide minerals if sufficient sulfur is available).

Although both the mantle and the crust are composed largely of lithophile elements, their specific compositions are quite different (Table 1). The differences result from the way in which the crust originates by melting of the Earth's interior. Partitioning between mantle minerals and the melts produced from them determines the crustal inventories of most elements. The major mineral phases in the

portion of the mantle that melts to form crustal rocks are olivine $[(\text{Mg,Fe})_2\text{SiO}_4]$, pyroxene $[(\text{Ca,Mg})\text{SiO}_3]$, spinel $(\text{MgAl}_2\text{O}_4)$, and garnet $(\text{Mg}_3\text{Al}_2\text{Si}_3\text{O}_{12})$. These are idealized chemical formulas; however, virtually every element in the periodic table occurs in the mantle, many dissolved in very small amounts in these mantle minerals, and some concentrated in other very rare minerals which may exist in some parts of the mantle. When melting occurs, those elements that do not fit easily into the structure of the major mantle minerals tend to be partitioned strongly into the melt. These are the elements that are enriched in the crust. In practice, they typically have very different ionic radii or valences than magnesium, iron, calcium, aluminium, or silicon, which are the major constituents of the mantle minerals. An example is the alkali element rubidium, which has a comparatively large ionic radius and therefore cannot fit into the structural sites normally occupied by magnesium, iron, or calcium in the mantle minerals. Estimates suggest that between one-third and two-thirds of the Earth's total inventory of rubidium is contained in the crust which, as mentioned, constitutes only 0.4% of the Earth. Elements such as rubidium that are strongly partitioned into melts are commonly termed incompatible elements by geochemists. The siderophile elements, as already noted, are concentrated in the core, and their abundances in the crust are very low regardless of their valence or ionic radius.

Earth crust. Although geochemists have a good general knowledge of the overall distribution of elements in the core and mantle, much more detailed information is available about the chemical composition of the crust, which is accessible. The crust is actually composed of two major parts with quite different compositions, thickness, and average age: the continental crust and the oceanic crust (Table 1).

The oceanic crust is continually renewed at the mid-ocean ridges and destroyed at subduction zones. It is a direct product of large-scale melting of the upper mantle, and its composition reflects this fact. Compared to the continental crust, it has higher concentrations of magnesium, iron, and calcium and lower concentrations of silicon, potassium, and the highly incompatible elements. Its overall composition is basaltic. *See* BASALT; MID-OCEANIC RIDGE; SUBDUCTION ZONES.

The continental crust is older than the oceanic crust and is more chemically fractionated relative to the mantle. There is good evidence that the continental crust itself is chemically zoned, with the upper parts being more highly enriched in the incompatible elements than the lower parts.

Although all elements are present, the crust is made almost entirely of just nine chemical elements: oxygen, silicon, aluminum, iron, magnesium, calcium, sodium, potassium, and titanium. Oxygen and silicon are by far the most abundant. The most common minerals in the crust are those of the silicate family, in which the basic building block is a silicon atom surrounded by four oxygen atoms in the form of a tetrahedron. The crust is essentially a framework of oxygen atoms bound together by the common cations. *See* CHEMICAL BONDING; SILICATE MINERALS.

A variety of processes act to make the crust chemically heterogeneous on many scales. Many of these processes involve liquid water. Running water physically sorts particles depending on size and density, which are ultimately related to chemical composition. It is also a superb solvent, carrying many elements in solution under different conditions of temperature and pressure, and depositing them when these conditions change. Processes involving water account for many ore deposits, in which extreme concentrations of some elements occur relative to their average abundance in the crust. One example is the circulating hydrothermal solutions in volcani-

cally active parts of the crust, which can leach metals from their normally dispersed state in large volumes of volcanic rocks and deposit them in concentrated zones as the solutions cool and encounter different rock types. Another example is the action of weathering in tropical regions with high rainfall, which can leach away all but the least soluble components from large volumes of rock, leaving behind mineral deposits rich in aluminum or, depending on the original composition of the rocks being weathered, metals such as iron and nickel. *See* ELEMENT (CHEMISTRY); GEOCHEMISTRY; ORE AND MINERAL DEPOSITS; WATER; WEATHERING PROCESSES.

Rivers. The chemical composition of water in rivers depends on many factors, such as the geology of the drainage basin, the climate and its seasonal variation, the amount and composition of the suspended sediment, the acidity of the water, and the amount of biomass that the water contains. In addition, the chemical composition of rivers is affected by the discharge of industrial, municipal, and agricultural wastewater, all of which vary seasonally and enter rivers at different sites along the channel. Therefore, it is difficult to make an accurate determination of the chemical composition of the water of a particular river. However, when the available data are averaged by continent, eight different constituents are found to be dominant in all rivers: calcium, magnesium, sodium, potassium, chloride (Cl^-), sulfate (SO_4^{2-}), bicarbonate (HCO_3^-), and silica (SiO_2) [Table 2]. The concentrations of other elements are comparatively low and vary widely among the rivers of the world depending on the factors identified above. *See* RIVER; SILICATE MINERALS.

Seawater. The chemical composition of seawater depends on the amount of each element that enters the ocean each year and the amount that is removed from it. Chemical elements enter the oceans by the discharge of river water and by deposition from the atmosphere either dissolved in rain or as dry particulate fallout. Certain chemical elements (sodium and chloride ions) remain in the oceans for long periods of time and therefore have high concentrations in seawater. Other elements are removed from seawater because they form insoluble compounds (calcium and magnesium), because they are adsorbed on the surfaces of sediment particles (beryllium and potassium), or because they are nutrients that are

TABLE 2. Average chemical compositions of rivers in milligrams per liter*

Continent	Ions in solution							SiO_2^\dagger
	Ca^{2+}	Mg^{2+}	Na^+	K^+	Cl^-	SO_4^{2-}	HCO_3^-	
N. America	21.2	4.9	8.4	1.5	9.2	18.0	72.3	7.2
S. America	6.3	1.4	3.3	1.0	4.1	3.8	24.4	10.3
Europe	31.7	6.7	16.5	1.8	20.0	35.5	86.0	6.8
Africa	5.7	2.2	4.4	1.4	4.1	4.2	26.9	12.0
Asia	17.8	4.6	8.7	1.7	10.0	13.3	67.1	11.0
Oceania	15.2	3.8	7.6	1.1	6.8	7.7	65.6	16.3
World	14.7	3.7	7.2	1.4	8.3	11.5	53.0	10.4

* After E. K. Berner and R. A. Berner, *The Global Water Cycle*, Prentice-Hall, 1987.

† Silicic acid (H_4SiO_4) is reported as SiO_2 .

TABLE 3. Concentrations of major ions in seawater of salinity 35‰ in milligrams per kilogram*

Ion	Concentrations
Ca ²⁺	412
Mg ²⁺	1,290
Na ⁺	10,770
K ⁺	399
Cl ⁻	19,354
SO ₄ ²⁻	2,712
HCO ₃ ⁻	120
SiO ₂ [†]	6

*After E. K. Berner and R. A. Berner, *The Global Water Cycle*, Prentice-Hall, 1987.

[†]Silicic acid (H₄SiO₄) is reported as SiO₂.

absorbed by marine plants (nitrogen and phosphorus).

The chemical composition of seawater is described in terms of the major ions whose concentrations depend only on the salt content of the water, and in terms of the other ions whose concentrations vary regionally within the oceans, or with increasing depth, or seasonally, or even on a daily basis (Table 3). The concentrations of the major ions in seawater are significantly higher than in average river water. Silica is the only exception because it is actively removed from seawater by certain sponges and diatoms which have siliceous skeletons. See SEAWATER.

Organic matter. The diversification of the chemical compositions of the rocks of the continental crust and of water that interacts with them is evident also in the composition of organic matter. The tissues of plants and animals are composed primarily of proteins, carbohydrates, lipids, and lignin, all of which contain carbon, oxygen, hydrogen, and minor amounts of other elements such as nitrogen, phosphorus, and sulfur. Carbohydrates and lignin are dominant in plant matter, whereas proteins, carbohydrates, and lipids are abundant in animal tissues. Plants derive certain nutrient elements from soil and from the atmosphere. The essential nutri-

TABLE 4. Concentrations of essential nutrient elements in plants expressed in weight percent dry tissue

Element	Concentration
Carbon	45
Oxygen	45
Hydrogen	6
Nitrogen	1.5
Potassium	1.0
Calcium	0.5
Phosphorus	0.2
Magnesium	0.2
Sulfur	0.1
Chlorine	0.01
Iron	0.01
Manganese	0.005
Zinc	0.002
Boron	0.002
Copper	0.0006
Molybdenum	0.00001

ent elements in plants (Table 4) are propagated to herbivorous and carnivorous animals with only minor shifts in composition. See PLANT MINERAL NUTRITION.

Some chemical elements are toxic to plants and inhibit their growth even when present in soil or water at low concentrations. The toxic elements include beryllium, arsenic, antimony, selenium, mercury, nickel, copper, cobalt, lead, zinc, silver, cadmium, and tin. These elements may be present in the soils either because they occur naturally or because of environmental contamination by industrial or municipal waste products. Certain toxic elements, such as selenium, accumulate in plants tolerant of the element and enter the food chain when these plants are consumed by herbivorous animals. See SELENIUM; ELEMENTS, GEOCHEMICAL DISTRIBUTION OF.

G. Faure; J. D. MacDougall

Bibliography. E. K. Berner and R. A. Berner, *The Global Water Cycle*, Prentice-Hall, 1987; G. Faure, *Principles and Applications of Geochemistry*, 1998; A. E. Ringwood, *Composition and Petrology of the Earth's Mantle*, 1975; B. J. Skinner and P. B. Barton, Jr., Genesis of mineral deposits, *Annu. Rev. Earth Planet. Sci.*, 1:183-211, 1973; S. R. Taylor and S. M. McLennan, *The Continental Crust: Its Composition and Evolution*, 1985.

Elephant

The common name for three living species of mammals in the family Elephantidae, one of several families included in the order Proboscidea. The remaining families contain extinct animals, such as the mammoth. Two of the living species (*Loxodonta africana* and *L. cyclotis*) are indigenous to Africa, and the other (*Elephas maximus*) ranges throughout Southeast Asia.

These animals are terrestrial and entirely herbivorous. They have 26 teeth and the dental formula is I 1/0 M 0/0 Pm 3/3 M 3/3. The nostrils and upper lip are elongated into a proboscis, the trunk, which is a powerful and sensitive organ specialized to form a prehensile, food-gathering structure. The upper incisors protrude from the mouth on either side of the trunk as tusks, which continually grow throughout the life of the animal. The head is massive, and the leathery ears, which are large, especially in the African species, are used as cooling organs as well as for hearing. The hard, thick skin of the elephant is sparsely covered with hair and serves as an insulator, much as does hair or fat in other animals. The eyes are tiny but vision is keen. The tail is short, large columnar legs support the massive body, and the feet may measure 20 in. (50 cm) across.

Elephants live and travel in herds which were originally composed of several hundred individuals, but the usual number is now around 20 animals, with a mature bull, a number of cows and calves, and some younger bulls. Elephants are at ease in water, and they bathe and roll in the mud as a protective



Loxodonta africana, African bush elephant. (Photograph by Robert Thomas and Margaret Orr; © 2001 California Academy of Sciences)

measure since, despite the thickness of the skin, they are sensitive to intense sun and insects. These animals are constantly on the move for pasturage, and when they travel they walk in single file with an experienced female in the lead and the males in the rear. During the heat of the day they take refuge in the cooler forested areas; however, there is insufficient herbage in these locations, and later in the day they move out to the tall-grass areas. Bulls reach sexual maturity at the age of 15, while cows mature earlier. The gestation period averages 20–22 months with a single calf being born. Twins are rare. The newborn, which may be 3 ft (1 m) tall and weigh 200 lb (90 kg), grows rapidly, and one elephant was reported to weigh 650 lb (293 kg) at 15 months. The female that is about to give birth is attended by other females, which in effect act as midwives and remain with her until the calf can follow its mother. Lactation continues for months, and the young uses its mouth and not the trunk for suckling.

The African elephant genus *Loxodonta* comprises two species, the larger African bush elephant (*L. africana*) (see **illustration**) and the smaller African forest elephant (*L. cyclotis*). The African bush elephant is the largest living land mammal. An adult bull may weigh as much as 6 tons (5.4 metric tons) and stand over 11 ft (3.3 m) high, while the cow usually weighs about 4 tons (3.6 metric tons). Both the bull and the cow have tusks, although they are small in the cow. The trunk appears annulated, the back is concave, and the tip of the trunk has two triangular lips. They have a life-span of 60–70 years and have not been domesticated.

The Indian elephant is smaller and rarely reaches 10 ft (3 m) in height. It weighs about 4 tons (3.6 metric tons), and the cow is tuskless. The trunk appears to be smooth, the ears are much smaller than those of the African elephant, the back is convex, and between the ear and the eye there is an opening to a gland, the musth gland, which may be associated with sexual activity. This animal has been domesticated and is a valuable beast of burden in much of the Far East. See MAMMALIA; PROBOSCIDEA.

Charles B. Curtin

Bibliography. R. M. Nowak, *Walker's Mammals of the World*, Johns Hopkins University Press, 1999.

Eleutherozoa

One of the two subphyla into which the phylum Echinodermata had been customarily divided. The Eleutherozoa are now best considered as comprising at least two distinct subphyla: (1) the Echinozoa, spherical-bodied forms with meridional symmetry; and (2) the Asterozoa, star-shaped forms with radially divergent axes of symmetry. The Echinozoa would seem to have arisen from ancestors similar to cystoids, whereas the Asterozoa seem to be derived from crinoidlike forms. See CRINOIDEA; ECHINOZOA; PELMATOZOA.

At one time the Eleutherozoa were thought to comprise the free-living echinoderms, namely, the sea stars, the sea urchins, and sea cucumbers, in which the mouth is directed downward and the anus (if present) is usually placed on the upper surface. The Eleutherozoa were contrasted with the so-called Pelmatozoa, in which the body was attached to the substrate and the mouth and anus opened together on the upper surface, therefore making the gut U-shaped. However, further research showed that both groups are polyphyletic and that the eleutherozoan and pelmatozoan habit (with consequential morphological change) has been adopted independently more than once by various groups of echinoderms. Therefore the original definitions of the subphyla are abandoned because they described unnatural groupings of unrelated animals. See ECHINODERMATA.

Howard B. Fell

Bibliography. H. B. Fell, Phylogeny of sea stars, *Phil. Trans. Roy. Soc. London*, ser. B, 735, 246:381–435, 1963; H. B. Fell, The evolution of the echinoderms, *Smithson. Inst. Annu. Rep. 1962*, pp. 457–490, 1963; S. P. Parker (ed.), *Synopsis and Classification of Living Organisms*, 2 vols., 1982.

Elevating machines

Materials-handling machines that lift and lower a load along a fixed vertical path of travel with intermittent motion. In contrast to hoisting machines, elevating machines support their loads instead of carrying them suspended, and the path they travel is both fixed and vertical. They differ from vertical conveyors in operating with intermittent rather than continuous motion. Industrial lifts, stackers, and freight elevators are the principal classes of elevating machines.

Industrial lifts. A wide range of mechanically, hydraulically, and electrically powered machines are classified simply as lifts (**Fig. 1**). They are adapted to such diverse operations as die handling and feeding sheets, bar stock, or lumber. In some locations with differences in floor level between adjacent buildings, lifts take the form of broad platforms to serve as floor levelers to obviate the need for ramps. They are also used to raise and lower loads between the ground and the beds of carriers when no loading platform exists. Lifting tail gates attached to the rear of trucks are similarly used for loading or

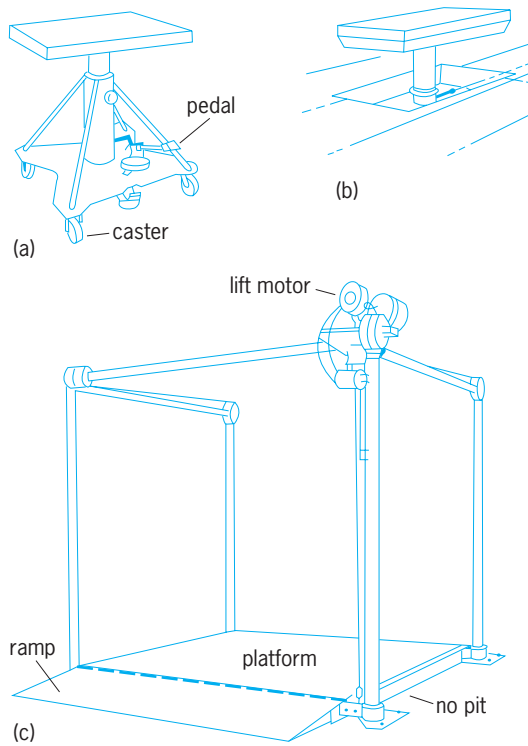


Fig. 1. Examples of industrial lifts. (a) Hydraulic elevating work table. (b) Hydraulic lift floor leveler. (c) Motor-driven floor leveler.

unloading merchandise on sidewalks or roads and at points where the lack of a raised dock would make loading or unloading difficult. These units are usually driven by battery-operated motors on a power takeoff from the drive transmission of the vehicle. Adjustable loading ramps are necessary because heights of truck and trailer beds vary. Advances in mechanized loading and unloading of vehicles have made necessary sturdier, more efficient dock boards or bridge plates, mechanically or hydraulically operated.

Stackers. Tiering machines and portable elevators used for stacking merchandise are basically portable vertical frames that support and guide the carriage,

to which is attached a platform, pair of forks, or other suitable lifting device (Fig. 2). The operation of these units varies in the sense that the carriage can be raised and lowered by hand, by an electrically driven winch, or by a hydraulic cylinder, which actuates the system of chains or cables, operative by hand lever, pedal, or push button. Early electric motors used on stackers were plugged into adjacent power lines to receive current. This limited the flexibility of the stackers. Since the trend is to make the machines independent of this source, there are now models powered by either storage batteries or by small gasoline or gas engines. Horizontal movement is effected by casters on the bottom of the vertical frame, and can be accomplished manually, or mechanically, by using the same power source as the lifting mechanism. These casters are usually provided with floor locks bolted in position during the elevating or lowering operation.

The basic type of stackers can be varied in several particulars. Masts, which are part of the frame, can be hinged or telescopic, and the platforms can be plain, equipped with rollers, or constructed specially to handle a specific product. Some stackers have devices for tilting barrels and drums or for lifting and dumping free-flowing bulk materials. Used in conjunction with cranes, they are widely applied to the handling of materials on storage racks and die racks. Stackers have a significant place in the development of materials-handling equipment. They are the prototypes of completely powered noncounterbalanced platform and forklift trucks. See INDUSTRIAL TRUCKS.

Industrial elevators. Examples of industrial elevators range from those set up temporarily on construction jobs for moving materials and personnel between floors to permanent installations for mechanized handling in factories and warehouses. Dumbwaiters are a type of industrial elevator, having capacities up to approximately 500 lb (227 kg), with a maximum floor space of 9 ft² (0.8 m²); they carry parts, small tools, samples, and similar small objects between buildings, but are not permitted to carry people.

Oil hydraulic plunger electric elevators are designed for low-rise, light- or heavy-duty freight handling. Although they can be installed without special building alterations, they are restricted to buildings with only a few floors because of the limitations of the plunger length and design.

The most common and economical elevator employs electric motors, cables, pulleys, and counterweights (Fig. 3). Powered machines impose a severe operating condition on elevators. Elevator platforms and structures are subjected to impact loading, off-balance loading, and extra static loading. To meet these forces, in addition to load forces, freight elevator design and construction provide greatly increased ruggedness over that of passenger units.

Special-purpose freight-handling elevators are equipped with platforms or arms for carrying specific articles such as rolls of paper, barrels, or drums. Some of these elevators load and discharge

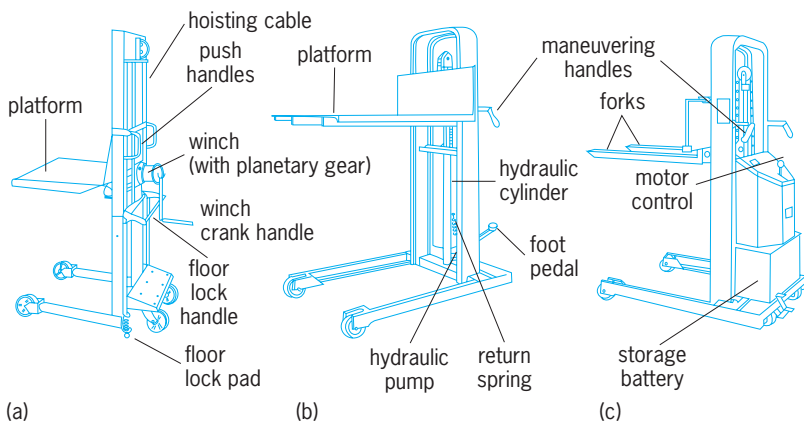


Fig. 2. Three types of electric and hydraulic stackers. (a) Hand type. (b) Hydraulic foot type. (c) Electric lift type.

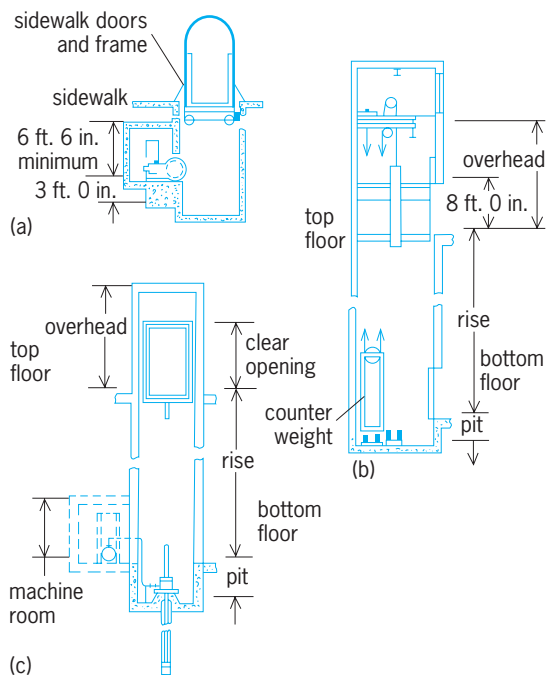


Fig. 3. Three types of industrial elevator. (a) Sidewalk elevator type. (b) Heavy-duty freight elevator type. (c) Hydraulic electric elevator type.

automatically and are arranged so that they can operate at any selected floor by means of remote control. See MATERIALS-HANDLING EQUIPMENT.

Arthur M. Perrin

Bibliography. American National Standards Institute, *Safety Code for Elevators and Escalators*, ANS A17.1-1981, ASME, 1981; American Society of Mechanical Engineers, *Safety Code for Elevators and Escalators: Handbook on A17.1*, 1997; E. A. Donoghue (ed.), *Safety Code for Elevators and Escalators: Handbook on A17.1*, 1993.

Elevator

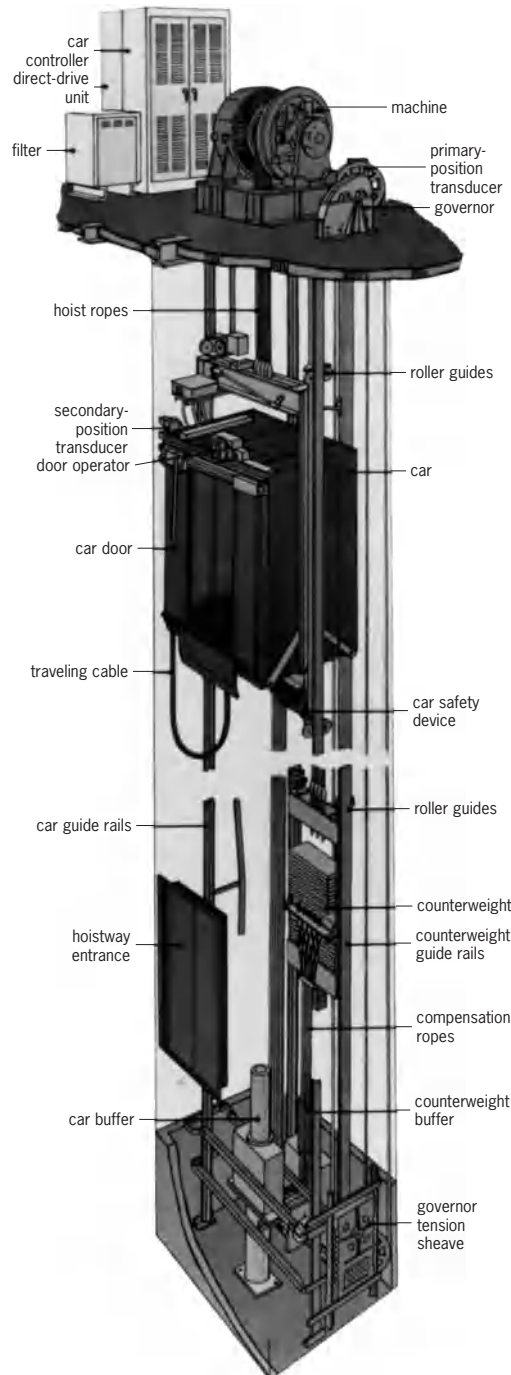
A platform or enclosure that is raised and lowered in a vertical hoistway to transport freight or people. The term elevator can also encompass all of the hoisting equipment, motor, cables, and accessories (see *illus.*). See ELEVATING MACHINES.

Operation. The closed passenger car of a modern elevator rests inside a steel frame. The car and the car frame ride up and down on steel rails in an elevator shaft or hoistway. Guide shoes or rollers on the frame keep the car in place on the rails. Most elevators also have a heavy weight, called a counterweight, attached to the other end of the steel hoisting ropes that pass over the driving machine pulley. The counterweight offsets much of the weight of the car and passengers, thereby reducing power requirements.

Control devices. The typical elevator control system is made up of a speed-sensing device known as a governor, a clamping device (safety) mounted under

each end of the car frame that grips the guide rail when tripped, a tension sheave (pulley) in the pit, and a steel rope. The governor rope makes a complete loop around the governor sheave and the tension sheave in the pit. Since the rope is fastened to and travels with the car, the governor sheave rotates at a speed directly proportional to the speed of the car.

The governor rope is connected to the safeties through mechanical linkages and lift rods. If the hoist ropes break or the car overspeeds, the governor



Typical installation of a traction elevator. (Otis Elevator Co.)

trips, releasing a clutching device that grips the governor rope. As the car continues its descent, the pull of the rope on the operating lever sets the safety. The safety, in turn, applies sufficient force against the guide rail to bring the car to a controlled stop with frictional force. *See* GOVERNOR.

Control devices are also built into the door and its control circuit. When the doors open, control circuits prevent the car from moving away from the landing, but permit releveling if the car moves as passengers enter or leave the elevator (load changes). Either a safety shoe or electronic detector, mounted on the doors, prevents the doors from continuing to close if a passenger or object is in the doorway. The safety shoe touches the person or object, retracts, and causes the door to reopen. Light-ray devices, often used with safety shoes, cause reversal of door movement whenever the light ray is broken by a passenger entering or leaving the car.

Elevator controller. Sophisticated elevator controllers began to replace human elevator operators in the 1950s. In modern elevators, microprocessor computer systems control elevator position, direction of travel, speed, door operation, passenger waiting time, flight time, energy consumption, and system diagnostics.

Systems are being developed that learn from past building traffic patterns and predict future patterns, assigning elevator cars to destinations in advance of actual demand and reducing passenger waiting times. In addition, in the future elevators will be equipped with laser devices that scan a floor for waiting passengers and signal the elevator accordingly to stop or continue. Modern elevators include Braille buttons and voice announcements of the floors to help the sight-impaired.

Types. There are three major types of elevators: gearless traction, geared traction, and hydraulic.

Gearless traction elevators are used in high-rise buildings over 10–12 stories, and travel at speeds from 400 to 2000 ft/min (120 to 610 m/min). They use large, slow-speed electric motors directly connected to a large, grooved drive sheave (pulley). The “hoist ropes” (steel cables) are attached at one end to the top of the elevator car, pass over the drive sheave and are attached at the other end to a counterweight that slides up and down the shaftway on its own guide rails. The full weight of the car and about half of its passenger load is balanced by the counterweight, which goes down as the car moves up. Thus, the electric motor does not have to lift the full weight of the car.

The principal difference between geared and gearless traction elevators is speed. A geared elevator usually travels at speeds from 25 to 450 ft/min (8 to 140 m/min), carries loads up to 30,000 lb (13,500 kg) or more, and uses a high-speed motor to drive the hoisting sheave through a gear reduction unit.

Hydraulic elevators are used extensively in low-rise buildings, usually up to five stories. With speeds rarely exceeding 150 ft/min (46 m/min), the hydraulic elevator does not need overhead hoisting machinery. The elevator is mounted on a piston inside

a cylinder that extends into the ground to a depth equal to the height the elevator will rise. Relatively simple in design, the system uses an electric pump to force oil into the cylinder to give the elevator a controlled ascent. Electrically controlled valves release the oil for a controlled descent. Another form of hydraulic elevator is the so-called holeless model. A plunger that slides up and down on the side of the elevator is used, so that no hole is required beneath the standard shaftway space. Cynthia Di Tallo

Bibliography. A. Goetz (ed.), *Up, Down, Across: Elevators, Escalators, and Moving Sidewalks*, 2003; G. R. Strakosch (ed.), *The Vertical Transportation Handbook*, 3d ed., 1998.

Elevator (aircraft)

The hinged rear portion of the longitudinal stabilizing surface or tail plane of an aircraft used to obtain longitudinal- or pitch-control moments. The angular setting of the elevator is controlled by the human or automatic pilot through the flight-control system. A typical elevator control surface is shown in **Fig. 1**. Both leading-edge and horn-type aerodynamic balances and trailing-edge tabs are illustrated. These features reduce or eliminate the hinge moments required to deflect the elevator during flight. *See* FLIGHT CONTROLS.

General principles. The operating principles of elevator control surfaces are typical of all trailing-edge hinged control devices. Deflection of the elevator changes the camber of the entire surface. With the trailing edge down, high local flow velocities are obtained on the upper airfoil surface and low relative flow velocities are produced on the lower surface (**Fig. 2**). By Bernoulli's law, reduced pressure results on the upper surface, with increased pressures on the lower surface, creating lift. *See* BERNOULLI'S THEOREM.

The variation of elevator surface control effectiveness with percentage of total chord is shown in **Fig. 3**. Narrow-chord elevators are more effective than wide-chord elevator surfaces. Operating hinge moments vary as the square of the chord, and are thus greatly reduced on narrow-chord surfaces.

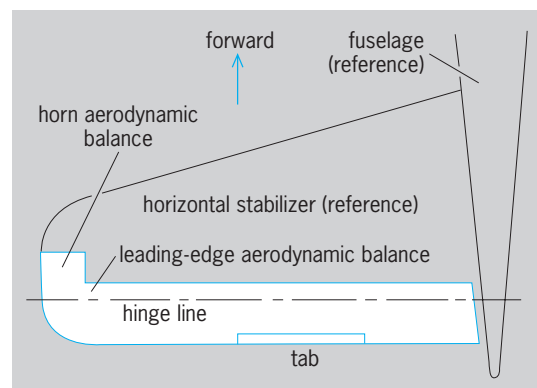


Fig. 1. Elevator control surface (left-hand side).

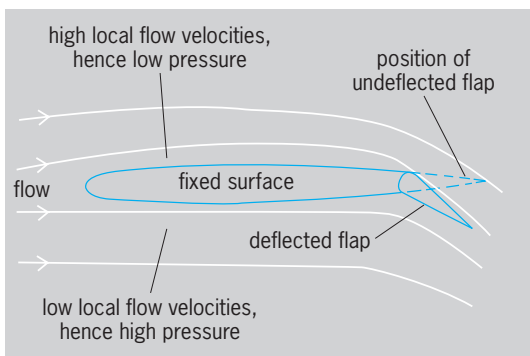


Fig. 2. Principle of operation of trailing-edge flap.

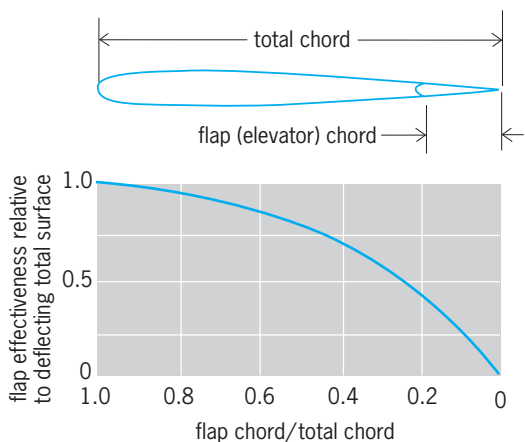


Fig. 3. Variation of trailing-edge flap effectiveness with flap chord to total chord ratio.

Flutter prevention. Elevator flutter is a divergent oscillation involving one or more degrees of freedom, such as rotation about the hinge line and flapping of the main surface. Hydraulic dampers, mass balancing, and structural stiffness are employed to prevent elevator flutter. See AEROELASTICITY; FLUTTER (AERONAUTICS).

Flight maneuvers. The elevator is used to perform pitching maneuvers, or maneuvers in which the aircraft's plane of symmetry is not disturbed. The maneuvers include airspeed adjustments and acceleration normal to the flight path (pull-ups or push-downs). The elevator also serves to adjust the aircraft's attitude with respect to the ground for takeoff and landing.

Special applications. All-moving horizontal stabilizers replace elevator control surfaces on supersonic aircraft and missiles to avoid large losses in effectiveness. The elevator may be geared to move at a fixed ratio to the deflections of an all-moving stabilizer. On aircraft without horizontal stabilizers (tailless aircraft), the elevator and aileron surfaces may be combined in surfaces that operate together as elevators and differentially as ailerons, for example, elevons. See ELEVON; STABILIZER (AIRCRAFT).

Malcolm J. Abzug

Bibliography. J. D. Anderson, *Introduction to Flight*, 4th ed., 1999; P. P. Wegener, *What Makes Airplanes Fly?*, 2d ed., 1996.

Elevon

A movable surface at the trailing edge of a tailless airplane that provides pitch and roll control. Elevons hinged on each side of the rear wing surface are shown in Fig. 1. They nose the airplane up or down, and roll one wing up and the other down. The term elevon is derived from elevator and aileron, and, in effect, elevons provide the same control as conventional elevators and ailerons. An example of an airplane employing elevons is the orbiter vehicle for the space shuttle. See AILERON; ELEVATOR (AIRCRAFT); SPACE SHUTTLE; TAIL ASSEMBLY.

The action of an elevon is illustrated in Fig. 2. A symmetrical airfoil is placed at an angle to the airflow, resulting in a lift force that acts approximately at a point located a quarter of the chord from the nose (Fig. 2a). The chord is defined as the distance from the nose, or leading edge, of the airfoil to the trailing edge. In Fig. 2b a trailing-edge flap on the same airfoil is deflected downward. This results not only in an increase in the lift on the airfoil but in a nose-down pitching moment about the quarter-chord point. Deflected upward, the flap produces a nose-upward pitching moment and a decrease in lift. See AIRFOIL.

When the elevons on both sides of the wing are deflected upward, they combine to produce a nose-up pitching moment on the wing and, hence, on the tailless airplane. If the elevon on the right wing is deflected upward and the one on the left wing downward, the differential movement does not result in any pitching moment, since the nose-up and nose-down pitching moments produced by the two elevons cancel each other. However, the lift on the

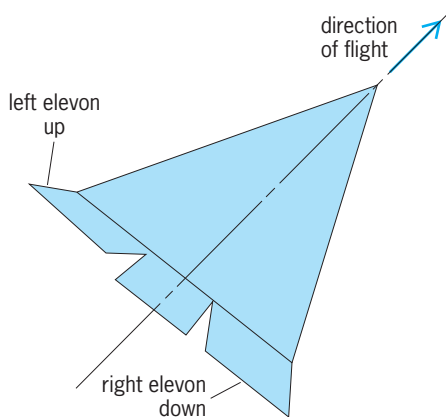


Fig. 1. Elevons on the trailing edge of a delta wing.

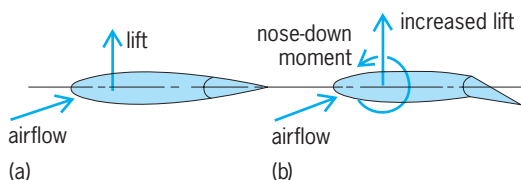


Fig. 2. Effect of flap on lift and moment of an airfoil at an angle of attack. (a) Symmetrical airfoil. (b) Symmetrical airfoil with flap deflected.

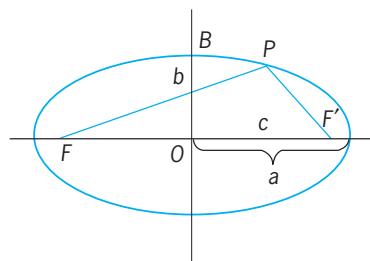
right wing decreases, and that on the left wing increases; the combined effect produces a moment on the airplane tending to roll it to the right. Deflecting the elevons in the opposite directions causes the wing to roll to the left, and deflecting both elevons downward produces a nose-down moment.

Pitching the nose up or down by elevon deflection changes the total lift on the wing, thereby accelerating the airplane vertically. Thus if the wing is in level flight, a pilot wishing to climb moves both elevons upward. Momentarily, the decreased lift causes the wing to drop, but the pitching moment quickly increases the angle of the wing relative to the airflow, the angle of attack; the result is a net increase of wing lift and an upward inclination of the flight path. See AIRPLANE; WING. Barnes W. McCormick, Jr.

Bibliography. J. D. Anderson, *Introduction to Flight*, 4th ed., 1999; B. Etkin and L. D. Reid, *Dynamics of Flight Stability and Control*, 3d ed., 1995.

Ellipse

A member of the class of curves that are intersections of a plane with a cone of revolution. The ellipse is obtained when the plane cuts all the elements of one nappe, and does not go through the apex. In the **illustration**, denote the distance between two points, F, F' of a plane by $2c, c > 0$, and let $2a$ be a constant, with $a > c$. The ellipse with foci F and F' and major axis $2a$ is the locus of points P of the plane such that $PF + PF' = 2a$, where PF denotes the distance of P and F . This suggests the following construction of an ellipse. Put pins at F and F' , and slip over them a loop of thread of length $2a + 2c$, pulling the thread taut with a pencil. If the pencil is moved, keeping the thread taut, its point traces an ellipse. Another way to construct an ellipse is to drill a hole in a stick (at any point other than the midpoint) and move the stick so that its ends slide along two mutually perpendicular lines. The point of a pencil inserted in the hole will trace an ellipse. Limiting forms of the ellipse are (1) a circle, as the two foci approach coincidence, and (2) the segment FF' , as c approaches a . If a circle is projected orthogonally on a plane not parallel to the plane of the circle, an ellipse is obtained, and every ellipse may be so obtained. Lines joining the foci to a point P of an ellipse make equal angles with the tangent to the ellipse at P , and consequently light or sound that emanates from one focus is reflected to the other focus. This property is used in construction



An ellipse; see text for explanation of symbols.

of “whispering galleries.” See CONIC SECTION.

The midpoint of F, F' is the center O of the ellipse, and the chord through O perpendicular to the major axis is the minor axis, whose length is denoted by $2b$. If B is a point in which the minor axis intersects the ellipse, then $BF = BF' = a$, and so $c^2 = a^2 - b^2$. The ratio $c/a = \epsilon < 1$ is the eccentricity of the ellipse. If the half chords perpendicular to the major axis are multiplied by a/b , their extremities will lie on the circle whose diameter is the major axis. Hence, if Δ denotes the area of the ellipse, $\Delta \cdot a/b = \pi a^2$; that is, $\Delta = \pi ab$. The determination of the length L of an ellipse leads to the integral

$$\int_0^{\pi/2} [1 - \epsilon^2 \sin^2 \phi]^{1/2} d\phi$$

which is the complete elliptic integral of the second kind. It follows that

$$L = 2\pi a \left\{ 1 - (1/2)^2 \frac{\epsilon^2}{1} - \left(\frac{1 \cdot 3}{2 \cdot 4} \right)^2 \frac{\epsilon^4}{3} - \left(\frac{1 \cdot 3 \cdot 5}{2 \cdot 4 \cdot 6} \right)^2 \frac{\epsilon^6}{5} - \dots \right\}$$

The volume bounded by the surface is obtained by revolving the ellipse with major axis $2a$ and minor axis $2b$ about its major axis $4\pi ab^2/3$. The area of the surface is

$$2\pi b^2 + 2\pi ab \frac{\sin^{-1} \epsilon}{\epsilon}$$

See ANALYTIC GEOMETRY; ELLIPTIC FUNCTION AND INTEGRAL. Leonard M. Blumenthal

Ellipsometry

A technique for determining the properties of a material from the characteristics of light reflected from its surface. The materials studied include thin films, semiconductors, metals, and liquids.

Principles. When an electromagnetic wave passes through a medium, it causes the electrons associated with the atoms of the medium to oscillate at the frequency of the wave. As a result, the wave is slowed so that its velocity v in the medium is less than its velocity c in empty space. Another result may be a transfer of energy from the wave to the electrons, thereby causing the amplitude of the wave to decrease as it penetrates into the material. These two processes are described phenomenologically by the complex refractive index, whose real part equals c/v and whose imaginary part equals $4\pi/\lambda$, where λ is the wavelength in vacuum. See ABSORPTION OF ELECTROMAGNETIC RADIATION; REFRACTION OF WAVES.

When an electromagnetic wave is incident on a medium, only part of it is transmitted into the medium (**Fig. 1**). The fraction reflected depends on the complex refractive index, the angle of incidence, and the polarization state of the wave. For multilayers with different complex refractive indices, the fraction also depends on the layer thicknesses. The two basic types of polarization are parallel, designated

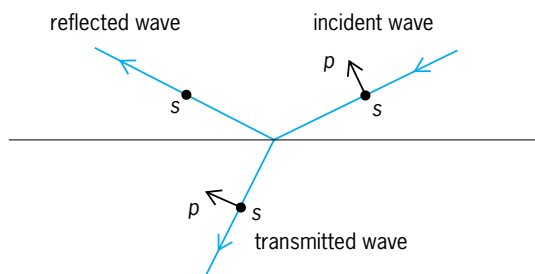


Fig. 1. Incident, transmitted, and reflected electromagnetic waves at the surface of a material. The angle of incidence shown is the Brewster angle, for which the p -polarized wave in the lower medium points along the reflection direction. The s wave is polarized perpendicular to the plane of the page.

p , and perpendicular, designated s . These terms refer to the orientation of the electric vector with respect to the plane of incidence, which is defined by the directions of the incident and reflected waves. The (intensity-independent) ratios of the amplitudes and phases of the reflected and incident p - and s -polarized electric fields are described by the complex reflectances r_p and r_s . See POLARIZATION OF WAVES; POLARIZED LIGHT; REFLECTION OF ELECTROMAGNETIC RADIATION.

The (also intensity-independent) ratio of the p - to the s -polarized component of such a wave is termed the polarization state. Simple examples include linear polarization, where the p - and s -polarized components are in phase, and circular polarization, where the p - and s -polarized amplitudes are equal but the phases differ by 90° . The geometric terms refer to the locus of the p and s (or y and x) components of the electric field when plotted in the complex plane (Fig. 2). The general polarization state is elliptical.

If an incident wave has electric-field components along both p and s , the components are independent during reflection and reassemble afterward with the changes in amplitude and phase described by r_p and r_s . That p - and s -polarized waves are affected differ-

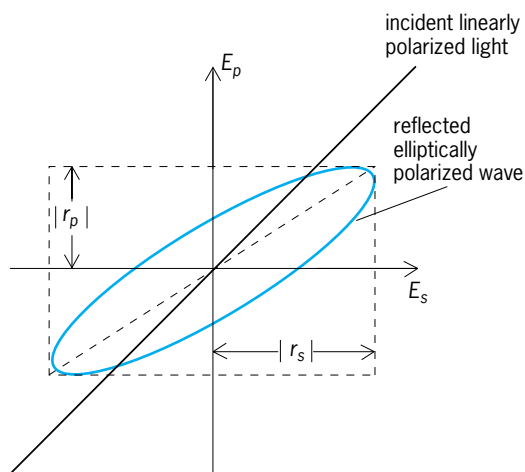


Fig. 2. Linearly polarized incident and elliptically polarized reflected waves. These represent the trajectories of the respective electric field vectors in the complex plane. E_p and E_s are the p and s components of the electric field; $|r_p|$ and $|r_s|$ are the magnitudes of the complex reflectances.

ently upon reflection is readily apparent. The incident s wave causes the electrons in the medium to oscillate at right angles to the reflection direction for any angle of incidence (Fig. 1). These oscillating electrons generate a reflected s wave with its electric field transverse to the direction of propagation. However, an incident p wave causes the electrons in the medium to oscillate partly along the reflectance direction, a motion that does not contribute to the reflected wave. Thus the magnitude of the complex reflectance r_p is always less than that of r_s . In fact, for transparent materials there exists a particular angle of incidence, called Brewster's angle, where the oscillation is entirely parallel to the reflectance direction and the p -polarized component is suppressed completely.

Since the complex reflectances r_p and r_s depend on the properties of the medium, the medium can be investigated by determining its reflectance for either p - or s -polarized light, that is, by determining the ratio of reflected and incident intensities $I_{\text{refl}}/I_{\text{inc}} = R = |r|^2$. This is the objective of reflectometry. Alternatively, because r_p and r_s are different, the complex reflectance ratio $\rho = r_p/r_s$, which is equal to the ratio of reflected and incident polarization states, can also be determined. This is the objective of ellipsometry. The ratio ρ is traditionally expressed in terms of angles ψ and Δ as in the equation below.

$$\rho = \frac{r_p}{r_s} = \tan \psi e^{i\Delta}$$

Because it deals with complex, intensity-independent quantities, an ellipsometric measurement is analogous to an impedance measurement. This gives ellipsometry certain advantages relative to reflectometry, such as higher accuracy and higher information content in a single measurement. A standard experimental approach, now used almost exclusively in spectroscopic applications, is to determine ρ by establishing a known state of polarization for the incident beam, for example, by passing it through a fixed linear polarizer, then determining the polarization state of the beam after reflection by passing it through a rotating linear polarizer, called an analyzer. The rotating analyzer essentially unrolls the polarization ellipse, allowing the azimuth angle of its major axis and its minor-major axis ratio to be determined by the phase and amplitude of the alternating-current component of the detected intensity. See ALTERNATING-CURRENT CIRCUIT THEORY; ELECTRICAL IMPEDANCE.

Applications. The primary application of ellipsometry is materials analysis, particularly the nondestructive analysis of thin films in semiconductor technology. Developing applications include the real-time monitoring and control of dynamic processes such as material deposition and etching. Spectroellipsometry, where the complex refractive index is measured and analyzed as a function of wavelength, has almost exclusively replaced reflectometry in materials analysis. Deposition and etching involve kinetic ellipsometry, where single-wavelength data are monitored as a function of time. See FILM (CHEMISTRY); INTEGRATED

CIRCUITS; NONDESTRUCTIVE EVALUATION; SEMICONDUCTOR HETEROSTRUCTURES.

Both applications rely on the fact that the complex refractive index depends directly on the ease by which electrons in a material can be set in motion by the incident electric field, which in turn depends directly on the properties of the material itself. Thus the optical properties of a metal such as chromium or aluminum, where the electrons may move freely, are quite different from those of an insulator such as plastic or glass, where the electrons are highly constrained. Further differences are observed if the material is crystalline or disordered or if two or more materials are mixed to form a composite. Other types of differences are observed in layered materials where the layers are thin enough to generate optical interference from back reflections within the sample. Such composites can be analyzed quantitatively for properties such as compositions and layer thicknesses, given an accurate database of optical responses for the pure materials. See AMORPHOUS SOLID; COMPOSITE MATERIAL. David E. Aspnes

Bibliography. F. Abeles (ed.), *Proceedings of the 1st International Conference on Spectroscopic Ellipsometry*, Paris, 1993; *Thin Solid Films*, vols. 233 and 234, 1993; R. M. A. Azzam and N. M. Bashara, *Ellipsometry and Polarized Light*, 1977, paper 1986; K. Riedling, *Ellipsometry for Industrial Applications*, 1988; H. Tompkins, *A User's Guide to Ellipsometry*, 1992.

Elliptic function and integral

In a certain sense, elliptic integrals are the simplest integrals not expressible in terms of elementary functions; elliptic functions arise as the inverse functions of certain elliptic integrals.

Let R be a rational function of x and y , and set $I = \int R(x,y) dx$. I can be expressed in terms of elementary functions if y^2 is a polynomial of degree 2 or less in x . If y^2 is a polynomial of degree 3 or 4 in x , I cannot in general be expressed in terms of elementary functions and is called an elliptic integral. (If y^2 is a polynomial of degree 5 or higher, I is called a hyperelliptic integral, and if y satisfies an algebraic equation whose coefficients are polynomials in x , I is called an abelian integral.) The standard elliptic functions are analogous to trigonometric functions. Trigonometric functions may be defined as the inverse functions of certain integrals of the form I ; they satisfy differential equations, are periodic functions, and may alternatively be obtained as the "simplest" periodic functions. The standard elliptic functions are the inverse functions of certain elliptic integrals; they satisfy differential equations of order 1 and degree 2, are doubly periodic functions, and may alternatively be obtained as the "simplest" doubly periodic functions.

Applications. In geometry elliptic functions or integrals arise in determining the length of an arc of an ellipse, hyperbola, or lemniscate, the surface of an ellipsoid, geodesics on quadrics of revolution,

parametric representations of plane cubic curves or, more generally, curves of genus 1, conformal representation problems, and other problems. In analysis there are applications to differential equations (Lamé's equation, diffusion equation, and others); in number theory, in a great variety of problems. In the physical sciences elliptic functions or integrals appear in potential theory both through conformal representations and in the potential of an ellipsoid, in the theory of elastica, the pendulum, in rigid body motion, in Green's functions in heat conduction and diffusion theory, and many other problems.

Reduction of elliptic integrals. By suitable homographic substitution, $x' = (ax + b)/(cx + d)$, $ad - bc \neq 0$, the elliptic integral I can be reduced to an elliptic integral in which the polynomial y^2 appears in normal form. The two customary normal forms are: Legendre's normal form, $y^2 = (1 - x^2)(1 - k^2x^2)$ where k , the modulus, is a real or complex number, $|k| \leq 1$ and $k^2 \neq 1$, and it is usual to set $x = \sin \phi$; and the Weierstrass canonical form, $y^2 = 4x^3 - g_2x - g_3$ where g_2 and g_3 , the invariants, are real or complex numbers. I can be expressed as a linear combination of an integral of rational functions and the elliptic integrals of the first, second, and third kinds defined in Legendre's normal form, respectively, as Eqs. (1)–(3), with $\Delta(t,k) = (1 - k^2 \sin^2 t)^{1/2}$, and in the Weierstrass canonical form as expressions (4) with $y = (4x^3 - g_2x - g_3)^{1/2}$.

$$F(\phi, k) = \int_0^\phi dt / \Delta(t, k) \tag{1}$$

$$E(\phi, k) = \int_0^\phi \Delta(t, k) dt \tag{2}$$

$$\Pi(\phi, n, k) = \int_0^\phi \frac{dt}{(1 + n \sin^2 t)\Delta(t, k)} \tag{3}$$

$$\int \frac{dx}{y} \int \frac{x dx}{y} \int \frac{dx}{(x - c)y} \tag{4}$$

In Legendre's theory, $\mathbf{K} = \mathbf{K}(k) = f(\pi/2, k)$ and $\mathbf{E} = \mathbf{E}(k) = E(\pi/2, k)$ are called the complete elliptic integrals of the first and second kinds, respectively, $k' = (1 - k^2)^{1/2}$ is the complementary modulus, and $\mathbf{K}' = \mathbf{K}'(k) = f(\pi/2, k')$, $\mathbf{E}' = \mathbf{E}'(k) = E(\pi/2, k')$. Complete elliptic integrals as functions of k satisfy linear differential equations of the second order and are hypergeometric functions of k^2 . They also satisfy Legendre's relation, $\mathbf{K}\mathbf{E}' + \mathbf{K}'\mathbf{E} - \mathbf{K}\mathbf{K}' = \pi/2$ identically in k .

Periods and singularities. Elliptic integrals are many-valued functions. Any two determinations of I differ by a sum of integral multiples of certain real or complex numbers, the periods. E , F , and Π are many-valued functions of the complex variable $x = \sin \phi$. All three functions have branch points at $x = \pm 1$, $\pm k^{-1}$, and Π has branch points also at $x = \pm in^{-1/2}$. The periods of F are $4\mathbf{K}$ and $2i\mathbf{K}'$, those of E are $4\mathbf{E}$ and $2i(\mathbf{K}' - \mathbf{E}')$. Since the complete elliptic integrals are real when $0 \leq k \leq 1$, the first (second) of these periods is called the real (imaginary) period. Although E and F are many-valued functions of x , E is a

TABLE 1. Properties of jacobian elliptic functions

Function	Primitive periods	Zeros	Poles
sn u	$4\mathbf{K}, 2i\mathbf{K}'$	$2m\mathbf{K} + 2ni\mathbf{K}'$	
cn u	$4\mathbf{K}, 2\mathbf{K} + 2i\mathbf{K}'$	$(2m + 1)\mathbf{K} + 2ni\mathbf{K}'$	$2m\mathbf{K} + (2n + 1)i\mathbf{K}'$
dn u	$2\mathbf{K}, 4i\mathbf{K}'$	$(2m + 1)\mathbf{K} + (2n + 1)i\mathbf{K}'$	

single-valued function of F provided that corresponding values of E and F are obtained by integration over the same path and using the same determination of $\Delta(t, k)$.

Inversion of elliptic integrals. Jacobian elliptic functions are determined by inversion of the functional relation $u = f(\theta, k)$. It is usual to write Eqs. (5).

$$\begin{aligned} \phi &= \text{am } u = \text{am } (u, k) \\ \sin \phi &= \text{sn } u = \text{sn } (u, k) \\ \cos \phi &= \text{cn } u = \text{cn } (u, k) \\ \Delta(\phi, k) &= \text{dn } u = \text{dn } (u, k) \end{aligned} \tag{5}$$

With the additional conditions $\text{sn } 0 = 0, \text{cn } 0 = \text{dn } 0 = 1$, it turns out that $\text{sn } u, \text{cn } u, \text{dn } u$ are single-valued analytic functions of the complex variable u , and that they are doubly periodic and also regular except for simple poles. Nine additional functions are introduced by the notations $1/\text{pn } u = \text{np } u, \text{pn } u/\text{qn } u = \text{pq } u$ where p and q stand for any of the letters s, c, d . Similarly, the Weierstrass \mathcal{P} -function is introduced by writing relation, Eq. (6), in the form

$$z = \int_{\infty}^w (4t^3 - g_2t - g_3)^{-1/2} dt \tag{6}$$

$w = \mathcal{P}(z) = \mathcal{P}(z; g_2, g_3)$, and $\mathcal{P}(z)$ turns out to be single-valued, doubly periodic, and regular except for double poles.

Doubly periodic functions. The term p is called a period of a single-valued analytic function $f(z)$, regular except for isolated singularities, if $f(z + p) = f(z)$. A nonconstant periodic function is either simply periodic, when all its periods are integral multiples of a single period, or else doubly periodic, when its periods are $2m\omega + 2n\omega'$ where m and n are integers, 2ω and $2\omega'$ are primitive periods, and $\text{Im } \omega'/\omega > 0$. Two points of the z plane are congruent if they differ by a period. A parallelogram in the z plane is a period parallelogram if every point in the plane is congruent to exactly one point of the parallelogram. If no singularity or zero of $f(z)$ lies on the boundary of the period parallelogram, the parallelogram is called a cell. An elliptic function is a doubly periodic function which is regular except for poles. An elliptic function has a finite number of poles in every cell; the sum of the residues at these poles is zero, and the sum of the orders of these poles is the order of the function.

Every elliptic function of order 0 is a constant. There is no elliptic function of order 1. An elliptic function of order $r > 1$ assumes in every cell each complex value r times (counting multiplicity). The difference of two elliptic functions with the same periods, same poles, and the same principal parts at each pole is a constant. The quotient of two elliptic

functions with the same periods, poles, and zeros (including multiplicities) is a constant. All elliptic functions with a given pair of primitive periods form an algebraic field, and any two such functions are connected by an algebraic relation. Every elliptic function satisfies an algebraic differential equation of the first order. Every elliptic function possesses an algebraic addition theorem, that is, an algebraic relation connecting $f(u), f(v)$, and $f(u + v)$. Conversely, any single-valued analytic function of z which is regular except for poles and possesses an algebraic addition theorem is either a rational function of z , or a rational function of e^{az} for some a , or else an elliptic function.

The simplest nontrivial elliptic functions are those of order 2. Choice of a basic function with two simple poles in a cell leads to Jacobi's functions, choice of a function with a double pole, to Weierstrass's functions.

Jacobian elliptic functions. Write $s = \text{sn } u, c = \text{cn } u, d = \text{dn } u, s' = ds/du$, and so on. the periods zeros, and poles are given in Table 1 in which m and n are integers. These functions possess symmetry properties around the points $u = 0, \mathbf{K}, i\mathbf{K}'$, which are set out in Table 2 and on account of which it is sufficient to study the functions in the parallelogram whose vertices are $0, \mathbf{K}, \mathbf{K} + i\mathbf{K}',$ and $i\mathbf{K}'$. There is a very large number of identities for these functions. Some of the basic ones are shown in Eqs. (7). The addition theorem for s is given in Eq. (8). There are similar addition

$$\begin{aligned} s^2 + c^2 &= 1, \quad k^2s^2 + d^2 = 1, \quad d^2 - k^2c^2 = k'^2 \\ s' &= cd, \quad c' = -sd, \quad d' = -k^2sc, \\ s^2 &= (1 - s^2)(1 - k^2s^2), \text{ etc.} \end{aligned} \tag{7}$$

$$\begin{aligned} \text{sn } (iu, k) &= i \text{sc } (u, k'), \quad \text{cn } (iu, k) = \text{nc } (u, k'), \\ \text{dn } (iu, k) &= \text{dc } (u, k') \end{aligned}$$

$$\text{sn } (u + v) = \frac{\text{sn } u \text{cn } v \text{dn } v + \text{sn } v \text{cn } u \text{dn } u}{1 - k^2 \text{sn}^2 u \text{sn}^2 v} \tag{8}$$

theorems for c and d . These, in combination with the formulas for $\text{sn } (iu)$, for example, serve to express $\text{sn } (u + iv)$ in terms of elliptic functions of u and v they provide formulas for $\text{sn } (2u), \text{sn } (u/2)$, and so on. By means of these formulas the values of the jacobian functions at the points $m\mathbf{K}/2 + ni\mathbf{K}'/2$

TABLE 2. Symmetries of jacobian elliptic functions

v	$-u$	$2\mathbf{K} - u$	$2i\mathbf{K}' - u$
sn v	$-\text{sn } u$	sn u	$-\text{sn } u$
cn v	cn u	$-\text{cn } u$	$-\text{cn } u$
dn v	dn u	dn u	$-\text{dn } u$

TABLE 3. Degenerate elliptic functions

K	K'	k	k'	sn u	cn u	dn u
∞	π/2	1	0	tanh u	sech u	sech u
π/2	∞	0	1	sin u	cos u	1
∞	∞			0	1	1

(m, n integers) and at the points mK + inK' + u can be obtained.

The elliptic functions reduce to elementary functions if one or both of the periods become infinite (degenerate elliptic functions, see Table 3).

Weierstrass's functions. With $w = 2m\omega + 2n\omega'$, sums and products running over all pairs of integers m, n except m = n = 0, there are these functions: Weierstrass sigma function, Eq. (9), which is

$$\sigma(z) = z\prod \left\{ \left(1 - \frac{z}{w}\right) \exp \left[\frac{z}{w} + \frac{1}{2} \left(\frac{z}{w}\right)^2 \right] \right\} \quad (9)$$

an entire function; Weierstrass zeta function, $\zeta(z) = \sigma'(z)/\sigma(z)$, which is a meromorphic function; and Weierstrass P-function, $\mathcal{P}(z) = \zeta'(z)$, which is an elliptic function of order 2 with double poles at $z = 0$ and congruent points. The invariants are $g_2 = 60\Sigma w^{-4}$ and $g_3 = 140\Sigma w^{-6}$. Legendre's relation becomes $\omega'\zeta(\omega) - \omega\zeta(\omega') = \pi i/2$. The P-function satisfies the differential equation, Eq. (10), and possesses the addition theorem given as Eq. (11). It is a

$$\mathcal{P}'^2(z) = 4\mathcal{P}^3(z) - g_2\mathcal{P}(z) - g_3 \quad (10)$$

$$\mathcal{P}(u+v) = \frac{1}{4} \left[\frac{\mathcal{P}'(u) - \mathcal{P}'(v)}{\mathcal{P}(u) - \mathcal{P}(v)} \right]^2 - \mathcal{P}(u) - \mathcal{P}(v) \quad (11)$$

homogeneous function of degree -2 in z, ω, ω'. Every elliptic function with primitive periods 2ω, 2ω' can be expressed in the form $R_1[\mathcal{P}(z)] + \mathcal{P}'(z)R_2[\mathcal{P}(z)]$, where $R_1(w)$ and $R_2(w)$ are rational functions of w; and there are also representations in terms of zeta and sigma functions. Degenerate cases lead to elementary functions.

Theta functions. The function in Eq. (12), with a

$$\theta(v|\tau) = \sum_{r=-\infty}^{\infty} e^{i\pi r^2\tau + i\pi r(2v+1)} \quad (12)$$

fixed τ and Im τ > 0, is an even entire function of v. It has period 1, it is multiplied by $e^{-i\tau(2v+\pi)}$ when v is increased by τ, and it has simple zeros at the points $v = m + (n + 1/2)\tau$ (m, n integers). It is usual to consider four theta functions, Eqs. (13). $\theta(x/2, i\pi t)$

$$\begin{aligned} \theta_1(v) &= -ie^{i\pi(v+\tau/4)}\theta\left(v + \frac{\tau}{2}\right) \\ \theta_2(v) &= e^{i\pi(v+\tau/4)}\theta\left(v + \frac{1+\tau}{2}\right) \\ \theta_3(v) &= \theta(v + 1/2) \\ \theta_4(v) &= \theta(v) \end{aligned} \quad (13)$$

satisfies the partial differential equation $\partial^2 y/\partial x^2 =$

$\partial y/\partial t$ and has a simple Laplace transform. Elliptic functions and elliptic integrals can be expressed in terms of theta functions: $\tau = \omega'/\omega$ in the case of Weierstrass functions and $\tau = iK'/K$ in the case of jacobian functions or Legendre's normal form of elliptic integrals.

Transformation theory. The set of periods of an elliptic function may be described by various pairs of primitive periods. The change from one pair of primitive periods to another pair is called a transformation of the elliptic function or integral. The quotient of the primitive periods, τ, undergoes a homographic substitution, $\tau = (a\tau + b)/(c\tau + d)$, where a, b, c, d are integers and $D = ad - bc$ is positive and is called the degree of the transformation. All transformations of degree 1 form the modular group. The study of these transformations is of great theoretical interest, has applications to number theory, and is also used for the numerical computation of elliptic functions. It is also connected with the study of elliptic modular functions, or analytic functions, f(τ), with the property that f(τ) and f(τ̄) are algebraically connected whenever τ and τ̄ are connected by a transformation of the modular group. See FOURIER SERIES AND TRANSFORMS.

A. Erdelyi

Bibliography. M. Abramowitz and I. A. Stegun (eds.), *Handbook of Mathematical Functions*, 10th ed., 1972; N. Akhiezer, *Elements of the Theory of Elliptic Functions*, 1990; P. F. Byrd and M. D. Friedman, *Handbook of Elliptic Integrals for Engineers and Scientists*, 2d ed., 1971; K. Chandrasekharan, *Elliptic Functions*, 1985; A. Erdélyi et al., *Higher Transcendental Functions*, vol. 2, 1953, reprint 1981; S. Lang, *Elliptic Functions*, 2d ed., 1987; D. F. Lawden, *Elliptic Functions and Applications*, 1989.

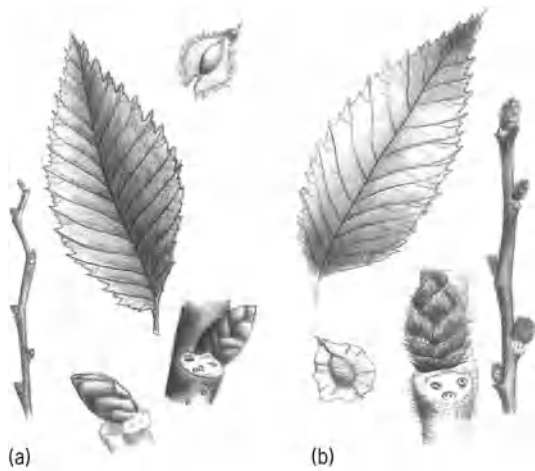
Elm

Any species of *Ulmus*, a genus of hardwood trees in the Northern Hemisphere, with simple, serrate, deciduous leaves. The American or white elm (*U. americana*) is the most important species (illus. a). It is a large, typically vase-shaped tree which usually grows to about 80 ft (24 m) but may reach a height of 140 ft (42 m). The leaves are unequal at the base and doubly serrate; that is, the leaf teeth are themselves toothed. The winter buds are brown, scaly, and usually at one side of the leaf scar. The fruit, ripening in late spring, is a small, flat, winged, elliptical samara, fringed with hairs on the edge.

It ranges from the eastern half of the United States westward as far as the base of the northern Rockies and southward through central Texas to the Gulf of Mexico. The tree is also found in southern Canada.

The tree is common in moist soil, usually at low elevations. It is very popular as a shade tree and is perhaps better adapted as a street tree than any other species because the upper branches spread joining with elms across the street to form an arch.

The wood is heavy, hard, and difficult to split. It is used for wagon parts, barrel staves, shipbuilding, furniture, flooring, sporting goods, boxes, and baskets,



Diagnostic features of elm species. (a) Twig, leaf, fruit, and terminal and lateral buds of American elm (*Ulmus americana*). (b) Leaf, fruit, lateral bud, and twig of the slippery elm (*U. rubra*).

and to some extent, veneer and wall paneling.

Once abundant, the elm is being severely attacked by the lethal Dutch elm disease imported from Europe. Shade trees are particularly susceptible and have been virtually eliminated in many towns and cities. The future of the species is uncertain. See PLANT PATHOLOGY.

Slippery elm (*U. rubra*), a smaller tree with larger, rough leaves (illus. *b*) and mucilaginous inner bark, is commercially important in the east-central United States. The uses of the wood are similar to those of American elm.

Rock or cork elm (*U. thomasi*) of the northeastern United States is not so common. It can be recognized by its usually corky branches, small, spherical winter buds, and flowers and fruit in elongated clusters. The wood is similar in character and use to that of the American elm. See FOREST AND FORESTRY; TREE.

Arthur H. Graves; Kenneth P. Davis

Elopiformes

An order of teleost fishes in the superorder Elopomorpha. All elopomorphs share the leptocephalus-type larvae. Elopiforms are further characterized by a slender body with abdominal pelvic fins and deeply forked caudal fin supported by seven hypurals (flattened bones along the ventral side of the urostyle); a large mouth, terminal or superior, bordered by premaxillae and toothed maxillae; teeth also present on parasphenoid and mesopterygoid bones, and tongue; mesocoracoid and postcleithra bones and a well-developed gular plate; cycloid scales; and small leptocephali, about 5 cm (2 in.) in length, with well-developed, forked caudal fin and strong teeth. Eel leptocephali lack a caudal fin (Fig. 1). There are two families of elopiforms, Elopidae and Megalopidae.

Elopidae (ladyfishes/tenpounders). Ladyfishes, also known as tenpounders, are elongate with a slightly compressed body, a large terminal mouth, large pseu-

dobranchiae (small gill-like structures, developed on the underside of the gill cover), and no swim bladder (Fig. 2). The last ray of the dorsal fin is not elongate, the scales are small, and the lateral-line scales have unbranched tubes (pores). Members of the family are primarily coastal marine species in tropical and subtropical oceans, rarely entering brackish or freshwater. Maximum length is 1 m (3.3 ft). This family consists of one genus and about six species.

Megalopidae (tarpons). Tarpons have a rather strongly compressed body; a large terminal or superior mouth; no pseudobranchiae; and a highly vascular, lunglike swim bladder (Fig. 3). The last ray of the dorsal fin is elongate, the scales are very large, and the lateral-line scales have branched tubes. Tarpons occur in tropical and subtropical oceans, often entering freshwater, especially as juveniles. Maximum length is about 2.4 m (7.9 ft). There are two species: *Megalops cyprinoides* of the Indo-West Pacific and *M. atlanticus* (*Tarpon atlanticus*) of the western Atlantic and tropical seas off West Africa.

Classification. The leptocephalus larval stage is a commonality of the Elopiformes, Albuliformes, Anguilliformes, and Saccopharyngiformes, groups that comprise the Elopomorpha (Fig. 1). The elopomorphs as well as all the remaining extant teleosts are thought to have descended from amiiforms. In

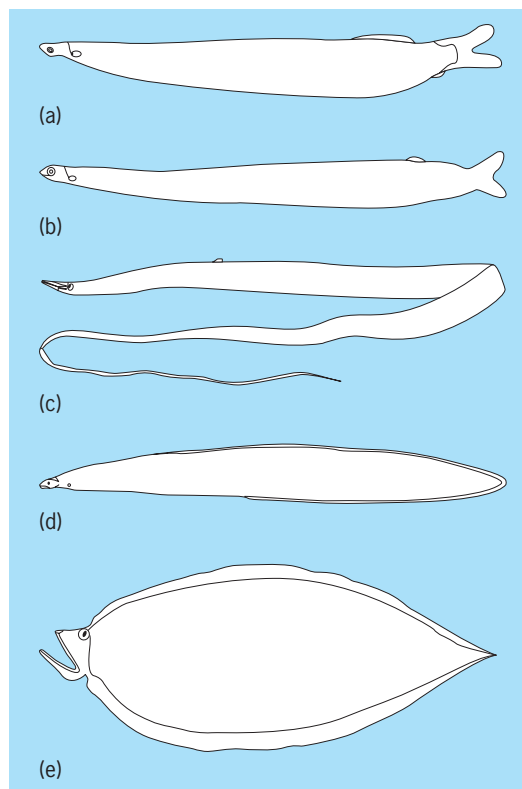


Fig. 1. Representative leptocephalus larvae of the major groups of Elopomorpha. Typical maximum length is in brackets. (a) Elopiformes [5 cm (2 in.)]; (b) Albuliformes [6 cm (2.4 in.)]; (c) Notacanthoidei [2 cm (0.8 in.)]; (d) Anguilliformes [12 cm (4.8 in.)]; and (e) Saccopharyngiformes [5 cm (2 in.)]. (Courtesy of J. S. Nelson, *Fishes of the World*, Wiley, 2006)

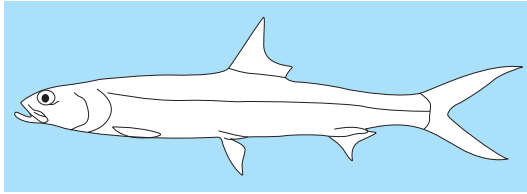


Fig. 2. Example from the family Elopidae (ladyfishes, or tenpounders). (Courtesy of J. S. Nelson, *Fishes of the World*, Wiley, 2006)

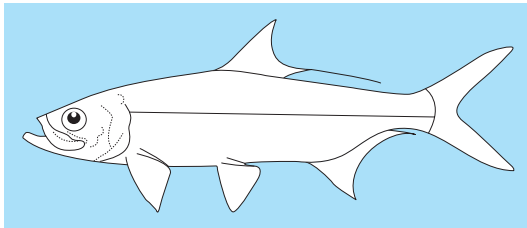


Fig. 3. Example from the family Megalopidae (tarpons). (Courtesy of J. S. Nelson, *Fishes of the World*, Wiley, 2006)

turn, albuliforms, anguilliforms, and saccopharyngiforms are believed to be descendants of elopiforms. The former order Notacanthiformes (spiny eels and halosaurs), which has leptocephalus larvae, is now classified in the order Albuliformes as the suborder Notacanthoidei. Some systemists reject the notion that monophyly (referring to a taxon that includes all the evolutionary descendants of the taxon's common ancestor and only those descendants) of the Elopomorpha is justified because the leptocephalus larval condition is believed to be a primitive rather than a derived condition. See ALBULIFORMES; ANGUILLIFORMES; SACCOPHARYNGIFORMES; TELEOSTEI.

Herbert Boschung

Bibliography. E. B. Böhlke (ed.), *Fishes of the Western North Atlantic*, Part 9, vol. 2: *Leptocephali*, pp. 657-1055, Sears Foundation for Marine Research, Memoir (Yale University), New Haven, 1989; J. S. Nelson, *Fishes of the World*, 4th ed., Wiley, New York, 2006; C. H. Wang et al., Molecular phylogeny of elopomorph fishes inferred from mitochondrial 12S ribosomal RNA sequences, *Zool. Scr.*, 32:231-241, 2003.

Embedded systems

Computer systems that cannot be programmed by the user because they are preprogrammed for a specific task and are buried within the equipment they serve. The term derives from the military, where computer systems are generally activated by the flip of a switch or the push of a button. For example, a military aircraft pilot may wish to turn on the countermeasures equipment with a switch. There is no need for the pilot to be involved with the computer. The same holds true for a soldier who may direct a ground-to-ground missile against a target tank by the push of a button. In both cases, an embedded computer quickly goes to work.

Embedded systems are also used in industrial, automotive, consumer, and medical applications. Specific machine-tool routines can often be initiated by the machine operator by a switch or a button. The cruise-control feature in automobiles, which is activated by a front-panel or steering-column switch, is the result of an embedded microcomputer. Automatic scoring and unattended bowling lanes depend on the operation of embedded microcomputers. Many types of intensive-care and health-care equipment often operate unattended because of embedded computer systems.

The continual increase in the densities of ever-smaller microprocessors, the brains and core of microcomputers, on silicon chips that fit on a thumbnail, and the attendant decreases in their costs, has pushed the concept of embedded systems well beyond the original military applications. Embedded microprocessors provide value-added features to a wide range of products. They are used in consumer electronics, automotive vehicles, telecommunications, computer peripherals, process control, factory automation, robotics, and intelligent instruments. Microwave ovens, stereo equipment, washer-dryers, refrigerators, blenders, videocassette recorders, and even pressing irons use embedded microcomputers based on small microprocessors. Engine and drive-train controllers are two of the largest automotive applications. Embedded systems are also widely used in computer printers and modems. The first microprocessor chip, designed in 1971, can be considered an embedded computer, since it was originally created for an electronic calculator application.

CISC and RISC architectures. Most embedded microprocessors are of the CISC (complex-instruction-set computer) type, and most of these are used in applications where low cost is paramount and performance is secondary, such as consumer products. The later-generation microprocessors have wider bus widths, up to 64 bits, and thus can do more computations.

Since about 1990, microprocessors of the RISC (reduced-instruction-set computer) type have appeared, with much greater computational capability and at greater cost. RISC processors are used mostly in those embedded applications where performance is primary and low cost is secondary, such as in the military. They are also used in engineering workstations, where the computational burdens of high-resolution graphics require such processors. See COMPUTER-AIDED ENGINEERING; COMPUTER GRAPHICS; COMPUTER SYSTEMS ARCHITECTURE.

Software development. It is generally more difficult to develop software for RISC processors than for CISC processors. Indeed, modern high-performance processors, whether RISC or CISC, require considerably more software development than their predecessors, and placing such processors in embedded applications further complicates the task of software development. The chief issue in the design of embedded processors is now the nature of the end application rather than the architecture of the hardware, as

had been the case in earlier-generation devices. See INTEGRATED CIRCUITS; MICROCOMPUTER; MICROPROCESSOR; SOFTWARE ENGINEERING. Roger Allan

Bibliography. A. Berger, *Embedded Systems Design: An Introduction to Processes, Tools and Techniques*, 2001; Q. Li and C. Yao, *Real-Time Concepts for Embedded Systems*, 2003; T. Noergaard, *Embedded Systems Architecture: A Comprehensive Guide for Engineers and Programmers*, 2005; R. Zurawski (ed.), *Embedded Systems Handbook*, 2005.

Embioptera

A peculiar order of silk-spinning, orthopteroid insects related to termites, commonly called the embiids or web spinners. This order comprises about 1000 species which are chiefly tropical in distribution. The body is linear and supple. Adults vary in length from 0.12 to 0.8 in. (3 to 20 mm). The legs are short with three-segmented tarsi. The forelegs are adapted for spinning silk, and the hindlegs for reverse locomotion. Cerci are short, two-segmented, and tactile. The head is prognathous, with orthopteroid mouthparts. Metamorphosis is incomplete. The females are neoteinic and wingless (apterous) [Fig. 1]. Males are usually winged (alate), but in certain genera and species they are apterous. The wings are subequal and elongate, with the vannal fold obsolete. Venation is simple, with the veins centered in pigment bands and separated by hyaline stripes. The wings are flexible when in repose and folded over the back, but are stiffened when extended for flight by the blood pressure in the saclike R_1 (radial) veins. Flight is a poorly directed, whirling flutter.

Plantar surfaces of the mid and basal foretarsal segments bear many hollow, silk-spinning setae, each connected by a duct to a globular, syncytial silk gland; many such glands are massed in the greatly swollen basal segment. Hundreds of silk strands are emitted in unison as the tarsi brush a surface. A labyrinth of silk galleries is produced rapidly, which constitutes a safe shelter for all embiid activities except adult dispersal (Fig. 2). All embiids spin silk regardless of the species, developmental stage, or sex.

The galleries radiate on or in the food supply which also constitutes the habitat and consists of bark, lichens, moss, dead leaves, or grass. Many individuals, usually the brood of one female, may



Fig. 1. Body form of a typical embiid (*Pararhagadochir trachelia*, female). (From E. S. Ross, *Insects Close Up*, University of California Press, 1953)

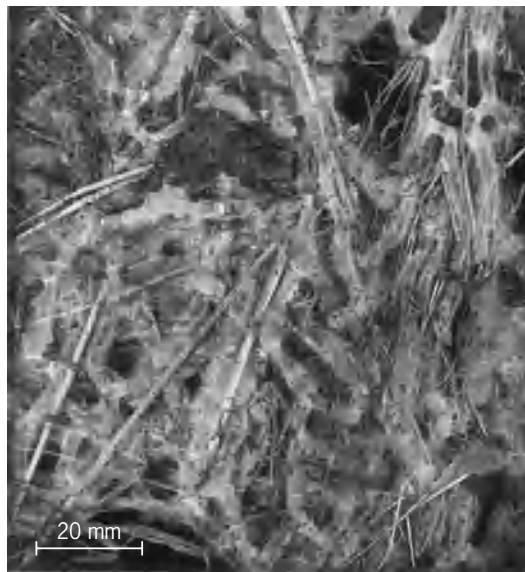


Fig. 2. Characteristic silk gallery system of the embiid *Haploembia solieri*.

occupy one gallery system. See INSECTA; ORTHOPTERA. Edward S. Ross

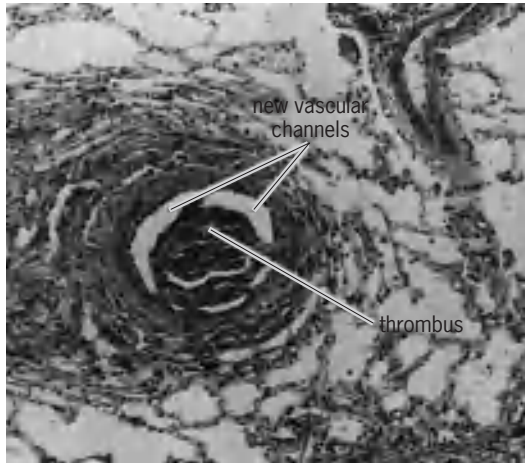
Bibliography. S. P. Parker (ed.), *Synopsis and Classification of Living Organisms*, 2 vols., 1982; E. S. Ross, The Embioptera of Europe and the Mediterranean region, *Bull. Brit. Mus. Entomol.*, 17:275-326, 1966; E. S. Ross, A revision of the Embioptera or web-spinners of the New World, *Proc. U.S. Nat. Mus.*, 94:401-504, 1940.

Embolism

The sudden blocking of an artery or vein by a clot or other substance which has been brought to its place by the blood current. The material carried in the circulation in this process is an embolus. Emboli may be composed of thrombi, fat, air, tumor cells, masses of bacteria or parasites, bone marrow, amniotic fluid, or atheromatous material from the vessel wall. A thrombus which has formed in the heart or one of the vessels is the usual form of embolus (see *illus.*). Emboli from the right side of the heart or the great venous system of the body come to lodge in the lungs; those from the portal system, in the liver; and those from the pulmonary veins or left side of the heart, in some segment of the peripheral arterial tree. Pulmonary emboli can result in infarction of the lung. However, if they are large enough to occlude the main pulmonary artery or one of the major branches, the individual may die of shock.

Embolization is a common method of spread of tumor cells, which is one reason for the frequency of metastatic tumors in the liver and lungs.

Injury to bone or adipose tissue may result in fat embolism. Following its escape from the injured fat cell, the fat gains access to the venous system. It may then come to rest in the capillaries of the lung, or it may pass through the lungs to find its way to some other tissue such as the brain or kidney.



Thrombus, a type of embolus, appearing in the pulmonary artery of a dog's lung, with the subsequent formation of new vascular channels.

Air embolism may be a complication of a surgical procedure, particularly those about the neck. It can also result from rapid decompression, as in caisson disease, with the formation of bubbles of nitrogen in the blood. Infected emboli can form new foci of infection at their sites of lodgement. See DECOMPRESSION ILLNESS.

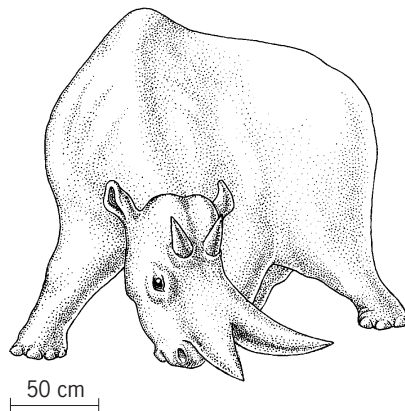
With complete obstruction of a vessel by an embolus an infarct may result. See INFARCTION.

Romeo A. Vidone

Bibliography. S. Z. Goldhaber, *Pulmonary Embolism and Deep Venous Thrombosis*, 1985; J. Hirsh et al., *Venous Thromboembolism*, 1981; G. Holm and M. Bjorkholm (eds.), *Biology of Vascular Disease*, 1987; O. Horwitz et al. (eds.), *Diseases of the Blood Vessels*, 1985; R. C. Little, *Physiology of the Heart and Circulation*, 4th ed., 1988.

Embrithopoda

An order established for the unique mammal *Arsinoitherium*, which has been found only in early Oligocene deposits in northern Egypt. This animal was of rhinoceros size with a large body and short pil-



Arsinoitherium, the early Oligocene embrithopod from Egypt.

larlike legs (see *illus.*). There were two huge, scimitarlike horns over the nose and two much smaller, peglike horns over the eyes. As in many herbivorous mammals, the cusps of the high-crowned teeth were connected to form two transverse lophs on each tooth. The exact relationship of this exotic order to other orders is not clear, but similarities to both the uinatheres, or Dinocerata, and Hyracoidea have been interpreted by some authorities as evidence of kinship. See DINOCERATA; MAMMALIA.

Malcolm C. McKenna

Embrittlement

A general set of phenomena whereby materials suffer a marked decrease in their ability to deform (loss of ductility) or in their ability to absorb energy during fracture (loss of toughness), with little change in other mechanical properties, such as strength and hardness. Embrittlement can be induced by a variety of external or internal factors, for example, (1) a decreasing or an increasing temperature; (2) changes in the internal structure of the material, namely, changes in crystallite (grain) size, or in the presence and distribution of alloying elements and second-phase particles; (3) the introduction of an environment which is often, but not necessarily, corrosive in nature; (4) an increasing rate of application of load or extension; and (5) the presence of surface notches. This article is restricted to metals and alloys, and describes embrittlement in terms of the above factors. See BRITTLENESS.

Effect of test temperature. In body-centered cubic metals (for example, iron, tungsten) and hexagonal close-packed metals (for example, zinc, magnesium), a critical temperature exists below which the metal exhibits limited toughness. Fracture is usually brittle in nature, occurring either through the crystal lattice (cleavage) or along the boundaries separating the crystallites or grains (intergranular fracture). In simple terms, low-temperature embrittlement results from a competition between deformation and brittle fracture, with the latter becoming preferred at a critical temperature. For a material to be useful structurally, it is desirable that this critical temperature be below the minimum anticipated service temperature; in most cases, this is room temperature. At high temperatures, internal structural changes that lead to intergranular embrittlement can occur. Embrittlement usually occurs in the range of creep temperature at which deformation can occur under very low stresses; and the two processes are believed to be connected. See CRYSTAL STRUCTURE; PLASTIC DEFORMATION OF METAL.

Effect of heat treatment. In many metals, particularly structural steels, annealing or heat treating in certain temperature ranges sensitizes the grain boundaries in such a way that intergranular embrittlement subsequently occurs during service. In steels, for example, high strength is achieved by quenching from the stable high-temperature face-centered cubic phase (called austenite). This can

produce a room-temperature structure called martensite, a distorted, body-centered cubic phase which is extremely hard and strong, but brittle. The transformation of austenite to martensite is controlled by several metallurgical variables, including the alloy composition and the quench rate. To reduce the brittleness, the steel undergoes an annealing treatment called tempering, which, while decreasing the strength, usually increases the toughness. The exception to this trade-off occurs when the steel is tempered at 1000°F (538°C). This can lead to a mode of intergranular fracture called temper embrittlement; such a process has led to catastrophic failures in turbines, rotors, and other high-strength steel parts. It is associated with the preferential segregation of undesirable "tramp" elements such as tin and antimony to the grain boundaries; these elements should be kept to as low a level as possible to help alleviate the problem. Tempering at 550°F (287°C) can also lead to embrittlement in steels. In other metals, there are less specific but similar types of embrittlement resulting from critical heat treatments. *See* HEAT TREATMENT (METALLURGY); STEEL; TEMPERING.

Effect of environment. Metals can fracture catastrophically when exposed to a variety of environments. These environments can range from liquid metals to aqueous and nonaqueous solutions to gases such as hydrogen. The phenomenology of these processes and some of the corrective procedures used or envisaged are described below.

Liquid metal embrittlement. If a thin film of a liquid metal is placed on the oxide-free surface of a solid metal, the tensile properties of the solid metal will not be affected, but the fracture behavior can be markedly different from that observed in air. In some cases, specimens stressed above a critical value will fracture catastrophically. Although the embrittlement is usually intergranular, transgranular cleavage has also been observed. Such fractures can be induced in highly ductile metals such as aluminum and copper. Although many different liquid metals are capable of inducing embrittlement in a variety of solid metals, some of the more common couples, many of which have important engineering and design consequences, are mercury embrittlement of brass (an alloy of copper and zinc), lead embrittlement of steel, and gallium embrittlement of aluminum. The mechanism of liquid metal embrittlement is thought to result from the reduction by the liquid metal atoms of the force necessary to break apart two solid metal atoms at a crack tip.

Stress-corrosion cracking. If a metal is stressed and simultaneously exposed to an environment which may be, but is not necessarily, corrosive in nature, cracking and fracture can occur. Both stress and environment are required; if only one of these elements is present, the metal usually displays no embrittlement. *See* CORROSION.

The analogy of stress-corrosion cracking to liquid metal embrittlement is obvious, but in this case, the environment can be any aqueous or nonaqueous solution. In specific metals cracking has been observed in such diverse solutions as high-purity distilled

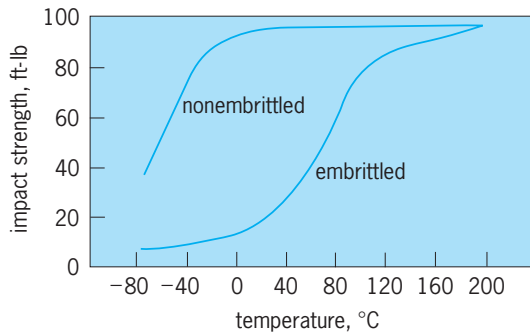
water, salt water, caustics (for example, lye), ammonia, molten anhydrous salts, and organics. Stress-corrosion cracking can occur over a wide temperature range, and can be a very serious problem in many service applications. For example, brass cartridge cases which were internally stressed by a forming operation cracked in storage because of the presence of trace amounts of ammonia in the air; and titanium alloys used in aerospace applications cracked when cleaned with methyl alcohol. There are a host of situations in which moisture or salt water led to cracks or total fracture of structural materials.

A mechanism which has been proposed as an explanation for many of these cases has been associated with the stress-assisted preferential dissolution of metal along specific metallurgical features, leading to intergranular or transgranular cracks. This process can be likened to the action of a battery. In this case, the cracking region is the anode; when the anode is coupled to a cathodic region, the circuit is completed. This electrolytic behavior has been exploited as a protective device. Critical parts are made to serve as cathodes in local cells, allowing other noncritical parts to act as sacrificial anodes. This procedure, termed cathodic protection, is used to protect not only against stress-corrosion cracking but also against general corrosion; for example, zinc sacrificial anodes are used to minimize corrosion of ship propellers. Other mechanisms, such as corrosion, have also been put forward to explain the stress-corrosion cracking phenomena. *See* ELECTROCHEMICAL PROCESS.

Hydrogen embrittlement. This form of embrittlement is often considered to be a type of stress-corrosion cracking. Hydrogen atoms can enter a metal, causing severe embrittlement, again with little effect on other mechanical properties. This phenomenon was originally observed, and is most critical in, steels, but it is now documented to occur in titanium and nickel alloys, and may lead to cracking in other alloy systems as well.

The source of hydrogen can be extremely varied, making this type of embrittlement difficult to control. For example, hydrogen can be retained internally during the melting and casting of alloys; it can be discharged at cathodic areas in electrolytic cells (it is apparent that, at both the anode and the cathode, events which may lead to embrittlement occur); it can enter the metal during a plating operation; and it can come from an external molecular gas environment, even at very low partial pressures of hydrogen. When it is realized that coal gasification plants produce hydrogen or methane (a hydrogen-containing gas) as a fuel, that long-distance transport of hydrogen is envisaged as a means of alleviating the energy crisis, that liquid hydrogen is used as a rocket fuel, and that hydrogen is a by-product of many chemical reactions, it becomes apparent that structural materials must be resistant to hydrogen embrittlement. *See* COAL GASIFICATION; HYDROGEN.

In all forms of environmental embrittlement, the best remedial practice is to physically isolate the



Effect of temper brittleness on the impact resistance of a low-alloy steel. 1 ft-lb = 1.36 J. $^{\circ}\text{F} = (^{\circ}\text{C} \times 1.8) + 32$. (After R. E. Reed-Hill, *Physical Metallurgy Principles*, Van Nostrand, 1973)

metal from the dangerous environment. This is not always practical or possible. Studies have indicated that the susceptibility of a given material to environmental embrittlement can be reduced by controlling metallurgical factors, such as alloy composition and internal structure. This type of alloy design should help control environmental embrittlement, even where there is uncertainty as to the controlling mechanism involved in the process. See ALLOY.

Other factors. Factors such as notches (which both raise and modify the local stress state) and the rate of application of stress can modify the response of a material to a specific type of embrittlement. In general, notches or surface flaws always enhance embrittlement, both by acting as a stress raiser and by providing a preexisting crack. An increasing loading rate (impact) enhances low-temperature embrittlement, and possibly thermal embrittlement, but, interestingly, suppresses environmental embrittlement. For environmental embrittlement, species such as hydrogen must keep up with the moving crack, and they do so by moving through the crystal lattice at a rate that is largely determined by the temperature. If the crack is moving too quickly, as under impact conditions, hydrogen can not keep pace, and the severity of embrittlement decreases. For normal rates of loading, the maximum effect of hydrogen is unfortunately at, or near, ambient temperatures. The combined effects of various embrittlement phenomena are shown in the **illustration**. Temper embrittlement increases the impact ductile-to-brittle transition temperature from about -58 to 176°F (-50 to 80°C), making this steel unacceptable for room-temperature applications. See BRITTLENESS.

I. M. Bernstein

Bibliography. Z. Ahmad, *Principles of Corrosion Engineering and Corrosion*, 2006; N. E. Dowling, *Mechanical Behavior of Materials: Engineering Methods for Deformation, Fracture, and Fatigue*, 2d ed., 1998; M. G. Fontana, *Corrosion Engineering*, 3d ed., 1986; G. Y. Lai, *High Temperature Corrosion of Engineering Alloys*, 1990; P. R. Roberge, *Handbook of Corrosion Engineering*, 1999; H. H. Uhlig and R. W. Revie, *Corrosion and Corrosion Control*, 3d ed., 1985.

Embryobionta

One of the two plant subkingdoms, the other being the Thallobionta. The Embryobionta (often also called Embryophyta) form a well-marked group with many morphological, anatomical, and physiological features in common, and taxonomists are agreed that the Embryobionta as a group are derived from the division Chlorophyta (green algae) in the subkingdom Thallobionta. The Embryobionta are here considered to include eight divisions, the Rhyniophyta, Bryophyta, Psilotophyta, Lycophyta, Spenoophyta, Polypodiophyta, Pinophyta, and Magnoliophyta. The Rhyniophyta are represented only by Paleozoic fossils, but the other seven divisions have both modern and fossil representatives. See the separate articles on each division. See CHLOROPHYCOTA.

Characteristics. The pigments of the chloroplasts of the Embryobionta are chlorophylls *a* and *b*, β - and usually α -carotene, and several xanthophylls, generally including lutein, cryptoxanthin, and zeaxanthin. Chlorophyll *a* is the principal chlorophyll, β -carotene is the principal carotene, and lutein is the principal xanthophyll. The carotenes and xanthophylls, collectively called carotenoid pigments, are masked by the more abundant chlorophylls. The most common carbohydrate reserve of the Embryobionta is starch. The primary cell wall consists largely of cellulose. Flagellated cells, when they are present (for example, the sperms of several divisions), have only the whiplash type of flagellum. In all these respects the Embryobionta resemble the green algae. See CAROTENOID; CHLOROPHYLL.

Life cycle. The Embryobionta differ from the green algae and from most Thallobionta in that the normal life cycle of the Embryobionta shows a well-marked alternation of generations in which the sporophyte (spore-producing, typically diploid) generation always begins as a parasite on the gametophyte (gamete-producing, typically haploid) generation. The young sporophyte is called an embryo. The only Thallobionta which have a definite sporophyte attached to the gametophyte are some of the red algae, and these are so different from the Embryobionta in other respects that no close relationship seems possible.

Reproductive structures. The more primitive divisions of Embryobionta have the gametes produced in multicellular sex organs (archegonia and antheridia), which is in contrast to the unicellular oogonia and antheridia of the Thallobionta in general. In the more advanced divisions (Pinophyta and Magnoliophyta) of Embryobionta, the antheridia and archegonia are highly modified or entirely suppressed, in conformity with the general reduction of the gametophyte generation.

Specializations. All divisions of Embryobionta except the Bryophyta have specialized conducting tissues (xylem and phloem) in the sporophyte. With the exception of most bryophytes, they also commonly have a characteristic stomatal apparatus which controls the opening and closing of numerous tiny pores (stomates) in the leaves and stems in response to

environmental conditions. These specializations, together with the progressive reduction of the gametophyte in the more advanced divisions, reflect the progressive evolutionary adaptation of the Embryobionta to life on dry land instead of in the ancestral water. The Embryobionta are therefore often called the land plants, in spite of the fact that many of them, such as the water lilies, have returned to an aquatic habitat. The seven divisions of Embryobionta which characteristically have xylem and phloem have sometimes been treated as a single comprehensive division under the name Tracheophyta. See PLANT KINGDOM; THALLOBIONTA.

Arthur Cronquist

Bibliography. A. Cronquist, *Introductory Botany*, 1961; A. Cronquist, A. Takhtajan, and W. Zimmerman, *Taxon*, 15:129-134, 1966.

Embryogenesis

The formation of an embryo from a fertilized ovum, or zygote. Development begins when the zygote, originating from the fusion of male and female gametes, enters a period of cellular proliferation, or cleavage. Cells of the embryo subsequently give rise to the tissues and organs of the body in a temporal and spatial pattern that creates a functional, multicellular organism. See CLEAVAGE (EMBRYOLOGY).

Animals

The animal zygote divides by mitosis to produce cells called blastomeres. Since there is little or no growth during cleavage, the blastomeres become smaller with each succeeding division, and the embryo maintains the relative size and shape of the zygote. In a cleaving cell, the spindle and asters of the mitotic apparatus are essential for division of the nucleus. The mitotic apparatus also determines the plane for division of the cytoplasm. A cleavage furrow arises on the cell surface along a line that is equidistant from the paired asters. As this furrow deepens, it gradually partitions the cytoplasm between the daughter blastomeres. The mitotic apparatus tends to situate in the central cytoplasm and orient parallel to the long dimension of a cleaving cell. The long dimension can, of course, change from one cleavage to the next, shifting the orientation of the mitotic apparatus and the plane of the cytoplasmic division. If cytoplasmic constituents, such as yolk granules, prevent the mitotic apparatus from assuming its customary position, the cleavage furrow will unequally distribute the cytoplasm, making the daughter blastomeres dissimilar in size and composition. See BLASTULATION; CELL DIVISION; MITOSIS.

Following cleavage, the cells of the animal embryo rearrange into three germ layers: an outer ectoderm, a middle mesoderm, and an inner endoderm. Cells, responding to intrinsic and extrinsic factors, eventually segregate from the germ layers and organize into the rudiments of the tissues (epithelia, connective tissues, muscles, and nerves) and the organs of the body. These rudiments alter both the size and the

shape of the embryo. Cellular growth and differentiation are the principal processes that transform the rudiments into functional tissues and organs. Once the organs and organ systems are formed, further development consists primarily of growth. See ANIMAL MORPHOGENESIS; GERM LAYERS.

Michael J. Cavey

Plants. The zygote in higher plants develops into a multicellular embryo with apical meristems that ultimately produce the adult plant. Embryogenesis in both plants and animals involves a combination of growth, morphogenesis, and cellular differentiation. Significant differences exist between patterns of embryogenesis in higher plants and animals: (1) cell movement does not play a major role in plant embryogenesis because plant cells are surrounded by cell walls; (2) programmed cell death does not occur during plant embryogenesis; (3) plant embryos do not contain a germ line that functions only in the formation of gametes; and (4) the mature plant embryo does not resemble the adult because after germination, groups of undifferentiated cells known as meristems continue to produce in an accretionary manner the cells, tissues, and organs that make up the adult plant. See CELL WALLS (PLANT).

Major features of embryogenesis in flowering plants include the formation of root and shoot apical meristems; differentiation of primary vascular tissue; the transition from a heterotrophic zygote to an embryo capable of independent growth and development; and preparations for desiccation, dormancy, and germination, all of which appear to be associated with differential gene expression. Embryogenesis in flowering plants is accompanied by development of the endosperm tissue produced after double fertilization involving the fusion of a second male gamete with additional cells of the female gametophyte located within the immature ovule. This nutritive endosperm tissue may be consumed during embryogenesis or be included as part of the dormant seed. The endosperm does not appear to serve a regulatory function during the final stages of embryogenesis, because immature embryos from a wide range of plants have been removed from developing seeds and cultured to maturity on a simple nutrient medium. This process of embryo culture has been used to study the regulation of embryogenesis and to rescue aborted embryos produced by interspecific hybridization.

Under appropriate conditions somatic cells of higher plants can produce embryos in culture through a process known as somatic embryogenesis. Differentiated cells from a wide range of plants have been shown to produce somatic embryos that resemble zygotic embryos in both morphology and patterns of gene expression. Somatic embryogenesis has been used to demonstrate the totipotency of plant cells and to produce a large number of genetically similar plants, which have been used in fundamental research and hold promise for crop improvement. See APICAL MERISTEM; CELL DIFFERENTIATION; DEVELOPMENTAL BIOLOGY; EMBRYOLOGY; EMBRYONIC DIFFERENTIATION; EMBRYONIC INDUCTION.

David Meinke

Bibliography. L. W. Browder et al., *Developmental Biology*, 3d ed., 1997; S. F. Gilbert, *Developmental Biology*, 8th ed., 2006; R. Grossinger, *Embryogenesis*, 1986; V. Raghavan, *Embryogenesis in Angiosperms: A Developmental and Experimental Study*, 1986.

Embryology

The study of the development of an organism, commencing with the union of male and female gametes. Embryology literally means the study of embryos, but this definition is restrictive. An embryo is an immature organism contained within the coverings of an egg or within the body of the mother. Strictly speaking, the embryonic period ends at metamorphosis, hatching, or birth. Since developmental processes continue beyond these events, the scope of embryology is customarily broadened to encompass the entire life history of an organism. Embryology may, in this wider context, consider the mechanisms of both asexual reproduction and regeneration.

Animals

Development of a multicellular organism is a continuum of processes. It is useful, however, to discuss development from the perspective of several overlapping phases, such as gametogenesis, fertilization, cleavage, gastrulation, and histogenesis and organogenesis.

Gametogenesis. The production of male and female gametes occurs in the gonads. The differentiating gametes arise from diploid stem cells. Cell division by meiosis reduces the number of chromosomes carried by a mature gamete to one-half that present in the stem cell. Following meiosis, the male gamete achieves a significant reduction in weight by compacting its nucleus and shedding excess cytoplasm, and it also acquires the organelles needed to approach and fuse with the female gamete.

Specializations of the female gamete take place during, rather than after, meiosis. The female gamete often enters a protracted period of growth, accumulating the yolk that will meet the early nutritional needs of the embryo. A protective covering, or envelope, is deposited around the female gamete before it leaves the gonad, and other envelopes may be added along the female reproductive tract. The envelope(s) and the gamete constitute the egg. *See* GAMETOGENESIS; MEIOSIS; OOGENESIS; OVUM; SPERMATOGENESIS.

Fertilization. The union of spermatozoon and ovum creates a diploid zygote with the potential to form an entire organism. Two events must occur for successful fertilization: the ovum must respond to contact with the spermatozoon by making preparations for further development, a process called activation; and the haploid nucleus of the spermatozoon must combine with the haploid nucleus of the ovum, a process called amphimixis.

The zygotes of different species can be very dissimilar in size, shape, and composition. Because the

spermatozoon contributes negligible cytoplasm to the zygote, these dissimilarities arise from the ovum. The cytoplasmic elements of an ovum are frequently localized in specific regions, and they are subsequently redistributed in consistent patterns during fertilization. Compositional differences, from one side of a zygote to the other and from the surface of a zygote to the interior, can be profound; disturbances of this organization can lead to serious developmental abnormalities.

Fertilization is the typical method to initiate development, but it is not the only way. In a few animals, the ovum develops independently by parthenogenesis, that is, without the participation of a spermatozoon. Parthenogenesis can sometimes be induced experimentally in organisms that would ordinarily opt for fertilization. Several physical and chemical agents are known to mimic the action of a spermatozoon in triggering the activation of an ovum. In such cases of artificial parthenogenesis, the ability of an ovum to develop rests largely on the response of its haploid nucleus. A normal sequence of developmental events is possible if the nucleus can, in some aberrant manner, become diploid. *See* FERTILIZATION.

Cleavage. A period of cell proliferation converts the unicellular zygote into a multicellular embryo. Cleavage is a modified form of cell division by mitosis, distinguished by little or no growth between the divisions. The cells of the embryo, or blastomeres, become smaller at the end of each division, so the embryo maintains the relative size and shape of the zygote. Small, fluid-filled spaces form between the cleaving blastomeres, and these spaces eventually coalesce to create an internal cavity, or blastocoele. Upon the appearance of a blastocoele, the cells of an embryo are referred to collectively as the blastoderm.

There is overwhelming evidence that the cytoplasm, not the nucleus, directs cleavage. The earliest demonstration of this control was achieved with subdivided ova, or merogones, produced by centrifugation. One merogone receives the haploid nucleus of the ovum, and its development can be initiated by fertilization. The development of the other merogone, which lacks a nucleus, must be initiated parthenogenetically. Both merogones show similar patterns of cleavage, implying cytoplasmic control over the process. The enucleated merogone cannot pass from cleavage to the next phase of development, suggesting that nuclear control is necessary in making the transition.

The isolation of blastomeres from cleaving embryos has also provided insights into early development. In amphibians, for example, the two blastomeres formed by the first cleavage can be separated and reared individually. Each blastomere continues to develop and produces a small but perfectly normal larva. The same experiment on the four blastomeres formed by the second cleavage gives different results. Two of the blastomeres develop into normal larvae, and two develop abnormally, forming larvae with conspicuous defects. Since the nuclei of the early blastomeres are identical, cytoplasmic

differences must specify whether a blastomere has the capacity to form an entire individual or only certain parts of that individual. Cleavage is, therefore, more than a period of rapid cell proliferation; it is also a time when cytoplasmic differences begin to guide the embryonic cells into specific pathways of development. *See* BLASTULATION; CLEAVAGE (EMBRYOLOGY); MEROGONY.

Gastrulation. Gastrulation is poorly delineated from cleavage because the cells of the embryo continue to divide. Gastrulation is distinguished from cleavage by extensive cell rearrangements that lead in most animals to the establishment of three germ layers: an outer ectoderm, a middle mesoderm, and an inner endoderm. Endodermal and mesodermal cells of the blastoderm migrate to the inside of the embryo, while ectodermal cells remain on the surface, where they spread to completely cover the body. Mapping experiments, using carbon particles, vital dyes, or radioisotopes to mark small areas of the blastoderm, illustrate that the cell movements during gastrulation are executed in a precise, orderly fashion. *See* FATE MAPS (EMBRYOLOGY).

Control of development passes from the cytoplasm to the nucleus immediately prior to gastrulation. Responding to cytoplasmic cues, the nuclei begin to specify the production of proteins that make the cells qualitatively different from one another. In a few invertebrates, the transfer of control from cytoplasm to nucleus actually fixes the developmental fate of a cell. In most other organisms and particularly in vertebrates, the determination of cell fate is not finalized until the blastoderm has rearranged into the three germ layers. *See* CELL LINEAGE; GASTRULATION; GERM LAYERS.

Histogenesis and organogenesis. The organization of cells into tissues (histogenesis) and tissues into organs (organogenesis) is closely allied with gastrulation. During creation of the germ layers, blastodermal rearrangements shift cells into new positions and bring about new intercellular relationships. To a considerable degree, the developmental fate of a cell can be the consequence of its new position. The influence exerted by one group of cells over the developmental fate of a neighboring group is called induction. Induction occurs by the transmission of chemical substances known as inducing agents.

An inducing agent presumably elicits a cytoplasmic change that, in turn, modifies the output of the nucleus. As cells specialize by accumulating different kinds of proteins, they begin to segregate from the germ layers and organize into four rudimentary tissues. These tissues, comprising epithelia, connective tissues, muscles, and nerves, subsequently combine to establish the rudiments of the organs. The spatial distribution of the organ rudiments alters both the size and shape of the embryo. Inductions between the rudimentary tissues are often required for conversion of the nascent organs into functional units. Growth and terminal differentiation are the principal processes in this conversion. Growth is possible in several ways: by augmenting the size of individual cells, by increasing the total number of cells, or by

accumulating materials in the intercellular spaces. Terminal differentiation is dependent upon the acquisition of all necessary proteins by which a cell maintains its structure and function; this acquisition is manifest by changes in cell number, size, shape, motility, and adhesion. *See* EMBRYOGENESIS; EMBRYONIC INDUCTION.

Differentiation, determination, and morphogenesis. Differentiation, the process by which a cell becomes specialized, correlates to a reduction in the amount of genetic information that is expressed. Determination, the fixation of a developmental fate, occurs when a cell has such a limited amount of usable genetic information that it must commit to a terminal pathway of differentiation. *See* CELL DIFFERENTIATION.

Differentiation and determination are deliberate processes set in motion at the time of fertilization. The heterogeneous cytoplasm of a zygote, when distributed among the blastomeres by cleavage, is undoubtedly responsible for initial constraints placed upon the nuclei. If the nucleus of a blastodermal cell is transplanted to an enucleated zygote, development proceeds normally, indicating that constraints may limit the functions of a nucleus but do not alter its potential. The early restrictions on utilization of genetic information account for the establishment of the three germ layers and the origins of the lowly differentiated ectodermal, mesodermal, and endodermal cells. Further restrictions are imposed partly by the composition of the cytoplasm and partly by the influence of inducing agents. Inductions between the cells of the germ layers give rise to the moderately differentiated cells of the embryonic tissues, and subsequent inductions between the rudimentary tissues provide the highly (and terminally) differentiated cells of the functional organs. *See* GENE; GENE ACTION.

Cellular differentiation is just one aspect of morphogenesis, or the development of form. Morphogenesis must be considered at all levels of organization, ranging from the individual cell to the whole organism. Such a broad perspective complicates the formulation of general theories of development. Presently, no comprehensive theory exists, but there are some embryologists who anticipate that a theory is possible once activities of the DNA molecule have been fully integrated into the topic of development. *See* ANIMAL MORPHOGENESIS; REPRODUCTION (ANIMAL).

Michael J. Cavey

Plants

Reproductive development in multicellular plants is generally divided into three phases: gametogenesis, fertilization, and embryogenesis. The zygote produced by the fusion of male and female gametes divides to form a multicellular embryo with meristematic regions that ultimately produce the adult plant. Many different patterns of reproductive development have been described in plants; emphasis here will be placed on the most common pathways of gametogenesis and embryogenesis in angiosperms.

Gametogenesis. Development of the cell in flowering plants begins with a diploid megasporocyte located within the nucellar tissue of an immature ovule. This megasporocyte undergoes meiosis to form a tetrad of four haploid megaspores. In the most common pattern of development, three of these megaspores degenerate, leaving a single functional megaspore that undergoes several postmeiotic mitoses to form a mature megagametophyte (embryo sac) composed of seven cells and eight haploid nuclei. One of these haploid cells is the egg cell. Variations on this common pattern of development include differences in the number of functional megaspores and the number and arrangement of haploid nuclei present in the mature megagametophyte.

Development of the male gametes begins with numerous diploid cells (microsporocytes) located within the anthers of an immature flower. Each microsporocyte undergoes meiosis to form a tetrad of four haploid microspores, which then separate and enlarge to form mature pollen grains. Each microspore divides unequally to form a large vegetative cell, and a small generative cell located within the cytoplasm of the vegetative cell. The generative cell divides again, in either the maturing pollen grain or the elongating pollen tube, to form two genetically identical male gametes, the sperm cells.

Fertilization and embryogenesis. The zygote is produced as part of a unique process known as double fertilization. One of the male gametes fuses with the egg cell to form the diploid zygote, while the other male gamete fuses with two polar nuclei, located near the center of the embryo sac, to form a triploid endosperm nucleus. Following double fertilization, the zygote develops into an embryo composed of two parts, the embryo proper and the suspensor. The embryo proper ultimately differentiates into the mature embryo, whereas the suspensor degenerates during later stages of development and is not usually present at maturity. The suspensor appears to play both a passive role in attaching the embryo proper to surrounding maternal tissues, and an active role in regulating development of the embryo proper during early stages of embryogenesis. The triploid endosperm nucleus in most angiosperms undergoes a series of free nuclear divisions to form a nutritive endosperm tissue that usually becomes cellular during later stages of development.

Flowering plants can be divided into two groups, monocots and dicots. In most dicots, the endosperm tissue is gradually absorbed by the developing embryo and is not present in the mature seed. Nutrients required for the germination of dicot seeds are generally stored in the embryonic leaves known as cotyledons. In contrast, most mature monocot seeds contain a significant amount of starchy endosperm tissue that serves as a source of nutrients for the germinating seedling.

Two important regions of the mature embryo are the root and the shoot apical meristems. The entire shoot system (stems, leaves, and flowers) of the adult plant forms from cells that are located in the shoot apical meristem of the mature embryo. In some species of angiosperms, this embryonic meristem

is composed of only a few cells, whereas in other species several true leaves are produced by the shoot apical meristem before the completion of embryogenesis. The root apical meristem that is formed during embryogenesis becomes active during the early stages of germination and ultimately produces the entire root system of the adult plant. *See* APICAL MERISTEM; ROOT (BOTANY).

The final stages of embryogenesis in angiosperms include maturation, desiccation, and preparation for seed dormancy. Lipids and storage proteins are produced at high levels during embryonic maturation, stored in membrane-bound vesicles, and utilized as a source of nutrients for the young seedling. Immature embryos are often capable of precocious germination in culture but become dehydrated within the fruit and usually germinate only when exposed to the appropriate environmental conditions. Release from seed dormancy may require abrasion of the seed coat, exposure to light, or cold treatment, in addition to sufficient moisture. *See* EMBRYOGENESIS.

Different patterns of embryo development are found in gymnosperms and in the more primitive vascular and nonvascular plants. Double fertilization and the development of a nutritive endosperm tissue are features unique to the angiosperms. The haploid microgametophyte (germinating pollen grain) in most gymnosperms contains two male gametes, but only one of these participates in fertilization. The nutritive function served by the endosperm tissue in angiosperms is served in gymnosperms by the large haploid megagametophyte. Early divisions of the zygote are also different in gymnosperms; the zygote typically undergoes a series of free nuclear divisions during the earliest stages of embryogenesis, and multiple embryos often arise from a single zygote through a process known as polyembryony. Even more striking differences in embryogenesis are found in ferns and mosses, where the haploid or gametophytic phase of the life cycle is much more extensive.

Several major differences also exist between embryogenesis in plants and animals. Plant cells are surrounded by a cell wall that limits the contact and movement between adjacent cells. Embryogenesis in plants therefore proceeds without the morphogenetic movements that are characteristic of animal development. Morphogenesis in plants is also not limited to embryo development, but occurs throughout the life cycle. The mature plant embryo is therefore not simply a miniature version of the adult plant. *See* PLANT MORPHOGENESIS.

Methods of study. Embryogenesis in plants has been approached through a combination of morphological, experimental, biochemical, and, more recently, genetic and molecular studies. Most of the basic anatomical features of reproductive development in angiosperms were elucidated through the observations of botanists in the nineteenth and early twentieth centuries. Some of the most detailed histological studies have dealt with *Capsella bursapastoris* (Cruciferae). Early stages of embryo development in *Capsella* are characterized by rapid growth of the endosperm tissue and the formation of a

proembryo which is composed of a filamentous suspensor and a small embryo proper. The suspensor is metabolically active during early stages of embryogenesis and reaches a mature size of 6–8 cells before degenerating during later stages of development. The embryo proper begins as a small sphere of cells that rapidly divides and becomes differentiated as it progresses through the globular, heart, torpedo, early cotyledon, and mature cotyledon stages of embryogenesis. The mature embryo contains two large cotyledons, a single hypocotyl, and both root and shoot apical meristems.

One of the most common experimental approaches to the analysis of embryo development in *Capsella* and other angiosperms has been the use of embryo culture. Immature embryos at the appropriate stage of development can be removed from the developing seed and cultured to maturity on a defined nutrient medium. This technique has been used not only to study the regulation of embryogenesis in plants but also to rescue aborted embryos produced by interspecific hybridization. It has also been possible in certain species of angiosperms to produce somatic embryos by growing cells derived from roots, leaves, or other sporophytic tissues on a defined nutrient medium with the appropriate growth regulators and then transferring the undifferentiated mass of cells onto a medium that supports the formation of nonzygotic embryos. This technique has demonstrated conclusively that the zygote is not the only plant cell capable of forming a normal embryo. Haploid plants have also been produced by culturing microspores on the appropriate nutrient medium. The immature microspores, in this instance, alter their pattern of development to form embryos instead of pollen grains.

The regulation of gene expression during embryogenesis in flowering plants has been approached through a variety of molecular and genetic techniques. Mutants defective in different aspects of embryo development have been analyzed, and many genes that are known to be expressed during embryogenesis are now being cloned and sequenced through the use of recombinant DNA technology. The genes that code for the major seed storage proteins have been studied in detail because they are transcribed into messenger ribonucleic acid at very high rates during later stages of embryogenesis. When more is learned about how the expression of these genes is regulated, it may be possible to alter the concentration and chemical composition of these storage proteins to improve the nutritional quality of seeds produced by important crop plants. It should also be possible to learn much more about how genes control the complex series of interactions that occur as the fertilized egg develops into a mature embryo. See DEVELOPMENTAL BIOLOGY; DEVELOPMENTAL GENETICS; REPRODUCTION (PLANT).

David W. Meinke

Bibliography. J. D. Bewley and M. Black (eds.), *Seeds: Physiology of Development and Germination*, 2d ed., 1994; L. W. Browder et al., *Developmental Biology*, 3d ed., 1997; B. M. Carlson, *Plant's Foundations of Embryology*, 6th ed., 2000;

R. E. Coalson, *Embryology*, 2d ed., 1992; E. H. Davidson, *Gene Activity in Early Development*, 3d ed., 1986; G. M. Edelman, *Topobiology: An Introduction to Molecular Embryology*, 1993; V. Raghavan, *Embryogenesis in Angiosperms*, 1986; G. C. Schoenwolf and W. W. Mathews, *Atlas of Descriptive Embryology*, 6th ed., 2003; T. A. Steeves and I. M. Sussex, *Patterns in Plant Development*, 2d ed., 1989; F. H. Wilt and S. C. Hake, *Principles of Developmental Biology*, 2004; S. J. Wright, *A Photographic Atlas of Developmental Biology*, 2005.

Embryonated egg culture

Embryonated eggs are among the most useful and available forms of living animal tissue for the isolation and identification of animal viruses, for titrating viruses, and for quantity cultivation in the production of viral vaccines. The embryo proper, chorioallantoic membrane, yolk sac, allantoic sac, or amniotic sac may be inoculated in hens' eggs of various ages, so that a wide choice of types of tissue is available to fit the characteristics of the virus under study or for special studies. The chorioallantoic membrane is frequently used; in some infections, such as smallpox, vaccinia, and herpes simplex, characteristic lesions are produced which in some cases may resemble those in the natural host. For example, smallpox virus when cultured on the chorioallantoic membrane produces pocks and typical inclusions within the infected cells. When the embryo is inoculated, characteristic skin eruptions appear. Influenza virus, however, when inoculated into the amniotic cavity, does not give rise to pathology like that of the natural infection. See HERPES; INFLUENZA; SMALLPOX.

The growth of influenza virus in the allantoic cavity of embryonated eggs is probably the best-known use of eggs for animal virus culture. After inoculation into the allantoic cavity, the number of virus particles decreases exponentially, with almost 50% adsorbed to cells within 1 h. After adsorption the virus particles disappear (presumably by penetration of the cells), and an "eclipse period" of about 4 h follows. Evidence of viral multiplication first appears in the middle of the eclipse period, when soluble complement-fixing antigen (noninfective and non-hemagglutinating) is detectable. Hemagglutinating antigens can be found 3–4 h following inoculation, but in no greater amount than in the original inoculum; increases in hemagglutinating activity are seen only after 4 h or more. Infective virus is detectable 5–6 h after the allantoic cavity is inoculated; it is released almost at once into the allantoic fluid. Infective virus particles thus released then infect other cells of the allantoic sac, and virus continues to accumulate in the allantoic fluid for about 48 h, after which the system is exhausted. See ANIMAL VIRUS; ANTIGEN; COMPLEMENT-FIXATION TEST; FETAL MEMBRANE.

Joseph L. Melnick

Bibliography. D. O. White and F. Fenner, *Medical Virology*, 4th ed., 1994.

Embryonic differentiation

The process by which specialized and diversified structures arise during development of the embryo. The process involves (1) an increase in the number of cell types, and (2) an increase in morphological heterogeneity through the arrangement of cells into increasingly complex structural patterns in the form of tissues and organs. *See* HISTOGENESIS.

Differentiation begins in most organisms with fertilization of an egg with a sperm, after which the relatively large egg divides into many smaller cells called blastomeres. The blastomeres receive unequal portions of the cytoplasmic materials of the egg and are therefore initially somewhat different from each other. At the end of cleavage, the blastomeres are organized into a blastula, commonly either a hollow ball of cells or a flattened two-layered disk of cells. The cells of the blastula lie in different relative positions from those that will be occupied by their descendants in the adult organism. By a process known as gastrulation, they move to their approximate final positions and are arranged into three basic layers, called germ layers.

However, only two layers form in the simpler multicellular organisms. The outer layer is the ectoderm, from which arise the nervous system and the epidermal layer of the skin. The innermost germ layer, the endoderm, forms the epithelial lining of the digestive tract and contributes the essential tissue of associated organs. In all but the most primitive animals a third germ layer, the mesoderm, is formed by cells which come to lie in the area between the other two layers. In higher animals the mesoderm gives rise to most of the cells of the organism, such as those found in the muscles, skeleton, blood, connective tissue, kidneys, gonads, and certain other organs. The molding of groups of embryonic cells into such diverse tissues and organs proceeds through a variety of morphogenetic processes, such as migration, aggregation, dispersion, delamination, folding, and differential local growth of cells. *See* BLASTULATION; GASTRULATION; GERM LAYERS.

Cellular differentiation. Underlying the visible structural diversification of the embryo is the more fundamental and concomitant process of cellular differentiation (chemodifferentiation), by which embryonic cells are transformed into the highly specialized cells of the adult. A characteristic feature of this process is the production of manifold kinds of cells, all derived from a single precursor cell. The appearance of these new types or strains of cells is not abrupt but is the consequence of a long chain of progressive transformations in the molecular composition and organization of the cell. During these transformations each cell gradually loses its capacity to develop in alternative directions. At the same time it acquires characteristics essential for further differentiation along its prospective pathway.

The variety of adult cells produced by the divergent pathways of differentiation varies enormously in different kinds of animals, but in each individual the cells cooperatively discharge, in an ordered man-

ner, the manifold functions of the organism. Simpler organisms are constructed of fewer kinds of cells, some of which seemingly perform a wider variety of functions than do the cells of more complex organisms. However, even the most specialized cells carry on a multitude of different chemical activities. Depending upon the criteria used, hundreds of different cell types may be recognized in the more complex animals. As methods of observation and analysis become more refined and penetrating, groups of apparently similar cells become resolvable into distinct types, in terms either of their functions or of their constituents.

Special synthetic products commonly distinguish many cell types, such as melanin in melanocytes, actomyosin in muscle cells, and hemoglobin in erythrocytes. But even these cell types may be further subdivided. Melanocytes of the skin are readily distinguished from those of the retina; and muscles of the heart, limbs, and digestive system each constitute distinctive cell strains. Moreover, the erythrocytes of the embryo have hemoglobin which is distinguishable from that found in adult erythrocytes. *See* CELL DIFFERENTIATION; CELL SENESCENCE; HEMATOPOIESIS.

Cell structure. A cell is an exceedingly complex structure, and presents numerous possibilities for alteration and diversification. Many of these possibilities have been realized in the specialized cells of adult organisms. In the center of the cell is the nucleus containing the chromosomes, which are responsible for determining the hereditary properties of the cell and for setting the limits of its developmental capacity. During cell division the chromosomes are duplicated and distributed in equal sets to the daughter cells; thus the cells of an organism are initially equipped with identical sets of chromosomes, one set derived from the egg, the other from the sperm. In each cell only a limited portion of this genetic endowment functions at any one time, and the most fundamental differences among cells stem from different sets of active genes. The genes are composed of deoxyribonucleic acid (DNA) that carries the code for synthesizing the proteins and enzymes that largely determine the properties of the cells of an organism. The more complicated the organism, the more DNA it requires: mammalian cells, for example, contain about 1000 times as much DNA as do bacterial cells. However, there are notable exceptions: among fish species the DNA content varies more than 150-fold and among amphibians more than 10-fold, yet these organisms are very similar in complexity. The significance of the vastly different amounts of DNA in the cells of these organisms is not understood at present. *See* DEOXYRIBONUCLEIC ACID (DNA).

During the differentiation of the cells of any organism the genetic code inscribed in the structure of active genes is transcribed into messenger ribonucleic acid (mRNA) molecules in the nucleus of the cell. These mRNA molecules then pass into the cytoplasm where they are translated into specific proteins or enzymes. The translation process requires two additional kinds of ribonucleic acid (RNA):

ribosomal RNA organized together with protein into small organelles called ribosomes, and transfer RNA responsible for carrying the amino acids from which proteins are made to the site of protein synthesis on the ribosomes. *See* GENETIC CODE; RIBOSOMES.

The activity of genes is dependent upon the chemical composition and physical structure of the chromosome. Morphological expressions of gene activity have been observed in some organisms: in certain dipteran insects, enlarged puffs appear at sites along the giant chromosomes during intense gene activity, while in amphibian oocytes, metabolically active loops extend laterally from the chromosomal axis. The molecular events underlying these phenomena are presently unknown, although at least two kinds of protein appear to be involved: the basic histones and the more acidic proteins. The histones are present in relatively constant amounts along the chromosomes and have been shown to inhibit gene-directed synthesis of ribonucleic acid. Chromosomes containing active genes also contain increased amounts of acidic or nonhistone proteins. These proteins may be involved in removing the histone inhibition of genes. *See* CHROMOSOME; MITOSIS; RIBONUCLEIC ACID (RNA).

In addition to the chromosomes, other constituents of the nucleus, the nucleoli, nucleoplasm, and the enclosing membrane also exhibit differences in various cell strains. However, the most conspicuous distinguishing characteristics of cells are found in the cytoplasm. In the cytoplasmic matrix are embedded various organelles, such as mitochondria, ribosomes, endoplasmic reticulum, centrosome, the Golgi complex, and a host of enzymes and other chemical substances, which by their proportions or type characterize the cell containing them. *See* CELL (BIOLOGY).

These components of the cell are organized into dynamic integrated patterns that confer on the cell the capacity to function and multiply. Moreover, once a terminal state of differentiation in the cell has been reached, it is conserved and perpetuated. Differentiated cells commonly do not multiply; but when division occurs under proper environmental conditions, the differentiated characteristics are preserved throughout descendant cell generations. Diverse cell strains are organized into still more complex patterns in the form of various tissues, which in turn are built up into organs. Tissues and organs vary greatly in the complexity of organization of their constituent cells but are generally composed of populations of heterogeneous cells held in rigid patterns by their mutual affinities, antagonisms, or both. The structure and function of tissues and organs undoubtedly are dependent upon the properties of the component cell; likewise, each cell is regulated in its function and course of differentiation by the cell population of which it is a part.

Mechanisms of differentiation. The mechanisms by which the course of cellular differentiation is realized are not precisely known. The factors involved may, however, be divided into two classes: (1) intrinsic, those operating within the cell, and (2) extrinsic,

those brought to bear upon the cell from outside. Both classes of factors play a role in the differentiation of every cell. However, the relative importance of these factors varies considerably from one cell strain to another and also within the same cell at different stages in its development.

The fertilized egg begins development with a rich endowment, consisting of a nucleus with a set of paternal and maternal chromosomes together with a complexly organized cytoplasm. The activation of the egg by the sperm sets off a chain of actions and reactions that progressively transform the physical and chemical constitution of each descendant cell. The emergence of new cell characteristics may be attributed to an oscillating interaction between the intrinsic gene makeup of the cell and the surrounding cytoplasm. The dynamic imbalance existing between these interacting components drives the cell along its path of differentiation. In certain kinds of invertebrate embryos, interactions within each separate cell seem sufficient for guiding differentiation to its terminal state. Such embryos exhibit mosaic development. By contrast, in the embryos of vertebrates and certain invertebrates such as echinoderms, influences from adjacent cells are an essential part of the differentiation process. These embryos show regulative development.

Embryos showing mosaic development, in which cells differentiate autonomously, seem to be the rule among such groups as tunicates, mollusks, and annelids. In these organisms the destruction of a blastomere during cleavage results in a corresponding defect in later stages of development; that is, the injury is not repaired. Moreover, isolated blastomeres of such embryos tend to continue development as if they were still a part of an intact embryo. Throughout differentiation, the cells of these embryos are fixed in their developmental capacity and are able to differentiate in only one direction even when placed in abnormal surroundings. This type of differentiation is the product of an unfolding sequence of reactions and segregation of activities that follows a course predetermined in the structure of the egg; each step leads to the next without the intervention of influences from adjacent cells.

On the other hand, in regulative development the differentiation of cells and tissues is initiated and directed by inductive influences emanating from adjacent cells or tissues. In addition to the guiding influence of one tissue upon another, the general physicochemical environment established within a population of cells serves to regulate the further differentiation of the constituent cells. These population effects are referred to as embryonic fields or gradients. In contrast to the cells of mosaic embryos, regulative blastomeres exhibit great developmental plasticity. When transplanted to new locations in the embryo, the cells respond to their new environment by developing along pathways appropriate to their new location. Such developmental plasticity declines as differentiation proceeds, and the variety of pathways open to a cell is continuously restricted until the terminal state of differentiation is reached. It

should be emphasized that all gradations between the extremes of mosaic and regulative development exist and that even the tissues of highly regulative embryos in later stages of development exhibit increasingly autonomous differentiation. In all cases, however, differentiation is a gradual, progressive process resulting in the formation of specialized cells from generalized precursor cells. Cells, when fully differentiated, are very stable, and only under exceptional conditions of growth in tissue culture or during regeneration of injured tissues do cells lose their differentiated attributes and acquire the characteristics of other mature cells. *See* CELL LINEAGE; DEVELOPMENTAL BIOLOGY; EMBRYONIC INDUCTION.

Clement L. Markert

Embryonic induction

In the early development of many tissues and organs of complex, multicellular organisms, the action of one group of cells on another that leads to the establishment of the developmental pathway in the responding tissue. The groups of cells which influence the responding cells are termed the inducing tissue. Since specific inducing tissues cannot act on all types of cells, those cells which can respond are referred to as competent to react to the action of a specific inducer stimulus.

Embryonic induction is considered to play an important role in the development of tissues and organs in most animal embryos, from the lower chordates (for example, *Branchiostoma*) to the higher vertebrates (for example, mammals).

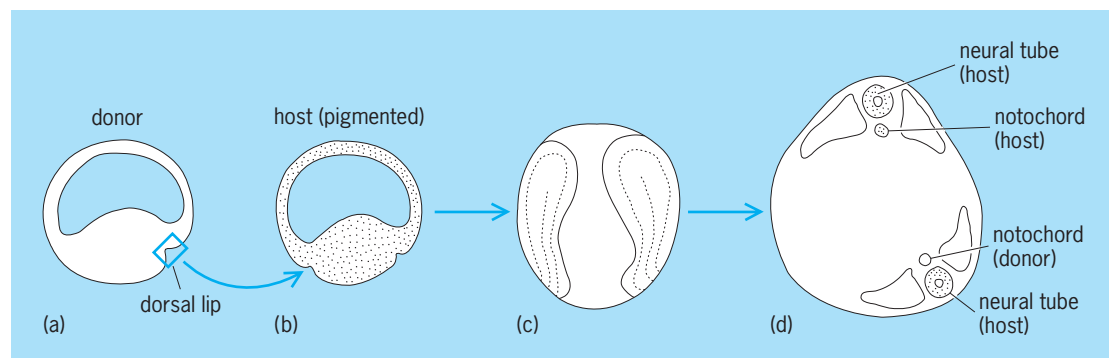
Primary induction. Perhaps the first major induction phenomenon occurs during the final stages of gastrulation of most animal embryos. Following fertilization, the egg divides to form a multicellular blastula-stage embryo. The cells of the blastula then undergo a series of movements which generate a more complex embryo, the gastrula, which contains three major groups of cells: ectoderm, mesoderm, and endoderm. The mesoderm actually arises

as cells move from the surface of the embryo to the inside. Once inside, they induce the cells which reside over them, the surface ectoderm cells, to develop into the neural tube. The neural tube eventually forms the central nervous system. The first induction event of early embryogenesis is called primary embryonic induction (see **illus.**). The migratory cells which invaginate from the surface and induce the development of the neural tube are termed the embryonic organizer. The first step in the sequence of events termed primary embryonic induction is the acquisition by the mesoderm of neural inducing activity. Proteins such as fibroblast growth factor and activin, which belong to a category of so-called peptide growth factors, play key roles in programming the mesoderm cells to induce overlying ectoderm to differentiate into neural structures. *See* GASTRULATION.

The development of a large number of tissues and organs is influenced by embryonic inductions. Various eye structures (lens, optic cup, and so on), internal ear structures, as well as several tissues (for example, vertebral cartilage) emerge from cells which were acted upon by inducer tissues.

Sequential inductions. Following primary embryonic induction, the developing neural tube is influenced by the surrounding tissues, especially mesoderm, to undergo regional specialization. Noggin is a protein which is secreted by some cells, and can be purified from amphibian embryos. When added to organ cultures of ectoderm cells that have been separated from associated mesoderm, noggin induces the explant to develop neural characteristics.

Once brain tissue is induced, one region, the hindbrain, induces the overlying ectoderm tissue to develop into the ear placode. Later the ear placode induces surrounding tissue to develop into other structures of the ear. Likewise, one region of the forebrain makes contact with another area of the ectoderm and induces it to form a lens rudiment. The lens in turn exerts an inductive influence on the forebrain, causing it to change shape to form an optic cup. Finally, the optic cup and lens induce overlying tissue to



Experimental demonstration of primary embryonic induction in amphibian embryo. (a) Mesoderm tissue surrounding dorsal lip of donor embryo was grafted to (b) ventral side of pigmented host embryo. (c) External development of host reveals two neural plates. (d) Histological examination of later-stage embryo reveals two notochords as well—one developed from host cells, the other from donor cells. Both neural tubes were, however, of host origin. The donor tissue (for example, notochord) apparently induced the development of the second neural tube.

develop into the transparent cornea. Should any single step be omitted, or should actual cell-cell contact between inducing and responding tissue be prevented, subsequent events are adversely affected. Thus, the development of many organs depends upon a set of sequential, and in some cases (such as lens and optic cup) reciprocal, induction. *See* NERVOUS SYSTEM (VERTEBRATE).

Genetic influences. The inducing stimulus, whether from an intact tissue or a cell extract, can act only within the constraints of the responding tissue's genetic constitution. For example, the tissue which induces hair development in mammals can induce feather (but not hair) development in birds, as well as scale development in reptiles. *See* DEVELOPMENTAL GENETICS.

Underlying mechanism. Limbs, kidney, nasal structures, salivary glands, pancreas, teeth, feathers, and hair are organs which require inductive stimuli. It is not known whether a single common mechanism underlies each of those inductions. In some cases, such as lens induction, an instructional role for the inducing tissue (forebrain) is postulated. In other cases (such as the pancreas) a permissive role for the inducing tissue is imagined. In this role the inducer is required only to promote the successful self-differentiation of the responding tissue; that is, the inducer does not directly influence the choice of the developmental pathway.

Many scientists believe that inductive interactions are mediated by cell-cell contacts; that is, the developmental information which is transferred from the inducing tissue is thought to reside at the cell surface of that tissue. Perhaps the surface of the responding tissue recognizes the signal molecules present on the surface of the inducing tissue.

In other instances, a secreted protein such as noggin might move among various cells or tissues and exert its effects on competent cells.

Plant morphogenesis. The principles of animal development also apply to plants. A greater role is, however, usually played by the diffusion of small-molecular-weight signal molecules rather than cell-cell contacts or protein growth factors. The earliest stages of plant embryo development involve groups of cells acquiring the competence to respond to inductive signals. Later in development, inductive signaling also becomes important. For example, in flowering plants the distance between nodes along the stem elongates, and lateral buds form below the shoot apex. The buds are believed to develop in response to a concentration gradient of signal molecules which exists along the stem. Thus, a process which is analogous to embryonic limb bud formation in animals is played out, and both plant and animal inductions can be conceptualized in similar terms. *See* CELL DIFFERENTIATION; DEVELOPMENTAL BIOLOGY; EMBRYOLOGY; PLANT MORPHOGENESIS.

George M. Malacinski

Bibliography. S. F. Gilbert, *Developmental Biology*, 6th ed., 2000; J. P. Trinkaus, *Cells into Organs: The Forces That Shape the Embryo*, 2d ed., 1984.

Emerald

The medium- to dark-green gem variety of the mineral beryl, $\text{Be}_3\text{Al}_2\text{Si}_6\text{O}_{18}$, crystallizing in the hexagonal system. A flawless emerald with good color is one of the most sought after and highly prized of all precious gems. Emerald is restricted in its occurrence, and seldom are exceptional stones found; most emeralds are flawed and cloudy, and few stones command high prices. What constitutes a true emerald is well defined by the gem dealers, and the dividing line between emerald and mere green beryl is sharp.

Occurrence. In contradistinction to beryl and its other gem varieties, aquamarine,morganite, and goshenite, which almost exclusively occur in granite pegmatite druses, emeralds have only been found in mica schists or metasomatized limestones. The most outstanding occurrences include the Muzo and El Chivor mines in Colombia. Here emeralds occur in calcite veins associated with a dark limestone which has been metasomatized by magmatic solutions. Noteworthy occurrences in mica schists include Tokovoja in the Ural Mountains, where emerald occurs with the beryllium minerals chrysoberyl (and its gem variety alexandrite) and phenakite; Habachtal, Austria; Transvaal, South Africa; and Kaliguman, India. The ultimate source of an emerald can often be assessed by a study of its inclusions.

Synthesis. Exceptional emeralds have been synthesized that rival the finest natural stones in quality and color. Broadly speaking, there are two main synthesis techniques. The molten-flux technique involves the compound oxides or gels of approximate beryl composition mixed with fluxes such as lithium molybdate and vanadium pentoxide. Fusion followed by slow cooling results in the growth of small prismatic crystals. The hydrothermal synthesis technique is capable of producing exceptional stones; emerald is hydrothermally grown on the plate of cut aquamarine which is used as a seed. The grown emerald is then sawed and used in turn as a seed for other emeralds. In all techniques 0.05–1.4% Cr_2O_3 is added. As in the natural emeralds, the Cr^{3+} chromophore is responsible for the green color. Detailed studies have been undertaken to define means of distinguishing natural emeralds from synthetics. Careful study (involving determination of birefringence and infrared absorption spectrum) will not only reveal whether a stone is natural or synthetic but will also distinguish the mode of synthesis if the emerald proved to be a synthetic stone. *See* BERYL; GEM.

Paul B. Moore

Bibliography. P. C. Keller, *Gemstones and Their Origins*, 1989; F. Ward, *Emeralds*, 2d ed., 2001.

Emergency medicine

The medical specialty that comprises the immediate decision making and action necessary to prevent death or further disability under emergency

conditions. It is based primarily in hospital emergency departments, but with extensive responsibilities for supervising emergency medical systems outside the hospital (paramedics).

Prehospital care. Prehospital care is professional emergency care delivered before the patient reaches the hospital. Before the 1960s there was almost no such care and ambulances merely provided transport. Care is provided by paramedics and emergency medical technicians. Paramedics are technicians trained in techniques such as electric shock for defibrillation of the heart, insertion of a tube into the trachea to assist respiration, and the use of intravenous drugs normally administered by physicians. Although possessing less training than registered nurses, they are authorized to carry out more procedures, and often function with less direct supervision. They operate under the supervision of an emergency physician, usually by radio communication. Emergency medical technicians (EMTs) are trained in first aid, give no medication, seldom communicate with a doctor, and often receive very little medical supervision. Paramedics are now found throughout the United States and are usually associated with ambulances and prehospital emergency care, but most such care is actually provided by emergency medical technicians.

The Federal Emergency Medical Services Act of 1974 established emergency medical service (EMS) districts throughout the United States. These districts set up regulations for paramedics and emergency medical technicians, administer the base hospitals that supervise the paramedics, and help establish rules for the management of trauma, burn, pediatric, and other emergencies. The control of emergency care is vested in different agencies in different regions; it may be the health department, the fire department, a private hospital, an ambulance company, or another body. Very few emergency medical service districts actually monitor the quality of care in a scientific manner, as hospitals are legally required to do. Nor does any impartial organization set standards for the districts' performance, audit the districts, or collect statistics about them. As a result, the quality of prehospital care varies widely and frequently is low.

There is little scientific knowledge about what constitutes effective prehospital care. Rapid transport of a victim of major trauma to a trauma center has been well studied and proven beneficial, as has prompt electrical defibrillation in cases of cardiac arrest. Whether these elaborate and expensive systems provide any benefit for the overwhelming majority of prehospital patients who suffer from neither of these conditions requires further research.

Trauma care. Trauma is physical injury, whether minor or major; it does not include medical emergencies due to other causes. In the United States, trauma constitutes only about 15% of all emergencies, and that includes minor trauma such as broken bones or lacerations. Life-threatening trauma needing emergency surgery constitutes only about 1-2% of all emergencies. In major trauma (usually caused

by vehicular or industrial accidents, stabbing or gunshot wounds, or falls) the best care requires rapid transport to the hospital. However, such rapid transport is the wrong approach to the far more common medical emergencies, such as cardiac arrest, for which emergency treatment must begin at the site of occurrence.

Professional organizations, such as the American College of Emergency Physicians and the American College of Surgeons, have set standards for trauma centers, which are hospitals that treat a minimum number of serious trauma cases, and that meet the standards of these organizations in terms of staffing, experience, and medical outcome. A trauma center should belong to a systematic plan implemented by an emergency medical service region, which determines under what conditions an injured patient should be transported to a trauma center rather than the closest hospital emergency department. The purpose of such a plan is to ensure that only serious trauma is referred to trauma centers, which will provide the best care, but which should not be overloaded with the more common minor injury cases that can be well cared for in any properly staffed hospital emergency department.

Proper treatment of trauma requires prompt assessment by qualified paramedics, who may have to immobilize the neck to prevent spinal cord injury, protect the breathing passages, or administer intravenous fluid to replace blood lost by hemorrhage. Further evaluation and treatment by a qualified emergency physician must take place immediately upon arrival at the emergency department or trauma center. Head injuries account for a large percentage of trauma deaths in the United States; ensuring an adequate airway is crucial since the oxygen supply to the brain and other tissue must be maintained. Also, efforts must be made to control edema (swelling) of the injured brain, which leads to increases in intracranial pressure that can cut off blood supply and result in brain death. Fatalities in other types of trauma are due to conditions such as internal bleeding and punctured lungs (pneumothorax), which usually require surgery within hours. The administration of large quantities of intravenous fluid and blood has made it possible to save most patients who reach a trauma center. Before properly cross-matched blood is available, it has been traditional to maintain blood volume with 0.9% saline solution, but recent research has shown that 5% saline is much more effective; it is even helpful in the small volumes that can easily be administered before the patient reaches the hospital. Trauma care has also been revolutionized by the advent of computerized tomographic x-ray, which reveals details of intracranial and intraabdominal injury that previously could be determined only by exploratory surgery. Magnetic resonance imaging is an even more sophisticated noninvasive technique that can reveal subtle details of tissue density, blood flow, and anatomy that previously could not be discerned. *See* COMPUTERIZED TOMOGRAPHY.

Burn care. Burns are a very serious injury because the protective barrier of the skin is destroyed,

leading to massive loss of body fluids, development of shock, and invasion by bacteria. Patients with major burns are best treated in specialized burn centers and should be taken there promptly. Burn mortality has dropped significantly in recent decades because of prompt intravenous administration of very large volumes of fluid to replace body fluid losses, and the use of antibiotics. See BURN.

Many burns are caused by smoking in bed, and could be avoided by the marketing of available self-extinguishing cigarettes and by more widespread installation of smoke detectors.

Management of cardiac arrest. Cardiac arrest occurs when the heart suddenly goes into an abnormal rhythm so that it does not pump blood; it is the chief cause of death in heart attack, and causes 400,000 deaths a year in the United States. Cardiac arrest tends to occur in the first few hours of heart attack, often before the patient has even decided to go to the hospital, and sometimes is not preceded by any symptoms. If treated by defibrillation (electric shock to the heart) within the first 6 min, it is reversible. Otherwise, death is certain. Therefore, to save lives, emergency medical service programs must deliver paramedics to the patient's side within 6 min of the arrest, not an easy task. Nonetheless, in cities where such programs succeed, and particularly where many citizens are trained in basic cardiopulmonary resuscitation (CPR) to keep the heart and brain tissues alive until professional help comes, about 25% of patients are resuscitated and return to a normal life. This fact spurred the standardization and widespread teaching of cardiopulmonary resuscitation and certification of nurses and physicians in this skill by the American Heart Association. Resuscitation of persons who have suffered cardiac arrest is a major advance that requires a public trained in cardiopulmonary resuscitation as well as public experienced and well-supervised paramedics. In most areas of the United States, too few of the key elements are in place, and survival rates from cardiac arrest are 10% or less.

Cardiopulmonary resuscitation includes the procedure of compressing the chest rhythmically 80 times a minute (in adults) to produce blood flow. Unfortunately, the amount of blood flow is a small fraction of normal and can keep heart and brain tissues alive only for a short time. Rescuers also must breathe for (ventilate) the patient; this can be done with a mask attached to a ventilating bag, or ideally, by placing an endotracheal tube in the patient's lungs and ventilating with 100% oxygen. The most successful cardiopulmonary resuscitation, and the usefulness of defibrillation, is limited to those patients who are in the abnormal rhythm known as ventricular fibrillation. However, about half of cardiac arrest victims are in other rhythms (asystole or electromechanical dissociation) for which no effective treatment is known.

Second to defibrillation, administration of epinephrine (to improve perfusion of the heart tissue during cardiopulmonary resuscitation) is the most useful treatment of cardiac arrest. This drug

may return perfusion of the heart to adequate levels long enough to allow the heart to start beating spontaneously again.

A revolutionary development was the production of automatic defibrillators that interpret heart rhythms electronically and administer electric shocks when (and only when) needed. There are about 33,000 paramedics in the United States who are trained to use defibrillators. However, one of the 450,000 emergency medical technicians can usually reach a victim much sooner than a paramedic. Automatic defibrillators allow fire fighters, emergency medical technicians, and others with very little training to treat patients in cardiac arrest.

Most heart attack victims have a blockage of the arteries supplying blood to the heart, which until recently was not promptly treatable. However, two new drugs—streptokinase and tissue plasminogen activator (tPA)—can permit rapid treatment. Both drugs dissolve clots rapidly and can be administered as soon as the patient arrives in the emergency department, saving heart tissue that would otherwise die from lack of blood. The earlier these drugs are given, the better the results, and researchers are studying the possibility of administering them before the patient reaches the hospital.

Unfortunately, there is no way to resuscitate patients who have heart rhythms that do not respond to defibrillation. An even greater problem is how to prevent the brain damage that often occurs from lack of blood supply during cardiac arrest or after resuscitation. Brain resuscitation is a serious problem because in many instances the heart is resuscitated but within hours a relatively healthy brain will swell from widespread cellular damage. Then pressure increases inside the skull, the blood supply is further worsened, and irreversible brain damage eventually occurs.

Treatment of poisoning. Innumerable poisonings occur every year in the United States, involving thousands of products. Many are accidental (particularly in children and confused elderly), but many result from deliberate attempts at suicide or abuse. The physician cannot stay abreast of all toxic syndromes and their constantly evolving treatment. Therefore poison centers have been established in every state, where trained physicians and pharmacologists equipped with electronic and on-line references provide emergency telephone consultation. Poison centers have greatly improved the quality of care in poisoning cases.

It used to be traditional in poisoning to administer emetics such as syrup of ipecac. However, this does not empty the stomach completely. For most poisonings it is preferable simply to administer activated charcoal, which binds the poison in the gut and prevents its absorption. This has simplified and improved the treatment of poisoning. Michael Callaham Bibliography. American College of Emergency Physicians, Guidelines for emergency department physician staffing, *Ann. Emerg. Med.*, 13:1165-1166, 1984; M. Callaham (ed.), *Current Therapy in Emergency Medicine*, 1987; P. Rosen et al. (eds.),

Emergency Medicine: Concepts and Clinical Practice, 2d ed., 1987.

Emery

A natural mixture of corundum with magnetite or with hematite and spinel. Emery has been used for centuries as an abrasive or polishing material. Because the mixture is very intimate and appears to be quite homogeneous, it was considered to be a single mineral species until the middle of the nineteenth century. The aggregate has a gray-to-black color and is extremely tough and difficult to break. The specific gravity varies from 3.7 to 4.3, depending upon the relative amounts of the constituent minerals. The hardness is about 8 (Mohs scale), less than that of pure corundum which is 9, and is more dependent upon the physical state of aggregation than on the percentage of corundum. *See* HEMATITE; MAGNETITE; SPINEL.

Since early times emery has been recovered from Cape Emery on the island of Naxos and from other islands in the Grecian archipelago. Here it occurs as irregular beds and lenses and in loose blocks associated with crystalline limestone and schists. It is also found at several localities in Asia Minor under similar conditions, notably at Gumach Dag, east of Ephesus, and at Kula, near Alashehr. Emery was worked during the latter part of the nineteenth century at Chester, Massachusetts, where it was associated with diaspore, margarite, and chloritoid. Because of its magnetic properties, resulting from the admixed magnetite, the Chester material was first worked unsuccessfully as an iron ore. Only after the similarity of the associated minerals with those of the Naxos emery was noted was its true nature determined.

Although synthetic abrasives have replaced emery in many of its earlier uses, it is still used as an abrasive and polishing material by lapidaries and in the manufacture of lenses, prisms, and other optical equipment. Emery wheels, emery paper, and emery cloth are used not only by lapidaries but also by machinists in the grinding and polishing of steel. *See* ABRASIVE.

Cornelius S. Hurlbut, Jr.

Emission spectrochemical analysis

An instrumental analytical technique used in qualitative or quantitative chemical analysis. It is conducted by monitoring and measuring the spectrum of light emitted by the material being analyzed. *See* ATOMIC SPECTROMETRY; QUALITATIVE CHEMICAL ANALYSIS; QUANTITATIVE CHEMICAL ANALYSIS.

In general, there are many ways in which to conduct an emission spectrometric measurement. The differences among approaches result from the choice of location within the electromagnetic spectrum at which to observe emitted radiation. Emission spectrochemical analysis, however, is the name given traditionally to those analytical determinations

based on radiation in the visible through vacuum ultraviolet region of the electromagnetic spectrum. The included wavelengths are approximately 800–150 nanometers. The technique is used principally to detect (qualitative analysis) and determine (quantitative analysis) concentrations of metals. A few non-metals can also be determined. Under optimum conditions, as little as 10^{-10} g of an element per gram of sample can be determined. Routine concentration ranges in which emission spectrometry is used are approximately 10^{-8} – 10^{-2} g per gram of sample (1% by weight).

The steps in emission spectrochemical analysis are vaporization and atomization of the sample, excitation of the atomic vapor produced, dispersion of the emitted radiation, and observation and measurement of the dispersed radiation.

A number of approaches to emission spectrochemical analysis have been developed, all of which follow the four steps listed but which differ fundamentally in the first two. These steps are involved with the transfer of energy to the sample so that the constituents of the sample will emit characteristic radiation. The application, control, and transfer of different forms of energy pose significantly different equipment requirements. Therefore, different versions of emission spectrometry have been developed. The principal modern versions of vaporization, atomization, and excitation are alternating-current (ac) spark, direct-current (dc) arc, inert-gas electrical plasmas [dc plasma (DCP), inductively coupled plasma (ICP)], and chemical flame.

Vaporization and atomization. First it is necessary to produce a vapor (vaporization) in which all the compounds of the sample are broken down into their constituent atoms (atomization). It is the presence of the individual atoms that is detected and quantitated. If the sample is originally a solid, it can be made part of an anode-cathode pair. When an electric discharge is struck across these electrodes, producing an ac spark or dc arc, the sample is vaporized into the electrode gap. There are numerous ways to configure an electrode so that the sample is an integral part of it. The most desirable designs make it possible to reproduce the procedure simply and easily, so that the technique is routinely applicable.

A more recent development for the production of a cloud of atoms from a solid sample is to fire a high-power pulse of laser radiation at the sample. The material that is vaporized when the laser creates a small crater is in just the physical state required for an atomic emission experiment. This approach has become more extensively used as the sample vapor generated has become more reproducible. The method is favorable for qualitative or semiquantitative analyses or in situations where it is desired to do sampling and analyses in the field.

Dissolved solids or liquid samples can be broken down into microscopic droplets in a process called nebulization, and a stream of these droplets can be injected directly into a discharge of hot gas. The hot gas then evaporates the solvent in the droplet, vaporizes the resulting solid microparticle, and atomizes

the compounds of which the particle is made.

A variety of techniques have been developed in which sample nebulization systems and components for generation of a stable hot gas work together. The earliest and most successful of these techniques involves nebulizing the sample and mixing it with the fuel and oxidant in a burner manifold; as a flame is generated, the sample is coincidentally vaporized and atomized. Many approaches for nebulizing the sample, mixing the droplets and gases, and configuring the burner manifold and burner head have proven successful.

Excitation of atomic vapor. In the majority of spectrochemical instrument systems, the same process is used both to produce the atomic vapor and to excite the atoms to cause emission of characteristic radiation. Energy is provided to a hot gas, the hot gas and sample are mixed to cause vaporization, and more of the energy of the hot gas is used to excite the atoms of vapor.

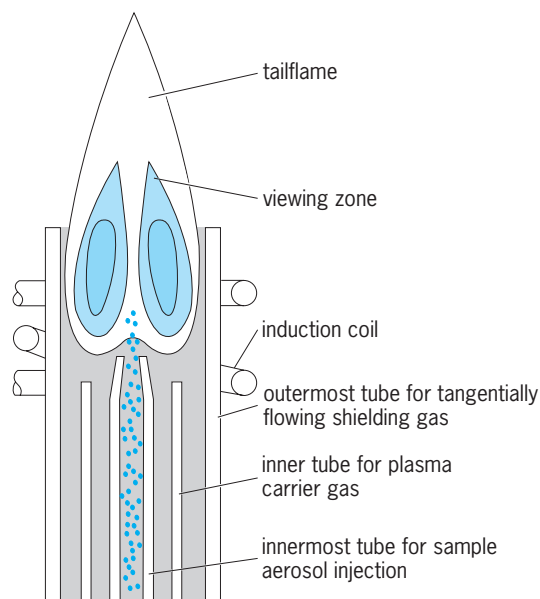
The hot gas is produced in a number of different ways. A chemical flame produces a hot gas from the combustion of some sort of fuel, for example, acetylene or nitrous oxide, with the release of many thousands of calories of heat. Chemical flames normally have temperatures between 1000 and 4000 K (1300 and 7000°F).

An electrical discharge, which is caused by applied potentials or applied currents across an electrode gap (ac spark or dc arc), produces a hot gas of atmospheric species, electrode material, and electrons. Temperatures up to about 7500 K (13,000°F) result.

More advanced types of electrical discharges, known as inert-gas electrical plasmas, produce hot gases of much higher temperature. A plasma is a cloud of vapor in which more than 1% of the atoms are ionized; that is, they have their outer electrons stripped away. In general, the greater the extent of ionization, the greater the amount of energy resident in the plasma.

Modern electrical plasmas are produced in two principal ways. Direct-current plasmas are electrode-based plasmas in which a conventional dc arc is ignited in an inert gas (usually argon) in an electrode gap by applying a current. Unlike most dc arcs, the arc in a direct-current plasma has a highly controlled flow of cooling gas (also argon) surrounding it. The cooling flow decreases the electrical conductivity of the outer regions of the arc, causing the applied current to be carried in a path of reduced volume. This is called a thermal pinch of the arc. As a result, the current density of the arc rises significantly, and this causes the temperature of the gas through which the arc is carried to increase. Analytical direct-current plasmas generate plasmas of argon with temperatures between 6000 and 8000 K (11,000 and 14,000°F).

Inductively coupled plasmas (a generic term for non-electrode-based plasmas) also are generated using argon. The inert gas reaches a plasma state through the interaction of electrons, carrier gas atoms, and an applied radio-frequency (kHz to MHz)



Inductivity coupled radio-frequency plasma.

or microwave-frequency (MHz to GHz) field of energy. Electrons will gain energy from the electric field established at an induction coil surrounding a flow of carrier gas (see **illustration**) or within a microwave-frequency resonant cavity through which the carrier gas flows. The electrons eventually gain extremely high kinetic energies. Although the electrons that gain this high energy have low mass individually, they are able to impart this energy through numerous collisions with the more massive argon carrier gas atoms. The carrier gas atoms eventually obtain a sufficient amount of energy to reach excited states, which takes only microseconds. At this point the plasma has ignited, and the carrier gas can release its energy by emitting radiation, losing electrons (ionizing), or colliding with other atoms (thermal transfer of energy). Temperatures routinely attained in induced plasmas are between 6000 and 8000 K (11,000 and 14,000°F). See MICROWAVE; PLASMA (PHYSICS).

Once the hot gas is produced, the vaporized and atomized sample is injected into it; or alternatively, the hot gas functions as both the vaporizer and atomizer, that is, cold sample is injected into it. In either case, the plasma is used to impart energy to the analyte atoms so that they may reach any number of excited states. Predominantly, collisional mechanisms of energy transfer contribute to the excitation of atoms. Other mechanistic pathways, such as resonance with applied fields and radiative transfer, probably occur. While more than one mechanism is followed as the analyte releases this energy, the most desirable route is that of radiational deactivation; that is, the analyte loses the energy by ejecting photons (emitting light).

Because the analyte atoms are allowed to attain a large number of excited states, a wide range of radiations (different wavelengths of light, photons of different energies) are emitted. The energy gained

and lost relates directly, through Planck's equation, $E = hc/\lambda$, to wavelengths between the infrared and vacuum ultraviolet regions of the spectrum. See PHOTOCHEMISTRY.

Dispersion and resolution. Identification of the element of interest and quantitative measurement of the radiation it emits requires that the light emitted be dispersed into its component wavelengths; that is, a spectrum of the light must be produced. This is done with a spectrograph using photographic detection or with a spectrometer using photoelectric detection. The instrument dispersion, its ability to separate wavelengths from one another, must be on the order of 0.01–0.1 nm/mm so that individual wavelengths are adequately separated at the detector. Prism spectrometers are no longer used since they do not have adequate dispersion over the entire wavelength range of interest. Diffraction grating instruments have mostly uniform dispersion and give quite high resolution. In general, spectrochemical analysis requires high spectral resolution—the ability to distinguish one dispersed wavelength from another. The reason is that the emission experiment can produce a larger number of spectral lines than other types of spectrometric approaches, such as fluorescence or absorption. See DIFFRACTION GRATING; SPECTROSCOPY; SPECTRUM.

Observation and measurement. Visual observation is not used in routine spectrochemical analysis, because the human eye does not possess any resolving power to sort out the complex spectra produced.

Photographic methods. Photographic observation has long been the method of choice for spectrochemical detection. It is used mainly for qualitative and quantitative analysis at the minor and trace constituent levels. Photographic glass plates are usually employed. Spectrograms are viewed with a projection comparator, which presents an image of both the sample and standard emulsion plates on a split-field screen, permitting any spectrum on one plate to be brought adjacent to and in register with any spectrum on the other. Any photographic emulsion is suitable for the 230–430-nm wavelength range. At higher wavelengths, special sensitization is needed. Emulsion sensitivity varies directly with the grain size in the emulsion. High sensitivity is obtained at the cost of resolving power, which varies inversely with grain size. Sensitivity and contrast vary with wavelength and age for all emulsions.

Measurement of line intensities requires calibration of the emulsion by constructing a relationship involving the blackening of the emulsion as a function of the intensity of the incident light. Lines whose intensities are to be compared photographically should be as close as possible in wavelength. Densitometers or microphotometers measure spectral-line blackening by scanning the illuminated spectrum with a fine slit or by projecting a fine illuminated line onto and through the spectrogram and onto a photocell.

The major advantage of photographic observation is that it is both the least expensive and the simplest method for obtaining the entire spectrum simultaneously and storing it nearly indefinitely. It is not the

most sensitive method and is not useful for ultratrace-element analysis.

Photoelectric methods. Photoelectric observation affords extremely sensitive and rapid detection. The electric current output of this detecting device is directly proportional to the rate at which photons strike a radiation-sensitive anode. This current is collected and processed by various electronic means, and is amplified to give an electrical signal that is proportional to the amount of analyte present in the sample vapor. Photoelectric detection is best suited to trace and ultratrace determinations because of its high sensitivity. However, it presents significant problems for simultaneous multielement analyses, since it needs one detector for each wavelength to be measured. Arrays of photomultiplier tubes have been constructed for simultaneous observation of 60 or more wavelengths. See PHOTOELECTRIC DEVICES; PHOTOMULTIPLIER; PHOTOTUBE.

Another approach to multielement detection is to use only one photomultiplier tube and to adjust the spectrometer precisely so that successive wavelengths pass by the detector in time. This is known as spectrometer scanning. Its principal disadvantage is that each wavelength is viewed sequentially, and so it takes longer to monitor more wavelengths. Its principal advantage is that there is no limit to the wavelengths which can be viewed; all of them are accessible.

Solid-state arrays. Another approach to multielement, simultaneous detection is to mimic the photographic plate with an electronic imaging device, such as a photodiode array or charge-transfer device. The simultaneity and simplicity of photographic detection is obtained in conjunction with the high sensitivity and speed of photoelectric detection. See CHARGE-COUPLED DEVICES; PHOTODIODE ARRAYS.

Solid-state imaging devices have become useful for spectrochemical detection. These devices are semiconductor circuits that are responsive to radiation. Most of them are made from silicon. If the silicon is held at a discrete electrical potential (a voltage is impressed across it), the absorption of an incident photon by the silicon crystal will cause an electron to migrate out of position (this is called forming a hole). The hole-electron pair formation is a charge-separation process. The positively charged holes are collected at the negative end of the crystal of silicon or stored at a capacitor. More incident photons make more holes and lead to increased charge collection. The amount of external current needed to neutralize the holes produced is measured and is proportional to the number of photons absorbed, constituting a measure of the radiation emitted by the sample.

Physically discrete photon-detecting units are formed when the semiconductor devices are fabricated. These units are called picture elements, or pixels. The action of each pixel is similar to that of a single photomultiplier. The pixel responds to the radiation it receives when it is located at a specific wavelength position on the focal plane of a spectrometer. Combinations of pixels, called arrays, can be configured to match the shape of the focal plane.

The pixels, then, are equivalent to a photographic plate, only with very large grain size.

One advantage of using solid-state arrays for detection is that large numbers of pixels can be fabricated on a single array. Detector arrays with many millions of pixels can come close to mimicking the exact performance of a photoplate. Another advantage is that arrays can be produced by cost-effective and rapid semiconductor processes, since the same technology used to make computer chips can be used to make photodetectors. Yet another advantage of solid-state array detectors is that they are addressed electronically, and this can be done rapidly with a computer. While arrays of photomultiplier tubes can also be addressed electronically, there are rarely more than a few dozen detectors used at any one time. Photographic plates have many more equivalent pixels on them than solid-state arrays, but it is orders of magnitude slower in retrieving their recorded signals. Solid-state arrays also have an advantage of photoplates in the ability to integrate signal; that is, the more photons that strike the pixels, the more signal is collected, so that large signals can be accumulated over time. *See* INTEGRATED CIRCUITS.

A disadvantage of solid-state arrays is that their pixel density determines the resolving power of the spectrometer. Since each pixel has a finite size, on the order of tens of micrometers, greater resolving power has to come through the use of larger and more expensive spectrometers. Finally, solid-state arrays of silicon cannot compete completely with photomultipliers with respect to radiation sensitivity. The signal range over which silicon responds to radiation will always be smaller than that of a phototube, making the arrays less useful for ultratrace determinations.

Andrew T. Zander

Bibliography. J. A. C. Broekaert, *Analytical Atomic Spectrometry with Flames and Plasmas*, 2002; M. Cullen, *Atomic Spectroscopy in Elemental Analysis*, 2004; J. R. Dean, *Practical Inductively Coupled Plasma Spectroscopy*, 2005; J. D. Ingle and S. R. Crouch, *Spectrochemical Analysis*, 1988; R. Kellner et al., *Analytical Chemistry: A Modern Approach to Analytical Science*, 2d ed., 2004; A. Montaser and D. W. Golightly (eds.), *Inductively Coupled Plasmas in Analytical Atomic Spectrometry*, 2d ed., 1992; J. Nolte, *ICP Emission Spectrometry: A Practical Guide*, 2003; F. A. Settle (ed.), *Handbook of Instrumental Techniques for Analytical Chemistry*, 1997; R. Thomas, *Practical Guide to ICP-MS*, 2004; H. H. Willard et al., *Instrumental Methods of Analysis*, 7th ed., 1988.

Emissivity

The ratio of the radiation intensity of a nonblackbody to the radiation intensity of a blackbody. This ratio, which is usually designated by the Greek letter ϵ , is always less than or just equal to one. The emissivity characterizes the radiation or absorption quality of nonblack bodies. Published values are readily available for most substances. Emissivities vary with tem-

perature and also vary throughout the spectrum. For an extended discussion of blackbody radiation and related information *see* HEAT RADIATION

There are several methods by means of which the emissivity can be determined. The one most commonly used is the cavity method. In this technique a fine hole is provided in a radiating surface, and the ratio of the radiation intensity from the surface to the radiation intensity from the hole yields the emissivity directly. This method is quite accurate. One can also use an optical pyrometer to determine the emissivity from the brightness temperatures of the hole and the surface in conjunction with Wien's law of radiation.

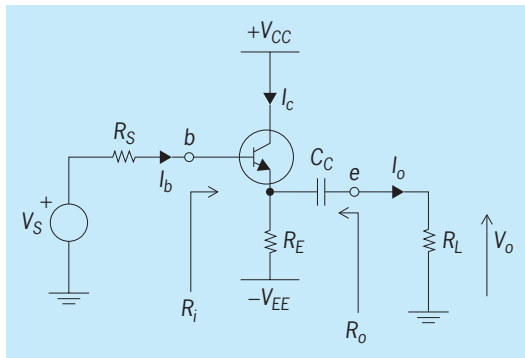
The total emissivity when introduced into the Stefan-Boltzmann law gives the total radiated energy W in joules per square centimeter of the real heat radiator as $W = \epsilon \sigma T^4$. Here T represents the absolute temperature and σ , the radiation constant, has the value 5.67×10^{-12} joule $\text{cm}^{-2} \text{K}^{-4}$. This energy is always smaller than the energy radiated by the blackbody, since ϵ is less than 1. For example, the total emissivity for tungsten is 0.32 at 2500°C (5032°F), which means that at the same temperature tungsten radiates approximately one-third the energy of a blackbody.

The spectral emissivity ϵ_λ (the subscript λ denotes the wavelength) provides information on the energy distribution. Any spectral emissivity value is valid only for a narrow wavelength interval. The wavelength at which ϵ_λ has been determined is indicated by a subscript, for instance, $\epsilon_{0.655}$. A spectral emissivity of zero means that the heat radiator emits no radiation at this wavelength. Strongly selective radiators, such as insulators or ceramics, have spectral emissivities close to 1 in some parts of the spectrum, and close to zero in other parts. Carbon has a high spectral emissivity throughout the visible and infrared spectrum, exceeding 0.90 in certain portions; thus carbon is a good blackbody radiator. Tantalum is the only metal with a spectral emissivity greater than 0.5 in the visible spectrum. All other metals have a lower spectral emissivity. Tungsten is a relatively good emitter, with a spectral emissivity of 0.43–0.47 within the visible region of the spectrum. *See* BLACKBODY.

Heinz G. Sell; Peter J. Walsh

Emitter follower

A circuit that uses a common-collector transistor amplifier stage with unity voltage gain, large input resistance R_i , and small output resistance R_o (see **illus.**). In its behavior, the emitter follower is analogous and very similar to the source follower in metal-oxide-semiconductor (MOS) circuits. Many electronic circuits have a relatively high output resistances and cannot deliver adequate power to a low-resistance load, or do suffer unacceptable voltage attenuation. In these cases, an emitter follower acts as a very simple buffer. Widely used, it is often found as the last stage of a multistage amplifier so that the circuit is better able to drive a low-resistance load. *See* AMPLIFIER.



Schematic diagram of an emitter follower circuit.

Operation. In the emitter follower, the collector is connected to the positive power supply V_{CC} , that is, to signal ground, so that the collector is common to alternating-current input and output (see illus.). The emitter is connected through a (large) resistor R_E to the negative power supply, $-V_{EE}$. The resistor R_E helps set the direct-current bias of the transistor. In integrated circuits, the function of R_E is normally replaced by a current source. The implementation makes no difference in principle in the following discussion: R_E simply represents the effective resistance of the current source. The emitter follower circuit, between input terminal b at the base and output terminal e at the emitter, is inserted between a load R_L and a source V_S with its source resistor R_S . The coupling capacitor C_C , if used, prevents R_L from influencing the transistor bias. The capacitance C_C is chosen large enough to act as a short circuit at signal frequencies. See CURRENT SOURCES AND MIRRORS; INTEGRATED CIRCUITS.

Characteristics. Under some simple practical assumptions, the input resistance R_i of the emitter follower is approximately $\beta + 1$ times the load resistor, $R_i \approx (\beta + 1)R_L$, where the parameter β is the transistor's current gain I_c/I_b . A typical value for β is 100. The same assumptions lead to the output resistance R_o being approximately equal to the source resistor divided by $\beta + 1$, $R_o \approx R_S/(\beta + 1)$. The voltage gain V_o/V_S is approximately equal to, but always slightly less than, unity (in practice 0.9 to 0.99) whereas the current gain is large: the output current I_o delivered by the emitter follower is approximately $\beta + 1$ times the input current I_b , $I_o \approx (\beta + 1)I_b$. In summary, the emitter follower is an impedance transformer. It effectively divides the source resistor and multiplies the load resistor by $\beta + 1$. It has approximately unity voltage gain but a large current gain, $\beta + 1$. Consequently, the emitter follower provides power gain. The power $V_o I_o$ delivered to the load is approximately equal to $\beta + 1$ times the power $V_S I_b$ drawn from the source, $V_o I_o \approx (\beta + 1) V_S I_b$.

Performance limits. An emitter follower has a very large bandwidth so that the simple equations in the previous paragraph, especially the constant near-unity voltage gain with no phase shift, describe the circuit's functioning well until very high frequencies. The ultimate limits are given by the parameters of the chosen transistor. The negative feedback caused

by R_E tends to reduce distortion in the emitter follower, but the signals must not get so large that the output current becomes equal to the bias current. Signal clipping will then result. When using an emitter follower, it must be observed that the forward-biased base-emitter diode imposes a 0.7-V direct-current voltage difference between the input (base) and output (emitter) terminals. Care must be taken to limit the current, should the output terminal be accidentally shorted to ground. See DISTORTION (ELECTRONIC CIRCUITS); TRANSISTOR.

Rolf Schaumann
Bibliography. P. R. Gray and R. G. Meyer, *Analysis and Design of Analog Integrated Circuits*, 3d ed., Wiley, New York, 1993; A. S. Sedra and K. C. Smith, *Microelectronic Circuits*, 4th ed., Oxford University Press, New York, 1998.

Emotion

An umbrella concept in the common language, typically defined by instantiation by reference to a variety of mental and behavioral states. These range from lust to a sense of liking, from joy to hostile aggression, and from esthetic appreciation to disgust. Emotions are usually considered to be accompanied by some degree of internal, frequently visceral, excitement, as well as strong evaluative components. Emotions are also often described as irrational, that is, not subject to deliberative cogitation, and as interfering with normal thought processes.

These latter qualities are often exacerbated in the emotional behavior and expression seen in clinical cases. The expression of strong emotions is typically considered to be symptomatic of some underlying conflict, and even the positive emotions are used as indices of unusually strong attachments and atypical earlier experiences. Sigmund Freud introduced the concept of repression to describe a defense mechanism against the occurrence of strong emotional experiences. From the psychoanalytic point of view, what is repressed is not the emotion itself, since the very concept of emotion implies conscious experience, but rather the memory of an event which, if it became conscious, would lead to strong conflicts and emotional consequences. Many other defense mechanisms, such as rationalization and compulsive or obsessive neurotic symptoms, are also seen as serving the purpose of avoiding conscious conflict and emotional sequelae. See PSYCHOANALYSIS.

Classical views. The conventional view of emotion in classical psychological and biological theory, at least to the end of the nineteenth century, was embodied in Charles Darwin's classic study of emotional expression in animals and humans. Emotional expression was described as dependent on some prior emotional mental state on the one hand, and on the evolutionary history of both the expression and the prior mental state on the other. Some situation "elicited" the mental state of fear, for example, and this state was then manifested in transspecies patterns of facial (and other bodily) expression. The cumulative evidence speaks for the existence of

universal facial and muscular patterns which are in fact recognized, across societies and cultures, as communicating similar evaluative and “emotional” messages. However, whether these expressions are evolutionary products of prior emotional states or whether they are the remnants of primitive preverbal communicative systems is open to question. It certainly is the case that the facial and gestural acts can occur without an accompanying emotional state and can frequently serve purely communicative, evaluative purposes.

The first major break with the traditional position occurred in the late nineteenth century when William James and Carl Lange proceeded independently to turn the classical sequence upside down. As James proposed in 1884: “My thesis is . . . that the bodily changes follow directly the *perception* of the exciting fact, and that our feelings of the same changes as they occur *is* the emotion.” This position required that for each discrete emotional experience there should be a prior discrete and different pattern of bodily (visceral and skeletal) changes. The lack of evidence for this and other aspects of the James-Lange position led the physiologist W. B. Cannon to a series of critical expositions which culminated in the thalamic theory of emotion, a neurophysiological, rather than psychological, position which has been influential in physiological approaches to emotion.

Neurophysiological theory. The thalamic theory became one of the bases for neurophysiological speculations and experiments about emotion. It suggested that emotions arise from the stimulation or excitation of areas in the thalamus but that these thalamic tendencies are normally inhibited by the activities of the cerebral cortex. Strong, innate stimuli bypass this cortical inhibition and act directly on the thalamus.

On the other hand, acquired emotions depend in the first instance on the recognition of, and memory for, certain classes of events that have “emotional” meaning. The memory for these events is activated in the cortex, and the inhibition by the cortex is released. When the thalamus is activated, signals from the thalamus serve two functions: they activate peripheral autonomic effects as well as cortical ones; the former produce the bodily expression of emotion, the latter its subjective experience. Removal of the cortex in animals does in fact produce hyperemotionality, but this persists even when the thalamus is removed. These studies and others have implicated the hypothalamus in the production of emotional experience and behavior. There is evidence of even wider areas of the brain in the emotional complex, with much emphasis placed on the role of the limbic system and the neocortex of the frontal lobes. See NERVOUS SYSTEM (VERTEBRATE).

Peripheral physiology. The essential aspect of subjective emotions appears to be the occurrence of visceral arousal. It is the peripheral autonomic nervous system that is primarily responsible for the reactions that are called visceral, with the sympathetic division focused on energy-expending func-

tion and the parasympathetic system devoted to energy-conserving functions. Cannon was principally responsible for noting the homeostatic function of the autonomic nervous system which “keeps the internal environment constant and fit for continued . . . action.” However, the human homeostat was conceived as being a passive respondent, just like the physical thermostat. Subsequent evidence has shown that autonomic, particularly sympathetic, changes may have very active consequences on the organism. They mobilize not only internal balance but also the scanning of and attention to environmental changes. The autonomic system not only prepares the internal environment for action but also prepares the whole organism for action on the external environment. See AUTONOMIC NERVOUS SYSTEM.

Within the context of the homeostatic view of the autonomic nervous system, the events that produce autonomic arousal are usually listed as physical stressors, events that threaten the physical well-being of the organism, such as tissue injury, pain, or loss of support. With the additional emphasis on organismic action, it is possible to define a set of events that produce autonomic reaction, in addition to those that depend on physical insult. The interruption of ongoing action, the disconfirmation of current expectations, the violation of perceptual expectancies, all generally subsumed under the interruption of thought, perception, and action, produce visceral (autonomic) arousal. Such arousal, in turn, generates additional attention and action. This version of the function of the autonomic nervous system also permits a view of the emotions as being adaptive, rather than merely annoying or disruptive. Emotional (visceral) arousal will occur for events that are important for the organism, that is, events that diverge from its expectations, past experiences, and past habits. Thus, the emotional experience emphasizes important events in the organism’s history, events that are different from the run-of-the-mill of everyday life. What is different is important, and what is important is underlined by the emotional experience. At the same time, the action of the autonomic nervous system makes it possible to cope better with these important events.

Cognitive/autonomic interactions. Modern views of emotion were significantly influenced by a series of experiments by S. Schachter and J. E. Singer which demonstrated the important coaction of cognitive (knowledge-oriented) and physiological (autonomic) events. They showed that in the presence of gross sympathetic arousal of unspecified origin, individuals will report subjective emotional experiences in keeping with their evaluations (cognitions) of their environment; that when the arousal can be ascribed to some specific cause, environmental evaluations are relatively ineffective; and also that emotional states are reported only to the extent that the individual experiences a state of physiological (visceral) arousal.

The experience of emotion is, therefore, the result of two different but complementary systems:

the arousal experience brought about by the autonomic nervous system, and a cognitive/evaluative state depending on the individual's current experiences, past evaluations of similar states, and the general accumulation of beliefs, values, and attitudes. Emotions arise when the current situation is different from that usually experienced, when it is experientially salient. The degree of divergence from the usual and the expected determines the intensity of the arousal and of the emotion. On the other hand, the current evaluation of the situation and the self determines the quality of the emotional experience. Thus, both positive and negative emotional states arise out of the combination of arousal and evaluation.

Within the complex interactions of cognitive and physiological factors, a variety of theoretical views can and have been accommodated. Thus, whether expressive or evaluative/communicative, the role of facial and gestural aspects of emotion are seen as part of the more general complex of the emotional experience. Biochemical and neurological factors that affect both cognitive and autonomic variables come into play, as well as pharmacological agents that exacerbate or suppress the state of arousal.

The perception of visceral arousal and the attention to the environmental origins and conditions of the individual's emotional state are factors that make for the apparent disorganizing effects of emotional experiences. The span of human attention or consciousness is known to be limited, and it is also known to be reduced by internal or external events that require some of that attentional capacity. Just as a loud noise or an extraneous thought reduces attentional capacity during a demanding task, so too does the internal "noise" of visceral activity. In addition, the external demands of an "emotional" situation intrude on a limited attentional capacity. As the attentional capacity is reduced during emotional states, the remaining attentional ability is focused on the aspects of the situation that the individual considers central. The judgment of centrality or importance is of course a subjective evaluation and may or may not be objectively correct. In moments of high emotional excitement or stress, it is frequently the case that attention is focused on events that may or may not be adaptive. In the latter case, emotion becomes disorganizing; in the former it is beneficial. Contemporary analyses of emotional states thus prevent any simple generalizations about emotions as being "good" or "bad" for the organism; they are properly seen as a component of the total adaptive complex of human thought and behavior rather than as a peculiar evolutionary hang-over. See CONSCIOUSNESS.

George Mandler

Bibliography. J. Callwood, *Emotions: What They Are and How They Affect Us*, 1986; M. Clynes and J. Panksepp (eds.), *Emotions and Psychopathology*, 1988; C. Darwin, *The Expression of Emotions in Man and Animals*, 1873; N. H. Frijda, *The Emotions*, 1987; C. E. Izard, *The Face of Emotion*, 1988; J. G. Thompson, *The Psychobiology of Emotions*, 1988.

Emphysema

The abnormal enlargement of air spaces distal to the terminal bronchioles in the lung. In the United States, the term emphysema is often used imprecisely to refer to any member of the group of illnesses known as chronic obstructive pulmonary disease, including emphysema, chronic bronchitis, and nonallergic asthma. The collective term obstructive pulmonary disease arises because many patients have elements of all three diseases. Although histological examination of lung tissue from surgical specimens or at autopsy is the only definitive diagnostic method, advances in physiological testing and radiologic techniques usually make an earlier diagnosis possible.

Types. There are two major anatomic patterns of emphysema, centriacinar or panacinar. In centriacinar emphysema, which is the form associated with heavy smoking, the grossly enlarged air spaces tend to appear in the respiratory bronchioles in the proximal part of the acinus, the unit of lung tissue; the upper portions of the lungs are most heavily involved. Histologically the lesion appears as a dilated or punched-out hole surrounded by relatively normal or only mildly thickened alveolar septa. In the panacinar type, the enlarged air spaces involve the entire acinus, from respiratory bronchioles to alveoli. Panacinar emphysema either favors the lower part of the lung or is diffuse, and is associated particularly with a deficiency of the glycoprotein α_1 -antitrypsin.

Physiology. The destruction of alveolar walls results in a loss of the elastic recoil of lung tissue and an increase in lung compliance, that is, a greater than normal lung volume with any distending pressure developed by the respiratory muscles. With progressive loss of elasticity, the thorax enlarges and the distending or transpulmonary pressure falls. The drop in pressure across the lung results in a loss of support of the airways, narrowing of the airways, and progressive airflow limitation. With such limitation, the ability of the body to respond to the metabolic demands of exercise by increasing total ventilation is reduced, and the individual's tolerance for physical activity gradually declines. Enlargement of air spaces, with its concomitant loss of capillary bed, causes less efficient gas transfer. The "wasted" ventilation can be quantified by measuring the transfer of tiny amounts of carbon monoxide into the capillaries during breath holding. The destruction of blood vessels causes a rise in pulmonary vascular resistance and increased pumping demand from the right ventricle. Initially the right ventricle responds by increasing thickness and strength of its muscle wall, a condition known as cor pulmonale. Eventually it is unable to handle the load, and dilates and goes into right ventricular failure with systemic venous congestion. See HEART DISORDERS.

Pathogenesis. Emphysema is characterized by loss of elasticity of the lung due primarily to the destruction of elastin, an important structural protein in the lung's supporting connective tissue. The elastin is

destroyed when there is an imbalance between enzymes that can damage tissue and those that block this effect. In response to a threat to the airways, such as infection or cigarette smoke, alveolar macrophages and neutrophils move into the lungs to digest the foreign material and in so doing release some of their rich supply of proteases. Without the neutralizing effect of antiproteases, the proteases can damage alveolar tissue and create emphysematous lesions. See CONNECTIVE TISSUE.

Both animal and human research has demonstrated the protease-antiprotease action in emphysema. Human emphysema has been replicated in the lungs of laboratory animals by the instillation of papain. The active ingredient in papain is elastase, a proteolytic enzyme that destroys elastin. Human susceptibility to emphysema increases in individuals that are homozygously deficient in the serum antiprotease α_1 -antitrypsin, which makes up about 90% of serum α_1 -globulin and is the most important antiprotease in humans.

Tobacco smoking causes a low-grade pulmonary injury that attracts inflammatory cells into the lung, possibly setting off the chain of events that progresses to emphysema. Although virtually all smokers have increased inflammatory cells that can be recovered from their airways, only a minority develop emphysema. Since the level of α_1 -antitrypsin does not distinguish those who will develop emphysema, other factors must be involved. In the absence of a smoking history, emphysema is distinctly unusual. Homozygous deficiency of α_1 -antitrypsin—perhaps other, as-yet-unrecognized proteins—result in premature emphysema, especially in the presence of even modest smoking. Certain occupational exposures, such as contact with cadmium, can also cause emphysema in the absence of smoking.

Diagnosis and treatment. The diagnosis of emphysema is usually suspected when a smoker develops symptoms of chronic shortness of breath and cough. However, the disease may be present for many years before a person becomes symptomatic. Physical examination rarely leads to early diagnosis, although a decrease in the intensity of breath sounds and an increase in chest volume may be apparent. Pulmonary function tests may be normal in the very early stages but will eventually show airflow limitation, an increase in total volume of the lungs, and a loss of diffusing capacity for carbon monoxide. Measurement of lung compliance is technically feasible but rarely performed, because it requires passage of a balloon into the esophagus to measure the pressure that develops at different lung volumes. Chest x-rays may show flattening of the diaphragm and hyperinflation of the chest, but more significant is the appearance of regions of increased lucency, which reflects the enlargement of air spaces and loss of pulmonary blood vessels. Advances in computerized axial tomography of the chest allow radiologic diagnosis of emphysema at a stage of disease progression that would not be visible on ordinary chest x-rays. See MEDICAL IMAGING.

Because tissue destruction in emphysema is irreversible, the key to treatment is prevention, which usually involves cessation of smoking as early in the course of the disease as possible. Even when the disease is well established, however, smoking cessation may slow its progression, reduce cough and sputum production (from the accompanying chronic bronchitis), and prevent the recurrent respiratory infections that seem to hasten the negative course of the disease. Because muscle tone in the smooth muscle encircling the bronchi may be increased, drugs that relax this muscle tone (bronchodilators) may be helpful, especially when administered as an aerosol directly into the lungs. Drugs known as sympathomimetics may mimic the effects of stimulation of the sympathetic nervous system, and others referred to as anticholinergics can block the cholinergic nervous system. In addition, since some degree of inflammation often accompanies emphysema, corticosteroid medications may be of some use.

Even if the course of the disease cannot be altered, general measures such as avoidance of irritants, vaccinations against influenza and pneumococcus, prompt treatment of infections, adequate nutrition, family support, regular exercise, and participation in multidisciplinary rehabilitation programs may substantially improve well-being and the ability to function independently. See LUNG; RESPIRATORY SYSTEM DISORDERS.

Michael Stulberg

Bibliography. B. Burrows, An overview of obstructive lung disease, *Med. Clin. N. Amer.*, 65:455-471, 1981; N. J. Gross and M. S. Skorodin, Role of the parasympathetic system in airway obstruction due to emphysema, *N. Engl. J. Med.*, 311:421-425, 1984; M. Higgins, Epidemiology of COPD: State of the art, *Chest*, 85(suppl.):35-75, 1984; T. L. Petty (ed.), *Chronic Obstructive Pulmonary Disease*, 1985; T. L. Petty and L. M. Nett, *Enjoying Life with Emphysema*, 1983; G. L. Snider et al., The definition of emphysema (report of the National Heart, Lung and Blood Institute, Division of Lung Diseases Workshop), *Amer. Rev. Respir. Dis.*, 132:182-185, 1985; W. M. Thurlbeck, *Chronic Airflow Obstruction in Lung Disease*, 1976.

Empirical method

The empirical method is not sharply defined. It is generally characterized by the collection of a large amount of data before much speculation as to their significance, or without much idea of what to expect, and is to be contrasted with more theoretical methods in which the collection of empirical data is guided largely by preliminary theoretical exploration of what to expect. The empirical method is necessary in entering hitherto completely unexplored fields, and becomes less purely empirical as the acquired mastery of the field increases. Successful use of an exclusively empirical method demands a higher degree of intuitive ability in the practitioner.

Percy W. Bridgman; Gerald Holton

Emulsion

A system of small liquid droplets dispersed in a second, immiscible liquid. Emulsions are a type of colloidal dispersion. By classical definition, emulsion droplets have diameters between 1 nanometer and 1 micrometer. However, for practical purposes, the principles of colloid science can be usefully applied to emulsions whose droplets are as large as tens or even hundreds of micrometers. The droplets in an emulsion are large enough that they do not behave like the atoms and molecules of classical chemistry. For example, emulsions do not generally behave like true solutions and may have undetectable freezing-point depressions. On the other hand, the droplets are small enough that they do not behave like the macroscopic particles of classical physics. For example, emulsion droplets may rise (cream) or fall (sediment) extremely slowly in apparent violation of Stokes' law. *See* COLLOID; SEDIMENTATION (INDUSTRY).

In most emulsions, one of the liquids is aqueous (that is, water containing dissolved substances), while the other is an oil (usually a hydrocarbon liquid). The two simplest kinds of emulsion are those in which oil droplets are dispersed in water (denoted oil-in-water, or O/W) and those in which water droplets are dispersed in oil (denoted water-in-oil, or W/O). In practice, emulsions can be much more complex and can consist of droplets within droplets. A double emulsion might contain aqueous droplets dispersed in oil droplets that are in turn dispersed in a continuous aqueous phase. Such an emulsion is termed water-in-oil-in-water, or W/O/W. The opposite kind of double emulsion is O/W/O. Double emulsion droplets can be quite large (tens of micrometers), and each can contain many droplets within it.

Preparation. If a small amount of vegetable oil is added to a cup of tap water and vigorously agitated, the oil will disperse into droplets: an emulsion has just been made. This emulsion will not be very stable, and in a short time the oil droplets will rise to the surface of the water and fuse (coalesce) together to produce a layer of oil floating on the water. Emulsions having a reasonable degree of stability contain an additional, emulsifying (stabilizing) agent. The emulsifying agents are usually surfactants or mixtures of surfactants and may include polymers or even fine solids. An ideal emulsifier has several properties. First, it makes the emulsion easy to create by reducing interfacial tension between the oil and aqueous phases. Second, it promotes creation of the desired kind of emulsion, whether O/W or W/O. Finally, it stabilizes the dispersed droplets against coalescence by providing an electrical and/or physical barrier at the droplets' surfaces. Beyond the oil and aqueous phases plus the emulsifier, chemical components may be added to an emulsion formulation to provide greater stability and to adjust other properties such as viscosity, color, and appearance. Properly formulated emulsions can remain stable for months or even years. *See* SURFACTANT.

Emulsions are prepared by the breaking-up of one liquid phase into droplets and then the further breaking-up, or comminution, of these droplets into smaller droplets. Typically, the emulsifying agent is dissolved into the phase where it is most soluble, after which the second phase is added and shear is applied to the mixture using either high-speed mixing or vigorous agitation. A variety of high-shear mixing/agitating devices are available for this purpose, including paddle, propeller, or turbine mixers, rotor-stator homogenizers, and ultrasound generators.

An important aspect of emulsion preparation is selecting an appropriate emulsifying agent. Although some experimentation will always be required, there are some rules of thumb that provide useful starting points for predicting the type of emulsion that will form under a given set of circumstances. For example, Bancroft's rule predicts that if an emulsifying agent is preferentially dissolved in, or wetted by, one of the liquid phases, that phase will likely form the external (continuous) phase. A related rule of thumb, the oriented wedge theory, predicts that soaps of monovalent metal cations, such as sodium, tend to produce O/W emulsions while those of polyvalent metal cations, such as calcium, will tend to produce W/O emulsions. The same basic concepts underlie the dimensionless hydrophile-lipophile balance (HLB) scale, in which the relative tendency of an emulsifying agent to prefer to dissolve in oil versus water is quantified. On the HLB scale, oil-soluble (lipophilic) emulsifiers have low values, usually less than 9, while water-soluble (hydrophilic) emulsifiers have high values, usually greater than 11. A very polar nonionic surfactant might have an HLB value of close to 20, while an ionic surfactant, such as sodium dodecyl sulfate, has an HLB of about 40. Good emulsifiers for W/O emulsions tend to have HLB values of about 3 to 8, while good emulsifiers for O/W emulsions tend to have HLB values of about 8 to 18. These rules of thumb provide indications of the kind of emulsion a potential emulsifying agent may promote, but not its efficiency. In addition, mixtures of emulsifiers tend to be more efficient at emulsifying than pure compounds having the same HLB. Laboratory emulsification tests are virtually always required in order to select an appropriate emulsifying agent (or mixture) and to determine the optimum amount of emulsifier to use.

Properties. Although emulsions can be prepared that will remain stable for weeks, months, or even years in some cases, most emulsions are not truly stable. Eventually, the dispersed droplets will fall or rise due to gravity (sedimentation/creaming). As they come into contact with each other, some will begin to stick together (aggregation), and some of those droplets that stick together will eventually fuse into larger droplets (coalescence). Over time, these processes will cause an emulsion to "break" or separate into two distinct liquid layers. A special kind of emulsion, called microemulsion, is actually thermodynamically stable; microemulsions do not break on standing or with centrifuging. Besides the emulsifying agent common to all conventional

emulsions, microemulsions contain an additional stabilizing agent (cosurfactant). The combined emulsifying agents promote the formation of very small size droplets, often 10 nm or less, and change the properties of the droplet surfaces so that the droplets behave almost as though they were truly dissolved. The droplets in a microemulsion can be so small that they scatter very little light and appear to be transparent.

Emulsions have some special properties. Having a system of small liquid droplets dispersed in a second, immiscible liquid means that emulsion droplets have high surface area-to-volume ratios. This means that, by being present in an emulsion, a very small volume of (dispersed) liquid can be made to present a very large amount of surface area, such as for reactions between species at the surface and species in solution.

Another consequence of having a dispersion of small droplets is an increased resistance to flow, or viscosity. An emulsion with a small volume of liquid (less than about 10% by volume) making up the dispersed droplets will have, essentially, a small number of small droplets so that individual droplets will not encounter each other very often and will have little effect on viscosity. However, as the volume of liquid making up the dispersed droplets increases above about 10% by volume, the number of droplets increases and so does the frequency of droplet-droplet encounters. As droplets encounter and interact with each other more often, there is a "crowding" effect, and it becomes more difficult to make the emulsion flow. As the total volume occupied by the droplets increases, so does the emulsion viscosity. In some emulsions the surfaces of the dispersed droplets carry electric charges, while in others the droplet surfaces have long surfactant or polymer molecules extending outward from them. In these cases, the electric fields and extended molecules cause additional forces between droplets when they encounter each other, which further acts to increase emulsion viscosity. *See* VISCOSITY.

In theory, the maximum volume possible for the liquid making up the dispersed droplets is 74%—the maximum volume for close-packed rigid spheres. In special cases, it is possible to make emulsions in which the droplets distort and occupy a greater volume. Generally when this happens, the emulsion inverts. For example, an O/W emulsion pushed to have a total oil droplet volume of 80% by volume will usually quickly invert to become a W/O emulsion having a 20% dispersed water droplet volume.

Characterization. Emulsions may be characterized in a number of ways.

1. The nature of the emulsion can be characterized as O/W, W/O, O/W/O, W/O/W, or O/W/O/W. The nature of the dispersed phase can often be determined using dyes and/or microscopy, whereas the nature of the external phase can be determined by electrical conductivity. The existence of multiple dispersed phases can usually be determined using optical or electron microscopy.

2. Droplet size and droplet size distribution can be

determined by such techniques as optical or electron microscopy and light scattering.

3. Composition is usually determined by a titration method for simple emulsions. For more complex emulsions, composition is usually determined by separating the oil and water, and any gas or solids, and measuring their respective quantities.

4. Stability of the emulsion against breaking is usually determined using standardized bottle or jar tests, in which emulsion samples are allowed to stand for a specified period of time, with or without centrifuging, and then examined. *See* ELECTRON MICROSCOPE; ELECTROLYTIC CONDUCTANCE; LIGHT-SCATTERING TECHNIQUES.

Occurrence. Systems containing colloidal-sized particles, droplets, or bubbles are important because they feature prominently, in both desirable and undesirable contexts, in a wide variety of practical disciplines, products, and industrial processes (see **table**). The problems associated with colloids are usually interdisciplinary in nature, and a broad scientific base is required to understand them completely.

Emulsions are commonly used in many industries. Emulsion products include some foods, insecticides and herbicides, polishes, drugs, biological systems, personal-care creams, ointments, and lotions, greases, inks, paints, and varnishes. Double emulsions have applications in cosmetics, agriculture, food, photography, leather, and drug delivery. The use of emulsions to provide permanent and transient antifoams is an application that is important to an even broader variety of products and processes, including foods, cosmetics, pharmaceuticals, pulp and paper, water treatment, and minerals beneficiation. Emulsions may be encountered in virtually all of the process industries, including mineral processing, petroleum production and refining, asphalt paving, and paper de-inking. Emulsions also occur in the environment. Examples include the difficult-to-treat "chocolate mousse" type of emulsion that may be formed when crude oil is spilled on the ocean, and soil remediation, in which soil contaminants are removed in the form of an emulsion during a surfactant flushing process. A currently emerging area is that of nanoemulsions, having droplet sizes of about 1 nm or less, which are of interest in medical diagnostics and drug treatments because of their potential ability to cross the endothelial cells that protect the brain.

Demulsification. Demulsification refers to emulsion breaking. Some emulsions occur where they are not desired, such as in the process industries, and these emulsions have to be broken so the two liquid phases can be separated. The first step is to understand the type of emulsion being dealt with in terms of its nature (O/W, W/O, or multiple emulsion) and preferably in terms of the mechanism by which the emulsion was stabilized. Emulsion breaking involves two steps: causing the droplets to aggregate, or come into contact; and causing the droplets to coalesce, or fuse together. If these can be accomplished, the two liquid phases can be separated.

Some examples of emulsion occurrence		
Application area	Emulsion examples	Emulsion types*
Environment	Insecticide and herbicide formulations	O/W
	Water and sewage treatment emulsions	O/W
	Oil spill emulsions	W/O, O/W
	Sulfide liquids in magma or lava	L/L
Mineral processing	Emulsion flotation media	O/W
	Oil flotation froth	O/W
Petroleum production	Oil well stimulation emulsion	O/W, W/O
	Emulsion drilling fluid	O/W, W/O
	Produced (well-head) emulsions	W/O
	Enhanced oil recovery in situ emulsions	O/W
Petroleum transportation	Transportation fuel emulsion (70% heavy oil)	O/W
	Heavy oil pipeline emulsions	O/W
	Oil spill mousse emulsions	W/O
	Tanker bilge emulsions	O/W
Manufacturing and materials	Asphalt (paving) emulsion	O/W
	Aqueous polymer-type paints	O/W
	Microemulsion metalworking oil	O/W
Food products	Milk, ice cream, creams, whippable toppings	O/W
	Mayonnaise, hollandaise, béarnaise sauces	O/W
	Butter, margarine, spreads, processed cheese	W/O
	Salad dressings	O/W, W/O
Agriculture	Pesticide, herbicide and fungicides	O/W, W/O/W
Biology and medicine	Blood	O/W
	Vesicles, vacuoles	W/W, W/W/W
	Emulsion encapsulated drugs	O/W, W/O, W/O/W
Personal care	Cosmetic and skin care creams	W/O, O/W, W/O/W
	Microemulsion hair dyes	W/O

*Liquid (L), water (W), oil (O).

Sometimes an emulsion can be broken by simply changing the temperature or by applying mechanical shear. A variety of emulsion-breaking devices and vessels are available under names such as separators, settlers, heaters, treaters, desalters, and free-water knockouts. All are designed to induce droplet aggregation and coalescence. In most cases, a chemical demulsifier is also needed.

Chemical demulsifiers generally act by neutralizing the stabilizing effect of the original emulsifier. For example, a W/O emulsion can sometimes be broken by adding just the right amount of a surfactant that would normally promote an O/W emulsion. Dosage is critically important. Adding too much of such a surfactant will usually create the reverse emulsion, replacing one problem emulsion with another. Similarly, if an emulsion is stabilized by electrical forces, demulsification could be brought about by adding a chemical agent that can overcome or reduce these forces. Following addition of an appropriate demulsifier, mechanical agitation is often all that is needed to cause the separation of the oil and water phases. A wide range of potential chemical demulsifiers is available, but selecting an appropriate demulsifier, and the optimum dosage, can involve considerable testing due to the large number of factors that can influence emulsion stability and demulsifier performance.

Laurier L. Schramm

Bibliography. P. Becher, *Emulsions: Theory and Practice*, 3d ed., American Chemical Society, 2001;

P. Becher (ed.), *Encyclopedia of Emulsion Technology*, vols. 1–3, Marcel Dekker, 1983, 1985, 1988; B. P. Binks (ed.), *Modern Aspects of Emulsion Science*, Royal Society of Chemistry, Cambridge, 1998; L. L. Schramm, *Emulsions, Foams, and Suspensions, Fundamentals and Applications*, Wiley-VCH, Weinheim, 2005; A. L. Smith (ed.), *Theory and Practice of Emulsion Technology*, Academic Press, 1976.

Emulsion polymerization

A variety of free-radical polymerization—a method where individual molecules of high-molecular-weight polymer are rapidly generated by a free-radical chain reaction—in which a monomer with low water solubility is dispersed in small droplets (~10 micrometers in diameter) in water, and stabilized with surfactant to prevent separation of these two phases. A dilute “seed” latex of minute polymer particles in water may also be added. The reactants are agitated and heated, and a free-radical initiator is added, leading to the formation of a latex—a dispersion of polymer particles of colloidal dimensions in water (**Fig. 1**). See COLLOID; EMULSION; FREE RADICAL; POLYMER; POLYMERIZATION; SURFACTANT.

The term “emulsion polymerization” is a misnomer, as polymerization does not take place to any significant extent in the emulsified monomer droplets. The misnomer arose because the process was

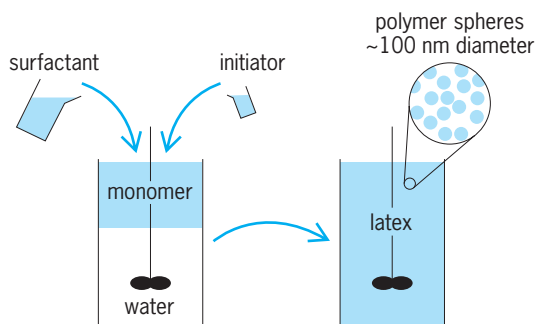


Fig. 1. Process of emulsion polymerization.

originally designed to polymerize in emulsion droplets, and the pioneers in the field did not realize that a much more complex process actually occurred. The first patents describe attempts to reproduce products similar to the natural polyisoprene latex produced by the rubber tree (*Hevea brasiliensis*), by dispersing monomers in water. Emulsion polymerization became an industrial process by the 1930s, with production of neoprene and of butadiene-acrylonitrile rubbers in Germany and the United States. Development and understanding of the process dramatically accelerated during World War II, when a research effort second only to the Manhattan Project was made in the United States for the industrial preparation of synthetic rubber.

Mechanism. A typical initiator of emulsion polymerization is an aqueous-phase thermal initiator (for example, ammonium persulfate, $\text{NH}_4\text{S}_2\text{O}_8$), although a redox couple (for example, potassium persulfate/sodium metabisulfite, $\text{K}_2\text{S}_2\text{O}_8/\text{Na}_2\text{S}_2\text{O}_5$) may also be used. Decomposition of the initiator in the aqueous phase generates hydrophilic radicals, which rapidly add to the small but not negligible amount of monomer present in the aqueous phase [reaction (1)]. The

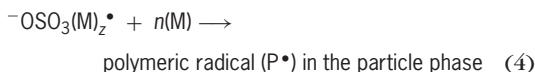
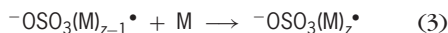


monomeric radical formed may then add a second monomer molecule [reaction (2)], and so on. Two



of these radical species may combine to form an inert product (bimolecular termination) in the water phase. As more monomer units are added, the radical is increasingly hydrophobic. At a critical degree of polymerization z , it becomes surface-active and rapidly adsorbs to the nearest hydrophobic interface. In the early stages of a typical emulsion polymerization, this interface is a micelle—an aggregate of surfactant molecules with a typical diameter of a few nanometers—or a preformed polymer seed particle. Entry to emulsion droplets is highly unlikely, because although emulsion droplets are large they are much less numerous than micelles. In the environment of a micelle or small particle, the reactive center immediately encounters a concentration of monomer much greater than in the aqueous phase and propa-

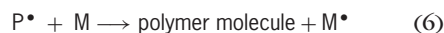
gates rapidly to form polymer [reactions (3) and (4)]



(Fig. 2). The growing polymer chain will terminate either by bimolecular termination [reaction (5)] with



an entering radical, or by abstraction of hydrogen from another species (typically monomer) to generate a radical capable of exiting to the aqueous phase [reaction (6)]. See INTERFACE OF PHASES; MICELLE.



As polymerization proceeds, the polymer particles will grow until no free surfactant remains to stabilize additional polymer/water interfaces. The monomer droplets merely serve as reservoirs of monomer, which will diffuse from the droplets to the growing polymer particles. Once the micelles have disappeared, new z -meric radicals in the water phase will all enter latex particles. (While this micellar nucleation mechanism is the most common way that particles are formed, there are other mechanisms, for example, homogeneous nucleation, involving precipitation of polymer chains of degree of polymerization $j_{\text{crit}} > z$, which occurs in surfactant-free systems.) Eventually, all the monomer droplets will disappear, and the polymerization will consume the monomer present in the particles. It is common in industry to use a controlled feed, whereby ingredients, such as monomer, are added during the polymerization process. Thus, it is possible to maintain an approximately constant monomer concentration during much of the polymerization process, or to add a different monomer after the original monomer has been consumed.

Polymer solid contents of 50% by weight are routinely obtained in industry, with negligible amounts of coagulation and typical particle diameters of

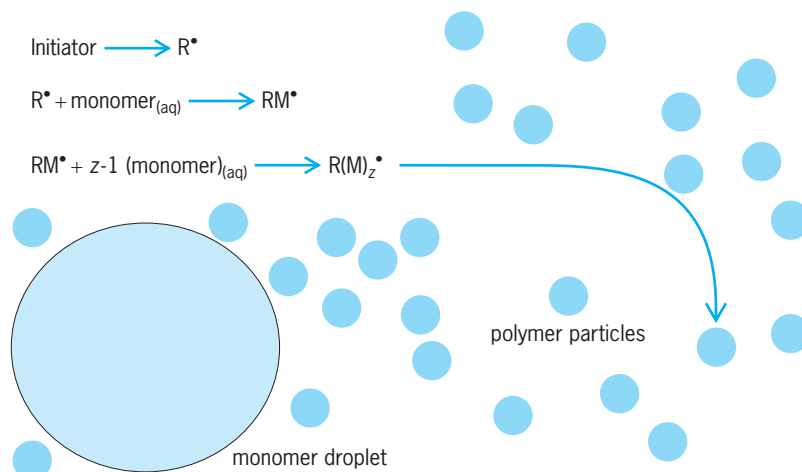


Fig. 2. Entry in emulsion polymerization.

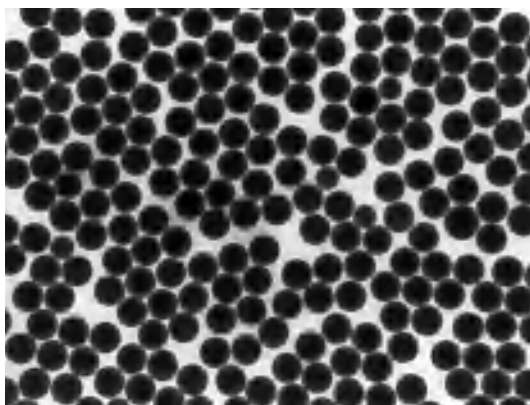


Fig. 3. Spontaneous packing of low-polydispersity poly(styrene) particles about 200 nm in diameter.

100 nanometers. Each particle contains a large number of polymer chains, but the size of a particle is not related to the molecular weight of the component polymer.

Benefits. Applications of emulsion polymerization include latexes for use as paints, adhesives, paper coatings, carpet backings, barrier products (for example, rubber gloves), and leather coatings. These latexes are frequently based on styrene, butadiene, acrylates, methacrylates, vinyl acetate, and acrylonitrile, with small amounts of water-soluble comonomers, such as acrylic acid, playing a crucial role in ensuring colloidal stability of the resulting latex. In addition to products used as a latex, there are many products where the polymer is isolated from the latex by coagulation: for example, styrene-butadiene rubber (SBR), used for tires, polytetrafluoroethylene (Teflon), and neoprene [poly(2-chlorobutadiene)].

Many advantages of emulsion polymerization arise from the compartmentalization of the growing chains, which reduces termination far below typical levels in bulk or solution polymerization to enable higher-molecular-weight products and faster reaction rates. Dispersion of the polymer in water keeps the viscosity low and makes the end product easy to handle, as opposed to the highly viscous and intractable materials produced by bulk or solution polymerization. As the process avoids use of flammable and toxic organic solvents, it has environmental and workplace safety benefits.

Recent developments. The mechanistic understanding of emulsion polymerization is now at a level where the size, composition, and topology of the colloidal particles can be controlled to achieve desired end-use properties.

The properties of a film cast from a latex are optimal when the particles can pack optimally, as is the case when they are all the same size (monodisperse) [Fig. 3]. While particle nucleation is now well understood, the process is sensitive to small disturbances (for example, changing amounts of oxygen present in the starting materials). Hence, the early stages of an emulsion polymerization in the absence of preformed latex may give an uncontrolled number of particles with a broad distribution of sizes. After nucleation, each particle grows at approximately the same rate in terms of volume, so the size distribution in terms of diameter narrows as the reaction proceeds. For this reason, industrial emulsion polymerizations are often seeded with a preformed latex of small diameter. Control of size and size distribution requires accurate measurement or prediction of particle-size distribution during polymerization, and this has been achieved by a number of groups in recent years by combining calorimetry or online spectroscopy with occasional sampling and novel process-control algorithms.

Most novel emulsion polymer products are likely to be made by combining, or more efficiently controlling, already common monomers, rather than by introducing new commodity monomers (Fig. 4). Possibilities to be considered are copolymers, polymer blends, and novel topologies. Indeed, the production of nanostructured particles has been routine in the industry for many years. See COPOLYMER.

In copolymers, different strategies must be used, depending on whether the two monomers should be distributed homogeneously within the particles or concentrated in separate domains. This can also be achieved by the online monitoring of the particle composition to optimize a mathematical model of the process that can be used for system control.

In some cases, it is desired to form two populations of particles with different copolymer compositions. For example, blocking of paints (the adhesion between two painted surfaces that can make a painted window frame difficult to open) can be minimized by generating a paint latex containing a population of small, relatively hard particles in addition to the larger, softer particles required for film formation. This may be done by varying the monomer feed

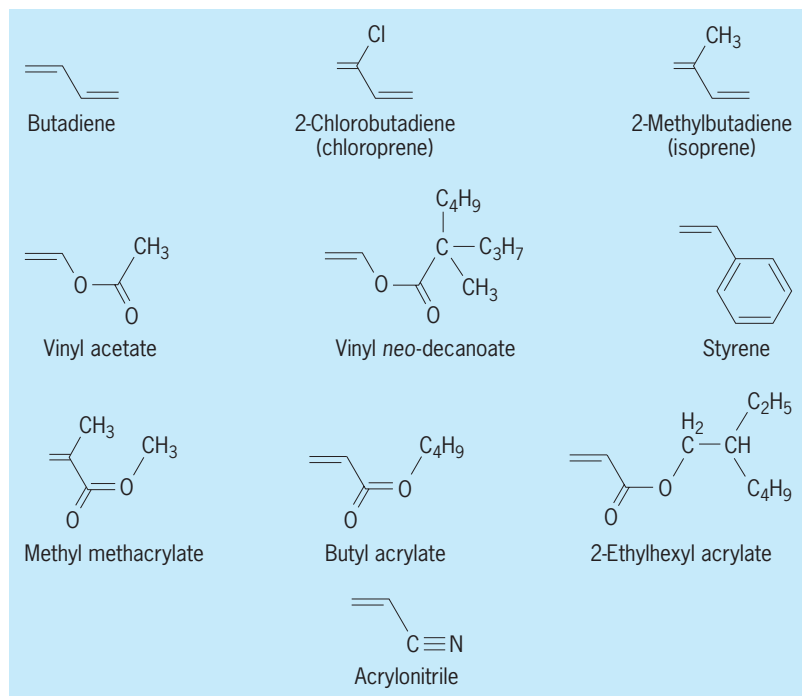


Fig. 4. Monomers used in emulsion polymerization.

composition in a single, continuous feed process. Water-soluble monomers are frequently incorporated in latexes to generate surfactant, during the course of the reaction, which will be physically grafted to the surface of the polymer particles. Unlike conventional surfactants, these will not separate from the particle surfaces on freezing or shearing to give particle coagulation, nor will they diffuse out of the formed films into the surrounding environment.

Core-shell particles, in which a softer outer shell [for example, poly(butyl acrylate)] may give good film formation, while a harder core [for example, poly(methyl methacrylate)] provides strength, can be prepared by a two-stage controlled monomer feed.

The most rapid advances in emulsion polymerization in the near future are likely to arise from application of the methods of controlled radical polymerization (CRP). Reversible addition fragmentation chain transfer (RAFT), atom transfer radical polymerization (ATRP), and nitroxide-mediated radical polymerization have been used recently on the laboratory scale to give control of molecular weight and allow formation of novel block copolymers in emulsion polymerization systems. These methods use a reversible capping of the growing polymer radical to greatly reduce termination processes. Rather than the few growing chains of very short lifetime of conventional free-radical polymerization, CRP gives many (slowly) growing chains of very long lifetime. This leads to a much narrower distribution of molecular weights and the ability to create a range of polymer such as blocks and combs. Recent work suggests the practicality of eliminating all surfactant from emulsion polymerization through use of an aqueous-phase RAFT agent, and being able to make polymer chains of any desired architecture (Fig. 5). In the first stage, a "living" (termination-free) aqueous-phase polymer is generated, which reacts with the

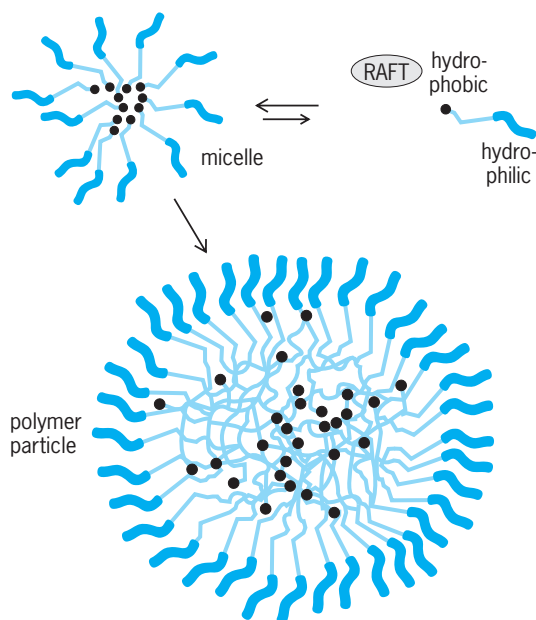


Fig. 5. Amphiphilic RAFT in emulsion polymerization.

hydrophobic monomer added in the second phase to form polymeric surfactant molecules. These form rigid micelles which grow into polymer particles by continued feeding of the hydrophobic monomer. It has been suggested that because the individual polymer molecules in such a system are irreversibly attached to the micelles, unlike conventional surfactants which undergo rapid exchange, the number of micelles originally present can be correlated directly to the number of polymer particles. Thus any desired polymer microstructure can be built up. These methods also provide a pathway for incorporating specific functionalities on the surface of a polymer particle, for instance, to attach antibodies for use in latex agglutination tests in biomedicine. See CONTROLLED/LIVING RADICAL POLYMERIZATION; MACROMOLECULAR ENGINEERING.

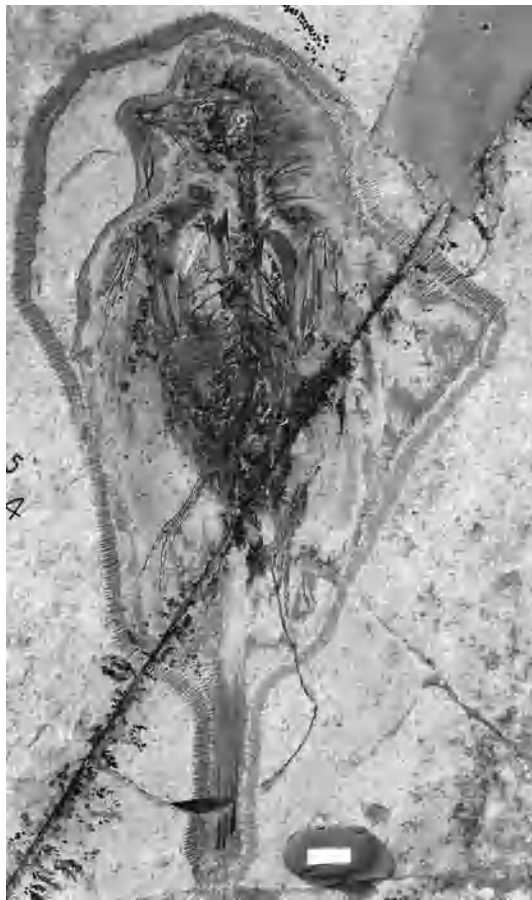
Robert G. Gilbert; Christopher M. Fellows

Bibliography. C. Burguiere, C. Chassenieux, and B. Charleux, Characterization of aqueous micellar solutions of amphiphilic block copolymers of poly(acrylic acid) and polystyrene prepared via ATRP: Toward the control of the number of particles in emulsion polymerization, *Polymer*, 44:509-518, 2002; C. J. Ferguson et al., Effective ab initio emulsion polymerization under RAFT control, *Macromolecules*, 35:9243-9245, 2002; R. G. Gilbert, *Emulsion Polymerization: A Mechanistic Approach*, Academic Press, London, 1995; K. Matyjaszewski, New (co)polymers by atom transfer radical polymerization, *Macromol. Symp.*, 143:257-268, 1999; M. Vicente et al., Control of microstructural properties in emulsion polymerization systems, *Macromol. Symp.*, 182:291-303, 2002.

Enantiornithes

An extinct and specialized group of land fossil birds that lived during the Cretaceous Period. Enantiornithes represents the most diverse avian group in the Mesozoic. This group was not discovered until 1981, at which time some very unusual, isolated bones from the Late Cretaceous of Argentina were recognized as belonging to a new group of birds. Enantiornithine birds are now known to be widely distributed, with remains from Argentina, North America, Mexico, Mongolia, Australia, Spain, and China. Because their anatomy includes a mixture of specialized and primitive features, they are thought to represent an evolutionary side branch in early avian evolution that became extinct at the end of the Cretaceous and left no living descendants. See AVES; CRETACEOUS; FOSSIL BIRDS; MESOZOIC.

Characteristics. Enantiornithine birds are readily distinguishable from living birds. The most notable characteristic is the articulation between the scapula (shoulder blade) and coracoid (a bone connecting the scapula and the sternum). In enantiornithines, the scapula has a fossa for articulation with a process on the coracoid, whereas modern birds have the opposite arrangement in which the coracoid bears the fossa and the process is found on the scapula.



Protopteryx, the most primitive known enantiornithine bird, collected from the Early Cretaceous of Hebei Province, northern China.

The group name refers to this difference: the Greek word “enanti” means directly opposite, and therefore Enantiornithes means “opposite birds.” Enantiornithes is a natural (monophyletic) group, which is united by the shared possession of several unique features, including the characteristic contact between the scapula and coracoid, the distal (lower) end of the third metacarpal of the hand extending markedly past that of the second metacarpal, and a furcula (wishbone) with a long process (hypocleidium) that projects forward.

Enantiornithine birds were good fliers, possessing a more advanced flight apparatus than *Archaeopteryx*, including a short skeletal tail (comprising a few free vertebrae and a pygostyle or “parson’s nose”), a sternum with a keel, and a wing that was generally similar to that of modern birds. Enantiornithines also possessed asymmetric flight feathers and the alula (bastard wing), indicating that they had acquired sophisticated flight capability. Conversely, some primitive avian characters, such as the presence of toothed jaws, were retained in many enantiornithes (although at least one form, *Gobipteryx* from Mongolia, completely lost the teeth). Primitive features were also retained in such structures as the pelvis and foot of enantiornithines.

Fossil record. Enantiornithine birds were initially described on the basis of often fragmentary Late

Cretaceous forms such as *Enantiornis*, *Alexornis*, and *Gobipteryx*. However, the recognition of more complete Early Cretaceous enantiornithines has allowed significant progress to be made in our understanding of these birds during the last decade. Early Cretaceous enantiornithines have been discovered in Spain, China, Australia, and Russia. These early forms are generally small—much smaller than other contemporaneous birds such as *Confuciusornis* and *Jeholornis*. The proportions of the hindlimbs and toes, and the presence of large, curved toe claws, indicate that enantiornithines were mainly arboreal forms. Most enantiornithines, such as *Protopteryx* (see **illustration**), probably fed on insects. However, there are some exceptions: for example, *Longipteryx* (from the Early Cretaceous of China) had an elongated beak, rather like a living kingfisher, which suggests that it ate fishes (although in the case of *Longipteryx*, densely packed teeth were also present). Another Chinese Early Cretaceous enantiornithine, *Longirostravis*, probably had a probing feeding behavior. Discoveries of hundreds of specimens of Early Cretaceous enantiornithines from northeastern China (referable to 12 genera) have shown that by this time Enantiornithes had significantly differentiated not only in morphology, size, and flight capability but also in terms of dietary adaptation. See CONFUCIUSORNITHIDAE.

Zhonghe Zhou

Bibliography. L. M. Chiappe and C. A. Walker, Skeletal morphology and systematics of the Cretaceous euenantiornithes (Ornithothoraces: Enantiornithes), in L. M. Chiappe and L. M. Witmer (eds.), *Mesozoic Birds above the Heads of Dinosaurs*, pp. 240–267, University of California Press, Berkeley, 2002; A. Feduccia, *The Origin and Evolution of Birds*, 2d ed., Yale University Press, New Haven, 1999; L. D. Martin, The origin and early radiation of birds, in A. H. Bush and G. A. Clark, Jr. (eds.), *Perspectives in Ornithology*, pp. 291–338, Cambridge University Press, 1983; C. Walker, New subclass of birds from the Cretaceous of South America, *Nature*, 292:51–53, 1981; Z. Zhou, The origin and early evolution of birds: Discoveries, disputes, and perspectives from fossil evidence, *Naturwissenschaften*, 91:455–471, 2004.

Encalyptales

An order of the true mosses (subclass Bryidae) that grow in dull, dark tufts on soil or soil-covered rock, generally in calcareous areas. The Encalyptales consist of a single family and two genera, the better known being *Encalypta*, the extinguisher moss, so called because of its long calyptra of candle-snuffer form. The leaves resemble some of the Pottiales in shape and papillosity, but the ladderlike thickenings of the basal cells and the peristome characters indicate a fairly distant relationship.

Encalyptales are characterized by broad, papillose leaves and erect capsules covered by very long calyptrae. The stems are erect and simple or forked with folded, incurved leaves. Leaves are broadly acute to rounded and often abruptly short-pointed to

hair-pointed with a strong midrib ending at or beyond the leaf tip. The leaf margins are often revolute below. The upper cells are generally short and opaque because of numerous C-shaped papillae; the basal cells are smooth and rectangular with thin side-walls and thickened, reddish crosswalls. The sporophytes are terminal, with elongate setae, and erect, cylindrical, and ribbed capsules. The operculum is long-rostrate, and the peristome variable in structure. The spores are often large and quite various in sculpturing. The calyptrae are very long-mitrate and long-beaked, with the base often lobed or fringed. The chromosome number is usually 13 or multiples of 13, although 12 and 14 have also been reported. See BRYIDAE; BRYOPHYTA; BRYOPSIDA; POTTIALES.

Howard Crum

Bibliography. D. G. Norton, A revision of the Encalyptaceae (Musci), with particular reference to the North American taxa, *J. Hattori Bot. Lab.*, 53:365–418, 1982, 54:357–532, 1983; D. H. Vitt and C. D. Hamilton, A scanning electron microscope study of the spores and selected peristomes of the North American Encalyptaceae (Musci), *Can. J. Bot.*, 52:1973–1981, 1974.

Endangered species

A species that is in danger of extinction throughout all or a significant portion of its range. “Threatened species” is a related term, referring to a species likely to become endangered within the foreseeable future.

Species diversity. The term “species” is defined in different ways. Most commonly, and for purposes of protecting from extinction, species are defined using the biological species concept: a species is a population or series of populations within which there is a significant amount of gene flow under natural conditions. Species extinction is a natural process. In fact, less than 1% of all species that ever existed are present today. Throughout the history of life, there have been certain periods of time with high extinction rates. For example, many dinosaur species became extinct within a short time, along with many other types of species, at the end of the Cretaceous Period. These periods of mass extinction have occurred five times. Currently, there is wide agreement among biologists that human activity has created a sixth mass extinction spasm, which continues today. According to some estimates, at least one in eight plant species and one in four mammal species are threatened with extinction.

The regions with the highest number of extinctions and endangered species are those with the greatest species diversity. Species diversity refers to the abundance of species in a given area. In general, larger areas contain greater species diversity. Species diversity tends to follow geographic patterns: there are more species at lower latitudes and altitudes. Tropical rainforests in countries such as Brazil and Indonesia harbor most of the world’s species. In addition, many unique species live on islands, particularly tropical islands. Island species often exist under spe-

cial conditions where they have comparatively few competitors and predators. They are restricted to a small range and commonly have small populations. These endemic species are particularly susceptible to extinction, because their ranges and population sizes are small.

Causes of species loss. The main factors that cause species to become endangered (and that have led to the current mass extinction spasm) are habitat destruction, invasive species, pollution, and overexploitation.

Habitat destruction. Habitat destruction is the single greatest threat to species around the globe. Natural habitat includes the breeding sites, nutrients, physical features, and processes such as periodic flooding or periodic fires that species need to survive. Humans have altered, degraded, and destroyed habitat in many different ways. Logging around the world has destroyed forests that are habitat to many species. This has a great impact in tropical areas, where species diversity is highest. Although cut forests often regrow, many species depend upon old-growth forests that are over 200 years old; these forests are destroyed much faster than they can regenerate. The threatened northern spotted owl, which lives in old-growth forests in the Pacific northwest of the United States, is a good example of a species that needs old-growth forests for the small mammals that they eat and for nest sites.

Agriculture has also resulted in habitat destruction. For example, in the United States, tallgrass prairies that once were home to a variety of unique species have been almost entirely converted to agriculture. Less than 2% of tallgrass prairies remain. Housing development and human settlement have cleared large areas of natural habitat. The construction of roads (which is associated with logging, mining, and housing developments) is a major threat to species, because animals are killed by vehicles or eventually avoid areas with too many roads and human activity. Mining has destroyed habitat because the landscape often must be altered in order to access the minerals (creating huge pits, removing mountain-tops, or stripping off the top layer of large swaths of land). Mining also creates pollution from procedures for extracting minerals (for example, cyanide is used to extract gold). Finally, water development, especially in arid regions, has fundamentally altered habitat for many species. Dams change the flow and temperature of rivers and block the movements of species up and down the river. Also, the depletion of water for human use (usually agriculture) has dried up vegetation along rivers and left many aquatic species with insufficient water.

Invasive species. The invasion of nonnative species is another major threat to species worldwide. Invasive species establish themselves and take over space and nutrients from native species. Because they are transplanted, these invasive species often live in the absence of their natural predators and competitors. Invasive species are especially problematic for island species, which often do not have defensive mechanisms for the new predators or competitors.

For example, the brown tree snake was introduced to Guam and nearby islands. The snake multiplied vastly and ate so many birds and other animals that most of the bird species on those islands became extinct. Habitat destruction and invasion of nonnative species can be connected in a positive feedback loop: when habitat is degraded or changed, the altered conditions which are no longer suitable for native species can be advantageous for invasive species. In the United States, approximately half of all endangered species are adversely affected by invasive species.

Pollution. Pollution directly and indirectly causes species to become endangered. In some cases, pesticides and other harmful chemicals are ingested by animals low on the food chain. When these animals are eaten by others, the pollutants become more and more concentrated, until the concentration reaches dangerous levels in predators and omnivores. These high levels cause reproductive problems (birds with high levels of the pesticide DDT cannot reproduce) and sometimes death. In addition, direct harm often occurs when pollutants make water uninhabitable. Agriculture and industrial production cause chemicals such as fertilizers and pesticides to reach waterways. Lakes have become too acidic from acid rain, which contains previously airborne chemicals. Other human activities such as logging, grazing, agriculture, and housing development cause siltation (another form of pollution) in waterways. Largely because of this water pollution, two out of three freshwater mussel species in the United States are at risk of extinction. Pollution has also led to other phenomena which present risks to species. Most scientists agree that pollution is also causing climate change (often called global warming), as well as greater exposure to ultraviolet radiation from ozone-layer depletion.

Overexploitation. Many species have become endangered or extinct from killing by humans throughout their ranges. For example, the passenger pigeon, formerly one of the most abundant birds in the United States, became extinct largely because of overexploitation. This overexploitation is especially a threat for species that reproduce slowly, such as large mammals and some bird species. Overfishing by large commercial fisheries is a threat to numerous marine and fresh-water species.

Strategies for protection. Habitat destruction and other threats to species worldwide increased in the twentieth century; steps are under way to reduce those threats. Many people have realized that endangered species can signal current or eventual threats to human health and safety. Efforts to save species focus on ending exploitation, halting habitat destruction, restoring habitats, and breeding populations in captivity.

In the United States, the Endangered Species Act of 1973 protects endangered species and the ecosystems upon which they depend. Internationally, endangered species are protected from trade which depletes populations in the wild, through the Convention on International Trade in Endangered

Species (CITES). Over 140 member countries act by banning commercial international trade of endangered species and by regulating and monitoring trade of other species that might become endangered. For example, the international ivory trade was halted in order to protect elephant populations from further depletion. In addition, the International Union for the Conservation of Nature identifies which species are in danger of extinction and initiates international programs to protect them.

Typically, the first step is identifying which species are in danger of extinction throughout all or part of their range and adding them to an endangered species list. In the United States, species are placed on the endangered species list if one or more factors puts it at risk, including habitat destruction or degradation, overutilization, disease, and predation. In addition, there is an assessment of whether the species is otherwise protected from those threats. Subspecies and distinct populations of species can be listed as endangered as well. For example, the Arizona population of the cactus ferruginous pygmy owl is endangered in the United States, even though this subspecies is more abundant across the border in Mexico. The International Union for the Conservation of Nature places species in categories of endangerment based on various factors, including population size, extent of population decline, and predicted population decline, either in population or area of occupancy.

From looking at patterns of where endangered species exist in the United States, it is clear that areas with high species diversity and high human population tend to have the most endangered species. Florida and California contain the most endangered species of all the contiguous 48 states. In addition, Hawaii has more endangered species than any other state. Hawaii, like other islands, has a diversity of unique species that occur nowhere else in the world. These species are also highly susceptible to endangerment because they tend to have small population sizes, and because they are particularly vulnerable to introduced competitors, predators, and disease.

Once a species is determined to be endangered, the U.S. Endangered Species Act protects and aims to restore endangered populations through two main strategies: halting exploitation, and habitat protection and restoration. It is prohibited to kill, harm, or harass an endangered species. In addition, habitat destruction is prohibited if it results in impairment of the animal's ability to forage, breed, or seek shelter. To complement prohibitions on harming endangered species, the act requires that critical habitat for a species be delineated and prioritized for protection. The act also requires a recovery plan for each species, which defines the status of the species, threats, measures to be taken to restore its numbers, and goals to be met so that delisting can occur. There are increasing numbers of cities and counties that develop conservation plans for one or more endangered species. These plans contain information on important habitat that should be preserved, as well as permits for housing development, agriculture, and

other activities to occur elsewhere. These protective efforts for endangered species also serve to protect other species in the same areas that might otherwise become endangered. All of these methods, from prohibitions to large-scale planning, are aimed at recovering species so that they are no longer in danger of extinction.

For many endangered species, a significant captive population exists in zoos and other facilities around the world. By breeding individuals in captivity, genetic variation of a species can be more easily sustained, even when the species' natural habitat is being destroyed. Some species exist only in captivity because the wild population became extinct. For a few species, captive individuals have been reintroduced into natural habitat in order to establish a population where it is missing or to augment a small population. Depending on the species, reintroduction can be very difficult and costly, because individual animals may not forage well or protect themselves from predators. For example, the golden lion tamarin has been nearly extirpated from Atlantic coastal rainforests in Brazil. These monkeys have been bred in captivity and trained for reintroduction at a nature reserve in Brazil. Although some individuals do not survive reintroduction, others are able to breed, and a viable wild population is being reestablished. See ECOLOGY; EXTINCTION (BIOLOGY).

Laura Hood Watchman

Bibliography. C. H. Flather, M. S. Knowles, and I. A. Kendall, Threatened and endangered species geography, *BioScience*, 48:365-375, 1998; G. K. Meffe and C. R. Carroll, *Principles of Conservation Biology*, 2d ed., Sinauer Associates, Sunderland, MA, 1997; E. O. Wilson, *The Diversity of Life*, Harvard University Press, Cambridge, MA, 1992.

Endocrine mechanisms

Those regulatory phenomena in animals or plants which involve, as intermediaries, one or more hormones.

Hormones

These are specific chemical entities, secreted by specialized cells, tissues, or organs and transported in solution in body fluids to other cells, tissues, or organs where they exert a specific physiological action at low concentrations (see **table**). The concept of a hormone was introduced in 1902 by W. Bayliss and E. Starling in a study of the coordination of digestive secretion; the term hormone was introduced 3 years later by Starling. The cells which form hormones are called endocrine because they pass their secretion into the blood or body fluid in contrast to exocrine cells, which secrete into the digestive tract or other regions considered outside the boundary of the body proper. See GLAND; HORMONE.

Regulation. Regulatory phenomena are those whereby the various activities of the component parts of the organism are modified so that they contribute to a coherent pattern of activity of the organism as a whole. The fact that hormones, like vitamins,

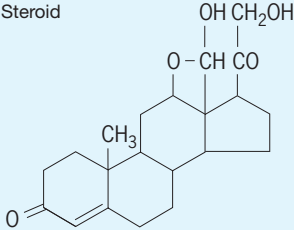
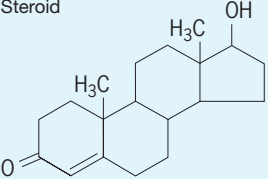
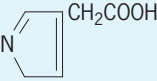
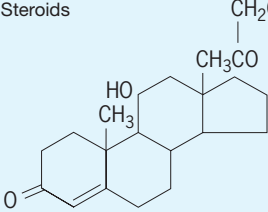
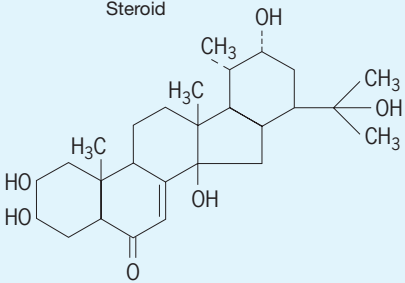
exert their actions at very low concentrations has led to the suggestion that both kinds of substances act as catalysts. See VITAMIN.

Action mechanisms. Several mechanisms of action have been demonstrated for hormones, although there is often uncertainty as to whether a particular action is primary or is a consequence of some other more fundamental action of the hormone. Hormones may change the properties of the plasma membrane or of some internal membrane of the cell, induce a change in the conformation of a protein molecule and thus alter its activity as an enzyme, serve as a coenzyme in an enzymatic reaction, induce the formation of a coenzyme or other substance modifying enzyme action, or induce or suppress the synthesis of a protein.

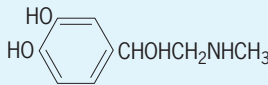
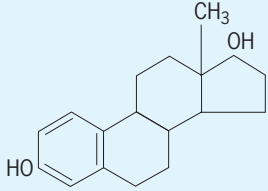
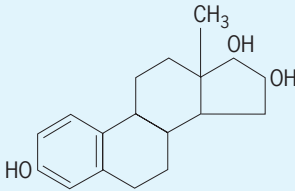
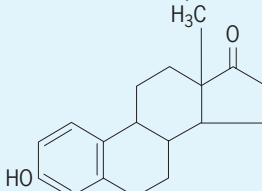
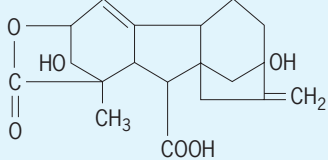
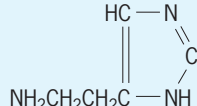
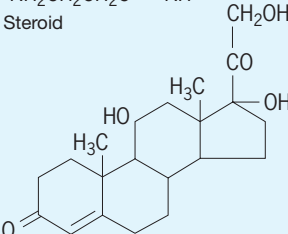
Many hormones have been shown to act on target cells through the "second messenger" mechanism. When the hormone molecule is bound to the surface of the target cell by combination with a specific receptor molecule on the surface, this combination stimulates an increase in activity of the enzyme adenylyl cyclase on the inner surface of the membrane. This enzyme catalyzes the conversion of adenylic acid (adenosine monophosphate, AMP) to cyclic 3',5'-adenylic acid by attaching two free sites of the phosphate radical through esterification to form a ring structure. Cyclic AMP (cAMP) is formed by a great variety of cells in response to specific stimuli, and activates certain enzymes within a cell, thus increasing its normal secretory or other activity. The response of the cells of the thyroid gland to the thyroid-stimulating hormone of the pituitary is also modulated by a prostaglandin, a member of a ubiquitous class of lipoid substances, first isolated from the prostate gland but since shown to occur in several different forms and to act as part of the regulatory mechanisms of cells.

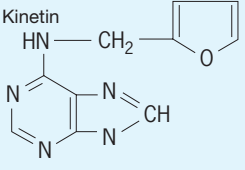
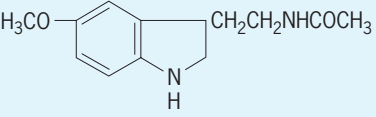
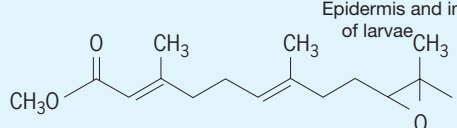
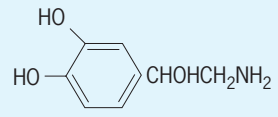
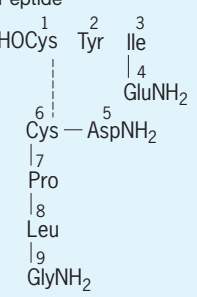
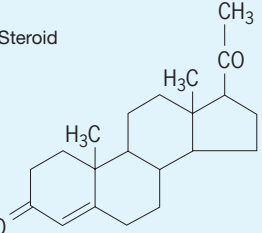
Any of these actions will alter the rate and direction of cellular activities and may thus result in widespread modification of physiological processes in the organism. Hormones appear to have more or less specific target organs or target cells; some hormones seem to act exclusively on one type of cell in a specific organ, while others act on many kinds of cells in many organs of the body. The specificity of target cells which react only to a certain hormone has been traced to the occurrence, at the cell surface or in the cytoplasm, of specific receptor molecules, each reacting only with a certain hormone. The hormone may be bound to the cell surface and exert its action there, as noted above, but other hormones, notably the steroids, enter the cell passively, and are then bound to specific receptors in the cytoplasm. In this latter case, they are carried, in combination with the receptor, to the cell nucleus, where they presumably effect their specific action by influence on the control function of the chromosomes.

There is a set of hormone-releasing factors formed by certain cells in the brain; each of these acts on a single type of cell in the anterior pituitary of a mammal, causing release from the pituitary cell of a specific hormone. On the other hand, hormones

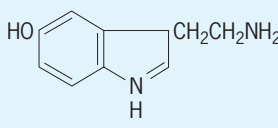
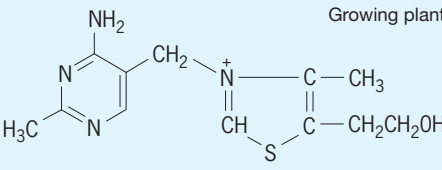
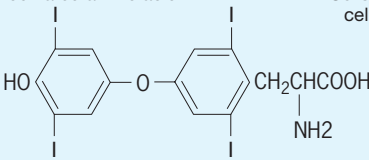
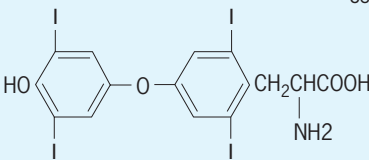
Properties of hormones of vertebrate animals, unless noted otherwise				
Names and abbreviations	Source; stimulus for secretion or release	Chemical nature	Target; receptor if known	Action or function
Acetylcholine (Ach)	Cholinergic neuron endings	Ester, $\text{CH}_3\text{COOCH}_2\text{N}(\text{CH}_3)_3$	Postsynaptic neuron or effector (muscle)	Transmitter; in mammals excites skeletal muscle, slows heartbeat, and stimulates intestinal contractions
Adrenocorticotrophic hormone (ACTH)	Anterior pituitary; CRF	Polypeptide	Adrenal cortex	Stimulates secretion of corticosterone and other steroids
Aldosterone	Adrenal cortex; angiotensin	Steroid 	Cells of walls of kidney tubules; receptors are intracellular	Stimulates active reabsorption (retention) of salt (Na^+)
Androgen (testosterone)	Interstitial cells of testes; ICSH	Steroid 	Male sex organs and ducts, muscle, and brain; receptors are intracellular	Stimulates development and maintains male secondary sex structures and behavior
Angiotensin	Blood, through enzymatic action of renin	Polypeptide	Muscle of arterioles, adrenal cortex	Stimulates contraction, increasing blood pressure; stimulates secretion of aldosterone
Auxins	Growing tips of plant shoots	Several known, including indoleacetic acid 	Growing plant cells and meristems	Stimulate growth of cells in roots and shoots, and formation of roots; inhibit development of buds in shoots
Calcitonin	Thyroid gland	Polypeptide	Bone	Inhibits reabsorption, promotes formation, decreasing Ca^{2+} in blood
Corticosterone, cortisone	Adrenal cortex	Steroids 	Many body cells; receptors are intracellular	Stimulate events of general adaptation syndrome in response to physical or emotional stress: breakdown of protein, synthesis of glycogen from amino acids, increased activity of immune mechanisms
Corticotropin-releasing factor (CRF)	Hypothalamus, posterior pituitary	Polypeptide	ACTH-secreting cells, anterior pituitary	Release of ACTH into blood
Corticotropin-release-inhibiting factor (CRIF)	Hypothalamus, posterior pituitary	Polypeptide	ACTH-secreting cells, anterior pituitary	Inhibits release of ACTH
Crustecdysone	Crustacean molt gland, Y organ	Steroid	Epidermis and other tissues concerned in periodic growth and development	Stimulates molting and associated processes
Ecdysone	Thoracic gland of larval and pupal stages of insects	Steroid 	Epidermis and other tissues concerned in periodic growth and development	Stimulates molting and associated processes

Properties of hormones of vertebrate animals, unless noted otherwise (cont.)

Names and abbreviations	Source; stimulus for secretion or release	Chemical nature	Target; receptor if known	Action or function
Enterocrinin	Cells of small intestine, in response to high concentrations of sugar and fat	Polypeptide	Musculature of stomach	Inhibits gastric contractions, depresses appetite
Epinephrin (adrenalin)	Chromaffin cells of adrenal medulla, in response to nervous stimulation		Liver, heart, blood vessels, other tissues	Increases conversion of glycogen to glucose, blood sugar, rate and force of heart; dilates blood vessels of skin; erects hairs
Estrogens: estradiol, estriol, estrone	Ovarian follicles, FSH	Steroids	Endothelium (inner lining) of uterus and female sex organs and ducts; source of FSHRF; intracellular receptors	Stimulate proliferation of mature endothelium; promote development of secondary female sex characteristics; promote estrus in some mammals; inhibit secretion of FSHRF
				
				
FSH-releasing factor (FSHRF)	Hypothalamus, posterior pituitary	Polypeptide	FSH-secreting cells of anterior pituitary	Promotes release of FSH
Gamma-aminobutyric acid (GABA)	Some presynaptic neuron endings	Amino acid, $\text{H}_2\text{NCH}_2\text{CH}_2\text{CH}_2\text{COOH}$	Postsynaptic neurons	Transmitter substance
Gastrin	Wall of stomach in response to filling	Polypeptide	Secretory cells of stomach	Stimulates production of gastric juice
Gibberellins	Plant tissues	Derivatives of gibberellic acid	Growing cells of plants, leaves, fruits	Stimulate cell elongation, growth of leaves and fruits
				
Glucagon	Alpha cells of pancreatic islet tissue, in response to low blood sugar	Protein	Liver, kidney, intestine	Stimulates formation of glucose from glycogen; stimulates elimination of salt; inhibits contraction
Histamine	Wall of stomach; injured cells in general	Decarboxylated amino acid	Stomach cells, arterioles, capillaries	Stimulates production of gastric juice; dilates, increasing blood flow; increases permeability and exudation of fluid
				
Hydrocortisone, cortisol	Adrenal cortex, ACTH	Steroid	Tissues generally	Effects are intermediate between those of aldosterone and corticosterone
				
Insulin	Beta cells of pancreatic islet tissue	Protein	Muscle, adipose tissue	Immediate effects include increased uptake of glucose, decreasing blood sugar (surface receptors); long-term increase in utilization of glucose and synthesis of lipids (intracellular receptors)
Intermedin	Intermediate pituitary	Polypeptide	Pigment cells (chromatophores) in skin of lower vertebrates	Disperses melanin, darkening body surface
Interstitial-cell-stimulating hormone (ICSH) or luteinizing hormone (LH)	Anterior pituitary, ICSHRF	Protein	Interstitial cells of testes, ovarian follicle	Stimulates production of androgen; promotes transition to corpus luteum after ovulation

Properties of hormones of vertebrate animals, unless noted otherwise (cont.)				
Names and abbreviations	Source; stimulus for secretion or release	Chemical nature	Target; receptor if known	Action or function
Isotocin	Posterior pituitary of fishes	Peptide: 4-serine, 8-isoleucine oxytocin	Testes, gills	Stimulates spawning; increases uptake of Na ⁺ in fresh water; turnover of Na ⁺ in sea water
Kinins (kinetin)	Plant tissues		Meristematic tissues	Increase rate of cell division
Melatonin	Pineal body	Dimethyl serotonin 	Pigment cells of skin of lower vertebrates	Concentrate melanin, lightening body surface
Mesotocin	Posterior pituitary of amphibians and reptiles	Peptide: 8-isoleucine oxytocin	Epidermis and internal tissues of larvae	Uncertain; may be similar to that of vasotocin Prevents development of adult characteristics in larval molts
Neotenin (juvenile hormone)	Corpora allata of insects			
Noradrenalin (Norepinephrin)	Postganglionic neuron endings of sympathetic nervous system		All structures with sympathetic innervation	Generally similar to those of epinephrin, antagonistic to those of acetylcholine
Oxytocin	Posterior pituitary	Peptide 	Kidney of lower vertebrates, uterus and mammary glands of mammals	Antidiuretic; stimulates contractions of uterus; stimulates flow of milk
Pancreozymin	Walls of duodenum, in response to HCl	Polypeptide	Exocrine cells of pancreas	Secretion and release of digestive enzymes
Parathormone	Parathyroid glands	Protein	Bone, kidney	Stimulates reabsorption of calcium; increases Ca ²⁺ of blood; increase elimination of phosphates
Progesterone	Corpus luteum of ovary	Steroid 	Endometrium of uterus	Maintains in secretory state, permitting implantation of embryo
Prolactin (luteotropic hormone)	Anterior pituitary	Protein	Corpus luteum	Maintains this structure and promotes secretion of progesterone
Renin	Juxtaglomerular cells of kidney	Protein	Precursor of angiotensin in blood	Liberates angiotensin by partial hydrolysis of precursor
Secretin	Walls of duodenum in response to HCl	Polypeptide	Exocrine cells of pancreas	Stimulates secretion of fluid containing salts, including NaHCO ₃

Properties of hormones of vertebrate animals, unless noted otherwise (cont.)

Names and abbreviations	Source; stimulus for secretion or release	Chemical nature	Target; receptor if known	Action or function
Serotonin (SHT)	Central nervous system neurons	5-Hydroxy tryptamine 	Postsynaptic neurons	Central transmitter substance; precursor of melatonin in pineal
Somatotropin (growth hormone, STH)	Anterior pituitary	Protein	Tissues in general, especially long bones and muscle	Promotes synthesis of protein and breakdown of fat, growth of young mammals; complemented by insulin
Somatotropin-releasing factor (STHRF)	Posterior pituitary	Polypeptide	Somatotropin-secreting cells of anterior pituitary	Promotes release of STH
Somatotropin-release-inhibiting factor (STHRIF)	Posterior pituitary	Polypeptide	Somatotropin-secreting cells of anterior pituitary	Inhibits release of STH
Thiamin	Leaves and seed germs of plants		Growing plant tissues	Stimulates growth of roots and shoots; thiamin phosphate is coenzyme for carboxylation, must be supplied as a vitamin for animals
Thyrotropin (thyroid-stimulating hormone, TSH)	Anterior pituitary, TSHRF	Protein	Thyroid gland; intracellular receptors	Maintains gland and stimulates secretion
Thyrotropin-releasing factor (TSHRF)	Posterior pituitary; low levels of thyroid hormones	Polypeptide	Thyrotropin-secreting cells of anterior pituitary	Stimulates release of TSH
Thyrotropin-release-inhibiting factor (TSHRIF)	Posterior pituitary; high levels of thyroid hormones	Polypeptide	Thyrotropin-secreting cells of anterior pituitary	Inhibits release of TSH
Thyroxin (tetraiodo-thyronine, T4)	Thyroid gland	Iodinated amino acid 	Cells generally (receptors in cell nuclei)	Stimulates metamorphosis in amphibians; stimulates metabolism in homeothermous animals (birds, mammals)
Triiodo thyronine (T3)	Thyroid gland	Iodinated amino acid 	Cells generally (receptors in cell nuclei)	Stimulates metamorphosis in amphibians; stimulates metabolism in homeothermous animals (birds, mammals); the principal circulating thyroid hormone of mammals
Vasopressin (antidiuretic hormone, ADH)	Posterior pituitary, inhibited by increased internal fluid volume due to intake of fluid	Peptide, 3-phenylalanine, 8-arginine oxytocin or 3-phenylalanine, 8-lysine oxytocin	Kidney	Decreases volume of urine formed by increasing permeability of tubule walls and osmotic reabsorption of water
Vasotocin	Posterior pituitary of lower vertebrates	Peptide, 8-arginine oxytocin	Kidney, skin, and bladder wall of amphibians	Increases osmotic uptake of water from urine and medium by increasing permeability of membranes to water

like thyroxin or insulin act on many, perhaps all, cells in the body, and their actions on different cells may be distinct; thus insulin causes muscles to take up glucose from the blood and convert it to glycogen, and it causes cells of fatty tissues to liberate less fat into the blood.

Control demonstrations. Evidence for hormonal control of a process may take the form of demonstrations that (1) the process is influenced by events occurring elsewhere in the organism when there is no nervous connection involved; (2) surgical removal or pathological change in a particular organ or tissue is followed by changes in another organ or tissue elsewhere in the body; and (3) extracts of an organ or tissue, injected or otherwise administered, have effects opposite to those of surgical removal of the organ or tissue or similar to those of pathological hypertrophy of the tissue. The final proof of hormonal control, after points 1, 2, and 3 above have been established, comes with isolation, purification, and chemical identification of the hormone and the demonstration that the pure substance has the same effects as the extract of the tissue of origin. In some instances isolation and purification have preceded proof of hormonal function. Development of tracer techniques using radioisotopic or fluorescent labels and electron or fluorescence microscopy, of techniques of cell, tissue, and organ culture, and of extremely sensitive radioimmune assay techniques has greatly expanded the range of methods available for study of endocrine mechanisms, and has added much to knowledge of them.

The types of processes in which hormonal control is involved may be classified conveniently as coordinative, conservative, progressive, and cyclical. Any classification of hormones on the basis of mechanism of action is premature until more is known about such mechanisms.

Coordinative Control

This phenomenon was first demonstrated in the coordination of digestive secretion in the vertebrates, and has been much extended since the original discovery, by W. N. Bayliss and Ernest Starling, of the secretin mechanism. When the acidic contents of the stomach, during digestion, are ejected into the first part of the small intestine (duodenum), the hydrochloric acid of the gastric juice elicits liberation of the hormone secretin into the blood. When this reaches the exocrine cells of the pancreas, it stimulates them to form, and release into the ducts, a juice poor in enzymes but containing salts, including sodium bicarbonate, which neutralizes the gastric juice. The duodenal mucosa also forms pancreaticozym, which stimulates liberation of digestive enzymes into the pancreatic secretion; cholecystokin, which stimulates the gallbladder to contract and eject bile into the intestine; and enterocrinin, which inhibits gastric secretion. Contact of food with the stomach causes secretion. The effect of these hormonal secretions is to ensure a supply of digestive fluids at the time when food arrives in the appropriate region of the digestive tract and, in the case of

enterocrinin, to terminate gastric secretion when intestinal digestion has begun. *See* DIGESTIVE SYSTEM.

Neurohumors. Hormonal substances formed by and liberated from special neurosecretory cells in nervous tissue (neurohumors) are often involved in coordinative control. In many and perhaps in all cases, the transmission of a state of excitation from nerve cell to nerve cell within the nervous system, or from nerve cell to muscle cell, involves a transmitter substance. The best known of these is acetylcholine, which has been identified as the transmitter between motor nerves and skeletal muscles, at parasympathetic nerve endings, between nerve cells of sympathetic ganglia in the vertebrates, and at some invertebrate nerve endings. At the vertebrate motor end plate in skeletal muscle, acetylcholine appears to act by increasing the permeability of the end plate membrane to salts; as a result, there is a flow of ions across the end plate which causes electrical depolarization sufficient to excite the muscle. The transmitter at postganglionic sympathetic fibers of vertebrates, such as those innervating the arterioles, is noradrenaline. Nerve endings at which adrenalinelike substances act as transmitters are called adrenergic. In some invertebrates, and possibly in vertebrates as well, certain nerve endings liberate gamma-aminobutyric acid (GABA) as a transmitter substance. There is growing evidence that GABA is a transmitter substance in the brains of vertebrates. Serotonin (5-hydroxytryptamine), first demonstrated as a transmitter substance in mollusks, is especially abundant in the pineal body in the roof of vertebrate brains, where it is the precursor of melatonin.

Color change. Control of color change in animals often involves neurohumors, and in crustaceans these substances act as true hormones, being liberated in one portion of the body and acting on structures in other parts of the body. Animals which can change color in response to visual and other stimuli, and thus match their body color or pattern to the background, do so by means of chromatophores. In crustaceans the movements of the pigments within the chromatophores are controlled by neurohumors formed in special groups of neurosecretory cells in the brain and eyestalk. Many of these cells have long extensions or axons terminating in the sinus gland, a small structure in the eyestalk. The secretions are apparently formed in the bodies of the neurosecretory cells and pass along the axon to be liberated from the ending in the sinus gland or elsewhere. None of the chromatophorotrophins, as these hormones are called, has been isolated, but partial purification has been accomplished, and it is clear that there are several substances, each active on one type of chromatophore. *See* CHROMATOPHORE.

Color change in the lower vertebrates is under a complex control which varies from species to species. In some, the only control is by means of intermedin, from the intermediate lobe of the pituitary gland. This substance has the effect of dispersing the black pigment melanin within the cells and thus darkening the body surface. Accumulating

evidence suggests that intermedin exerts this effect by changing the ionic balance of the cells, and the dispersion of pigment is a consequence of this change.

The pineal body, in the roof of a vertebrate brain, has a structure resembling the retina of the eye, and functions as an accessory photoreceptor, not quite a "third eye," in certain reptiles, where the skull over the pineal is relatively thin and translucent. A suspected endocrine function of the pineal—which Descartes considered, naively as it now seems, as the seat of the soul—has been verified by the demonstration that this body forms and releases into the blood the hormone melatonin. This substance, dimethyl serotonin, acts on melanophores (black pigment cells) in the skins of amphibians and reptiles to concentrate the melanin in the center of the cells, and thereby to lighten the body surface, acting thus in opposition to intermedin. These changes in color (or in shade) are adaptive in matching the shade of the animal to that of its background, and also in regulation of body temperature. A dark body surface absorbs sunlight, but also radiates heat, and a light surface absorbs and radiates less energy.

Intermedin and melatonin occur in birds and mammals, although these animals have no chromatophores and cannot change the color of their skin. Functions of these hormones in the higher vertebrates are not yet established. The secretory function of the pineal in mammals has been shown to be under the control of an adrenergic (with norepinephrine as transmitter) nerve tract originating in the hypothalamus, and the pattern of release of melatonin and of changes in the enzyme which forms it from serotonin suggests that melatonin may be concerned in regulation of the many circadian rhythms of activity and metabolism in mammals.

The effects of hormones are, in general, exerted slowly and in no pattern of local action other than that of the distribution of their targets. On the other hand, color change can be quite rapid in some lower vertebrates, and may occur in a variable pattern which matches the pattern of light and shade in the background. Professor Francis B. Sumner had a photograph of a flatfish which had matched the pattern of its skin to that of a checkerboard which had been placed underneath the fish. Rapid and patterned changes imply the action of the nervous system. *See* PROTECTIVE COLORATION.

In general, cholinergic nerves disperse pigment and adrenergic nerves concentrate it. Some vertebrates have only adrenergic nervous control; others have both.

Control of hormone secretion. Studies have demonstrated, in mammals, an elaborate system of controls centered in the hypothalamus, the lower portion of the midbrain, and acting primarily on the anterior pituitary.

Anterior pituitary. The anterior pituitary gland secretes six distinct hormones with varied actions. These are thyrotropin (TSH), which controls secretion of the thyroid hormones; adrenocorticotropin (ACTH), which controls secretion of steroids by the adrenal cortex; follicle-stimulating hormone (FSH);

luteinizing hormone (LH), also known as the interstitial cell-stimulating hormone; prolactin, also known as the luteotrophic hormone (LTH)—which three hormones control secretion of sex steroid hormones by the gonads; and somatotropin (STH), which, besides its action on growth, controls secretion of insulin by the islet tissue of the pancreas. Each of these trophic (feeding, nourishing, or maintaining) or tropic (turning toward or acting upon) hormones acts to influence secretion of another hormone by its target cells. The secretion of the six hormones of the anterior pituitary has been traced to a specific cell type, one type for each hormone. *See* PITUITARY GLAND.

Neurosecretory activity. Control of these cells, in their function of secreting hormones, resides in a group of specialized neurosecretory cells located in the hypothalamus. These neurosecretory cells send fibers or axons into the stalk that attaches the pituitary gland to the base of the brain, along a ridge called the median eminence. In the stalk and the median eminence there is a specialized set of blood vessels arising from arteries supplying the brain. These arteries branch to form capillaries; as the blood flows through these capillaries, it picks up the products liberated by the endings of the fibers from the neurosecretory cells and carries them through a small "portal vein" directly to the anterior lobe of the pituitary gland.

Andrew Schally and Roger Guillemin, in work that earned them the 1977 Nobel prize, with their coworkers, demonstrated a set of releasing factors. These are specific peptide hormones formed in the hypothalamopituitary neurosecretory system; each of these hormones, on arrival in the anterior pituitary via the portal route, stimulates the release of a specific hormone from the cell type in which it is formed. In addition, release-inhibiting factors have been identified for most of the hormones. The entire set of systems and hormones assures the interaction between the nervous system and the endocrine system through the pituitary. This allows for rapid and integrated control of many aspects of body function, as well as for psychosomatic disorders when disturbances in brain function disrupt the integrated operation of the endocrine system and consequently normal body functions. The neurosecretory connection also takes part in the feedback control mechanisms by means of which the concentrations of various hormones in the blood are kept constant. *See* NEUROSECRETION.

Conservative Control

It is well known that hormones exert a conservative control of intermediary metabolism and of salt and water balance in mammals and other vertebrates.

Carbohydrate metabolism. The primary effect of insulin is to increase the utilization of glucose by the cells. This involves increased ability to transport glucose into the cells and increased synthesis of glycogen and oxidation of glucose. In many tissues, and especially in fatty tissues, the synthesis of fats from carbohydrate is increased, and the breakdown

of fats with the liberation of fatty acids into the blood is decreased. *See* INSULIN.

The pancreas also forms the hormone glucagon in the alpha cells of the islets. This hormone acts primarily on the liver, where it stimulates the breakdown of glycogen with consequent liberation of glucose into the blood; it also stimulates the synthesis of glycogen from amino acids. The action on glycogen breakdown involves the activation of the enzyme phosphorylase, which is the principal catalyst for the breakdown of glycogen. Glucagon and insulin are antagonists, with respect to the concentration of glucose in the blood (blood sugar); insulin decreases and glucagon increases the blood sugar. The normal level of sugar in the blood is controlled primarily by insulin, with glucagon serving to raise the level if it should fall below normal. *See* GLUCAGON; PANCREAS.

Experimental diabetes can be alleviated by surgical removal of the pituitary, and the injection of extracts of the anterior lobe of the pituitary into normal animals will temporarily increase the blood sugar (diabetogenic effect). The active principles of pituitary extracts act as antagonists to insulin, decreasing the utilization of glucose by the tissues. The growth hormone, somatotropin (STH), is the major insulin antagonist in pituitary extracts; this hormone also stimulates protein synthesis in the presence of insulin. Adrenocorticotropin (ACTH) also acts as an insulin antagonist in the intact animal, but this action is largely indirect and is exerted through stimulation of the cortical cells of the adrenal glands to form and liberate steroid hormones such as corticosterone and hydrocortisone. These steroids stimulate protein breakdown, the synthesis of liver glycogen from protein (gluconeogenesis), and the breakdown of liver glycogen to glucose (glycogenolysis) and hence cause an increase in the blood sugar. Removal of, or damage to, the adrenals is followed by disappearance of glycogen from the liver and inability of the animal to maintain blood sugar levels in fasting or in vigorous activity. *See* DIABETES.

In the normal animal these factors are in balance, and the blood sugar is held within well-defined limits. Increased levels of sugar in the blood act directly on the beta cells of the pancreas, stimulating secretion of insulin. The level at which this action takes place, however, is determined in part by STH; in the presence of STH, the beta cells secrete insulin at lower levels of blood sugar. Glucagon inhibits secretion of insulin. Secretion of glucagon occurs as a direct response of the alpha cells to decreased levels of sugar in the blood. The levels of adrenal steroids in the blood are maintained relatively constant, in normal health, by the ACTH feedback control, in which the rate of secretion of ACTH is increased when the levels of cortical steroids in the blood falls, and is decreased when the level increases. *See* CARBOHYDRATE METABOLISM.

Water balance. In mammals water balance is in part controlled by the posterior pituitary.

Pituitary gland. The injection of extracts of the posterior lobe of the pituitary, or of vasopressin (antidiuretic hormone, ADH), causes a distinct decrease in

the volume of urine formed by a normally hydrated animal. Damage to certain neurosecretory cells in the hypothalamus, or section of the axons connecting these cells to the posterior pituitary, has the opposite effect, inducing the condition diabetes insipidus, in which large quantities of dilute urine are produced. This condition is entirely unrelated to the more common condition known as diabetes mellitus or sugar diabetes, which results from a deficiency of insulin secretion. ADH is normally released continually, and its release is inhibited through nervous action when the volume of the internal body fluid is increased, or the salt concentration of the fluid decreased, usually by intake of large quantities of liquids. The action of diuretic substances such as alcohol or caffeine or of some medications used to stimulate urine flow is exerted by inhibition of secretion of ADH. *See* EXCRETION.

Thyroid and parathyroid glands. The calcium and phosphate balance of blood, bone, and other tissues is controlled by the parathyroid and thyroid glands. Removal of the parathyroids is followed by a decrease in the urinary excretion of phosphate, increased calcium deposition in bone, and decreased levels of blood calcium. Nervous and muscular excitation are increased in consequence of the decreased calcium concentration, and tetany and death follow. Injection of parathormone decreases phosphate excretion and causes mobilization of calcium and phosphate from bone, increasing the blood calcium level. *See* PARATHYROID GLAND; PARATHYROID HORMONE.

The normal level of parathormone secretion is probably determined by the level of calcium in the blood, acting directly on the parathyroids. In addition to the well-known secretion of thyroxine, the thyroid gland of the vertebrates contains specialized cells that secrete a hormone known as calcitonin; this substance, when injected, causes a decrease in the calcium level of the blood and hence is antagonistic to parathormone. Calcitonin appears to act primarily by preventing the reabsorption of calcium and phosphate from bone into the blood and secondarily by increasing the loss of calcium in the urine. The stimulus for release of calcitonin seems to be an increase in levels of calcium in the blood. It acts rapidly, and for a short period of time, and consequently serves as a "fine control," with the long-term regulation being dependent upon the parathyroids. Other functions of the thyroid, in control of development and of metabolism, have been known much longer than that of control of calcium metabolism. Of two known products of secretion by this gland, thyroxine or tetraiodo thyronine (T₄) and triiodo thyronine (T₃), the latter appears to be the circulating form of thyroid hormone in mammals. Injected T₄ has been shown to be converted to T₃ in the body. *See* THYROID GLAND; THYROID HORMONES.

Total metabolic level. The primary effect of thyroid hormone in the warm-blooded or homeothermic animals—birds and mammals—is that of maintaining normal total metabolism, oxygen consumption, and heat production, and indirectly contributes to a constant body tempera-

ture. Thyroidectomy (removal of the thyroid gland) is followed by a decrease in metabolic rate, and injection of thyroid extract or of T₄ or T₃ increases this rate. Thyroidectomy has also been shown to decrease the resistance of animals to cold. *See* METABOLISM; THERMOREGULATION.

Integrity and normal function of the thyroid gland depends on the thyrotropic hormone (TSH), secreted in the anterior pituitary and released under the influence of the thyrotropin-releasing factor (TSHRF), from the hypothalamopituitary neurosecretory system. TSHRF is secreted when the level of circulating thyroid hormone falls below a set point. The response of the thyroid gland to TSH depends upon surface receptors and, through these, formation within the cells of cyclic AMP. The familiar analgesic aspirin has been shown to decrease blood levels of thyroid hormone, independent of TSH.

The receptors of target cells to T₃ have been demonstrated in cell nuclei. This suggests that the thyroid hormones exert their actions on developmental processes and on metabolic rate through effects on cellular control processes, rather than by a direct effect on oxidation processes as such, as was formerly supposed.

Sodium ion concentration. The level of sodium ion concentration in the blood, and the ratio of sodium to potassium, is decreased after removal of the adrenal glands from mammals. This change is largely a consequence of an increase in sodium chloride excretion and a decrease in loss of potassium ion through the kidneys. These changes are reversed by aldosterone. When frogs are acclimated to salt solutions, the uptake of sodium ion across the skin and from the kidney tubules is decreased, and this change is accompanied by decreased levels of aldosterone in the blood.

Secretion of aldosterone by the adrenal cortex is in general not controlled by ACTH, which controls secretion of the other adrenal steroids. Renin is a protein secreted by specialized juxtaglomerular cells in the kidney in response to decreased blood flow to, or decreased concentration of sodium ion in the blood flowing to, the glomeruli. Renin functions as a hormone, but is also an enzyme, and its target in both functions is a protein component of blood. Renin hydrolyzes this component partially, to produce the peptide angiotensin, which then acts in two ways: At low concentration, it stimulates the adrenal cortex to secrete aldosterone. At higher concentrations, it acts on the circular muscle of the arterioles, causing them to contract and constrict these vessels, increasing the blood pressure. Aldosterone stimulates reabsorption of sodium ion from the kidney tubules thus retaining this ion and its salts in the body, and increased blood pressure increases blood flow to the glomeruli of the kidney.

Progressive Processes

The progressive processes under hormonal control are those of growth and differentiation in plants and animals. For example, the growth and development of higher plants is under complex control by at least

three classes of hormones, known as auxins, kinins (kinetin), and gibberellins. Other examples are the periodic growth of arthropods (insects and crustaceans), the associated developmental processes, and the color changes of some crustaceans which are controlled by a neuroendocrine system having some interesting analogies with the vertebrate system.

Arthropods. The arthropod system is centered in the brain, with extensions into the eyestalks in some crustaceans. The thoracic gland of insects, and its homolog, the Y organ of crustaceans, is an ordinary endocrine gland (not neural in origin or connections) which secretes a steroid hormone (ecdysone in insects, and a modified version, crustecdysone, in crustaceans) that initiates and maintains the processes of molting, the shedding of the integument which is essential to growth and development in these animals. Secretion of ecdysone is initiated by a neurohumor formed in the brain. In larval insects the corpora allata, another pair of glands, secrete the juvenile hormone neotenin at each molt, and this hormone retards developmental changes leading to the adult condition until the final molt, when neotenin secretion ceases and the animal takes on adult characters.

Crustaceans have an androgenic gland associated with the testes which is responsible for the differentiation of male secondary sex characteristics, and the eyestalks form a hormone which restrains ovarian growth and delays the transformation of males into females, characteristic of the life history of certain crustaceans. In some crustaceans, also, the ovary secretes a hormone which induces development of a brood pouch and other temporary female characters. In insects the corpora allata appear to secrete one or more hormones concerned with development of female reproductive structures. In many other invertebrate animals there are correlations of neurosecretory activity with sexual maturation and the discharge of sexual products. In mussels and oysters, for example, neurosecretory cells show marked development before the maturation of the gametes and their discharge into the water; removal of the ganglia containing these cells initiates maturation and spawning. *See* ARTHROPODA; CRUSTACEA; INSECTA.

Plants. The growth and differentiation of plants is also under hormonal control. The growing tips of higher plants form auxins, which are then transported from the tips toward the base of the plant. The auxins have the effect of stimulating elongation of growing cells in proportion to auxin concentration, up to a limit; higher concentrations may inhibit growth. They also inhibit the growth of lateral buds and the abscission of fruits and leaves, and stimulate the enlargement of fruit cells and the initiation of new roots. The phenomenon of apical dominance, whereby the presence of a growing tip inhibits the growth of lateral buds in the basal regions of the stem, is attributed primarily to auxin production in the tip and polar transport toward the base. The phototropic responses of plants, whereby a growing stem tip turns toward the light and a growing root tip turns away, and the geotropic responses, in which

the stem grows upward and the root downward, are attributed to lateral polar transport of auxin in stem and root under the influence of light and gravity, respectively, and the inactivation of auxin by light.

The functions of two other plant-growth hormones remain uncertain. Kinetin is formed in many plants and has a stimulating effect on cell division. Its effects on intact plants and isolated plant parts depend partly on auxin concentration. Gibberellins, originally isolated from a fungus which attacks rice plants and other substances of similar nature, are also widely distributed. Like the auxins, the gibberellins stimulate cell elongation, and they are particularly effective in causing dwarf varieties of plants to grow to sizes characteristic of normal varieties. Unlike the auxins, gibberellins move in all directions in the plant and are not transported in polar fashion. *See* PLANT GROWTH; PLANT HORMONES; PLANT MORPHOGENESIS.

Cyclical processes. Cyclical processes of reproduction and molting have been noted above. The vertebrate female sexual cycle forms the best-known instance of hormonal control of such a cycle.

In mammals the anterior pituitary secretes a gonadotropin or follicle-stimulating hormone (FSH), which causes development of ovaries in the female and testes in the male. A second gonadotropin, the luteinizing (LH), or interstitial-cell-stimulating hormone (ICSH), stimulates growth of the interstitial cells of the testes and with FSH causes maturation of the ovarian follicle and ovulation. The interstitial cells of the testis secrete testosterone, which induces development of male sexual characteristics and male behavior. The ovarian follicle secretes estrogens, which induce development of female sexual characteristics. As the follicle grows, and estrogen secretion increases, the estrogen stimulates proliferation of the uterine lining and secretion of more LH by the pituitary. The LH, in turn, causes ripening of the follicle and ovulation. The ruptured follicle becomes a corpus luteum and begins to secrete progesterone, which further stimulates growth of the uterine lining. If fertilization occurs, the fertilized egg is implanted in the uterus and the placenta secretes estrogen, progesterone, and LH; these maintain the corpus luteum and the uterine lining in active condition and inhibit the secretion of FSH by the pituitary. If fertilization does not occur, the pituitary begins to secrete FSH, the corpus luteum breaks down, and in the absence of progesterone the uterine lining breaks down, completing the cycle. *See* REPRODUCTIVE SYSTEM.

In most mammals the maximum period of estrogen secretion is during heat or estrus, the only time the female is receptive to the male. In some mammals, such as the rabbit, ovulation occurs only following copulation, as a result of neurosecretory stimulation to the anterior pituitary, leading to LH secretion. In other mammals, and many other vertebrates as well, the female cycle is initiated by an effect of changing lengths of daylight period, mediated through neurosecretory cells acting on the pituitary. In oviparous vertebrates the cycle is essentially the

same up to the ovulation. Following deposition of eggs, or parturition in mammals, the anterior pituitary forms prolactin which induces maternal behavior and in mammals is essential to lactation. As the luteotropic hormone, it also maintains the ovarian corpus luteum in an active condition secreting progesterone. The flow of milk in lactating mammals is initiated by oxytocin. *See* ENDOCRINE SYSTEM (INVERTEBRATE); ENDOCRINE SYSTEM (VERTEBRATE); ESTRUS; LACTATION. Bradley T. Scheer

Bibliography. K. L. Becker et al. (eds.), *Principles and Practice of Endocrinology and Metabolism*, 1990; I. Chester-Jones and I. W. Henderson (eds.), *Fundamentals of Comparative Vertebrate Endocrinology*, 1987; L. J. DeGroot, *Endocrinology*, 3 vols., 3d ed., 1995; J. E. Griffin and S. R. Ojeda (eds.), *Textbook of Endocrine Physiology*, 4th ed., 2000.

Endocrine system (invertebrate)

The chemical integrating system in animals that lack a vertebral (spinal) column. An endocrine system consists of those glandular cells, tissues, and organs whose products (hormones) supplement the rapid, short-term coordinating functions of the nervous system. Evidence has been presented for hormones in a wide variety of invertebrates. However, the large majority of the published reports pertains to the more highly evolved groups that will be discussed below, the annelids, mollusks, and particularly two classes of arthropods, the insects and crustaceans. Many of the hormones in invertebrates are neurohormones; that is, they are produced by nerve cells. Just as in the vertebrates, a wide variety of functions are regulated by hormones in the invertebrates. *See* NEUROSECRETION.

Insects. The endocrine system of insects has been studied more intensively than that of any other invertebrate group. Increase in linear dimensions of an insect, as in all arthropods, can occur only at periodic intervals when the restricting exoskeleton is shed during a process known as molting. The orderly sequence of postembryonic molts that leads from the immature insect, which has hatched from the egg, to the adult is controlled by three major classes of hormones. The brain is the source of prothoracicotropic hormone, which is the starting component of the endocrine sequence that leads to a molt. This hormone activates a pair of glands in the prothorax, the prothoracic glands, causing them to secrete ecdysone (Fig. 1a). The ecdysone is then converted to 20-hydroxyecdysone in the Malpighian tubules and in the fat body, and released into the blood. 20-Hydroxyecdysone is the actual ecdysteroid that initiates the molt. Ecdysteroids are in fact the molt-inducing factors of all arthropods. The third class of hormones is the juvenile hormones (Fig. 1b), produced by the corpora allata, a pair of glands near the brain. Juvenile hormone functions during the juvenile molts to suppress the differentiation of adult characteristics, permitting growth but

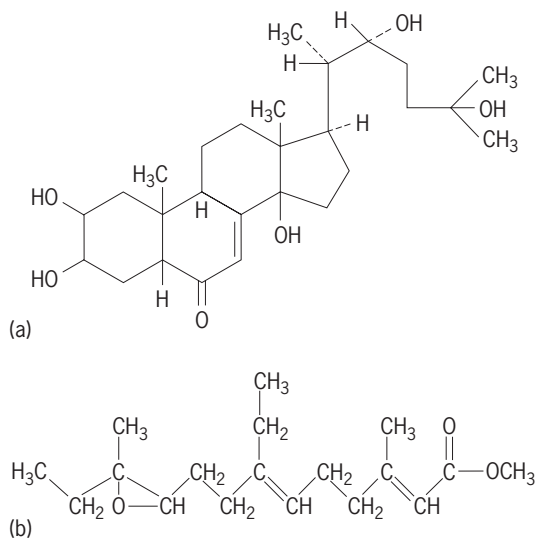


Fig. 1. Chemical structures of two insect hormones. (a) Ecdysone. (b) Juvenile hormone I.

preventing maturation. Surgical removal of the corpora allata early in development leads to the precocious development of an adult, which, however, is smaller than normal. Thus the character of the molt, whether juvenile characteristics are expressed or there is metamorphosis to the adult, is determined by the presence or absence of juvenile hormone in the blood. See MOLTING (ARTHROPODA).

Prothoracicotropic hormone appears to be a polypeptide or small protein. Ecdysone and 20-hydroxyecdysone are steroids. Juvenile hormone (JH I) was the first such isolated; at least three additional forms of juvenile hormone (JH 0, JH II, and JH III) that are structurally very closely related to JH I were isolated later. Prior to the isolation of JH I, a substance from the corpora allata was found to be necessary for the synthesis of yolk and its deposition in the developing eggs of an adult female. Now that pure juvenile hormone is available, it is clear that juvenile hormone, not another product of the corpora allata, has this role after the adult stage of development has been attained. See ECDYSONE; INSECT PHYSIOLOGY; INSECTA.

Water regulation. Two hormones with antagonistic actions are involved in regulating the water content of insects. One, the diuretic hormone, promotes water loss by increasing the volume of fluid secreted into the Malpighian tubules, the excretory structures, which are hollow threadlike structures that develop as outpocketings from the portion of the digestive tract anterior to the rectum, at the junction of the midgut and the hindgut. The second, the antidiuretic hormone, acts to conserve water by stimulating the wall of the rectum to increase the volume of water resorbed from its lumen.

Other hormones. Bursicon, a protein neurohormone, is responsible for the tanning and stiffening of a newly formed cuticle. The adipokinetic hormone is a peptide from the corpora cardiaca, which lie in the head near the corpora allata. This hormone produces an increase in the blood lipid concentration;

the lipid is released from the fat body. This increase occurs normally during flight, the lipid presumably serving as a source of energy for the flight muscles. A few species of insects, notably the stick insect *Carausius morosus*, have the ability to change color. This insect has a circadian rhythm of color change, becoming darker at night and lighter by day as a result of the rearrangement of pigment granules within its epidermal cells. Darkening is due to a neurohormone. See PROTECTIVE COLORATION.

Crustaceans. In crustaceans, as in insects, hormones have vital roles in a wide variety of physiological processes. Higher crustaceans (Fig. 2) have a structure, the sinus gland, that in most stalk-eyed species lies in the eyestalk and is the storage and release site for several hormones, including a molt-inhibiting hormone. The Y-organs, a pair of glands found in the anterior portion of the body near the excretory organs, are the source of ecdysone. The ecdysone is then converted in other tissues such as the ovaries and testes to 20-hydroxyecdysone, which as in insects is the actual molting hormone. The sinus glands are neuroendocrine structures, but the Y-organs are nonneural. Molt-inhibiting hormone controls the time of onset of molting activity by preventing the output of ecdysone by the Y-organs. When molt-inhibiting hormone release is suppressed, the Y-organs become active. See CRUSTACEA.

Color changes. Crustaceans exhibit some of the most spectacular color changes seen in the animal kingdom. The pigment cells (chromatophores) are under hormonal control, as in insects. The chromatophorotropic hormones are released from the postcommissural organs, which lie near the esophagus, and from the sinus glands. Neuropeptides have been found to cause dispersion of the pigment within the chromatophores, as well as substances with the opposite action (those which cause the pigment to become punctate). A pigment-concentrating hormone of 8 amino acids and pigment-dispersing hormones of 18 amino acids have been purified, and their chemical structures have been determined. See CHROMATOPHORE.

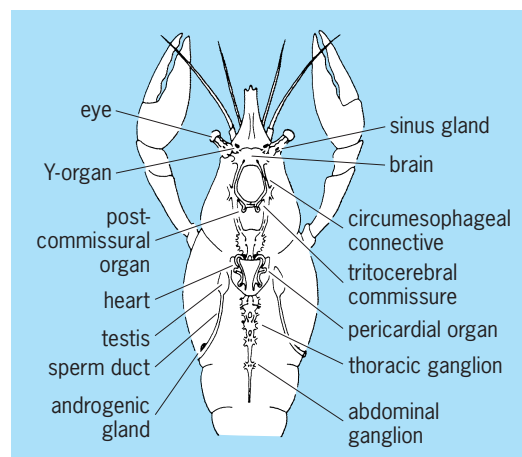


Fig. 2. Endocrine system of a male crustacean. (After A. Gorbman and H. A. Bern, *A Textbook of Comparative Endocrinology*, Wiley, 1962)

Reproduction. Although the sex of a crustacean is genetically determined, the structural and functional expression of the sex (male or female) is related to the presence or absence of functioning androgenic glands. These glands, which develop to a functional state only in males, are in all but a few species attached to the sperm ducts, and are required for the differentiation of a normal male. Not only are the androgenic glands essential for the normal development and functioning of the testes, but the androgenic glands (not the testes) are also responsible for the normal differentiation of a male's secondary sex characters, such as the copulatory appendages. In genetic females, however, in the absence of androgenic gland hormone, the gonadal rudiments autodifferentiate into ovaries; and the ovaries then produce a hormone that promotes the normal differentiation of the female secondary sex characters, such as those that are used to brood the developing young. A substance in the sinus gland inhibits the functioning of the gonads (the gonad-inhibiting hormone), whereas the brain and thoracic ganglia contain a gonad-stimulating hormone. The gonad inhibitor and stimulator act directly on the ovaries but only indirectly on the testes by acting directly on the androgenic glands. Through the actions of gonad-inhibiting and gonad-stimulating hormones, the reproductive cycles are properly timed.

Retinal pigments. The compound eye of crustaceans has three retinal pigments whose movements in accordance with changes in the illumination level control the amount of light impinging on the photosensitive cells. Of these pigments, the distal or outermost one has been studied most intensively. In the sinus gland is a neurohormone (the light-adapting hormone) that causes the migration of the distal pigment toward its light-adapted position, shielding the photosensitive cells, and another neurohormone (the dark-adapting hormone) that causes migration of this pigment toward its dark-adapted position, thereby uncovering the photosensitive cells. *See* EYE (INVERTEBRATE).

Other endocrine mechanisms. The pericardial organs, which are found near the heart, are neuroendocrine organs. The pericardial organ hormone produces an increase in the amplitude of the heartbeat. The sinus gland contains the hyperglycemic hormone, which causes a rise in the blood glucose concentration. Studies on the endocrine control of salt and water regulation revealed the existence of a neurohormone that prevents dilution of the blood of crustaceans placed into diluted seawater.

Annelids. Strong evidence exists for hormones in annelids, particularly from studies of nereids, which are polychaetes, and earthworms, which are oligochaetes. In nereids a hormone from the brain, that is, cerebral ganglia, is required for growth, which occurs normally by the formation of body segments posteriorly. This hormone is called the brain hormone. A substance from the nereid brain also inhibits precocious sexual maturation. These studies have been done mostly with nereid females, but the

few studies done with nereid males suggest the control of sexual maturation is the same in both sexes. Removal of the brain from an immature worm stops growth but leads to precocious sexual maturation. It is possible that the brain produces two hormones, one that stimulates growth and another that inhibits reproduction, but the consensus is that the evidence favors the interpretation that both actions are due to a single substance. For example, there is a parallel decline in the rates of secretion by the brain of both the growth-stimulating hormone and the maturation-inhibiting hormone. The brain hormone in its combined role would stimulate growth in immature animals while exerting its inhibitory action upon sexual maturation processes. Later, secretion of this hormone would slow and eventually stop, and this leads to sexual maturation.

In contrast to the nereids, in earthworms a stimulatory hormone from the brain is required for sexual maturation. Removal of the brain from a sexually mature earthworm stops gametogenesis and leads to the regression of the secondary sexual characteristics, such as the clitellum that secretes the cocoon. The earthworm brain produces hormones that regulate other functions also. One such helps to control salt and water balance, possibly by lowering the permeability of the body wall to water influx or by acting as a diuretic hormone to stimulate excretion of excess water; still another one is a hyperglycemic hormone that, as in crustaceans, elevates the blood glucose concentration. However, the possible hormonal control of growth in earthworms is a subject that needs further study. There are apparently conflicting reports that normal growth is dependent upon the brain; brain removal does not prevent the formation of new segments; and brain implants inhibit growth. *See* ANNELIDA.

Mollusks. As in the annelids and arthropods, in mollusks hormones are known to control a wide variety of functions. In cephalopods, such as the octopus and squid, there is a pair of optic glands which lie on either side of the central part of the brain, on the stalks of the optic lobes. These glands secrete a gonadotropic hormone that stimulates maturation of the ovary and testis and also stimulates development of the male and female accessory sex organs. The optic glands are under nervous inhibitory control. Severing the nerve supply to these glands results in precocious sexual maturation. Similarly, pulmonate snails have paired hormone-producing organs, the dorsal bodies, attached to the dorsal surface of their cerebral ganglia. Pulmonate snails are hermaphroditic. The dorsal bodies secrete a female gonadotropic hormone that appears to be proteinaceous. Removal of the dorsal bodies inhibits maturation of egg cells, whereas reimplantation of the dorsal bodies restores the normal condition. Spermatogenesis in at least some pulmonates, such as the land slug *Limax*, is stimulated by a substance called the maturation hormone that appears to be released from two sites: the cerebral ganglia and a group of cells (the collar cells) in the optic tentacles.

The hormonal control of egg-laying behavior in snails has been studied extensively. The caudo-dorsal cells in the cerebral ganglia of the fresh-water snail *Lymnaea* and the bag cells which lie at the rostral margin of the abdominal ganglion in the marine snail *Aplysia* produce neuropeptides that evoke egg-laying behavior. This egg-laying behavior is a complex series of events, including ovulation of the oocytes, extrusion of the egg string, and its fixation to the substrate. The egg-laying hormones from the caudo-dorsal cells and bag cells consist of a chain of 36 amino acids, but are not identical. In fact, *A. californica* and *A. parvula* even show amino acid sequence differences in their bag-cell egg-laying hormone.

A growth hormone, which appears to be a peptide, has been found in the cerebral ganglia of several snails. It stimulates growth of both the shell and the soft body parts. A hyperglycemic factor has been found in the nervous systems of the mussel *Mytilus* and of *Lymnaea*. Also, hormones controlling salt and water balance have been found. For example, the pleural ganglia of *Lymnaea* produce a diuretic hormone. Cauterization of these ganglia reduces the rate of urine formation, resulting in increased body weight due to water retention. The visceral ganglia of the oyster *Crassostrea* appear to produce such a diuretic hormone. Also, the cerebral ganglia of *Lymnaea* produce a substance, which seems to be a peptide, that stimulates influx of sodium ions. Removal of the cerebral ganglia results in a decrease in the blood sodium concentration. See MOLLUSCA.

Milton Fingerman

Bibliography. D. E. Bliss and L. H. Mantel, *The Biology of Crustacea*, vol. 9, 1985; M. Fingerman, The endocrine mechanisms of crustaceans, *J. Crust. Biol.*, 7:1-24, 1987; P. M. Hopkins, Ecdysteroid titers and Yorgan activity during late anecdyosis in the fiddler crab, *Uca pugilator*, *Gen. Comp. Endocrinol.*, 63:362-373, 1986; H. Laufer and R. G. H. Downer, *Endocrinology of Selected Invertebrate Types*, 1988; G. G. Lunt and R. W. Olsen, *Comparative Invertebrate Neurochemistry*, 1988; M. H. Schaffer, Functional and evolutionary relationships among the RPCH-AKH family of peptides, *Amer. Zool.*, 26: 997-1005, 1986.

Endocrine system (vertebrate)

A system of chemical communication among cells. The classical vertebrate endocrine system consists of a group of discrete glands that secrete unique products (hormones) into the bloodstream. These products travel in the blood to distant sites or targets where they cause specific physiological responses. Thus endocrine glands differ from exocrine glands, in that they lack ducts and deliver their secretions in the bloodstream. The classical definition of an endocrine system is harder to apply nowadays with the discovery of scattered cells rather than discrete glands that act as endocrine organs,

of endocrine cells that affect themselves (autocrine effect) or nearby targets (paracrine effect) by diffusion through extracellular fluids rather than the bloodstream, and of neurons that secrete hormones (neurosecretion). All of these mechanisms, however, allow for chemical intercellular communication and can be considered part of the endocrine system. See NEUROSECRETION.

Homeostasis. One important function of the endocrine system, along with the nervous system, is to maintain homeostasis, that is, a constancy of the internal environment of an organism. Thus an organism reacts and adjusts physiologically to changes in its external environment. For example, certain species of fish migrate from fresh water to seawater, or vice versa, and adjust to changes in salt concentrations by endocrine and neuronal responses. The roles of the endocrine and nervous systems in maintaining homeostasis are many, complementary, and overlapping. See HOMEOSTASIS; NERVOUS SYSTEM (VERTEBRATE).

Nature of hormones. Hormones are the products of endocrine cells. They are either proteinlike (amino acids, peptides, proteins, glycoproteins) or steroidal. Peptide hormones are produced by protein synthetic mechanisms directed by the genes of the endocrine cells. Peptide hormones are usually produced as precursors (prohormones). In some cases, two or more different hormones are enclosed within the larger prohormone. These prohormones are cleaved to the smaller peptide hormone or hormones just before leaving the endocrine cell or when they reach their target. One reason that prohormones are produced may be that they are biologically inactive and thus do not cause physiological changes when stored. Protein hormones are stored in endocrine cells in secretory granules that bud off the endoplasmic reticulum and Golgi membranes, where protein synthesis occurs. The granules leave the cell by endocytosis and enter the bloodstream. See PROTEIN.

Steroid hormones, on the other hand, are produced from cholesterol by a number of well-characterized, enzyme-catalyzed steps. Cholesterol is thus converted stepwise to various steroid families: the hormones of the adrenal cortex (cortisol and aldosterone), the estrogens (estradiol) from the ovary, and the androgens (testosterone) from the testis. Steroid hormones diffuse across the endocrine cell plasma membrane to enter the circulation. See HORMONE; STEROID.

Receptors. Hormones cause physiological responses in their targets. Since in most cases the blood carries hormones throughout the body, there must be a system by which only certain tissues respond to each hormone. This is accomplished by receptors, which are binding sites either on the surface of the target cell or within its nucleus. Receptors are high-molecular-weight proteins; the structure of some, such as the insulin receptor, is known. In general, peptide hormones cannot cross the plasma membrane, so their receptors are located there, whereas steroid hormones do pass through the plasma membrane of their targets and

bind to nuclear receptors, probably located in the deoxyribonucleic acid (DNA).

Second messengers. In order for peptide hormones to stimulate physiological changes within the target cell, the "message" must be passed from the hormone-receptor complex of the plasma membrane to the interior of the cell. This process of signaling across the plasma membrane is accomplished by so-called second messengers of which the best known is adenosine 3',5'-cyclic monophosphate (cyclic AMP). The breakdown of cell-membrane phospholipids in response to hormone-receptor interaction is another widespread mechanism of transmembrane signaling which mobilizes calcium ion within the target cells. Calcium ion, in turn, is an important second messenger. Steroid hormones, of course, have no need for second messengers since they are lipid-soluble and pass readily through the plasma membrane and into the cell. *See CELL MEMBRANES.*

Physiological responses. Once hormones are bound with their receptors and have stimulated their target cells, physiological responses occur. This may involve such biochemical processes as conversion of an inactive form of an enzyme into an active one, stimulation of critical enzymatic pathways, increase in transport of glucose or amino acids into cells, or synthesis of new proteins. These events may result in overall changes in cell or organ function, metabolism, growth, or even the behavior of the organism.

Feedback systems. The endocrine system is regulated by control mechanisms, the means by which homeostasis is achieved. The most common relationship between the hormone and its target is one of negative feedback, whereby the response to the hormonal stimulus turns off the original stimulus. For example, the endocrine pancreatic beta cells produce insulin in response to high blood sugar levels. Insulin is released into the blood, where it causes its target cells to take up glucose, thus reducing blood sugar. When blood glucose concentration falls, the secretion of insulin is turned off. The system is turned back on by the following process: when less insulin is secreted, the target cells take up less glucose, and thus the blood glucose content gradually increases again until insulin secretion is stimulated once more. Thus the blood level of hormone (insulin) and the substance responding to the hormone (glucose) are kept within normal ranges. *See CARBOHYDRATE METABOLISM; INSULIN.*

Most relationships between hormones and their targets involve negative feedback. However, positive feedback also occurs here and there, whereby the production of the hormone is stimulated rather than turned off by the product of the hormone-target interaction. *See ENDOCRINE MECHANISMS.*

Pituitary gland and hypothalamus. The pituitary gland, or hypophysis, is located near the base of the brain. It secretes many hormones and controls the function of other endocrine glands. The production and release of the various pituitary hormones are regulated in turn by small peptide-releasing hormones

from the hypothalamus of the brain. These factors are produced by neurosecretory neurons and travel to the adenohypophysis (anterior lobe of the pituitary) by way of a portal blood system. The releasing hormones stimulate specific cells of the adenohypophysis to produce and release their hormones. Generally speaking, each of the adenohypophysial hormones is affected by a separate releasing hormone. Thus, the hypothalamic thyrotropin-releasing hormone stimulates the synthesis and release of thyroid-stimulating hormone (thyrotropin) by the adenohypophysis.

Other adenohypophysial hormones include adrenocorticotrophic hormone, which stimulates the production of steroid hormones by the adrenal cortex; growth hormone, which stimulates protein synthesis and growth in many cells; prolactin, which stimulates the production of milk by the mammary glands and is important in salt and water balance and many other functions; follicle-stimulating hormone, which induces growth of the follicles of the ovary prior to ovulation; and luteinizing hormone, which induces ovulation in the ovary. The release of both follicle-stimulating and luteinizing hormones is governed by gonadotropin-releasing hormone.

Other hypothalamic hormones do not reach the pituitary by way of the bloodstream; instead they travel down the long axons of the neurosecretory cells into the neurohypophysis (posterior lobe of the pituitary). These hormones, oxytocin and vasopressin (or antidiuretic hormone), are released directly from the axon end bulbs into the blood. Oxytocin acts upon the mammary glands of female mammals to cause milk release in response to suckling by the young, and stimulates the uterus to contract at the end of pregnancy to aid in expulsion of the offspring. Vasopressin is important in water conservation by the kidney tubules (less urine excretion) and also produces an increase in blood pressure. Oxytocin and vasopressin are, therefore, hypothalamic hormones that are stored in the pituitary. *See PITUITARY GLAND.*

Thyroid gland. The thyroid gland lies in the neck region of mammals. It produces two closely related hormones, triiodothyronine and thyroxine. These both increase the metabolic rate of an organism, and increase enzyme activity and protein synthesis. The thyroid hormones act along with growth hormone to promote cell growth and development. These functions are well illustrated in cases of hypothyroid human infants, who are retarded both mentally and physically. Thyroid hormones are peptides but their three-dimensional structures may be similar to those of steroid hormones. Thus they are unusual in their ability to pass through the plasma membrane of their target cells and bind to nuclear receptors, directly affecting genes that control protein synthesis.

The control of hormone secretion by the thyroid (as well as by the adrenal cortex and gonads) involves more complex feedback relationships. These endocrine glands are affected by the levels of hormones from the adenohypophysis (which were affected by

releasing hormones from the hypothalamus). In the case of the thyroid gland, thyrotropin-releasing factor from the hypothalamus stimulates the release of thyroid-stimulating hormone by the adenohypophysis. In response, the thyroid secretes thyroxine and triiodothyronine. High blood levels of the thyroid hormones inhibit the secretion of both thyrotropin-releasing factor (long-loop feedback) and thyroid-stimulating hormone (short-loop feedback). *See* THYROID GLAND.

Calcium regulation. The parathyroid glands derive their name from the fact that in mammals they are embedded within the thyroids. These small glands are essential for life, as they regulate the concentration of calcium ion in blood and other extracellular fluids. If calcium is too low, the animal goes into tetanic convulsions and dies, whereas if calcium is too high, abnormal calcification and stone formation can occur. Thus the regulation of calcium ions is of utmost importance. Parathyroid hormone is a protein hormone that raises the blood calcium levels (hypercalcemia). The hormone acts upon bone to cause the release of calcium and phosphate, and upon the kidney to increase the reabsorption (conservation) of calcium and excretion of phosphate.

Vitamin D is now recognized as a steroidlike hormone, although it does not originate from an endocrine gland. It is synthesized from precursors present in the diet or produced after exposure of skin lipids to ultraviolet light. The two final steps in vitamin D biosynthesis involve the addition of hydroxyl groups to the molecule. The first hydroxylation occurs in the liver, the second in the kidney, thus forming the most active metabolite, 1,25-dihydroxyvitamin D₃. Vitamin D plays roles in calcium conservation by the kidney and in bone mineralization, but its most important function is to enhance calcium transport across intestinal cells and thus conserve dietary calcium. *See* VITAMIN D.

Calcitonin is a newly recognized peptide hormone produced by thyroid cells in mammals (different cells from those that produce thyroid hormones) and from the ultimobranchial glands of nonmammalian vertebrates. Calcitonin is hypocalcemic and acts by inhibiting calcium loss from bone. Of the three calcium-regulating hormones, it appears to be the least important. *See* PARATHYROID GLAND.

Carbohydrate regulation. Insulin and glucagon are peptide hormones produced by endocrine cells of the pancreatic islets. Insulin is a protein hormone produced by the pancreatic beta cells and is the only hormone that decreases blood sugar (glucose) levels. It acts on its target cells (skeletal muscle, fat cells) to increase the uptake of glucose, amino acids, and fatty acids. Once taken into cells, glucose is used in metabolic reactions or stored as glycogen, a large carbohydrate. Insulin also causes the conversion of amino acids to proteins and fatty acids to fats in the target cells. In the absence of insulin, as in diabetes mellitus, the target cells cannot take up glucose, and thus the body must utilize amino acids and fats as energy sources. These processes result in the accumulation of toxic metabolic products which eventually

disrupt the acid–base balance, leading to coma and death. *See* DIABETES.

Glucagon, in contrast, is a hyperglycemic hormone. It is a small peptide from the pancreatic islet alpha cells that acts upon liver cells to cause the conversion of glycogen to glucose by activation of key enzymes in a complex metabolic pathway.

Many other hormones elevate blood sugar levels. For example, epinephrine (adrenalin), an amino acid derivative from the adrenal medulla, acts by the same pathway as glucagon to convert glycogen to glucose, except that the targets of epinephrine are skeletal and heart muscle. Epinephrine is secreted in times of stress and serves to prepare the body for an emergency by increasing the availability of energy in the form of glucose and by increasing the heart rate and blood pressure.

Growth hormone, a large protein hormone from the adenohypophysis, is secreted in response to low blood sugar levels. This hormone elevates blood sugar by blocking the uptake of glucose by cells and by favoring the utilization of fats rather than glucose as an energy source.

Many of the adrenal cortical hormones, such as cortisol, are known collectively as glucocorticoids, because they also elevate blood glucose levels. These steroid hormones favor the production of glucose from proteins and fats rather than glycogen (gluconeogenesis). Glucocorticoids also exert an anti-inflammatory action, which makes them useful for treatment of arthritis and other diseases. *See* ADRENAL GLAND.

Salt and water regulation. Several hormones affect the ability of the kidney to conserve or excrete salts and water. The antidiuretic hormone (vasopressin) promotes water reabsorption by the kidney tubules, so that the organism excretes less water. The secretion of vasopressin is regulated by hypothalamic neurosecretory neurons that are sensitive to the concentration of salts in the extracellular fluids. In the absence of vasopressin, an individual excretes great volumes of dilute urine, leading to severe dehydration (diabetes insipidus).

Salt excretion is regulated mainly by two hormones that act in opposition. Aldosterone is an adrenal cortical steroid that promotes the reabsorption of sodium by the kidney tubules and thus decreases its excretion in the urine. In contrast, atriopeptin (atrial natriuretic factor), a peptide that originates in heart muscle, acts upon the kidney to increase the excretion of sodium in the urine. *See* OSMOREGULATORY MECHANISMS.

Reproductive hormones. Probably the best-studied endocrine glands are the gonads, the testes of the male and the ovaries of the female. The gonads are regulated by the follicle-stimulating hormone and luteinizing hormone from the adenohypophysis. In the male, follicle-stimulating hormone stimulates the initiation of sperm formation by the testis tubules, and luteinizing hormone acts on the nearby Leydig cells of the testis to produce testosterone, the principal male sex hormone. Testosterone acts by a paracrine mechanism to cause the final maturation

of sperm, and by way of the blood to stimulate development of the male reproductive system and secondary sex characteristics (body shape, beard, muscle growth). *See* TESTIS.

In the female, follicle-stimulating hormone stimulates the growth of the ovarian follicles at the beginning of each reproductive cycle. As the follicles grow, they produce estradiol, an important female sex hormone. Increasing levels of estradiol cause feedback inhibition of the released gonadotropin-releasing hormone. High levels of estradiol also have an unusual positive feedback effect upon the hypothalamus and adenohypophysis to cause a surge in the secretion of luteinizing hormone, which results in ovulation (the release of the ovum). The corpus luteum, a remnant of the ovulated follicle, produces both estradiol and the second major female sex hormone, progesterone. Progesterone is necessary for the maintenance of a quiescent uterus during pregnancy, and both estrogen and progesterone are important in the regulation of the female reproductive cycle (the rhythmic timing of the menstrual cycle in primates, for example). Estradiol is also essential for the growth and maturation of the female reproductive system and secondary sex characteristics (breast growth, body shape, patterns of fat deposition). In both males and females, the sex hormones affect reproductive behavior. *See* OVARY; REPRODUCTIVE BEHAVIOR.

Other hormones. There are many other factors that act in various ways to achieve homeostasis or intercellular communication. Of these, the existence of a large number of peptides found in both gastrointestinal cells and the brain is of interest. These were recognized for many years as gastrointestinal hormones which aid in secretion of digestive juices and motility of the gastrointestinal tract. Their function in the brain appears to be different, and there is evidence that they act in pain reception or analgesia, as factors that stimulate or curb appetite, or in memory or other functions. This field of neuropeptide hormones is in its infancy and serves to emphasize the close relationship between the endocrine and nervous systems in intercellular communication.

Nancy B. Clark

Bibliography. A. Gorbman et al., *Comparative Endocrinology*, 1983; G. Litwak and A. W. Norman, *Hormones*, 2d ed., 1997; D. Randall et al., *Eckert's Animal Physiology: Mechanisms and Adaptations*, 4th ed., 1997.

Endocrinology

The study of the glands of internal secretion, the endocrine glands, and the hormones which they synthesize and secrete. These glands are ductless; the hormones are secreted directly into the blood to be carried to the target tissue or organ. The hormones, or chemical messengers, are highly specific and their action may be selective or generalized. *See* ENDOCRINE SYSTEM (VERTEBRATE).

Most of the physiological processes of the body are under hormonal and nervous system regulation. The quantities of circulating hormones are maintained in proper balance by the delicate feedback mechanisms involving endocrine gland interactions, blood levels of hormones, and physiological activities of the target organ. The pituitary is considered the master gland because so many of its hormones are trophic in nature; that is, they regulate the activities of other endocrines. *See* ENDOCRINE MECHANISMS; PITUITARY GLAND.

The hormones carry out their biological activities by regulating cellular metabolism; they act as inhibitors or accelerators of intracellular biochemical reactions to serve three general functions. These are control of growth and development, determination of secondary sex characteristics, and maintenance of homeostasis. *See* HOMEOSTASIS.

Chemically, there are three kinds of hormones: steroids, amines, and peptides or proteins. However, many subtypes and variations may appear in an individual. *See* HORMONE; STEROID.

Most hormones have been purified, analyzed, and standardized so that a whole gamut of synthetic analogs are available for the highly selective actions desirable in therapy. *See* ENDOCRINE SYSTEM (INVERTEBRATE).

Sybil P. Parker

Endocytosis

The process by which animal cells internalize particulate material (such as cellular debris and microorganisms), macromolecules (such as proteins and complex sugars), and low-molecular-weight molecules (such as vitamins and simple sugars). The size difference among these different nutrient sources requires cells to develop a hierarchical plan for endocytosis. Thus, cells engage in at least three different types of endocytosis: phagocytosis where cells engulf particulate material, receptor-mediated endocytosis of macromolecules, and potocytosis of small molecules.

Some of the essential nutrients that a cell needs are scarce in the environment. The cells overcome this problem by expressing high-affinity receptors, or binding sites, on the membrane surface. Each type of receptor is specific for either macromolecules or molecules. These endocytic receptors are capable of concentrating their ligand at the cell surface before carrying it into the cell, thus increasing the efficiency of uptake so effectively that a malfunctioning endocytic receptor can cause a disease. For example, individuals who have a defective endocytic receptor for low-density lipoprotein (LDL), the major cholesterol-carrying molecule in the blood, develop a disease called familial hypercholesterolemia that causes premature heart attacks.

In all three endocytic pathways the internalization step begins with the invagination of plasma membrane and the conversion of this membrane into a closed vesicle called an endosome. Each of the pathways has its own set of molecules that control

internalization. These molecules assemble at the cell surface and physically deform the membrane into the shape of a vesicle. The vesicle, the endosome, then detaches and migrates to other locations within the cell. The same cell-surface assemblage of molecules also attracts endocytic receptors that are moving around on the cell surface, causing them to cluster over the site of internalization. Receptor clustering, which is essential for efficient uptake, is sometimes stimulated by ligand binding.

In general, molecules must be smaller than 1000 daltons to pass through a membrane barrier, because this size can easily fit within water-filled channels in the membrane. Therefore, particles and macromolecules can benefit the cell only after they are broken down into sugars, amino acids, and other molecules that can readily cross membranes. The breakdown of macromolecules is the major function of a membrane-bound digestive organelle called the lysosome. Therefore, endosomes that are generated by the phagocytic and receptor-mediated endocytic pathways often fuse with lysosomes that contain many different hydrolytic enzymes. Small molecules, by contrast, do not need further processing, so during potocytosis they are delivered directly to the cytoplasm. *See* CELL MEMBRANES; LYSOSOME.

Phagocytosis. Phagocytosis is a receptor-mediated process where the receptors function as adhesive elements that bond the plasma membrane to the particle. Inactivation of either the ligand on the particle or the receptor on the membrane inhibits phagocytosis.

The adhesive interaction of the phagocytic receptors with the membrane stimulates invagination. A critical molecule in this activity is actin, the same protein that provides power for muscle contraction. Surface membranes contain actin-binding proteins that link the phagocytic receptor to the actin cytoskeleton of the cell. Thus, when a particle binds to its endocytic receptor, a signal cascade is initiated that stimulates the recruitment of actin filaments to the site of phagocytosis. *See* CYTOSKELETON; MUSCLE PROTEINS.

Sometimes a cell initiates phagocytosis by extending lamellipodia that engulf the particle in a phagocytic cup. Lamellipodia are actin-rich ruffles of the plasma membrane that can assemble without a particle first binding to the membrane. They are particularly abundant in migrating cells, but in certain sessile cells, such as epithelial cells, they assemble in response to hormonal stimuli. The fusion of two lamellipodia that are next to each other can trap small amounts of extracellular fluid in an endocytic vesicle. This process is called pinocytosis, or cell drinking. All three endocytic pathways contribute to pinocytosis by nonspecifically trapping molecules and macromolecules present in the cell environment each time they generate an endocytic vesicle. *See* PHAGOCYTOSIS.

Receptor-mediated endocytosis. The clathrin-coated pit is the segment of cell membrane that is specialized for receptor-mediated endocytosis. Each pit can be recognized by the presence of a polygonal lattice

on the cytoplasmic surface of the membrane. This lattice shapes the plasma membrane into a coated vesicle that immediately uncoats and fuses with endosomes. The endosome functions as a switching area that directs membrane and content molecules to specific locations within the cell. Each of the structural units of the coated pit, as well as any associated endocytic receptors, recycles many times during its life.

Endocytosis by coated pits takes place in several steps. There are four major groups of molecules to consider: receptors that use coated pits to internalize molecules; clathrin triskelions that form the distinctive polygonal lattice of the coat; assembly protein complexes (adaptors) that link the triskelion lattice to the membrane; and a membrane receptor for the adaptors that controls coated-pit assembly. Coated-pit formation starts when an adaptor binds to its receptor. The activated adaptors then aggregate and bind to triskelions that are free in the cytoplasm. The triskelions polymerize into a polygonal lattice that contains both hexagons and pentagons. Receptors migrate into fully formed coated pits and are retained there as the lattice transforms into a deeply invaginated pit. A molecular marker on the cytoplasmic tail of each receptor is the signal for retention. The coated pits then detach from the membrane.

Shortly after budding, the interior of the new vesicle becomes very acidic. A low pH is essential for the proper separation of a receptor from its ligand. Once the ligand is free in the lumen, the receptor is shunted into a vesicle that buds from one end of the endosome. These recycling vesicles return to the cell surface, while the ligand contained in the endosome proceeds to its next cellular destination. *See* CYTOPLASM.

Potocytosis. Potocytosis uses membrane proteins that are anchored by lipid rather than protein as endocytic receptors. The lipid anchor causes the attached proteins to migrate in the plane of the membrane and cluster in a membrane specialization called a caveola. Clustering ensures that any ligand bound to these receptors will be concentrated in this location. When caveolae close, they create a tiny compartment of uniform size that is sealed off from the extracellular space. When the ligand dissociates from its receptor, it reaches such a high concentration that it naturally flows through water-filled membrane channels into the cell.

Like coated pits, caveolae also have a distinctive membrane coat on their cytoplasmic surface. Presumably this coat plays an essential role in controlling the curvature of the membrane. A protein component of the coat, called caveolin, appears to be a regulatory molecule that helps to maintain a high concentration of cholesterol in the caveolar membrane. Unlike other endocytic organelles, the function of caveolae is strictly dependent on cholesterol. For example, both the shape of the caveolar membrane and the clustering of lipid-anchored endocytic receptors are abnormal in cells that have been treated to lower cholesterol levels.

Although caveolae are found on the surface of most cells, they are particularly abundant in endothelial cells. Caveolae in endothelial cells appear to deliver molecules present in the blood directly to tissue cells by a process called transcytosis. This involves the formation of a plasmalemmal vesicle at the apical surface of a special endosome, and this vesicle moves across the cell and fuses with the basal surface.

Caveolae work somewhat differently in nonendothelial cells. In epithelial cells, they appear to create a transient compartment that remains associated with the cell surface. It is not clear whether this involves the actual formation of a plasmalemmal vesicle or a simple opening and closing cycle without actual detachment from the plasma membrane. Nevertheless, there is no evidence that vesicles derived from caveolae ever merge with vesicles created by other endocytic pathways.

The closed caveolar compartment appears to be a unique space for the cell. It is transient, does not merge with other organelles, and can selectively concentrate extracellular molecules or ions and deliver them to the cytoplasm. In addition to importing molecules, cells can also use this space to store and process incoming or outgoing messengers that affect cell behavior. It has been determined that caveolae contain the molecular equipment for sending signals. For example, isolated caveolae are enriched in nonreceptor tyrosine kinases and heterotrimeric guanosine triphosphate (GTP) binding proteins, two sets of proteins that are known to be intermediates in a variety of signaling cascades. Furthermore, some of the lipid-anchored endocytic receptors are binding sites for signaling molecules like cyclic adenosine monophosphate (cAMP), while others are enzymes that can act on nearby substrates. Since the lipid anchor directs all of these molecules to the same caveola, a complex array of processing reactions can occur when it closes. *See CELL (BIOLOGY).*

Richard G. W. Anderson

Bibliography. R. G. W. Anderson, Dissecting clathrin-coated pits, *TIC*, 3:177–179, 1993; R. G. W. Anderson et al., Potocytosis: Sequestration and transport of small molecules by caveolae, *Science*, 255:410–411, 1992; J. L. Goldstein et al., Receptor-mediated endocytosis: Concepts emerging from the LDL receptor system, *Annu. Rev. Cell Biol.*, 1:1–39, 1985; J. H. Keen, Clathrin and associated assembly and disassembly proteins, *Annu. Rev. Biochem.*, 59:415–438, 1990; C. J. Steer and J. A. Hanover (eds.), *Intracellular Trafficking of Proteins*, 1991.

Endodermis

The single layer of plant cells that is located between the cortex and the vascular (xylem and phloem) tissues (**Fig. 1**). It has its most obvious development in roots and subaerial stems. In the aerial stem and in leaves it is not always detectable except by histochemical tests. However, aerial parts of many plants may develop a characteristic endodermis

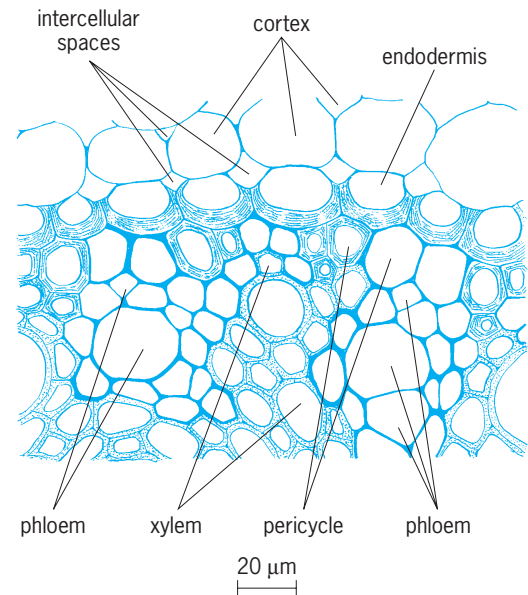


Fig. 1. Transection from part of *Zea* root illustrating endodermis in relation to contiguous tissues. (After K. Esau, *Plant Anatomy*, 2d ed., John Wiley and Sons, 1967)

when subjected to belowground growing conditions.

Initial development. The primary phase in the development of the endodermis is identified by a thin band (Casparian strip) of suberin or ligninlike deposition around each cell in the anticlinal (perpendicular to the surface) walls, that is, the radial and transverse walls (**Fig. 2a**). There are no intercellular spaces in the endodermis, and the anticlinal walls between the cells are blocked by the Casparian strip, giving a “watertight” appearance to this tissue (**Fig. 2b**). When plasmolyzed, the protoplasts of the endodermis appear to adhere to the Casparian strip. The contiguity of protoplast and the material of the Casparian strip (**Fig. 2c**) is supposed by some to prevent the movement of solute ions and water through the walls, or between the walls and the cytoplasm and thus restrict this movement to the cytoplasm. This relation, in turn, is thought to play a role in the selective absorption of ions and in maintaining

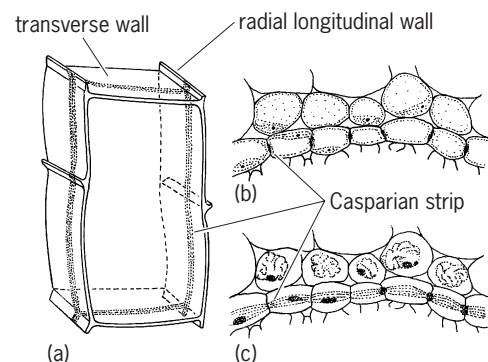


Fig. 2. Details of endodermal structure. (a) Diagram of cell showing location of Casparian strip. (b) Cells of endodermis and of ordinary parenchyma before treatment with alcohol. (c) Cells after treatment. Casparian strip is seen only in sectional views in band c. (After K. Esau, *Plant Anatomy*, 2d ed., John Wiley and Sons, 1967)

hydrostatic pressure. Such a function of the endodermis, however, has not been proved conclusively. See CELL WALLS (PLANT).

Later development. The secondary phase in development of the endodermis consists of a complete suberization of all the walls. This event may be followed by a tertiary phase in which cellulose deposition takes place on the inner tangential and radial walls. An oxidation of phenols, naphthols, and anthrols to quinones occurs in the endodermis of roots and subaerial stems. These substances possibly serve as a barrier against the entry of pathogens such as bacteria, fungi, and nematodes. During the tertiary phase of wall development polymerized and oxidized inclusion products are deposited in the cellulose walls.

Chromosome complement. The endodermis has been found to have extra sets of chromosomes as compared with cortical and other cells in the plant. In some plants the chromosome numbers may be so high in the endodermis that four sets of chromosomes may occur in each endodermal cell. The larger amount of nuclear material and nucleic acid in the cells of the endodermis may in part account for the great capacity of endodermal cells to produce large amounts of chemical substances high in caloric energy.

Acetylenic oils. Unsaturated fatty acids develop in the endodermis, resulting in acetylenic oils in hundreds of species of plants. Acetylenic oils, called polyynes, may have the highest caloric content of any organic substance in the world. Large amounts of energy are required to produce these oils. Because of their rare chemical structure, the oils are often toxic to many organisms. The functional value of the endodermis includes the making of these and many other kinds of chemical molecules, which may serve as protective agents against invasion by bacteria and fungi and may render the plant unpalatable or poisonous to insects and herbivores. See SECRETORY STRUCTURES (PLANT).

Functions. The endodermis has many apparent functions: absorption of water, selection of solutes and ions, and production of oils, antibiotic phenols, and acetylenic acids. See CORTEX (PLANT); HYPODERMIS; PLANT TISSUE SYSTEMS. Dick S. Van Fleet

Endophytic fungi

Fungi that live in the interior of plant host tissues without causing external symptoms. The term endophyte literally means "living within" and is contrasted with epiphyte, which means "living on the surface." These descriptive terms are used in several disciplines such as microbiology, mycology, and ecology, and at times their use is controversial. This article uses the term endophyte in a broad and topographical sense following the definition of O. Petrini (1991): "endophytes colonize asymptotically the living, internal tissues of their hosts, even though the endophyte may after an incubation or latency period cause disease."

Endophytes range from biotrophic parasites to transient facultative saprotrophs (organisms that live on dead and decaying organic matter), and as a result their associations with hosts span the continuum from biotrophic mutualists (two species whose mutual interactions are beneficial to both) to necrotrophic, antagonistic pathogens. The endophytic habit apparently has evolved independently several times and is represented by fungi in various orders of Ascomycetes and Basidiomycetes. See ASCOMYCOTA; BASIDIOMYCOTA; FUNGAL ECOLOGY; FUNGI; MYCOLOGY.

Grass endophytes. The evolution and ecological relevance of endophytes are dramatically illustrated in the case of the clavicipitalean grass endophytes. These fungi are widespread in the grass family in many natural populations. Grass endophytes produce systemic infection in the aboveground parts of the grass host, developing hyphae [filaments composing the vegetative body (mycelium) of fungi] in the intercellular spaces of leaf sheaths, culms (jointed, usually hollow grass stems) [Fig. 1], developing flowers, and seeds. With most infections of grasses, the seed becomes the dissemination and survival structure for both the grass plant and the fungal endophyte. Another characteristic of the clavicipitalean grass endophytes is the high specificity of compatibility with the plant host. In general, a specific endophyte is adapted to a particular plant species; rarely do two strains coexist in the same plant host. This is in stark contrast to non-clavicipitalean endophytes, where many species may be present in a single host individual.

Taxonomically, grass endophytes are mainly anamorphic (a sexual state) forms of genus *Neotyphodium* in the Ascomycota family Clavicipitaceae. *Neotyphodium* spp. are frequently interspecific hybrid strains derived from the sexual forms in genus *Epichloë*. *Epichloë* causes choke disease (Fig. 2) in numerous grasses, producing total or partial abortion of inflorescences at the time of sporulation.

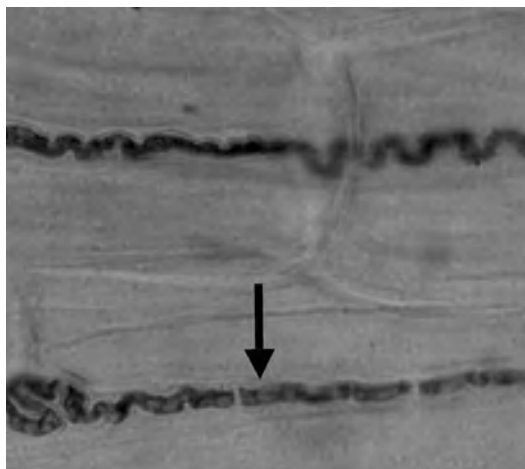


Fig. 1. *Neotyphodium coenophialum* in intercellular spaces of tall fescue grass. Arrow shows hyphae; Magnification 1000 \times .



Fig. 2. Choke disease caused by the clavicipitalean fungus *Atkinsonella texensis*.

The “sleepy grass” (*Achnatherum robustum*) common in western North America derives its name from the sleep-inducing alkaloid lysergic acid amide produced by an unidentified clavicipitalean endophyte inhabiting leaves, stems, and seeds of large natural populations of the grass. Horses that consume small quantities of this grass will sleep for up to 3 days, after which they recover gradually. See ALKALOID.

Similarly, in South America the clavicipitalean endophyte *Neotyphodium tembladerae* is widespread in several grasses—including *Poa buecu*, *Festuca hieronymi*, and *F. argentina*—where the fungus produces alkaloids that cause toxicity in grazing animals and result in muscular trembling and sometimes death. In northwestern China and Mongolia, an endophyte infects the grass commonly known as “drunken horse grass” (*Achnatherum inebrians*) and produces toxic alkaloids that cause a nervous disorder resulting in trembling and death. In South Africa, “drunk grass” (*Melica decumbens*) contains a clavicipitalean endophyte producing indole-terpenoid tremorgen (fungal mycotoxin that causes trembling) that produces nervous disorders in animals comparable to those of other endophytes.

Recognition that endophytes are responsible for many toxicoses of animal herbivores has developed only in the last three decades. C. W. Bacon et al. (1977) began the modern era of clavicipitalean endophyte investigation when they linked an endophyte (*Neotyphodium coenophialum*) in tall fescue (*Festuca arundinacea*) to widespread toxicities in cattle. The deluge of research that followed this historical epiphany included the pioneering ecological work of Keith Clay (1988) and the articulation of the “defensive mutualism hypothesis” that the clavicipitalean endophyte grass association is a mutualistic association where plants benefit from endophytes through production of fungal metabolites that are defensive against insect and mammalian herbivores.

It is now well supported that numerous clavicipitalean endophytes are defensive to their hosts. Recent research has extended endophyte defense of the hosts to include protection from fungal diseases of grasses. In certain grasses the endophytes exhibit an epiphytic stage where a network of mycelium extends over the epidermis of the grass leaf blades. These epiphyllous nets are believed to inhibit the colonization of leaf blades by potential pathogens of the host through a niche exclusion mechanism rather than through use of toxic metabolites. Endophyte-mediated resistance to drought in tall fescue grass has also been documented and may be interpreted as a defense against environmental challenges to plant survival. See FESCUE; GRASS CROPS.

Because of the many benefits of clavicipitalean endophytes in grasses, plant breeders have sought to take advantage of these endophytes to enhance performance of grasses for diverse applications. Seed containing endophytes can now be purchased from many seed suppliers. Turfgrass established with endophyte-infected seeds is hardier and more resistant to environmental extremes and pests requiring fewer chemical inputs to manage. For forage grasses, plant breeders have selected strains of endophytes that are low in production of alkaloids toxic to animals but high in alkaloids that impart insect resistance to grasses. See BREEDING (PLANT).

Although some fungi classified as endophytes may be early colonizing saprophytes, most endophytic fungi have arrived at endophytism generally from pathogenicity of the host through evolving reduced virulence or prolonged latency. In the Clavicipitaceae there is evidence that endophytes evolved initially from insect parasites. Clavicipitalean fungi are believed to have made a host shift to plants by infecting scale insects that are plant parasites themselves. This permitted clavicipitalean fungi to adapt to an epibiotic niche on plants prior to evolving endophytism. Once clavicipitalean fungi had adapted to plants as biotrophic pathogens (for example, sexually reproducing *Epichloë* spp.), evolutionary selection reduced virulence resulting in numerous asexual *Neotyphodium* endophytes. Because this process of ameliorating virulence occurred many times with different endophytes, asexual *Neotyphodium* endophytes are widespread in grasses in all habitats.

Curvularia endophytes. Non-clavicipitalean endophytes have shown similar enhancements in host resistance to environmental stresses. It has been demonstrated that a *Curvularia* species (fungal family Pleosporaceae) endophytic of rosette grass (*Dichanthellium lanuginosum*) substantially increases resistance of the infected plants to heat stress. Grass individuals containing the endophyte growing near hot springs were able to tolerate soil temperatures up to 65°C (149°F) for several days whereas endophyte-free grasses rapidly succumbed to heat stress. The mechanisms for endophyte-mediated tolerance to biological stresses such as drought and heat (thermo-tolerance) is currently unknown.

Fusarium endophytes. Another group of widespread endophytes and pathogens in many plant hosts from around the world are found in fungal genus *Fusarium* (Nectriaceae). One species, *Fusarium moniliforme*, is ubiquitous as an asymptomatic endophyte of corn wherever corn is grown. In some cases, *F. moniliforme* may exhibit a pathogenic phase (including ear, kernel, or stalk rots, and seedling blights), but generally it is entirely asymptomatic. It has been estimated that 90% of the corn grown worldwide may contain *F. moniliforme* endophytically and asymptotically. Like the clavicipitalean endophytes, *Fusarium* endophytes produce chemical metabolites, including fumonisins and other mycotoxins. Some of these compounds are extremely toxic to animals that consume infected plants. The associations between plant and *Fusarium* endophytes have been proposed to be a defensive mutualism in that infected corn plants show enhanced resistance to several disease-causing fungi.

It is generally hypothesized that most fungal species encountered as asymptomatic endophytes may be in some way mutualistic with the host plant. However, it has been difficult in numerous cases to document the beneficial effects to hosts of some endophytes.

Woody plants: foliar endophytes. In contrast with grass endophytes, foliar endophytes of woody plants are highly diverse within the plant host. Foliar endophytes are transmitted horizontally via spore fall from plant to plant, leaves are free of the endophytes when emerging, and they become infected shortly thereafter. Fungal spores are disseminated from older leaves by water, wind, and probably insects that germinate on the leaf surface, and they infect the leaves via cuticular penetration or growth through stomata. Endophyte-enhanced defense against plant pathogens was demonstrated for the chocolate tree (*Theobroma cacao*). Taxonomically, foliar endophytes of woody plants represent a broad range of taxa in several orders and families in the Ascomycota and Basidiomycota.

Sources of bioactive compounds. Endophytic fungi have been shown to be a source of bioactive compounds that may have potential as new medicines. For example, the Pacific yew endophyte *Taxomyces andreaneae* was shown to produce the anticancer drug taxol. This chemical previously had been believed to be produced only by the host plant. *Serratia marcescens*, the bacterial endophyte of the aquatic plant *Rhynchospora penicillata*, was found to produce the potent antifungal compound oocycin A. *Pestalotiopsis jesteri* has been shown to produce jesterone, another antifungal, apparently defensive compound. *Muscador albus* from the cinnamon tree (*Cinnamomum zeylanicum*) was found to produce volatile antibiotics that protect the host from other microbes, a phenomenon termed mycofumigation.

Common methods to study endophytes. Endophyte detection in grass endophytes can be carried out through microscopic observation of different plant



Fig. 3. Colony of *Neotyphodium lolii* isolated from *Lolium perenne*, perennial ryegrass, growing on potato dextrose agar.

parts using fresh plants or dried herbarium material. Sheaths of internal leaves in the tiller are appropriate in the vegetative stage, and flowering culms or seeds are suitable for endophyte detection. Leaf sheaths can be processed using a clearing agent followed by staining of the endophyte. For examination of endophytes in culm tissues, culms should be cut longitudinally and the parenchyma tissue scraped out with a scalpel. Then the tissue is transferred to a drop of aniline blue or rose bengal stain on a slide and observed under the microscope. To determine the presence of endophytes in seeds, seeds are soaked overnight in a 5% solution of sodium hydroxide, after which they undergo several washes with tap water for removal of sodium hydroxide followed by microscopic observation using aniline blue or rose bengal stains. For clavicipitalean endophytes there are commercially available immunoblot assay kits for detection of the fungi in plant leaves and seeds of tall fescue and perennial ryegrass, and also for detection of ergot alkaloids produced by the endophytes in plants.

For identification of endophytes in general, the common approach is direct isolation from the host plant. The procedure consists of surface sterilization of host tissues followed by culture in standard microbiological media (Fig. 3). Small pieces of host tissue are placed on the surface of agar medium (for example, potato dextrose agar). If the isolate sporulates, it is further characterized by its morphological features using a light microscope. If isolation results in sterile mycelia, the use of nucleic acid sequence approaches will add to strain characterization, allowing the researcher to determine the approximate phylogenetic position and assign a taxonomic category.

[Acknowledgments: The authors are grateful for support from the Fogarty International Center (NIH) under U01 TW006674 for International Cooperative Biodiversity Groups and the Rutgers Turfgrass Research Center.] Monica S. Torres; James F. White, Jr.

Bibliography. C. W. Bacon and J. F. White, Jr. (eds.), *Microbial Endophytes*, Marcel Dekker, 2000; K. Clay, Fungal endophytes of grasses: A defensive mutualism between plants and fungi, *Ecology*, 69:10-16, 1988; S. Faeth, Are endophytic fungi defensive plant mutualists?, *Oikos*, 98:25-36, 2002; O. Petrini, Fungal endophytes of tree leaves, in J. H. Andrews and S. S. Hirano (eds.), *Microbial Ecology of Leaves*, pp. 179-187, Springer-Verlag, New York, 1991; C. L. Schardl, A. Leuchtman, and M. J. Spiering, Symbioses of grasses with seedborne fungal endophytes, *Annu. Rev. Plant Biol.*, 55:315-340, 2004; J. K. Stone, J. D. Polishook, and J. F. White, Jr., Endophytic fungi, pp. 242-270 in G. M. Mueller, G. F. Bills, and M. S. Foster (eds.), *Biodiversity of Fungi: Inventory and Monitoring Methods*, Elsevier/Academic Press, 2004.

Endoplasmic reticulum

An intracellular membrane system that is present in all eukaryotic cells. In most cells the endoplasmic reticulum is thought to consist of only one continuous membrane enclosing only a single space. However, in protozoa, some unicellular algae, and possibly some fungi, the endoplasmic reticulum occurs as separate, multiple vesicles. See CELL MEMBRANES.

Subdomains. Several morphologically and functionally distinct domains of this continuous membrane system can be distinguished. At the level of the nuclear pores, the inner nuclear membrane is continuous with the outer nuclear membrane; both membranes together are referred to as the nuclear envelope. The outer nuclear membrane in turn is continuous with the rough endoplasmic reticulum, which contains specialized regions, termed transitional elements, and is continuous with the smooth endoplasmic reticulum. The two membranes of the nuclear envelope enclose the perinuclear space. The rough and smooth endoplasmic reticula and the transitional element enclose a space called the intracisternal space, or lumen. Both intracisternal and perinuclear spaces form a single compartment. All nucleated cells contain at least a nuclear envelope, but the amount of smooth and rough endoplasmic reticula varies greatly among different cell types. See CELL NUCLEUS.

This single-membrane system undergoes an extensive and reversible fragmentation during mitosis. During the first stage of mitosis (prophase), extensive fission of the nuclear envelope and the endoplasmic reticulum takes place, yielding numerous closed vesicles. At the end of mitosis (telophase), the separated vesicles fuse and a single continuous membrane is restored. Fragmentation ensures that the two daughter cells each end up with approximately half of the mother's nuclear envelopes and endoplasmic reticulum. See MITOSIS.

When cells are broken by mechanical homogenization, the endoplasmic reticulum and often part of the outer nuclear membrane are fragmented. To a large extent, the fragments close at their ends to form vesi-

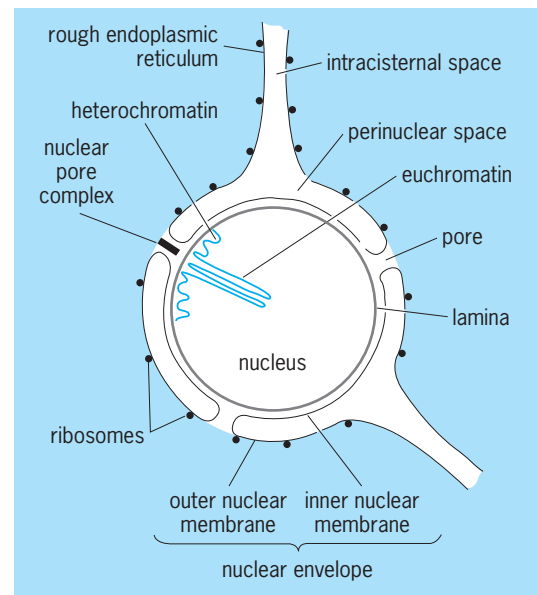


Fig. 1. Nuclear envelope and associated structures.

cles that can be separated on the basis of their size and density, yielding fractions that represent rough microsomes, which are derived from rough endoplasmic reticulum, and smooth microsomes, which are derived from the smooth endoplasmic reticulum. Biochemical analyses of these fractions have revealed many of their functions.

Nuclear envelopes. When cells are broken mechanically, most of the connections between nuclear envelope and rough endoplasmic reticulum (Fig. 1) are severed. The nuclei, still surrounded by their nuclear envelopes, can easily be separated from the homogenate because of their large size and high density. Following enzymatic digestion of the chromatin with nuclease, the nuclear envelope can readily be isolated from the nuclease digest with its principal associated organelles, the nuclear pore complexes, and the nuclear lamina still attached. Inner and outer nuclear envelope membranes have not yet been separated.

Associated with the nuclear envelope are the nuclear pore complexes (Fig. 1). An average mammalian nucleated cell contains approximately 5000 pore complexes, each with an estimated mass of 10^8 daltons. Pore complexes are thought to regulate macromolecular traffic into and out of the nucleus and to function in the three-dimensional organization of the genome.

Associated with the nuclear side of the inner nuclear membrane is the nuclear lamina, a fibrous meshwork 15-50 nanometers thick that is intercalated between inner nuclear membrane and chromatin. The nuclear lamina is composed of the lamins, which structurally belong to the family of intermediate filament proteins. The lamins (specifically lamin B) are physically connected to the cytoplasmic intermediate filament proteins at the level of the nuclear pore complex and are thought to function in the three-dimensional organization of the chromatin.

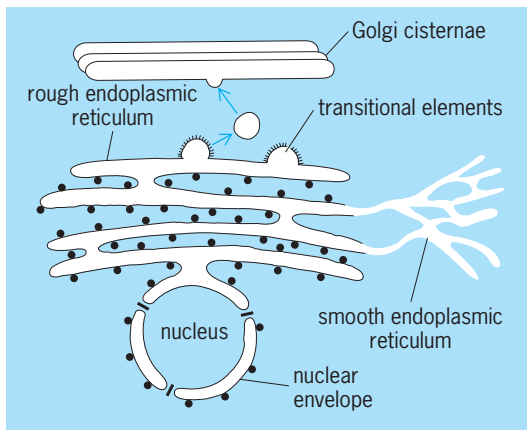


Fig. 2. Nuclear envelope, connected to the rough and smooth endoplasmic reticulum. The rough endoplasmic reticulum is linked to the cis cisterna of the Golgi complex by vesicles that shuttle between the two compartments.

Associated with the outer nuclear membrane are ribosomes. In fact, the outer nuclear membrane and the rough endoplasmic reticulum appear to be functionally equivalent.

The inner and outer nuclear membranes are likely to differ in their composition of integral membrane proteins. So far, only one such protein, the lamin B receptor, has been assigned exclusively to the inner nuclear membrane. This protein can therefore serve as a marker for the inner nuclear membrane. Another integral membrane glycoprotein of 190 kilodaltons has been assigned to the pore-wall region of the nuclear envelope; it is thought to function in anchoring nuclear pore complexes.

Rough endoplasmic reticulum. The term rough endoplasmic reticulum is based on the morphologic appearance of attached ribosomes, which are absent from smooth endoplasmic reticulum. Another morphologic distinction is the organization of the former in interconnected flattened sacs (called cisternae), whereas the latter forms a tubular network (**Fig. 2**). *See RIBOSOMES.*

The rough endoplasmic reticulum is the site of translocation of secretory and lysosomal proteins from the cytosol to the intracisternal space, and of integration into the membrane of integral membrane proteins. Except for integral membrane proteins of chloroplast, mitochondria, and peroxisomes, essentially all other integral membrane proteins are integrated into the endoplasmic reticulum and either remain there (resident endoplasmic reticulum membrane proteins) or are subsequently distributed to other cellular membranes.

The signal hypothesis was formulated to explain how these proteins are targeted to and then translocated across or integrated into the endoplasmic reticulum membrane. Its tenets are that all polypeptides targeted to this membrane contain a discrete sequence (termed the signal sequence), that a complex machinery recognizes this sequence, and that recognition triggers the opening of a proteinaceous channel (composed of several subunits) through which the polypeptide passes across the membrane. In

the case of membrane proteins, the existence of an additional topogenic sequence, the so-called stop-transfer sequence, was postulated. This sequence is thought to trigger opening of the channel to the lipid bilayer to abort translocation and thus integrate the protein into the lipid bilayer.

Many predictions of the signal hypothesis have been verified. The presence of a signal sequence in these proteins has been demonstrated. The critical information for targeting is encoded in about 10–15 amino acid residues of the signal sequence, which in most cases is cleaved during or shortly after translocation by a special endoplasmic reticulum-associated signal peptidase. The machinery for recognition of the signal sequence and targeting to the endoplasmic reticulum is complex. As soon as the signal sequence has emerged from within a channel of the large ribosomal subunit, it is recognized by the signal recognition particle, a ribonucleoprotein particle consisting of one 7S ribonucleic acid (RNA) molecule and six distinct polypeptide chains. The ensemble of ribosome, signal recognition particle, and nascent chain is then targeted to the signal recognition particle receptor (located exclusively in the endoplasmic reticulum). This receptor displaces the signal recognition particle and frees the signal sequence to interact with a second signal recognition system in the endoplasmic reticulum membrane. The second recognition step may be coupled to the opening of a proteinaceous channel; however, the existence of such a channel has not yet been demonstrated. The chain would then pass through that channel, or in the case of integral membrane proteins, the channel would be opened to the lipid bilayer by a stop-transfer sequence. After passage, the channel would close, terminating translocation.

The fact that the integration of integral membrane proteins requires proteins of the same type indicates that cells cannot assemble membranes *de novo*. Pre-existing membranes with an integrated translocation machinery are required to make more such machinery and more membranes. Thus, if the daughter cells would not inherit endoplasmic reticulum from the mother cell, they would be able to synthesize but not integrate membrane proteins.

The rough endoplasmic reticulum also contains numerous enzymes, most of which are involved in the modification of the nascent chain on the cisternal side. For example, all the enzymes (and the required sugar nucleotide transporters) for the synthesis of a lipid-linked oligosaccharide and for the transfer of the oligosaccharide moiety to an asparagine residue of a polypeptide are located in the rough endoplasmic reticulum.

Thus the main function of the rough endoplasmic reticulum and the outer nuclear membrane is to serve as a port of entry of secretory, lysosomal, and integral membrane proteins and as the site of their initial modification. *See CYTOCHEMISTRY.*

Transitional elements. Secretory and lysosomal proteins as well as those integral membrane proteins that are not residents of the endoplasmic reticulum

are next transported to the cis Golgi cisternae. The transitional elements represent sites of transport from the rough endoplasmic reticulum. Coated vesicles carrying proteins to be transported form at these sites and, after uncoating, eventually fuse with the cis Golgi cisternae. *See* GOLGI APPARATUS.

Smooth endoplasmic reticulum. Smooth endoplasmic reticulum contains enzymes for phospholipid biosynthesis, steroid biosynthesis, and drug detoxification. The principal phospholipid synthesized in the smooth endoplasmic reticulum is phosphatidyl choline. The enzymes for assembly from fatty acyl coenzyme A, glycerophosphate, and cytosine diphosphatecholine are integral membrane proteins with their active sites exposed to the cytosol. Phosphatidyl choline then becomes part of the cytoplasmic leaflet of the endoplasmic reticulum lipid bilayer. Transport to the luminal leaflet may be catalyzed by an enzyme called a flippase. Certain hydrophobic compounds that dissolve in the lipid bilayer can be hydroxylated and further modified by enzymes, thus becoming water-soluble and excretable. *See* CELL (BIOLOGY); CELL ORGANIZATION; ENZYME.

Günter B. Blobel

Bibliography. W. M. Becker, *The World of the Cell*, 4th ed., 2000; T. C. Cheng (ed.), *The Structure of Membranes and Receptors*, 1984; C. A. DeDuve, *A Guided Tour of the Living Cell*, 1985; E. Holtzman and A. B. Novikoff, *Cells and Organelles*, 3d ed., 1984; R. C. Warren, *Physics and the Architecture of Cell Membranes*, 1987; S. L. Wolfe, *Cell Ultrastructure*, 1985.

Endopterygota

A division (also known as Holometabola) of the subclass Pterygota, including those insects that undergo complete metamorphosis during their life cycle. That is, individual development goes through four distinct stages: egg, larva (trophic), pupa (reconstructive), and adult (reproductive). The larval and adult stages often live in very different adaptive zones and have very different life forms, necessitating a quiescent pupal stage in which extensive restructuring takes place. Typically growing through a sequence of immature forms (instars) punctuated by integumental molts, the larval instars correspond to those of nymphal Exopterygota in most ways, but are distinct in having wings, internalized as tiny wing buds. The buds are tucked out of the way of the larva, grow very rapidly in the pupal stage, and expand and quickly become functional in the newly hatched adult.

Internalization of the developing wings and the advent of the pupa combine to form a key general adaptation that breaks the constraining linkage between trophic and reproductive life forms, allowing each to evolve for its own special function. For example, most larvae have evolved life forms with legs that are reduced or absent, enabling them to live and feed cryptically within a substrate. This adaptation probably accounts for the rapid rise of the Endopterygota

to a dominant position among the insects from the Permian to the present.

The division Endopterygota comprises a large majority of all insects, distributed in the orders Coleoptera, Neuroptera, Strepsiptera, Mecoptera, Siphonaptera, Trichoptera, Lepidoptera, Diptera, and Hymenoptera. It is thought to be monophyletic. *See* EXOPTERYGOTA; INSECTA; PTERYGOTA.

William L. Brown, Jr.

Bibliography. O. W. Richards and R. G. Davies, *Imm's General Textbook of Entomology*, 2 vols., 10th ed., 1994; S. P. Parker (ed.), *Synopsis and Classification of Living Organisms*, vol. 2, 1982.

Endorphins

A family of endogenous, morphinelike peptides present within the central nervous system. The enkephalins, the first members of the endorphin family, were discovered by John Hughes and Hans Kosterlitz in 1975 and opened research into a wide number of centrally active neuropeptides. The term endorphin is generic, referring to all the opioid peptides, while specific peptides are given individual names, such as the enkephalins, dynorphins, and β -endorphin. Their isolation, shortly after the identification of the opioid receptors, brought together the concept of endogenous opioid peptide systems within the central nervous system which can modulate pain perception and through which opioids act. *See* OPIATES; PAIN.

Morphine, codeine, and their many synthetic and semisynthetic analogs are effective painkillers that act through specific recognition sites, or receptors, localized on the surface of neurons within selected brain regions. The existence of these highly specific receptors implied that, for example, morphine was mimicking endogenous compounds within the brain with morphinelike actions. This concept of an endogenous pain-modulating system provided an explanation for prior clinical and experimental observations. In a classic study, it was reported that soldiers wounded in World War II required fewer painkillers than civilians undergoing elective surgery, even though the soldiers had far more extensive injuries. Clearly, some aspect of the stress of combat affected the perception of pain. More than 20 years later, scientists reported that electrical stimulation of a specific brain region, the periaqueductal gray, produces pain relief, or analgesia, in animals; this was demonstrated later in humans. This pain relief was reversed by the drug naloxone, a potent and highly selective opiate antagonist that reverses morphine actions. The ability of naloxone to reverse this analgesia implied that the stimulation was releasing endogenous morphinelike, or opioid, materials, which have since been termed endorphins. Three families of opioid receptors involved with pain modulation have been identified that selectively interact with different classes of opioids and endorphins: mu (morphine and the endomorphins), kappa

TABLE 1. Structures of the opioid peptides

[Leu]enkephalin	Tyr-Gly-Gly-Phe-Leu
[Met]enkephalin	Tyr-Gly-Gly-Phe-Met
Dynorphin A	Tyr-Gly-Gly-Phe-Leu-Arg-Ile-Arg-Pro-Lys-Leu-Lys-Trp-Asp-Asn-Gln
Dynorphin B	Tyr-Gly-Gly-Phe-Leu-Arg-Arg-Gln-Phe-Lys-Val-Val-Thr
α -Neoendorphin	Tyr-Gly-Gly-Phe-Leu-Arg-Lys-Tyr-Pro
β -Endorphin	Tyr-Gly-Gly-Phe-Met-Thr-Ser-Glu-Lys-Ser-Gln-Thr-Pro-Leu-Val-Thr-Leu-Phe-Lys-Asn-Ala-Ile-Ile-Lys-Asn-Ala-Tyr-Lys-Lys-Gly-Glu
Endomorphin 1	Tyr-Pro-Trp-Phe-NH ₂
Endomorphin 2	Tyr-Pro-Phe-Phe-NH ₂

(dynorphins), and delta (enkephalins). See MORPHINE ALKALOIDS.

Structure and processing. The enkephalins were the first endogenous opioid peptides, or endorphins, to be isolated. These two pentapeptides have the same first four amino acids followed by either methionine or leucine. A number of opioid peptides have been isolated (Table 1). With the exception of the endomorphins, all contain the sequence of either methionine (Met) enkephalin or leucine (Leu) enkephalin at their NH₂ terminus. This common sequence among the various peptides raised questions as to whether the enkephalins might be breakdown products of larger peptides, such as β -endorphin and the dynorphins. However, three separate genes have been cloned that encode the precursor proteins for the enkephalins, the dynorphins, and β -endorphin, which clearly demonstrates the independence of each of these endorphin families of peptides.

The enkephalins are generated from a larger precursor peptide through the actions of enkephalin convertase, a carboxypeptidase. The pro-enkephalin gene codes for a peptide containing six copies of [Met⁵]enkephalin and one copy of [Leu⁵]enkephalin which is generated by enzymatic cleavage. The enkephalins are very sensitive to degradation by endopeptidases, particularly enkephalinase. Inhibitors of this enzyme have been synthesized and have modest analgesic actions, presumably by enhancing the enkephalin levels in the nervous system. The prob-

lem of metabolism has been greatly diminished by substituting an unnatural D-amino acid in the second position of a series of enkephalin analogs. A large variety of compounds have thus been synthesized, some of which are quite potent and active in animals and humans, even after systemic administration. Interestingly, alterations in the sequence of these peptides can change their selectivity toward the classes of opioid receptors.

The dynorphins and α -neoendorphin are derived from the pro-dynorphin gene, also termed the pro-enkephalin B gene, which is distinct from the one involved with enkephalin synthesis. Again, the presence of two basic amino acids directs processing to form dynorphin A, dynorphin B, and α -neoendorphin. These peptides have not been studied as extensively as the enkephalins, and few analogs are available. Like the enkephalins, all three are very sensitive to proteases and are broken down very rapidly in the brain.

β -Endorphin is perhaps the most interesting peptide. It is encoded by the pro-opiomelanocortin gene. This gene generates a large peptide, termed Big ACTH, which then is processed to form β -lipotropin and the stress hormone adrenocorticotropin (ACTH). β -Lipotropin is further processed to yield β -endorphin and β -MSH (melanocyte-stimulating hormone). Thus, β -endorphin is cogenerated with important, nonopioid hormones. See ADENOHYPOPHYSIS HORMONE.

Regional distributions and actions. The enkephalins are distributed unevenly throughout the brain with very high levels in the basal ganglia, the thalamus, and the periaqueductal gray (Table 2). In addition, there are high concentrations of enkephalins in the adrenal medulla, where they are coreleased with norepinephrine in response to stress, among other stimuli. The dynorphins and α -endorphin are located within the central nervous system with a distribution similar to that of the enkephalins. See STRESS (PSYCHOLOGY).

β -Endorphin has a unique distribution within the brain compared to the other opioid peptides. Unlike the enkephalins and dynorphins, which have high concentrations within a large number of brain

TABLE 2. Tissue distribution of opioid peptides

Location	Enkephalin	β -Endorphin	Dynorphin
Central nervous system			
Spinal cord	Dorsal horn (laminae I and II)	Absent	—
Brainstem	Cells and fibers	Fibers only	Fibers: medial forebrain bundle
Cerebellum	Absent	Absent	Absent
Diencephalon	High: thalamus, hypothalamus	Fibers: thalamus Cells: hypothalamus	Fibers: substantia nigra Cells: hypothalamus
Basal ganglia	High: globus pallidus	Low	Medium
Telencephalon	High: amygdala Low: hippocampus, cerebral cortex	Low	High: hippocampus Medium: cerebral cortex Low: amygdala
Peripheral tissue			
Pituitary gland	Low	Intermediate lobe	Posterior lobe
Gastrointestinal tract	Present	Present	Present
Adrenal medulla	High	—	—
Autonomic nervous system	Present	—	—
Placenta	—	High	—

regions, β -endorphin has been identified in only a single group, or nuclei, of cells within the hypothalamus that project to a number of regions. Its highest levels are in the pituitary gland. Within the pituitary, both ACTH and β -endorphin are derived from the same precursor protein and are located within the same cells. Stimuli that release ACTH, a stress hormone which in turn induces the adrenal gland to release steroids, also corelease β -endorphin at the same time. Thus, stressful stimuli that release ACTH and norepinephrine also release both β -endorphin from the pituitary and enkephalins from the adrenal into the blood. This is particularly intriguing in view of the decreased perception of pain reported under periods of stress. See ENDOCRINE SYSTEM (VERTEBRATE); PITUITARY GLAND.

All the endorphins are analgesic despite the fact that they act through different classes of opiate receptors. Unlike the stabilized synthetic enkephalin analogs, the natural endorphins are not active when given systemically due to their rapid degradation and their difficulty traversing the blood-brain barrier. The natural endorphins relieve pain only when given directly into the brain. However, the endorphins are likely to have a far more diverse range of actions beyond modulation of pain, as suggested by their presence in brain regions unrelated to pain perception, such as the striatum. Thus, many questions remain regarding the full physiological significance of these agents. See NERVOUS SYSTEM (VERTEBRATE).

Gavril Pasternak

Bibliography. B. M. Cox, A. Goldstein, and C. H. Li, Opioid activity of a peptide, β -lipotropin(61-91), derived from a β -lipotropin, *Proc. Nat. Acad. Sci. USA*, 73:1821, 1976; A. Goldstein et al., Porcine pituitary dynorphin: Complete amino acid sequence of the biologically active heptapeptide, *Proc. Nat. Acad. Sci. USA*, 78:7219-7223, 1981; J. G. Hardman and L. E. Limbird (eds.), *Goodman & Gilman's The Pharmacological Basis of Therapeutics*, 2001; J. Hughes et al., Identification of two related pentapeptides from the brain with potent opiate agonist activity, *Nature*, 258:577-579, 1975; D. J. Mayer and J. C. Liebeskind, Pain reduction by focal electrical stimulation of the brain: An anatomical and behavioral analysis, *Brain Res.*, 68:73-93, 1973.

Endotoxin

A biologically active substance produced by bacteria and consisting of lipopolysaccharide, a complex macromolecule containing a polysaccharide covalently linked to a unique lipid structure, termed lipid A. All gram-negative bacteria synthesize lipopolysaccharide, which is a major constituent of their outer cell membrane. One major function of lipopolysaccharide is to serve as a selectively permeable barrier for organic molecules in the external environment. Different types of gram-negative bacteria synthesize lipopolysaccharide with very different polysaccharide structures. See CELL MEMBRANES; LIPID; POLYSACCHARIDE.

Lipopolysaccharide. The polysaccharide component of lipopolysaccharide is often divided into two regions—the core polysaccharide and the O-antigen polysaccharide. The core region contains common hexose sugars, such as glucose and galactose (found in the outer core), and unique 7- and 8-carbon sugars, such as ketodeoxyoctulosonic acid (localized in the inner core). The O antigen consists of repeating oligosaccharide units, each containing four to six hexose sugars. The O antigen is different for each strain of microorganism, whereas the core polysaccharide is more consistent among species. Some pathogens, such as *Neisseria* and *Hemophilus*, synthesize a lipopolysaccharide which is devoid of O antigen. See GLUCOSE; GONORRHEA; MENINGITIS; OLIGOSACCHARIDE.

Covalently attached to the ketodeoxyoctulosonic acid of the inner-core polysaccharide is lipid A, with a structure that is very similar in almost all gram-negative microorganisms. This lipid consists of a backbone of two molecules of the sugar *N*-acetylglucosamine, to which are attached five to seven long-chain fatty acids. Sometimes attached to lipid A are additional phosphate or phosphorylethanolamine groups. Since lipid A is a chemically conserved structure among all species of gram-negative bacteria, the biological activity of endotoxic lipopolysaccharide resides almost entirely in the lipid A component

Host responses. When lipopolysaccharides are released from the outer membrane of the microorganism, either by natural cell death, by destruction of the bacterium following treatment with antibiotics, or by natural shedding, significant host responses are initiated in humans and other mammals. It is generally accepted that lipopolysaccharides are among the most potent microbial products, known for their ability to induce pathophysiological changes, in particular fever and changes in circulating white blood cells. For example, in humans as little as 4 nanograms of purified lipopolysaccharide per kilogram of body weight (0.0000008 ounce in a normal adult male) is sufficient to produce a rise in temperature of about 3.6°F (2°C) in several hours. This profound ability of the host to recognize endotoxin is thought to serve as an early warning system to signal the presence of gram-negative bacteria.

Unlike most microbial protein toxins (which have been termed bacterial exotoxins), endotoxin is unique in that its recognized mode of action does not result from direct damage to host cells and tissues. Rather, endotoxin stimulates cells of the immune system, particularly macrophages, and of the vascular system, primarily endothelial cells, to become activated and to synthesize and secrete a variety of effector molecules that cause an inflammatory response at the site of bacterial invasion. These effector molecules include a number of important cytokines, particularly tumor necrosis factor, interleukin-1 and -6, and interferon, as well as inflammatory mediators, including prostaglandins and leukotrienes. These mediator molecules promote the host response which results in elimination of the

invading microbe. One important consequence of endotoxin interaction with phagocytic cells is that these cells become more efficient at engulfing and killing bacteria. Endotoxin-stimulated macrophages also become more efficient at killing tumor cells. Thus, under these circumstances lipopolysaccharide is not a toxin at all, but serves an important function by helping to mobilize the host immune system to fight infection. *See* CYTOKINE; IMMUNOLOGY.

Even though endotoxin stimulation of host cells is important to host defense against infection, overstimulation due to excess production of endotoxin can lead to serious consequences. When excess mediator production occurs, the elicited pathophysiologic changes (termed the systemic inflammatory response syndrome) are expanded to include hypotension (abnormally low blood pressure), blood clotting, and changes in metabolism, heart rate, and respiration. These changes can severely damage many organs, including the lung, liver, and kidney. Endotoxin-induced multiple-organ failure continues to be a major health problem, particularly in intensive care; it has been estimated that as many as 50,000 deaths annually occur in the United States as the result of endotoxin-induced shock.

Immunization. Because of the extensive variability in chemical structure of the O antigen of various species, and the wide spectrum of gram-negative microorganisms that can cause disease, immunization of humans with endotoxin vaccines to protect against endotoxin shock has not been considered practical. In contrast, with bacterial protein exotoxins (such as diphtheria toxin, pertussis toxin, tetanus toxin) immunization leads to long-lasting immunologic protection against disease. Efforts to provide immunologic protection against endotoxin-related diseases have focused upon development of antibodies that recognize the conserved lipid A structure of endotoxin as a means of passive protection against the lethal effects of this microbial product. *See* BACTERIA; MEDICAL BACTERIOLOGY; VACCINATION.

David C. Morrison

Bibliography. D. C. Morrison and J. L. Ryan (eds.), *Bacterial Endotoxic Lipopolysaccharides*, vol. 1, 1992; E. T. Rietschel and H. Brade, Bacterial endotoxins, *Sci. Amer.*, 267:54–61, 1992; J. L. Ryan and D. C. Morrison (eds.), *Bacterial Endotoxic Lipopolysaccharides*, vol. 2, 1992.

Energy

One of the two fundamental constituents of the universe, interconvertible with the other constituent, matter. All physical processes involve exchange of energy, or conversion from one of its named forms (including relativistic mass-energy) into another. The total energy (including mass) of a closed system is a conserved quantity. *See* MATTER (PHYSICS).

Intuitive notions. On the human scale, energy is most familiar in the form of mechanical work and transfer of heat. The word “energy” is derived from Greek meaning “inner work.” This is intuitively rea-

sonable from the everyday notion of doing work because physical labor expends the body’s energy.

The causal relationships among energy, force, work, power, kinetic energy, potential energy, heat, and temperature are also intuitive. A person pushing a rock up a slope applies the force of his or her muscles. While the rock is in motion, it has kinetic energy (relative to the ground) that could, for example, dent a tree trunk that might be in the rock’s path. The amount of work done increases proportionally with the distance the rock is moved, and it results in a proportionate reduction of the person’s internal energy resources; hence the notion that work and energy spent are equivalent. To push the rock faster, the person must apply more power, or energy per unit time. The work done on the rock is converted into an increase of gravitational potential energy, because the rock is now higher up on the slope, and into heat energy, resulting from friction, that raises the temperature of the rock and proximate ground. The potential energy of the rock is so-named because it has the potential to be converted into other forms of energy such as kinetic energy in the event that the rock should be allowed to roll freely downhill. *See* FORCE; HEAT; POWER; TEMPERATURE; WORK.

Formal and quantitative concepts. The intuitive notions of energy and its related parameters were first quantified with Newton’s three laws of motion. Classical mechanics recognizes two types of energies:

1. Energy due to relative motion, or kinetic energy: A mass m moving with speed v relative to an observer has kinetic energy $T = \frac{1}{2}mv^2$.

2. Energy that is stored in restraining motion against forces, or potential energy: For example, lifting a mass m from ground level to a height b increases its gravitational potential energy by $\Delta U = mgb$, where g is the downward acceleration due to Earth’s gravitational field. *See* GRAVITATION; NEWTON’S LAWS OF MOTION.

With the development of atomic and subatomic physics came understanding of new types of forces with which potential energies could be associated. The electric and magnetic forces, with their associated potentials, act on electrical charges. In the nineteenth century, they were revealed to be two manifestations of the electromagnetic force. Twentieth-century research revealed that subatomic particles interact via the weak nuclear force and the strong nuclear force. *See* ELECTRIC FIELD; FUNDAMENTAL INTERACTIONS; MAGNETISM; MAXWELL’S EQUATIONS; STRONG NUCLEAR INTERACTIONS; WEAK NUCLEAR INTERACTIONS.

Most force fields of interest, such as gravity and electromagnetic force, are conservative. This means that taking a particle that “feels” that force and moving it through any round trip (starting from any point in the field and ending on the same point) will require no net energy input (the particle’s potential energy will be unchanged).

Units of energy. The unit of energy in the International System (SI) is the joule (J), defined as the force of one newton (N) acting through a distance of one meter. Power, the time rate of flow of energy, is

measured in watts (W): one watt of power is defined as the rate of energy flow of one joule per second. Electrical utilities use the kilowatt-hour (kWh), the energy delivered by one kilowatt over a period of one hour; one kWh equals exactly 3.600×10^6 joules. Energies of chemical interactions are measured in kilocalories per mole, where one calorie equals 4.187 joules. A kilocalorie of heat will raise the temperature of a liter of water by about one degree Celsius. *See* METRIC SYSTEM; PHYSICAL MEASUREMENT; UNITS OF MEASUREMENT.

In high-energy physics, subatomic-particle interactions frequently involve mass changing to radiant energy, and vice versa. It is therefore more convenient to express both energy and mass by the same unit of measure, the electronvolt (eV): the change in potential energy of one electron when moved through an electrical potential difference of one volt; 1 eV equals 1.602×10^{-19} joule. *See* ELECTRONVOLT.

Law of conservation of energy. The law of conservation of energy, also known as the first law of thermodynamics, is one of the most important laws of physics. It states that the total energy of closed, isolated systems is always constant. (More precisely, it is the mass-energy that is conserved, of major concern only for relativistic velocities and nuclear reactions.) This law is a result (according to Nöther's theorem) of the homogeneity of the universe with respect to time; that is, with all else being equal, the outcome of any experiment is independent of when it is performed. *See* CONSERVATION LAWS (PHYSICS); CONSERVATION OF ENERGY; SYMMETRY LAWS (PHYSICS).

Many periodic phenomena are due to the cyclic interchange between potential energy and kinetic energy, with their sum always constant. The classical example is a swinging pendulum whose potential energy is maximum at the highest part of its swing, its kinetic energy at the lowest. *See* PENDULUM.

Energy and special relativity. Albert Einstein's special theory of relativity introduced the notion of mass-energy $E = mc^2$, where c is the speed of light. In the mass's own frame of reference, within which the mass is stationary, this formula gives the rest-mass energy. For small (nonrelativistic) velocities, the total energy approximately equals the rest-mass energy plus the term $\frac{1}{2}mv^2$, exactly the classical formula for kinetic energy. *See* REST MASS.

When velocities are routinely relativistic, as in high-energy physics, it is convenient to express universally true physical relations as four-vectors. The momentum four-vector (also called the momentum-energy) embodies the conservation laws of both energy and momentum for all moving observers. *See* CONSERVATION OF MOMENTUM; MOMENTUM; RELATIVITY.

Energy in quantum mechanics. Quantum mechanics, and its offshoot quantum field theory, is founded on the principle of the quantum, the smallest amount of energy exchangeable in a given interaction. (Linear momentum, angular momentum, and mass are also quantized.) *See* QUANTUM FIELD THEORY; QUANTUM MECHANICS.

One interesting property of quanta of electromagnetic radiation is that they cannot have energies of zero even in empty space. This nonzero energy is called the zero-point energy and is experimentally confirmed by the Casimir effect: two parallel, closely spaced, uncharged conducting plates in empty space will attract each other.

Mechanical energy can be observed to be quantized: vibrational energy propagating through periodic media, such as silicon crystals, does so in quanta called phonons; essentially, these are quanta of sound energy. *See* LATTICE VIBRATIONS; PHONON.

A second principle of quantum mechanics, Heisenberg's uncertainty principle, limits the accuracy with which energy (as well as length, velocity, time, and position) may be measured. The error ΔE in measuring energy over time interval Δt is related by the inequality $\Delta E \cdot \Delta t \geq h/(2\pi)$, where h is Planck's constant. One consequence is that as masses (which, of course, are equivalent to energies) are examined at ever-shorter time scales, the uncertainties in the measured values increase. For short-enough particle interactions, particle physicists can invoke virtual particles that exist only briefly during the transformation of colliding particles into postcollision products. During their putative existence, they may violate the law of conservation of energy, a paradox resolved by the fact that they are never actually observed. *See* UNCERTAINTY PRINCIPLE.

System evolution, entropy, and the second law of thermodynamics. In a closed, isolated system with a certain distribution of energy within it there are many (usually practically uncountable) ways the energy could be redistributed. The second law of thermodynamics indicates what the ultimate distribution of energy within the system will be after a long time. It invokes the concept of entropy, a measure of the disorder within the system. Since the number of ways the water molecules can be arranged in a steam vessel is larger than the ways they can be arranged to form ice, a given quantity of water has a lower entropy as ice crystals than as steam. The second law states that a closed, isolated system will evolve over time so as to maximize the entropy. Thus, heat energy flows (if allowed) from hotter areas to colder ones, eventually equalizing the temperature. This maximizes the entropy; it also serves as an indicator of the flow of time (specifically, thermodynamic time). *See* ENTROPY; THERMODYNAMIC PRINCIPLES; TIME, ARROW OF.

All natural processes that locally decrease the entropy increase it elsewhere. When sodium and chlorine atoms migrate in aqueous solution to form a salt crystal, their entropy is decreased. However, the heat of crystallization that is released warms the solution, thereby increasing its entropy. Similarly, the biological processes of life create order but at the expense of the surroundings: consuming food creates highly ordered tissue but increases disorder via waste products and heat.

Types of energy. A number of named types of energy are found in scientific and technical literature. Most are macroscopic manifestations of

electromagnetic interactions at the atomic level.

Heat became understood, with the development of the atomic theory in the nineteenth century, as a manifestation of the collective kinetic energy of atoms in bulk materials. Temperature is a measure of their average kinetic energy. In solids, the motions are confined to vibrations of the atoms about equilibrium positions. In gases, the energy is equally partitioned into linear flight (between collisions), rotations, and vibration modes of the molecules. *See* KINETIC THEORY OF MATTER; STATISTICAL MECHANICS.

Elastic energy is potential energy stored when solids are deformed by various kinds of strains, such as compression or torsion. It is due to dislocation of the constituent atoms of the solid from their "relaxed" minimum-energy configurations. *See* ELASTICITY.

Surface-tension energy, which is responsible for the spherical shape of droplets, is the potential energy stored in the tangential cohesive forces between molecules at the surface of a liquid. *See* SURFACE TENSION.

The field of chemistry involves various kinds of forces, often expressed in terms of potential energies. Some are the ion-ion force, van der Waals force, and dipole-dipole force, of which the hydrogen bond is the strongest. The binding energy of a covalent or other bond between two atoms of a molecule is the energy required to pull them apart. Atomic cohesive energy is the energy required to separate the atoms of a solid far enough that they no longer feel cohesive forces. *See* CHEMICAL BONDING; COHESION (PHYSICS); DIPOLE-DIPOLE INTERACTION; INTERMOLECULAR FORCES.

Burning is a process of oxidation that converts fuel plus oxygen into combustion products. The total binding energy of fuel plus air is larger than of smoke plus ashes, and the difference is emitted as light and heat. Energy to power biological processes is stored in the binding energies of biomolecules such as adenosine triphosphate (ATP). This energy is also released through oxidation processes such as the Krebs cycle. *See* ADENOSINE TRIPHOSPHATE; BIOLOGICAL OXIDATION; CITRIC ACID CYCLE; COMBUSTION; METABOLISM.

Bending, stretching, and twisting of molecular bonds store and release energy. Analysis of molecules by infrared spectroscopy uses the fact that the various distortions of the bonds have telltale resonance frequencies. *See* INFRARED SPECTROSCOPY; MOLECULAR STRUCTURE AND SPECTRA.

The shapes that molecules assume are those that minimize the total potential energy due to the forces among their constituent atoms. This simple driving principle yields surprisingly complicated configurations in proteins. Deoxyribonucleic acid (DNA) codes for the one-dimensional sequence of amino acids, sometimes thousands of units long. These spontaneously crumple, by shedding potential energy, into precise and unique three-dimensional configurations that determine their highly specific biological functions. Predicting the ultimate shape of a given amino-acid chain is called the protein-folding

problem, which has not yet been solved. *See* PROTEIN; PROTEIN FOLDING.

Nuclear energy results from release of binding energies of atomic nuclei, orders of magnitude larger than the binding energies of molecules in chemical reactions. For example, solar nuclear fusion combines pairs of deuterium nuclei into single nuclei of helium having a smaller binding energy; the difference is released as heat. The binding-energy differences are large enough to be evident as a conversion of mass to energy. Much of the mass of nucleons is due to the binding energies confining the quarks that comprise them. Chemical reactions, too, result in changes of mass, but generally too small in magnitude to be measurable. *See* ELEMENTARY PARTICLE; NUCLEAR BINDING ENERGY; NUCLEAR FUSION; PROTON-PROTON CHAIN; QUARKS.

Energy as a commodity. The production of electricity and its consumption illustrate how energy can go through many forms. The energy in the coal burned in a power plant originated during the Earth's Carboniferous Period, 300 million years ago, when the Sun's energy drove weather systems to carry ocean water as rain to land-based plants. The plants' chlorophyll also used the Sun's energy to convert rainwater and airborne carbon dioxide (CO₂) to produce plant tissue, mainly cellulose and other polysaccharides. Over eons, the energies of geological processes, such as the folding of tectonic plates, concentrated the carbon content through prolonged heat and pressure. This resulted in the fossil fuels, of which coal has the highest carbon content. *See* CARBONIFEROUS; COAL; FOSSIL FUEL; PHOTOSYNTHESIS.

In one type of power plant, coal is burned to heat water in boilers (sending the carbon back into the atmosphere as CO₂). The heat converts water to high-pressure steam that is passed to steam turbines, where the potential energy of steam pressure is converted to kinetic energy of the turbines' rotors. The rotors are coupled mechanically to the electric generators, whose rotors carry conductive metal windings that are rotated within a strong magnetic field. This results in electromotive force (emf) that manifests itself as electrical power flowing out of the generators. The electrical energy is passed to the power grid, which distributes it to the power plant's customers. *See* ELECTRIC POWER GENERATION; ELECTRIC POWER SYSTEMS; ELECTROMOTIVE FORCE (EMF); GENERATOR.

In the home, the electrical energy undergoes various conversions depending on its use. A few examples are heat energy in electric ranges, radiant energy (light) in lamps, kinetic energy in motorized appliances, and chemical energy in the recharging of batteries. After performing its desired functions, most of the energy will be dissipated as heat. *See* BATTERY; FLUORESCENT LAMP; INCANDESCENT LAMP; MOTOR.

Potential misconceptions. Some common ways of discussing matters related to energy could be misleading.

Energy utilities are often described as "producing" energy. In fact, they are only converting one form of energy into another.

We speak of a glass vase as “gaining potential energy” when raised from tile floor to tabletop. However, this energy does not reside in the vase as the words imply; rather, it resides in the new spatial relationship between the vase and the Earth’s gravitational field. That the gained energy is real becomes evident if the vase should fall and shatter on the floor.

Similarly, an object’s kinetic energy is an attribute not of the object alone but of its mass and the reference frame used. Thus, a speeding car has no kinetic energy with respect to the driver, but relative to a pedestrian it has a deadly amount.

Energy and cosmology. A few minutes after the big bang, the universe was uniformly filled with baryonic matter, mostly hydrogen and helium. Unlike a closed system of gas, however, this configuration did not represent maximum entropy and therefore stability—gravitational attraction among the atoms caused them to clump into stars. Those that exploded as supernovae generated the full range of the natural elements that made planets and life possible. The thermonuclear energy radiating from our Sun is therefore a direct result of gravity, and the Sun’s radiated energy powers the entropy-lowering processes of life. *See* BIG BANG THEORY; SUPERNOVA.

Observations of star and galaxy dynamics imply that less than 5% of the mass of the universe is visible baryonic matter. The other 95% is hypothesized to be composed of invisible dark mass-energy, whose properties, other than its gravitational pull, are unknown. Since dark matter has mass, it represents a huge part of the universe’s total energy budget. Suffusing the universe, too, is dark energy.

Dark energy may be at least partly due to the cosmological constant in Einstein’s field equation for gravity that accounts for the expansion of the universe. Such energy could account for the recent evidence that the expansion of the universe is accelerating. Pressure itself, if large enough, can contribute to gravitational force, as occurs when stars collapse into black holes. The dark energy, however, can generate negative pressure, that is, tension, that would account for the acceleration. *See* ACCELERATING UNSTANT; COSMOLOGICAL CONSTANT; COSMOLOGY; DARK ENERGY; DARK MATTER.

Search for the “true” nature of energy. Although energy flows and transformations can be accurately measured and calculated, the true nature of energy—its essence, or what it is in the ontological sense—is not known. To define energy in terms of work, as is often done, does not help since work itself is defined in terms of energy transferred.

Understanding of energy’s true nature will likely come from a successor to the standard model that will successfully reconcile quantum mechanics with general relativity. In doing so, it will likely have to be valid for all stages of the evolution of the universe, including the very moment of the big bang itself at “time zero” where current theory breaks down.

Interestingly, the two main lines of research to solve this conundrum are founded on the conflicting world views of Isaac Newton and Gottfried Wilhelm Leibniz, archrival founders of classical mechanics.

String theory posits that the superstrings, of which subatomic particles are composed, exist in a fixed frame of time and space (Newton’s thesis); in contrast, loop quantum gravity is based on background independence, a concept that does away with absolute space and time (Leibniz’s thesis); it also yields the radical implication that spacetime itself must be quantized. *See* SUPERSTRING THEORY.

Energy, in the form of fire and the natural forces of nature, baffled and amazed prehistoric peoples. It is remarkable that thousands of years later energy is still a key part of the mystery of existence.

Andrej Tenne-Sens

Bibliography. R. P. Feynman, R. B. Leighton, and M. L. Sands, *The Feynman Lectures on Physics: The Definitive and Extended Edition*, vol. 1, Addison-Wesley, 2005; B. Greene, *The Fabric of the Cosmos*, Vintage Books, 2005; J. A. Peacock, *Cosmological Physics*, Cambridge University Press, 1999; W. K. Purves et al., *Life: The Science of Biology*, 7th ed., Sinauer Associates and W. H. Freeman, 2003; L. Smolin, *Three Roads to Quantum Gravity*, Basic Books, 2001; T. W. G. Solomons and C. B. Fryhle, *Organic Chemistry*, 8th ed., Wiley, 2003; E. F. Taylor and J. A. Wheeler, *Spacetime Physics: Introduction to Special Relativity*, 2d ed., W. H. Freeman, 1992; J. R. Taylor, C. D. Zafiratos, and M. A. Dubson, *Modern Physics for Scientists and Engineers*, 2d ed., Prentice Hall, 2004; F. Wilczek, Whence the Force of $F = ma$? I: Culture Shock, *Phys. Today*, 57(10):11–12, October 2004; F. Wilczek, Whence the Force of $F = ma$? II: Rationalizations, *Phys. Today*, 57(12):10–11, December 2004 .

Energy conversion

The process of changing energy from one form to another. There are many conversion processes that appear as routine phenomena in nature, such as the evaporation of water by solar energy or the storage of solar energy in fossil fuels. In the world of technology the term is more generally applied to operations in which the energy is made more usable, for instance, the burning of coal in power plants to convert chemical energy into electricity, the burning of gasoline in automobile engines to convert chemical energy into propulsive energy of a moving vehicle, or the burning of a propellant for ion rockets and plasma jets to provide thrust. *See* ENERGY.

There are well-established principles in science which define the conditions and limits under which energy conversions can be effected, for example, the law of the conservation of energy, the second law of thermodynamics, the Bernoulli principle, and the Gibbs free-energy relation. Recognizable forms of energy which allow varying degrees of conversion include chemical, atomic, electrical, mechanical, light, potential, pressure, kinetic, and heat energy. In some conversion operations the transformation of energy from one form to another, more desirable form may approach 100% efficiency, whereas with others even a “perfect” device or system may have a theoretical

limiting efficiency far below 100%. See BERNOULLI'S THEOREM; CONSERVATION OF ENERGY; THERMODYNAMIC PRINCIPLES.

The conventional electric generator, where solid metallic conductors are rotated in a magnetic field, actually converts 95–99% of the mechanical energy input to the rotor shaft into electric energy at the generator terminals. On the other hand, an automobile engine might operate at its best point with only 20% efficiency, and even if it could be made perfect, might not exceed 60% for the ideal thermal cycle. Wherever there is a cycle which involves heat phases, the limitation of the Carnot criterion precludes 100% conversion efficiency, and for customary temperature conditions the ideal thermal efficiency frequently cannot exceed 50 or 60%. See CARNOT CYCLE.

In the prevalent method of producing electric energy in steam power plants, there are many energy-conversion steps between the raw energy of fuel and the electricity delivered from the plant, for example, chemical energy of fuel to heat energy of combustion; heat energy so released to heat energy of steam; heat energy of steam to kinetic energy of steam jets; jet energy to kinetic energy of rotor; and mechanical energy of rotor to electric energy at generator terminals. This is a typical elaborate and burdensome series of conversion processes. Many efforts have been made over the years to eliminate some or many of these steps for objectives such as improved efficiency, reduced weight, less bulk, lower maintenance, greater reliability, longer life, and lower costs. For a discussion of major technological energy converters see POWER PLANT; ELECTRIC POWER GENERATION.

Efforts to eliminate some of these steps have been stimulated by needs of astronautics and of satellite and missile technology and need for new and superseding devices for conventional stationary and transportation services. Space and missile systems require more compact, efficient, self-contained power systems which can utilize energy sources such as solar and nuclear. With conventional services the emphasis is on reducing weight, space, and atmospheric contamination, on improving efficiency, and on lowering costs. The predominant objective of energy conversion systems is to take raw energy from sources such as fossil fuels, nuclear fuels, solar energy, wind, waves, tides, and terrestrial heat and convert it into electric energy. The scientific categories which are recognized within this specification are electromagnetism, electrochemistry (fuel cells), thermoelectricity, thermionics, magnetohydrodynamics, electrostatics, piezoelectricity, photoelectricity, magnetostriction, ferroelectricity, atmospheric electricity, terrestrial currents, and contact potential. The electromagnetism principle today dominates the field. Electric batteries are an accepted form of electrochemical device of small capacity. See BATTERY; ELECTRIC ROTATING MACHINERY; ENERGY SOURCES; FUEL CELL; MAGNETOHYDRODYNAMIC POWER GENERATOR; SOLAR CELL; THERMIONIC POWER GENERATOR; THERMOELECTRIC POWER GENERATOR.

Theodore Baumeister

Bibliography. E. A. Avallone and T. Baumeister III (eds.), *Marks' Standard Handbook for Mechanical Engineers*, 10th ed., 1996; H. A. Sorenson, *Energy Conversion Systems*, 1983; K. Weston, *Energy Conversion Engineering*, 1992.

Energy level (quantum mechanics)

One of the allowed values of the internal energy of an isolated physical system. This energy is not free to vary continuously above its minimum value, as predicted by classical mechanics, but is constrained to lie among a set or spectrum of particular values. This spectrum may consist of both an isolated discrete portion and a continuous component of restricted range. The term energy level usually refers to one of the allowed values in the discrete set.

Hydrogen spectrum. The primary indication for the existence of discrete energy levels came from the study of the spectrum of emissions of energetically excited atomic systems. Historically, the most important such spectrum is that of the simplest atom, hydrogen, a system of one proton and one electron bound together by their electromagnetic attraction. Within the framework of classical physics, the structure of the hydrogen atom poses fundamental problems. The first is the existence of a stable ground state: An electron in orbit around a proton is in constant acceleration, and therefore, according to Maxwell's classical electromagnetic theory, should continuously radiate away energy. Furthermore, the radiation emitted as the atom decays to a lower energy state should form a continuous spectrum of frequencies. However, the hydrogen atom both possesses a stable ground state and emits radiation at only a discrete set of frequencies. These frequencies, f , can be described, to a good approximation, by the formula (1) found by J. Balmer in 1885, where n_k and

$$f = \frac{R}{b} \left(\frac{1}{n_j^2} - \frac{1}{n_k^2} \right) \quad (1)$$

n_j are both positive integers. Here, R is Rydberg's constant, and b is Planck's constant. See ELECTROMAGNETIC RADIATION.

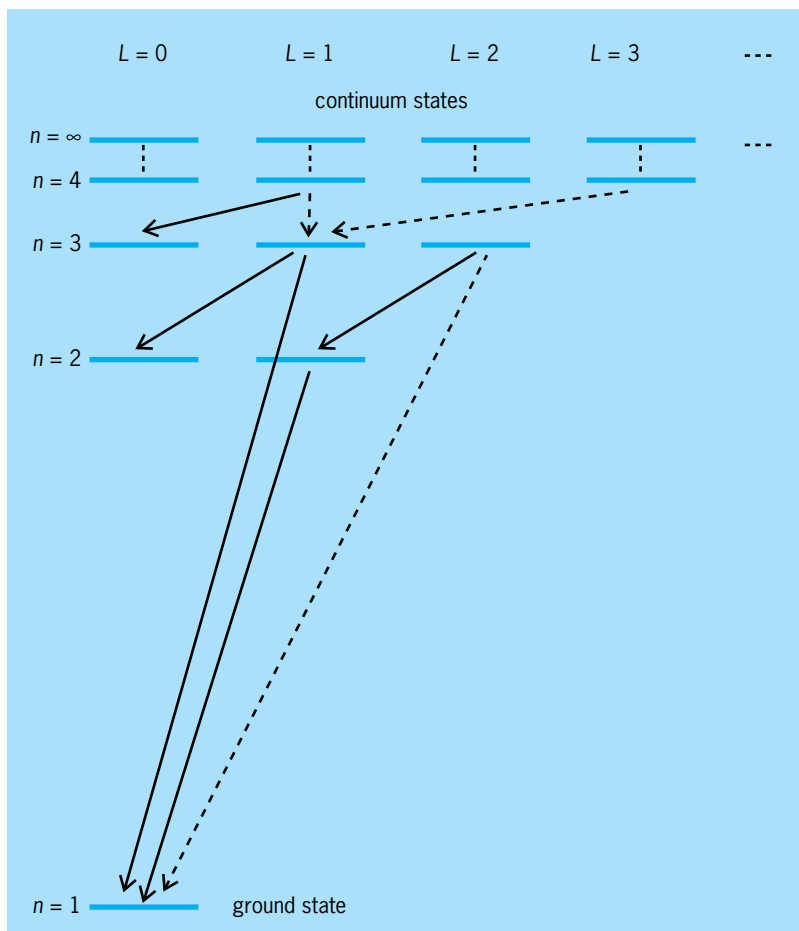
Bohr model. In 1913 N. Bohr made a fundamental advance by postulating that the angular momentum of the electron-proton system could take on only a discrete set of values given by Eq. (2), where \hbar is

$$L = (n - 1)\hbar \quad n = 1, 2, 3, \dots \quad (2)$$

Planck's constant divided by 2π . The angular momentum is said to be quantized. A consequence is that the excitation energies of the hydrogen atom also have a discrete spectrum, given by Eq. (3)

$$E_n = -\frac{R}{n^2} \quad (3)$$

[see **illus.**]. In addition, Rydberg's constant could be calculated in terms of more fundamental



Spectrum of hydrogen energy levels in the leading order of approximation. L = angular momentum. Some typical allowed transitions are indicated by solid arrows. Transitions that are less probable because they violate approximate selection rules are indicated by broken arrows. Above the zero of energy, there exists a continuum of unbound states for all angular momenta.

constants, $R = me^4/(2\hbar^3)$, where e and m are the electron charge and mass respectively. The integer $n = 1, 2, \dots$ now indexes the infinite number of discrete levels. They all have negative energy, relative to an infinitely separated electron-proton pair, corresponding to the bound nature of these states. The states with positive energy form a continuum and represent unbound states that can separate off to infinity. This is a generic feature of energy levels in quantum mechanics. See RYDBERG CONSTANT.

Bohr made the further postulate that the atom decays from an excited level, E_k , only by making a transition to a lower energy level, E_j , emitting a single light quantum (photon) in the process. The energy, E_γ , of this photon is given by the conservation of energy, $E_\gamma = E_k - E_j$. Then, by the Planck-Einstein relation, $E_\gamma = hf$, between the energy of the photon and the frequency, f , of the radiation, together with the expression (3) for the energies of the levels, the Rydberg formula (1) is derived. Although Bohr's postulates are in many ways without real foundation, they were later justified and extended by the development of quantum mechanics. See ATOMIC STRUCTURE AND SPECTRA.

Uncertainty relation. Crudely, the stability of the ground state of atomic systems in quantum mechan-

ics is a consequence of the Heisenberg position-momentum uncertainty relation $\Delta x \Delta p \geq \hbar$, which implies that as the electron more closely approaches the proton (and thus is localized to within a smaller uncertainty Δx) the potential energy is reduced but its momentum, p , and thus its kinetic energy, increases. The competition between these two effects leads to a state of lowest energy. See UNCERTAINTY PRINCIPLE.

Complications. Other quantities label the excited states of hydrogen besides the integer n appearing in Eq. (3). At the simplest level of approximation, the n -th energy level comprises n independent but degenerate states of total angular momentum $L = 0, \hbar, \dots, (n-1)\hbar$, each of which can have a component, m , of angular momentum along a given axis belonging to the set $m = L, L - \hbar, \dots, 0, -\hbar, \dots, -L$. Usually such a degeneracy is due to a symmetry of the system, possibly approximate. In the case of hydrogen, the exact rotational invariance of the system results in the exact degeneracy of the states with different values of m (for fixed n and L), while the special $1/r$ form of the Coulomb potential results in the degeneracy with respect to L . In fact, the potential is not exactly coulombic, and the levels with different L (for a fixed n) are found to be split by a small amount. The inclusion of the intrinsic angular momentum (spin) of the proton and electron further complicates the situation. See ANGULAR MOMENTUM; DEGENERACY (QUANTUM MECHANICS); SPIN (QUANTUM MECHANICS); SYMMETRY LAWS (PHYSICS).

Selection rules and broadening. An important modification of the simple picture presented above occurs when the decay of excited states is described in a more rigorous way considering the interactions of the hydrogen atom with the quantized electromagnetic field. Then, the excited levels are no longer infinitesimally narrow but are broadened by a calculable amount. This can be crudely understood by applying the Heisenberg energy-time uncertainty relation, $\Delta E \Delta t \geq \hbar$, which implies the shorter the lifetime of an excited state, Δt , the greater the intrinsic uncertainty in its energy, ΔE . Further, not all transitions are found to be equally probable. Since the intrinsic angular momentum of the photon is \hbar , conservation of angular momentum usually requires that the L values of the initial and final states differ by \hbar . Less frequently, the photon can also have some orbital angular momentum and mediate a transition with $\Delta L > \hbar$. Such differences in the relative probability of transitions are encoded in so-called selection rules. See LINEWIDTH; QUANTUM ELECTRODYNAMICS; SELECTION RULES (PHYSICS).

Quantization. The quantization of the allowed energy values that occurs in quantum mechanics has analogs for other physical quantities as well, such as angular momentum. The basic reason why such quantization occurs for bound systems of particles in quantum mechanics but not in classical mechanics is that in quantum mechanics particles have associated wavelike attributes, specifically a wave function which encodes the dynamical state of the

particle. (This is the content of wave-particle duality.) The wave function of a bound state satisfies an equation similar in many ways to the equation describing waves on a guitar string or drumhead of finite extent. Such musical instruments produce only certain specific notes, or frequencies, for a given length of string or size of drumhead. In other words, the frequencies are quantized. Similarly, the modes of oscillation of the wave function for a quantum system of finite extent are also quantized, leading to discrete energy levels, and so forth. An unbound quantum system, however, is analogous to a string of infinite length, which can play a continuous range of notes.

Eigenvalues. Formally, in quantum mechanics the energy spectrum of a physical system is given by the values of the energy, E_n , the energy eigenvalues, for which there exist nontrivial solutions of the time-independent Schrödinger equation (4). These

$$\left(-\frac{\hbar^2}{2m} \frac{d^2}{dx^2} + V(x)\right) \psi_n(x) = E_n \psi_n(x) \quad (4)$$

solutions, labeled $\psi_n(x)$, are the energy wave functions or eigenfunctions. For simplicity, the nonrelativistic Schrödinger equation has been chosen in one spatial dimension, x , for a single particle of mass m in a background potential $V(x)$. The energy eigenfunction, $\psi_n(x)$, encodes all possible information about the n -th quantum-mechanical state. In the general case, Eq. (4) becomes Eq. (5), where \hat{H} is the hamiltonian operator appropriate for the system under study. See EIGENFUNCTION; EIGENVALUE (QUANTUM MECHANICS).

Complex systems. Energy levels are of great importance for many systems other than simple atoms such as hydrogen. For instance, they determine the interactions and binding of molecules in chemistry and biochemistry, the stability or decay of nuclei, and the macroscopic properties of solids, such as the optical properties of dyes or semiconductors. The observed spectroscopy of the energy levels of a system can also elucidate the properties of a new force, just as the study of hydrogen led to the development of quantum mechanics and quantum field theory. An outstanding example is the discovery, in high-energy experiments, of the J/ψ and Υ families of particles, comprising bound states of charmed-anticharmed quarks and bottom-antibottom quarks respectively. However, the binding force is no longer electromagnetism as in the case of hydrogen, but the strong force that binds nuclei, namely quantum chromodynamics (QCD). The study of the excited states of these systems has yielded greater understanding of the properties of the extremely complicated QCD force. See ELEMENTARY PARTICLE; J/PSI PARTICLE; MESON; NONRELATIVISTIC QUANTUM THEORY; QUANTUM CHROMODYNAMICS; QUANTUM MECHANICS; QUARKS; UPSILON PARTICLES. John March-Russell

Bibliography. A. Bohm, *Quantum Mechanics: Foundations and Applications*, 3d ed., 1994;

A. Pais, *Inward Bound*, 1986; R. Shankar, *Principles of Quantum Mechanics*, 2d ed., 1994.

Energy metabolism

Energy metabolism, or bioenergetics, is the study of energy changes that accompany biochemical reactions. Animals, plants, and bacteria require energy to sustain life. Energy sustains the work of biosynthesis of cellular and extracellular components, the transport of ions and organic chemicals against concentration gradients (osmotic work), the conduction of electrical impulses in the nervous system, and the movement of cells and the whole organism. Sunlight is the ultimate source of energy for life. Photosynthetic cells use light energy to produce chemical energy and reducing compounds, used to convert carbon dioxide into organic chemicals such as glucose. The energy from the oxidation of carbohydrates, fats, and proteins sustains the biochemical reactions required for life.

Energy content of foodstuffs. The main sources of chemical energy for most organisms are carbohydrates, fats, and protein. Energy content is expressed in calories or joules. A calorie (cal) is equivalent to the amount of energy required to increase the temperature of 1 g (0.035 oz) of water by 1°C (1.8°F). One calorie is equivalent to 4.184 joules (J). The nutritional calorie, or kilocalorie (kcal), in foodstuffs is equivalent to 1000 calories. The energy content per gram of carbohydrate is 4 kcal (16 J); protein, 4 kcal (16 J); and fat, 9 kcal (36 J). The metabolism of foodstuffs yields chemical energy and heat.

Free energy changes. Energy is defined as the ability to do work, and metabolism represents the biochemical reactions that a cell can perform to produce energy. A chemical reaction is the process by which one or more substances are converted into other substances. For example, citrate (one biochemical compound) is converted into isocitrate (another compound) in the cell. The most important thermodynamic parameter in bioenergetics is the free energy change, ΔG , occurring at constant temperature and pressure (the usual conditions for chemical reactions inside the cell). The Gibbs free energy change is defined as the free energy content of the final state minus the free energy content of the initial state, as in Eq. (1). All feasible reactions occur with a nega-

$$\Delta G = G_{\text{final}} - G_{\text{initial}} \quad (1)$$

tive free energy change; the final state has less free energy than the initial state; that is, $\Delta G < 0$ (process is exergonic). If the free energy of the final state is more than that of the initial state, ΔG is positive and the reaction is not feasible without the input of energy; $\Delta G > 0$ (process is endergonic). When the free energy change is zero, the reaction or process is at equilibrium; $\Delta G = 0$ (process is isoergonic). See FREE ENERGY; THERMODYNAMIC PRINCIPLES.

Standard free energy change. For general chemical reaction (2), the free energy change is expressed by



Eq. (3), where ΔG° = standard free energy change;

$$\Delta G = \Delta G^\circ + RT \ln [C][D]/[A][B] \quad (3)$$

R = gas constant, or temperature-energy coefficient (8.314 J K⁻¹ mole⁻¹, or 1.98 cal K⁻¹ mole⁻¹); K = absolute temperature in Kelvin (273 + degrees Celsius); \ln = natural logarithm; $[C][D]/[A][B]$ = the product of the molar concentrations, or activities, of C and D divided by the product of the molar concentrations, or activities, of A and B.

The standard free energy change for reaction (2) is the free energy change that accompanies the conversion of one mole each of reactants A and B into products C and D under standard conditions. Most biochemical reactions of interest occur in water, and the concentrations of reactants and products under standard conditions are taken to be 1 M (1 mole/liter).

The standard free energy change is a logarithmic function of the equilibrium constant,

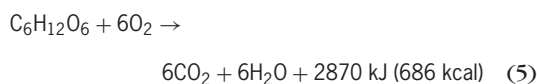
$$K_{\text{eq}} = [C][D]/[A][B]$$

where the concentrations of A, B, C, and D are those at equilibrium. The relationship of the standard free energy change and the equilibrium constant is given by Eq. (4). When the equilibrium constant is

$$\Delta G^\circ = -RT \ln K_{\text{eq}} \quad (4)$$

greater than 1, the standard free energy change is negative. Under standard conditions (1 M in reactants and products), the reaction will proceed toward products (C + D). The free energy of the products is less than that of the reactants; that is, ΔG is negative. The actual free energy change (ΔG) can be larger or smaller than the standard free energy change (ΔG°). Increasing the concentrations of reactants or decreasing the concentration of products produces a larger decrease in free energy (the reaction is more exergonic). Decreasing the concentration of reactants or increasing the concentration of products produces a smaller decrease in free energy (the reaction is less exergonic).

Adenosine triphosphate. The complete oxidation of one mole of glucose (C₆H₁₂O₆) to carbon dioxide (CO₂) and water (H₂O) is associated with the liberation of free energy, as in reaction (5). Energy is

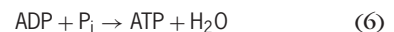


released in a stepwise fashion and is coupled to the biosynthesis of adenosine triphosphate (ATP) from adenosine diphosphate (ADP) and inorganic phosphate (P_i). The reaction of ATP with water to produce ADP and P_i results in the liberation of a large amount of energy (30 kJ, or 7 kcal per mole). Such compounds are said to be energy-rich and to possess a high-energy bond. Lipmann's law is the cornerstone of energy metabolism: ATP serves as the common currency of energy exchange in living systems (animals, plants, and bacteria). The ATP-ADP couple receives and distributes chemical energy in all living

systems. Creatine phosphate is an energy-rich compound found in vertebrate muscle and brain; it is a storage form of chemical energy and can energize the regeneration of ATP from ADP. Such a reaction occurs in vigorously exercising skeletal muscle when ATP is expended to produce contraction. See ADENOSINE TRIPHOSPHATE (ATP).

ATP synthesis. There are three general mechanisms for ATP generation in cells: substrate-level phosphorylation, oxidative phosphorylation, and photo-phosphorylation. Substrate-level phosphorylation is mediated by non-membrane-bound enzymes, and oxidative and photophosphorylation are mediated by an intricate membrane-associated ATP synthase. The energy for substrate phosphorylation is derived from exergonic chemical reactions. In the Embden-Myerhof glycolytic pathway for the conversion of glucose into two molecules of pyruvate, glyceraldehyde 3-phosphate is oxidized by nicotinamide adenine dinucleotide (NAD⁺) to form 1,3-bisphosphoglycerate, an energy-rich compound. 1,3-Bisphosphoglycerate reacts with ADP to form ATP. This two-step process is one prominent example of substrate-level phosphorylation.

The energy for oxidative phosphorylation [oxidation of NADH (reduced NAD) or of succinate coupled to the phosphorylation of ADP to produce ATP] is derived from a proton gradient. Exergonic electron transport reactions generate a proton gradient (outside positive) across the inner mitochondrial membrane or the bacterial plasma membrane. Next, protons move exergonically down their electrochemical gradient from the outer to the inner aspect of the membrane. This process provides the energy to drive the endergonic portion of the phosphorylation reaction (6).



$$\Delta G^\circ = +30 \text{ kJ/mole or } +7 \text{ kcal/mole}$$

Experiments with animal mitochondria show that the P:O ratio (number of ATPs formed from ADP + P_i per atom of oxygen consumed) using NADH as substrate is 2.5, and the P:O ratio with succinate as substrate is 1.5. There are three sites along the electron transport chain that are involved in proton translocation from the mitochondrial matrix to the outside of the inner mitochondrial membrane. Sites 1 and 3 are associated with the generation of a proton gradient sufficient for the formation of 1 ATP, and site 2 with 0.5 ATP. Transport of electrons from succinate dehydrogenase circumvents site 1 and results in the production of only 1.5 ATP molecules.

To recapitulate, there are four properties of oxidative phosphorylation. (1) Transport of electrons from reductant to oxygen along the respiratory chain is a very exergonic process. (2) Part of the chemical energy is conserved by the translocation of protons from the inside to the outside of the inner mitochondrial membrane or bacterial plasma membrane and the ensuing proton gradient. (3) The membrane is not freely permeable to protons. (4) Protons then move down their exergonic

electrochemical gradient and drive endergonic ATP formation in a process involving the ATP synthase of the inner mitochondrial membrane or plasma membrane (bacteria). Four protons need to be transported for the synthesis of each ATP in mitochondria: one for the translocation of phosphate and three for ATP synthase. See BACTERIAL PHYSIOLOGY AND METABOLISM.

In plant photophosphorylation, light energy fuels oxidation-reduction reactions that promote proton translocation in the chloroplast. Protons move down their electrochemical gradient to drive ATP biosynthesis by a chloroplast ATP synthase that is homologous to that in mitochondria and in bacteria. The generation of a proton gradient, or a proton-motive force, and its use to drive the generation of ATP biosynthesis is a cornerstone of bioenergetics, as first described by Peter Mitchell. The concept that osmotic energy derived from the proton gradient is used to drive the chemical reaction of ATP formation is called the chemiosmotic theory, or the Mitchell theory—a cornerstone of energy metabolism. See CYTOCHROME; ELECTRON-TRANSFER REACTION; MITOCHONDRIA; PHOTOSYNTHESIS.

Robert Roskoski, Jr.

Bibliography. D. G. Nicholls and S. J. Ferguson, *Bioenergetics 2*, 1992; R. Roskoski, Jr., *Biochemistry*, 1996; L. Stryer, *Biochemistry*, 4th ed., 1995.

Energy sources

Sources from which energy can be obtained to provide heat, light, and power. The term energy is used to describe an amount of work performed. There are two kinds of energy, kinetic energy, meaning work performed by the movement of matter, and potential energy, meaning work that is stored or at rest in matter. See ENERGY.

Energy Concepts and Needs

In the kinetic or potential state, energy takes on one of five forms: (1) Chemical energy results from changes in the chemical structure of substances, such as during combustion of fuel. (2) Electrical energy results from electrons and protons in motion in a stream called an electric current, or in temporary storage as in a battery or fuel cell. (3) Mechanical energy results from force applied or about to be applied to liquid, solid, or gaseous matter. (4) Thermal energy results from heat being applied to matter. (5) Nuclear fission is the splitting of the nucleus of an atom into two or more parts by collision with neutrons, with the consequent release of the force that binds protons and neutrons of the nucleus together. All living things on Earth depend on one or more of these forms of energy and must look to a wide variety of energy sources. See BATTERY; CHEMICAL ENERGY; FUEL CELL.

The Sun. The ultimate source of energy on Earth is the Sun, which produces energy sources in two ways on a daily basis and a stored-over-time basis (Figs. 1 and 2).

The Sun provides radiant heat on a daily basis to the Earth which drives many reactions. For example, solar heat evaporates water from the sea and the lakes, providing the moisture for cloud formations which break into rain in mountainous regions. The rain runs off into constructed reservoir lakes which feed water through hydroelectric dams, thus rotating electric generators that produce the kinetic electrical energy supplied through power lines to industrial plants, commercial buildings, and residential homes. Tidal power, wave power, solar power, and wind power are additional examples of daily energy sources.

On a stored-over-time basis, radiant heat from the Sun striking the Earth over millions of years provided the necessary energy input to convert vegetable matter into coal, petroleum, and natural gas. This phenomenon is still taking place, but the conversion process is so slow that it is meaningless in terms of its replacement capability, especially when compared to the enormous rate at which modern industrial society is consuming these resources.

Capital energy. Scientists refer to the stored-over-time energy sources collectively as capital energy, and this is subdivided into six categories: primary

Nonrenewable stored-over-time (capital)	Fossil:	coal peat crude petroleum natural gas	Combustion process						
	Nuclear:	uranium thorium deuterium lithium beryllium							
Renewable daily	Solar:	<table border="0"> <tr> <td>solar thermal conversion</td> <td rowspan="2">} direct</td> </tr> <tr> <td>photoelectric energy conversion</td> </tr> <tr> <td>photochemical conversion</td> <td rowspan="2">} indirect</td> </tr> <tr> <td>stored solar heat with heat pumps</td> </tr> </table>	solar thermal conversion	} direct	photoelectric energy conversion	photochemical conversion	} indirect	stored solar heat with heat pumps	Noncombustion process
	solar thermal conversion	} direct							
	photoelectric energy conversion								
	photochemical conversion	} indirect							
	stored solar heat with heat pumps								
	Hydro:	river-reservoir energy conversion							
	Tidal:	tidal energy conversion							
Wind:	windmill energy conversion								
Oceans:	ocean heat conversion ocean current conversion wave energy conversion								
Geothermal:	natural steam hot water hot dry rocks								
Biomass:	wood and other vegetation	Combustion process							

Fig. 1. Primary energy sources.

Nonrenewable	Electric:	electric power generation fuel cells	Noncombustion process
	Nuclear:	tritium plutonium	
	Fossil: (coal-derived)	coke char tar blast furnace gas water gas and carbureted water gas producer gas town gas briquetting coal slurries coal gasification coal methanol	Combustion process
	Fossil: (petroleum-derived)	gasoline kerosine petroleum coke oil shale petroleum fuel oils (No. 1, 2, 4, 5, and 6) liquefied natural gas (LNG) liquified petroleum gas (LPG) propane butane recycled lubricants and solvents	
Renewable	Biomass:	wood waste and bark bagasse hulls (grain, rice, cottonseed) peanut shells coffee grounds sugarbeets tobacco stems citrus rinds corncobs garbage and trash methane gas alcohol (ethanol, methanol)	

Fig. 2. Secondary energy sources.

energy, secondary energy, renewable energy, nonrenewable energy, combustion process, and noncombustion process.

Primary. This classification includes all forms of potential energy created mainly by the Sun in the Earth's crust that need no processing or treatment to transform them into usable energy (Fig. 1).

Secondary. This classification includes the forms of potential energy manufactured from primary energy forms generally by mechanical, chemical, thermal, or nuclear reaction means to transform them into usable energy.

Renewable. This term refers to forms of potential energy that constantly and rapidly renew themselves for steady, reliable use. This is a somewhat ambiguous definition, particularly when the word "rapidly" is added. Coal, crude petroleum, and natural gas are listed as nonrenewable in Fig. 1, but this is not totally accurate. These forms of potential energy are, in fact, constantly being created. However, modern industrialized societies are utilizing these energy sources so rapidly compared to the geological time period required for the formation of additional sources as to render the concept of renewable useless within any reasonable time frame. On the other hand, such

energy forms as solar and wind power are clearly renewed on a timely basis and are labeled as renewable energy.

Nonrenewable. Any form of potential energy that does not fall within the definition accepted for renewable energy is considered nonrenewable. For example, the fossil fuels may be defined in specific considerations as nonrenewable energy sources.

Combustion process. Many of the forms of potential energy shown in Figs. 1 and 2 must be utilized in a combustion process before they will release their stored energy into a work process. This is part of the process of conversion that will be discussed below.

Noncombustion process. There are also ways to release energy without a combustion process, such as river water turning a waterwheel. This was an important source of energy for nineteenth-century factories.

Energy needs. Human use of energy, particularly capital energy, has accelerated over time due to an increase in the human population and the discovery of new technologies for utilizing energy. This acceleration took on exponential proportions during the Industrial Revolution of the nineteenth and twentieth centuries (see **table**). As the utilization of, and therefore the demand for, energy rose, scientists and engineers discovered ways to utilize new forms of energy. The availability and technical feasibility of many energy forms yielded a proliferation of energy choices, so that decisions had to be reached on the selection of energy sources. The need to make energy decisions has resulted in the development of an energy choice system (**Fig. 3**). This system involves an economic analysis which compares the cost of input versus the benefit of output for two specific energy sources. Among the factors included in this analysis are depletion cost, the cost of using nonrenewable energy; social cost, the health and environmental issues; availability cost, the cost of utilizing energy sources subject to interruptions of supply; and switching or conversion cost, converting plant and equipment to other energy uses.

Technical Aspects

Certain technical aspects must be considered in order that energy users can make intelligent decisions regarding energy sources.

Fossil fuels. Crude petroleum, natural gas, and coal were formed in the Earth's crust over the course of millions of years and exist today in subsurface locations. For instance, coal usually exists in seams ranging from 3 to 6 ft (0.9 to 1.8 m), although one seam in Wyoming averages 100 ft (30 m) thick with a maximum thickness of 220 ft (67 m). Crude petroleum is found trapped in the pores of rock (sandstone, limestone, or dolomite) or sand overlain with some kind of impervious cap rock that prevents the liquid from dispersing. Natural gas is normally found trapped in the Earth alongside or with crude petroleum.

The specific ways in which fossil fuels were formed is not known. The theory most dominant since the 1920s, and the one accepted by the United States oil, gas, and coal industry, is that they were

Energy consumption in the United States*		
Item	1982	2000
<i>Key energy use factors</i>		
U.S. population, $\times 10^6$	232	268
Dwellings, $\times 10^6$	83	107
Passenger cars, $\times 10^6$	128	148
Other cars, $\times 10^6$	11	14
Buses and trucks, $\times 10^6$	20	25
Gross national product (1972\$), $\times 10^9$	1476	2426
<i>Fuel consumption</i>		
Petroleum (total), 10^3 bbl (m^3)/day	15,254 (2425)	12,940 (2057)
Gasoline, 10^3 bbl (m^3)/day	6538 (1039)	4890 (778)
Jet fuel, 10^3 bbl (m^3)/day	1009 (160)	1400 (220)
Aviation gasoline, 10^3 bbl (m^3)/day	26 (4.1)	10 (2)
Diesel, 10^3 bbl (m^3)/day	1298 (206)	1800 (290)
Other distillates, 10^3 bbl (m^3)/day	1506 (239)	500 (80)
Petrochemical chemical feed, 10^3 bbl (m^3)/day	507 (80.6)	940 (150)
LNG/LPG, 10^3 bbl (m^3)/day	1537 (244)	1000 (200)
Residual fuels, 10^3 bbl (m^3)/day	1694 (269)	1640 (260)
Still gas, 10^3 bbl (m^3)/day	554 (88.1)	360 (57)
Asphalts/road oils, 10^3 bbl (m^3)/day	343 (54.5)	200 (30)
Lubricating waxes, 10^3 bbl (m^3)/day	154 (24.5)	200 (30)
Petroleum coke, 10^3 bbl (m^3)/day	247 (39.3)	100 (20)
Miscellaneous, 10^3 bbl (m^3)/day	91 (14)	—
Fuel alcohol, 10^3 bbl (m^3)/day	3 (0.5)	300 (50)
Natural gas, 10^3 ft ³ (m^3)/day	17.9 (0.5)	17.9 (0.5)
Coal (not including exports), 10^6 tons (metric tons)	678 (615)	1160 (1053)
Solar, $\times 10^6$ units	0.1	3.0
Hydroelectricity, 10^9 kWh (megajoules)	310 (1122)	340 (1224)
Nuclear, 10^9 kWh (megajoules)	283 (1019)	589 (2120)
Electricity from wind-waste, 10^9 kWh (megajoules)	2 (7)	18 (65)
Geothermal, 10^9 kWh (megajoules)	4 (14)	52 (187)
Electricity from geothermal, quad (joules)	—	1 (1.055×10^{18})
TOTALS, quad (joules)	70.9 (74.8×10^{18})	84.7 (89.4×10^{18})

*Based on U.S. Department of Commerce statistics, July 1983.

formed as a result of the fossilization and carbonization of trees, ferns, and other vegetable matter under intense pressure and temperature over exceedingly long periods of time in the Earth's crust.

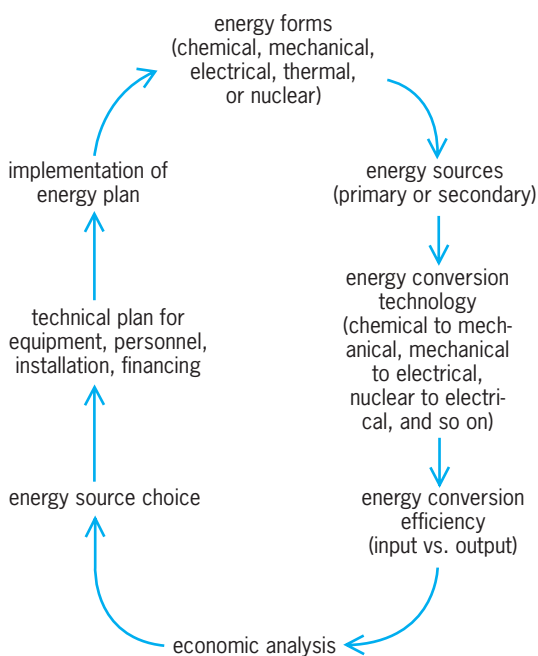


Fig. 3. Diagram of an energy choice system.

Another theory about the formation of fossil fuels is based on the deep-earth gas hypothesis formulated by T. Gold of Cornell University. According to Gold, incredibly huge quantities of methane gas (CH_4), the principal constituent of natural gas, may be trapped deep inside Earth. Tapping into this reserve would require wells drilled several miles below Earth's surface. According to this model, the known deposits of crude oil, natural gas, and coal were formed as a result of the mixing of the deep methane gas with deposits of vegetable and animal remains as the methane rose up through the earth.

Coal. This complex mineral substance is located principally in 29 of the 50 states of the United States and is actively and economically mined in 12 states in meaningful commercial quantities.

The principal chemical constituents of coal are water, carbon, hydrogen, sulfur, nitrogen, oxygen, and ash (noncombustible mineral residue). However, coal is not a uniform substance. It is almost infinitely variable in its composition from one location to the next, even within the same mine location. Nevertheless, 31 tests for defining coal have been established by the American Society for Testing and Materials. Different coals can be ranked according to the degree to which they have progressed from lignite through the bituminous stage to anthracite. This progression is roughly equivalent to the geologic time of development, although the other variables of depth-pressure,

heat, and vegetable matter constituents play an important part.

Coal is used in the United States for generating electric power, metallurgical production, general industrial processes, residential-commercial uses, and synthetic fuels. For the most part, coal is burned in fire-tube or water-tube boilers for raising steam that is used to generate electricity, provide heat for factories and buildings, and provide steam for production processes. *See* COAL; LIGNITE.

Peat. In the very early stages of coal formation, accumulations of decomposed and partially decomposed vegetation, trees, ferns, and mosses located in a wet, cold, and anaerobic (oxygen-deficient) environment are very likely to turn to peat at the rate of about 3 in. (7.5 cm) per 100 years. Only after hundreds or even millions of years would this material graduate to the status of coal.

Since the vegetable matter in peat, consisting of cellulose and other organic compounds, becomes only partially converted to carbon and hydrocarbons, peat has only one-third to one-half of the heating value of coals. It is used in very limited quantities for fuel after it is cut from the earth and formed into briquettes. Peat is harvested and sold mainly as a soil conditioner. *See* PEAT.

Crude petroleum. Oil is also a complex substance derived from the carbonized remains of trees, ferns, mosses, and other types of vegetable matter. Like coal there is doubt about the exact nature of its origin. The principal chemical constituents of oil are carbon, hydrogen, and sulfur. Crude petroleum extracted from the earth will burn and produce thermal energy, but virtually all crude oil is processed in refineries, where it is converted into several useful fuels and special products (for example, feedstock for chemicals, plastics, food products, medicines, and tires, plus tar, asphalts, and lubricating oils). The various fuels made from crude oil are jet fuel, gasoline, kerosine, diesel fuel (or No. 2 fuel oil), and heavy fuel oils (or No. 4, 5, and 6 fuel oils).

Major oil consumption in the United States is in the following areas: transportation, residential-commercial, industry boilers and other industrial uses, and generating electric power. *See* PETROLEUM; PETROLEUM PRODUCTS.

Natural gas. This energy source is 83–93% methane (CH₄), so that its principal chemical constituents are carbon and hydrogen. Natural gas is usually found in the immediate vicinity of crude petroleum, although some natural gas wells do not yield oil.

Of all the chemical or mineral sources of energy, natural gas may well be the most desirable because it can be pipelined directly to the customer, requires no storage vessels, is clean-burning, requires no air-pollution control equipment, produces no ash for disposal, and mixes with air easily to provide complete combustion at low excess air.

The principal uses for natural gas are residential, commercial buildings, industrial, transportation, and generating electric power. *See* NATURAL GAS.

Nuclear energy. The two ways to utilize radioactive fuel as an energy source are fission and fusion.

In fission, heavy atoms are split into two principal elements that form the nucleus of two new, smaller atoms. In fusion, the nuclei of two small atoms fuse together into a single, larger nucleus. In both cases large quantities of energy are released.

Nuclear fission. The splitting of atoms cannot be effectively accomplished with electrically charged matter such as alpha particles, beta particles, or protons, which tend to be diverted or slowed down as they approach other charged matter. Neutrons, however, are not deflected by positive or negative electrical charges, and this fact makes them ideal candidates for atom smashing. The fission process, therefore, involves the bombardment of atoms with neutrons so that a sufficient number of collisions will take place on a statistically predictable basis and split those atoms into two or more separate nuclei while releasing massive quantities of thermal energy. *See* NUCLEAR FISSION.

Nuclear fusion. The fusion process is the opposite of fission. Instead of splitting atoms into two or more pieces, the fusion process causes two atoms to collide with such force that their natural electronic repulsion is overcome and their nuclei fuse into one.

The most suitable fuels for fusion are hydrogen isotopes (hydrogen atoms with the same number of protons but different numbers of neutrons). In particular, deuterium and tritium are thought to be ideal for fusion reactions. The former is relatively abundant and can be found in ordinary seawater. The latter is scarce. *See* DEUTERIUM; TRITIUM.

Although this energy source offers tremendous potential for the twenty-first century, there are formidable problems that must be overcome before it can be utilized. The temperature in a fusion reactor should be between 1.0 and 2.5×10^8 K, and no solid material has been developed for a container for such a reaction. A solution to this problem may be found by setting up an extremely powerful magnetic field around the reacting materials. However, conventional magnets may consume more energy than the fusion reaction puts out. This could be solved by using superconducting magnets that operate at very low temperatures to decrease their resistance to electrical flow, except that it will be difficult to operate such magnets adjacent to the heat of a fusion reactor.

At any rate, fusion power is still an important potential source of continuous low-cost energy, provided that the technical problems can be solved. *See* NUCLEAR FUSION; SUPERCONDUCTING DEVICES.

Solar energy. By far the most attractive energy source is the Sun itself because it is free, is clean and nonpolluting, and does not involve the use of dwindling, finite reserves of capital energy. The problem with solar energy is its cost, particularly for large industrial uses. Roughly 50% of the sunlight approaching Earth is reflected or absorbed by the atmosphere and the other 50% strikes Earth's surface. In the tropical and temperate zones, the Sun provides the equivalent of, on average, about 622 MW of energy per square mile (240 MW/km²). Unfortunately, technology is able to convert only about 15% of the Sun's

radiant heat to usable work, so that a typical large electric generating station capable of producing 1000 MW would require an array of solar head collector-converters covering an area of at least 11 mi² (28 km²). See SOLAR ENERGY.

The collector-converters are made of silicon, which produces electricity when it is exposed to light. The high cost of these devices derives from the cost of growing and cutting large crystals into thin slices of silicon. The more this process becomes automated, the lower the price and the more the process is likely to be employed. See PHOTOVOLTAIC CELL; SOLAR CELL.

Another idea is the solar power tower. This type of technology was demonstrated by an Italian scientist, G. Francia, in 1976. A 10-MW electric generating station is in operation in Daggett, California; it utilizes a 72-acre (29-hectare) field of mirrors to concentrate sunlight at the top of a central tower. The thermal energy from the sunlight superheats steam that is used to power electricity-generating turbines.

Finally, the concept of solar ponds has received attention in many parts of the world. These are artificial, salt-gradient ponds. They are derived from the work of a Russian scientist, A. Von Kaleczinsky, who in 1902 found that the temperature of the water in Lake Medve in Transylvania a few feet below the surface was 185°F (85°C) due to the variation in salinity of the water. Experiments have been conducted with ponds in which layers of different salinity allow the water to trap and hold solar heat. The hot brine can be used for space heating or for producing electricity with a thermoelectric device or organic Rankine cycle engine.

Hydro energy. One of the oldest energy-producing mechanisms uses water flowing in a river or falling from a height to rotate work devices, ranging from the waterwheels of the past to the massive modern hydroelectric dams that employ gigantic electricity-generating turbines. Internationally only 7% of the total potential hydroelectric power estimated at 2.9×10^6 MW is being used. Much of that potential, unfortunately, is located in remote places from which lengthy transmission lines would be prohibitively expensive. Furthermore, environmental protectionists are resisting new projects that might threaten natural resources. See ELECTRIC POWER GENERATION; WATERPOWER.

Tidal energy. The only major tidal energy project in operation is the Rance River project in Brittany, France; however, estimates indicate a worldwide potential for producing 3×10^6 MW of electric power from tidal movement. The French project employs a barrage type of dam across the estuary of a river. Turbines located in this barrage pump water into the estuary when the tide is rising. When a sufficient head of water is built up in the estuary, the water is permitted to flow back through the turbines to produce electricity. See TIDAL POWER.

Wind energy. This is also a form of energy that can be used to generate electrical or mechanical energy. Rotating devices known as windmills can convert the mechanical energy of wind to useful work.

The largest known wind generator has a 300-ft (90-m) span to produce 2.5 MW of electricity for a local electric power grid in Goldendale, Washington. This installation employs huge aircraft propeller-type blades which must always be facing into the wind; a more efficient design might be the Darrieus vertical-axis windmill that rotates regardless of the direction of the wind. Such a device was erected at Sandia Laboratories, Albuquerque, New Mexico, with funds from the U.S. Energy Research and Development Administration. It measures 61 ft (19 m) in height and produces 60 kW of electric power in a 28 mi/h (45 km/h) wind.

In other parts of the United States (such as California and Hawaii), the concept of wind farms is being developed that involves the erection and operation of perhaps dozens of windmills for the production of significant amounts of electric power. See WIND POWER.

Ocean energy. Ocean power may become an energy source in the future, either as wave power or as ocean temperature differential. Both forms are highly experimental and may require many years to become significant energy sources.

Wave power. Many different devices have been proposed and tested for exploiting the wave motion of the sea. In Scandinavia, one device is designed so that ocean waves cause massive quantities of water to flow into the confines of the device which then generates electricity as the water tries to escape back to the sea. In England, S. Salter demonstrated his so-called nodding duck device, which is a floating device that rotates due to the motion of waves rolling over it. The float is capable of driving a hydraulic pump that, in turn, can drive an electrical generating device.

Ocean temperature. An energy-producing system known as ocean thermal energy conversion (OTEC) is based on the temperature differential in the oceans near the Equator where the surface water is about 40°F (20°C) warmer than the water a few thousand feet down. This temperature difference can be utilized to vaporize a working fluid (such as ammonia) that can be run through an electricity-producing turbine.

Geothermal energy. As with ocean power, the geothermal energy source may be exploited in either or both of two ways: hot rocks and hydrothermal.

Hot rocks. It is known that hot granite rock (up to 400°F or 200°C) exists almost everywhere on Earth. The heat is generated, for the most part, by a slow radioactive decay process deep within the earth. According to one estimate, a 40-mi³ (170-km³) chunk of granite at 350°F (177°C) would yield the equivalent energy output of 1.2×10^{10} barrels (2×10^9 m³) of oil, which is approximately the total yearly energy consumption in the United States. According to another estimate, the continental United States is underlain by hot rock with thermal energy amounting to 1.3×10^7 quad (1 quad = 10^{15} Btu) at a depth of just under 6 mi (10 km).

The technique (still experimental) for tapping this energy source involves conventional oil or gas

drilling expertise. The first hole drilled is used for injecting water pressurized to 5000 lb/in.² (35 megapascals) that will cause the hot granite to fracture vertically. Once the rock is fractured, the amount of pressure is decreased to a normal pumping pressure. The next hole drilled and any subsequent holes must follow the fracture line so that the water heated by the hot granite can be brought back up to the surface and utilized for raising steam and generating electricity for space heating or for industrial process steam.

Hydrothermal. The earth emits steam in many locations, such as Iceland, the United States, the Philippines, and New Zealand, where it is captured and employed for space heating, generating electricity, or industrial process steam applications. Hydrothermal steam now provides the equivalent of about 1200 MW of electric power. It is believed that only about 1% of the total potential hydrothermal energy can be utilized in the world and converted to electricity at a 25% efficiency factor for a total contribution of only 3×10^6 MW. See GEOTHERMAL POWER.

Biomass energy. As applied to the field of energy, the term biomass energy encompasses a broad selection of energy sources: Any and all types of living matter that can be converted to a form of energy can be said to be biomass. Hence, scientists tend to think of such items as wood, wood waste, coffee grounds, corn husks, peanut shells, rice hulls, garbage, animal and human waste, sugarcane waste (bagasse), and organic effluent from streams and ponds as biomass.

The biomass considered to be most significant in terms of energy sources is wood and wood waste. There are approximately 9.6×10^9 acres (3.9×10^9 hectares) of forest land in the world, of which about 4×10^9 acres (1.6×10^9 hectares) are economically accessible. Up until the end of the nineteenth century, when coal took over as the leading energy source, wood was a preeminent provider of energy to the world. Wood is again being used on a limited scale at the residential level. The lumber, furniture, plywood, and pulp-and-paper industries all utilize wood waste items (for example, bark, shavings, sawdust, slabs, and end pieces) for raising steam that, in turn, is employed for space heating and industrial processes. See BIOMASS.

Conservation

An energy source that is being increasingly considered to be significant is energy conservation. Better insulation of buildings and homes could slash heating and air-conditioning requirements in half. More people could travel by public transportation than by using private vehicles. There could be greater use of smaller, fuel-efficient automobiles, bicycles, or low-gas-consumption "mopeds," motorcycles, and such. Greater use could be made of fluorescent lighting in homes and industry. Government agencies, businesses, and individuals could lower their thermostats in winter and raise them in summer. In fact, this movement was already quietly under way in the late 1970s until its effects were dramatically felt in the early 1980s in terms of mass decreases in total energy consumption. Conservation, coupled with the

worldwide economic recession of 1982, resulted in the so-called oil glut in the mid 1980s. Therefore, conservation may acquire a major role in establishing energy source security. See CONSERVATION OF RESOURCES.

For further information on primary energy sources see BERYLLIUM; LITHIUM; NUCLEAR FUELS; OIL AND GAS, OFFSHORE; PETROLEUM RESERVES; THORIUM; URANIUM.

For further information on secondary energy sources see ALCOHOL FUEL; COAL GASIFICATION; COAL LIQUEFACTION; COKE; GASOLINE; KEROSENE; LIQUEFIED NATURAL GAS (LNG); LIQUEFIED PETROLEUM GAS (LPG); OIL SAND; OIL SHALE; PLUTONIUM.

William K. Fox

Bibliography. B. Aldridge, L. Crow, and R. Aiuto, *Energy Sources and Natural Fuels*, 1993; J. H. Harker and J. R. Backhurst, *Fuel and Energy*, 1981; J. T. McMullen, R. Morgan, and R. B. Murray, *Energy Resources and Supply*, 1976; D. Marier and L. Stoiaken (eds.), *Alternative Sources of Energy*, 1988; T. Ohta, *Energy Technology: Sources, Systems, and Frontier Conversion*, 1994; H. B. Smith, *Energy: Sources, Applications, Alternatives*, 1993.

Energy storage

The general method and specific techniques for storing energy derived from some primary source in a form convenient for use at a later time when a specific energy demand is to be met, often in a different location.

In the past, energy storage on a large scale had been limited to storage of fuels. For example, large amounts of natural gas may be stored under pressure in underground reservoirs during the summer and used to meet increased demands for heating fuel in the colder season. Petroleum and its products are stored at several points in the energy system, from the strategic petroleum reserve to the fuel tanks of automobiles. Since gasoline is a highly concentrated and readily portable form of energy, this method of energy storage makes the automobile independent of the supply system for appreciable distances and times. On a smaller scale, electric energy is stored in batteries that power automobile starters and a great variety of portable appliances. See OIL AND GAS STORAGE.

Energy storage in many forms is expected to play an increasingly important role in shifting patterns of energy consumption away from scarce to more abundant and renewable primary resources. For example, automobiles may be able to store transportation energy in the form of coal-derived synthetic fuels or as electricity in batteries charged with electric power from coal or nuclear power plants. Solar energy can already be made more usable by accumulating it during the day as warm water that can be stored for later use. Another example is the storage of electric energy generated at night by coal or nuclear power plants to meet peak electric loads during daytime periods. This storage is achieved by pumping water

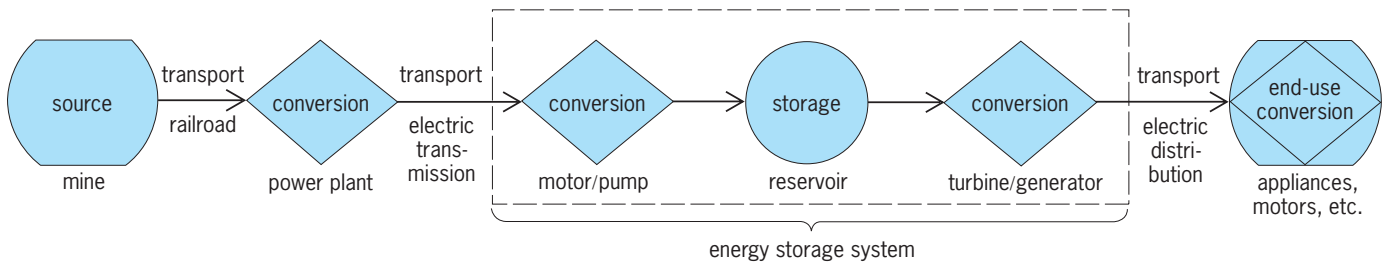


Fig. 1. Principal components of an energy storage system with specific reference to an electric power system.

from a lower to a higher reservoir at night and reversing this process during the day, the pump then being used as a turbine and the motor as a generator. As shown in Fig. 1, this example can be used to illustrate the conversion and storage functions of an energy storage system.

Broader application and use of new methods of energy storage could reduce oil consumption in the major energy-use sectors. Probably the largest practical potential for oil displacement is in the heating and cooling of buildings through storage of heat or coolness generated on site with solar energy or off-peak electricity (singly or in combination). In principle, the potential for oil displacement is largest in transportation, since highway vehicles consume large quantities of fuel. In practice, extensive displacement of conventional automobiles will be very difficult because of the very large energy storage capacity of fossil fuels compared to the most attractive alternative, electric storage batteries: a 20-gal (76-liter) tank of gasoline will give the average car a range of 300–600 mi (500–1000 km), while a lead-acid battery may provide a range of just 30–60 mi (50–100 km), weighs nearly a ton (0.9 metric ton), and costs between \$1200 and \$1500. The contrast between the two systems is illustrated in Fig. 2, which compares their specific ranges. Yet, because of their independence of oil and their more efficient as well as cleaner use of energy, electric vehicles powered by batteries are a possible alternative to conventional automobiles, especially for urban driving.

Energy storage in electric power systems has excellent potential for reducing the use of oil and gas for power generation, especially if pumped hydro storage can be supplemented by more broadly applicable techniques such as underground pumped hydro, compressed-air storage, and batteries. These techniques also could provide utilities with special operating advantages not found in conventional generating plants. Storage technologies that could find application in one or more of the major energy-use sectors are discussed below.

Pumped hydro storage. Until the late 1970s the only bulk energy storage method used by electric utilities, pumped hydro goes back to 1929 in the United States; the largest such facility in the world—with a capacity of 1.5×10^7 kWh—is the Ludington pumped storage plant on Lake Michigan. The construction of such facilities is limited by a variety of geographic, geologic, and environmental constraints, all of which may be considerably reduced

if the power generation equipment and discharge reservoir are placed deep underground (Fig. 3). The size of the surface reservoir required for such an underground pumped hydroelectric system will be several times smaller than that used in the conventional system because of the greater pressure available to run the generators. The difference in height between reservoirs in the underground system may be several thousand feet (1 ft = 0.3 m), compared to several hundred feet common in the conventional pumped hydroelectric system.

Compressed-air storage. Off-peak electric energy can also be converted into mechanical energy by pumping air into a suitable cavern where it is stored at pressures up to 80 atm (8 MPa). Compared to pumped hydro storage, compressed air offers several advantages, including a wider choice of suitable geologic formations, lower volume requirements, and a smaller minimum capacity to be economically attractive. The world's first commercial compressed-air energy storage system (Fig. 4) is in Huntorf, Germany. In this installation, expanding air must be heated by the combustion of natural gas before it can be used to drive the generator and turbines. In a more efficient concept, the heat evolved during air compression would be stored in a bed of pebbles and reused to heat the expanding air.

Batteries. The development of advanced batteries with characteristics superior to those of the

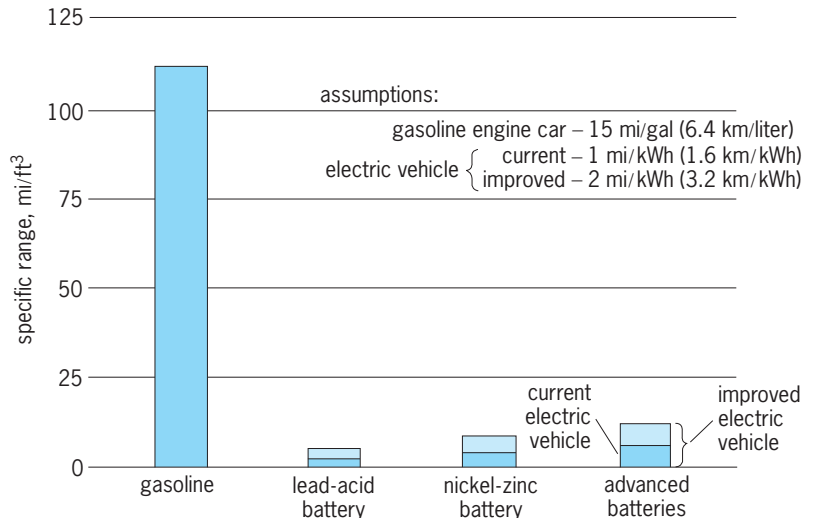


Fig. 2. Specific range of vehicles using energy stored in various forms. 1 mi/ft³ = 57 km/m³.

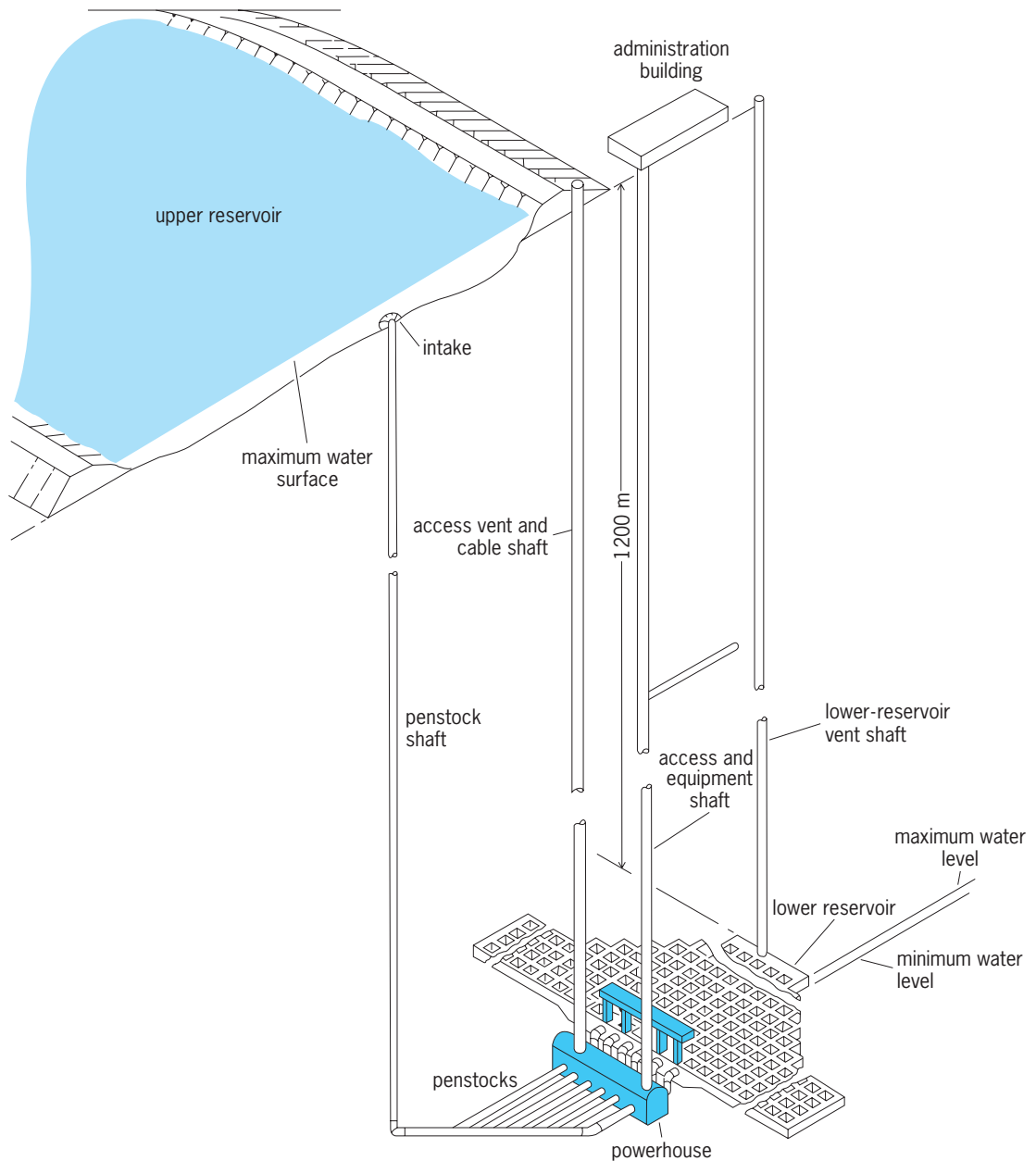


Fig. 3. Concept of an underground pumped hydroelectric storage system.

familiar lead-acid battery could result in use of battery energy storage on a large scale. For example, batteries lasting 2000 or more cycles could be used in installations with capacities of several hundred thousand kilowatt-hours in various locations on the electric power grid, as an almost universally applicable method of utility energy storage. Batteries combining these characteristics with energy densities (storage capacity per unit weight and volume) well above those of lead-acid batteries could provide electric vehicles with greater range, thus removing a major barrier to their broader use.

The development of several battery types has been undertaken to achieve systems with superior characteristics. Candidates include nickel-zinc and nickel-iron systems that could provide modest improvements over lead-acid. The zinc-chlorine,

sodium-sulfur, and lithium-iron sulfide systems could yield substantial improvements over lead-acid.

The zinc-chlorine battery (Fig. 5) operates near ambient temperature with an aqueous solution of zinc chloride as the electrolyte. On charge, zinc is deposited on a graphite substrate and chlorine is released as a slightly soluble gas at the opposite electrode, composed of porous graphite. Chlorine is stored in a separate vessel by freezing it out of the aqueous solution as a solid compound of water (ice) and chlorine at about 48°F (9°C); in this form, chlorine is relatively safe to handle. On discharge, zinc redissolves and chlorine is reduced at the porous graphite cathode, reforming the original zinc chloride aqueous solution. The active materials and other components offer the promise of a low-cost battery; however, the complexity of the system (which uses

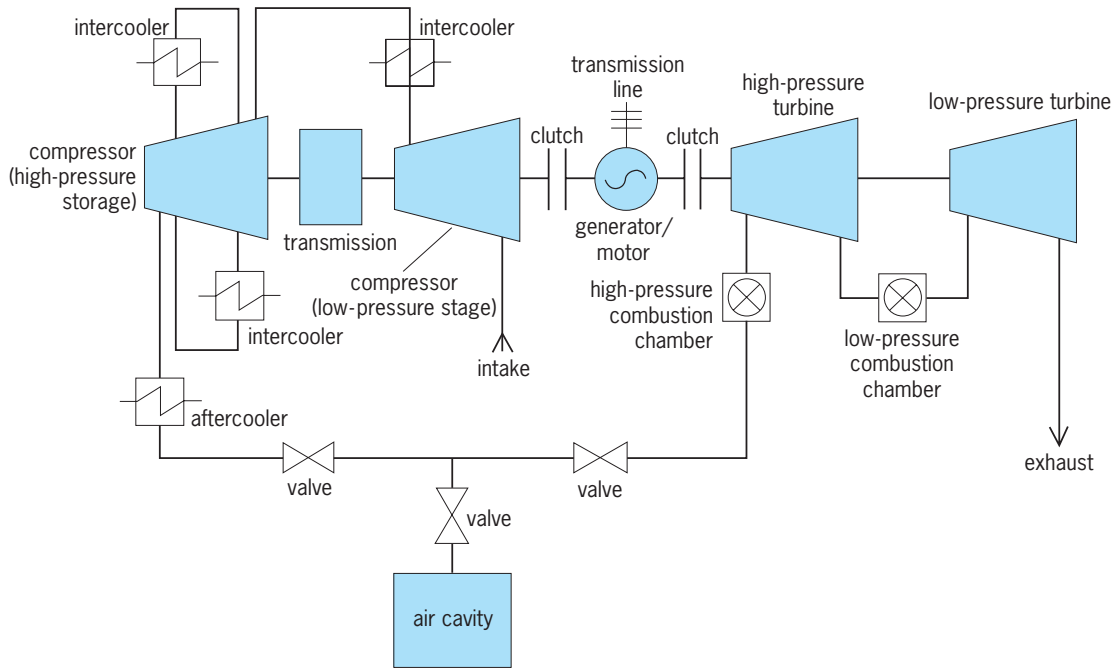


Fig. 4. Schematic diagram of compressed-air storage plant at Huntorf, Germany.

pumps, valves, and other auxiliaries) suggests that operating reliability could be a possible problem. See SOLID-STATE BATTERY.

The sodium-sulfur battery operates at temperatures above 570°F (300°C), where its electrodes exist as liquids separated by a unique solid ceramic electrolyte. Small experimental cells have been tested through 1000 charging cycles while showing no appreciable degradation. A disadvantage of the sodium-

sulfur battery, however, is that it must operate at a temperature of 570–660°F (300–350°C). The operating costs and safety characteristics of large batteries consisting of thousands of single cells in electric series and parallel connection are also potential drawbacks. See BATTERY.

Thermal storage. Ceramic brick “storage heaters” that store off-peak electricity in the form of heat are used for heating buildings in Europe; the barriers to

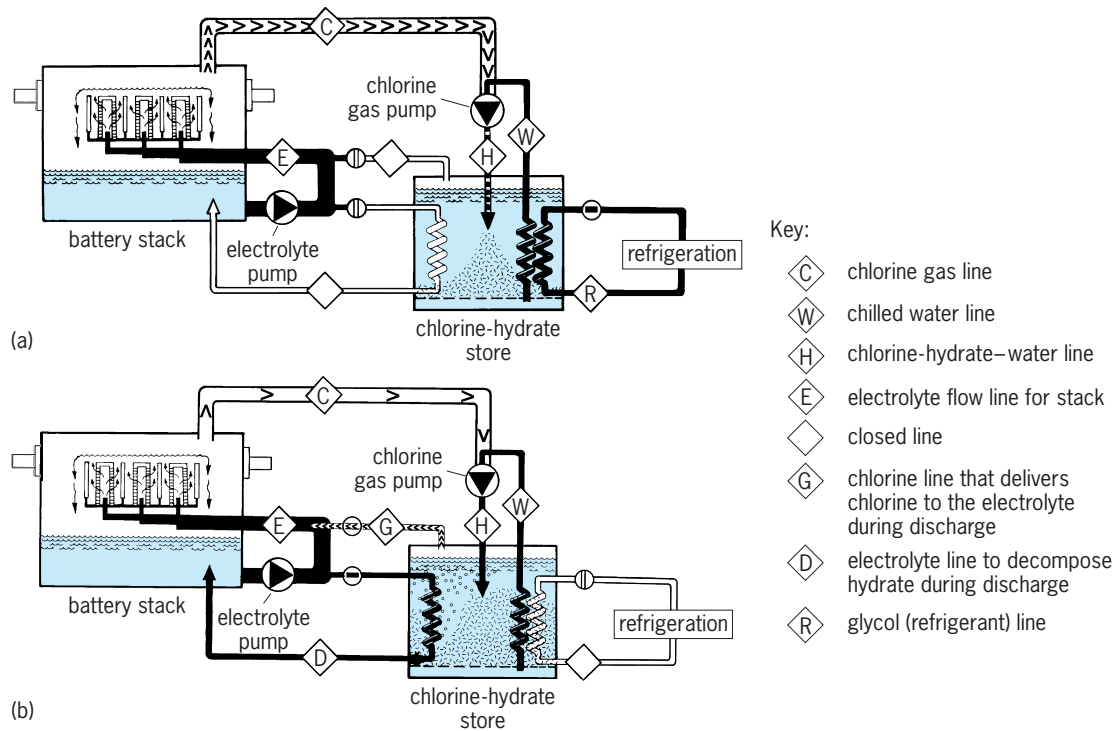


Fig. 5. Schematic diagram of zinc-chloride battery, showing two main compartments in the design, the battery stack and the chlorine hydrate store. (a) Charge cycle. (b) Discharge cycle. (Gulf and Western Co.)

their use in the United States are more institutional and economic than technological.

Testing has been done on prototype "coolness" storage systems, which use electric refrigeration to chill water or produce ice at night. Experimental installations indicate that daytime electric power demand for air conditioning could be reduced up to 75% by using such systems, but these systems are still relatively bulky and expensive.

Solar hot-water storage is technically simple and commercially available. However, the use of solar energy for space heating requires relatively large storage systems, with water or rock beds as storage media, and difficulties can arise in integrating this storage with existing buildings while keeping costs within acceptable limits. Innovative designs that may help overcome these problems include new types of heat transfer equipment, storage of heat in the walls and ceilings of buildings (passive solar heating), and the use of low-cost, roof-mounted water bags. See SOLAR ENERGY; SOLAR HEATING AND COOLING.

Chemical reaction systems. Heat or electricity may be stored by using these energy forms to force certain chemical reactions to occur. Such reactions are chosen so that they can be reversed readily with release of energy; in some cases the products can be transported from the point of generation to that of consumption, adding flexibility to the ways the stored energy can be used. For example, reactions which produce hydrogen could become attractive since hydrogen could be stored for extensive periods of time and then conveniently used in either combustion devices or in fuel cells.

Another possibility is to apply heat to a mixture of methane and water, converting them to hydrogen and carbon monoxide, which can be stored and transported to the end-use site, where a catalyst permits the reverse reaction to occur spontaneously with the release of heat. While not yet in practical

use, chemical reaction systems are under extensive investigation as economic and flexible ways of storing energy. See FUEL CELL.

Superconducting magnets. Electrical energy can be stored directly in the form of large direct currents used to create fields surrounding the superconducting windings of electromagnets. In principle such devices appear attractive because their storage efficiency is high, plant life could be long, and utilities would have few difficulties establishing the necessary conversion equipment. However, the need for maintaining the system at temperatures approaching absolute zero and, particularly, the need to physically restrain the coils of the magnet when energized require auxiliary equipment (insulation, vacuum vessels, and structural supports) which will represent a large expenditure. See SUPERCONDUCTING DEVICES.

Flywheels. Storage of kinetic energy in rotating mechanical systems such as flywheels is attractive where very rapid absorption and release of the stored energy is critical. However, even advanced designs and materials are likely to be too expensive for utility energy storage on a significant scale, and applications will probably remain limited to systems where high power capacity and short charging cycles are the prime consideration. Such applications do exist in pulse power supplies and in electric transportation for recovery of braking energy. See FLY-WHEEL.

Combined systems. The rising importance of energy storage comes from its potential for shifting demand from scarce to plentiful primary energy sources. As such, the most successful storage devices are likely to be those that are adopted as components of larger systems designed specifically for resource conservation. In electric power generation, for example, system-wide storage may be combined with solar-electric systems to flatten load curves and reduce oil consumption (Fig. 6). If such complex systems are eventually to be developed and commercially introduced, financial incentives and more coherent national energy policy planning must be adopted. Thus, while advanced energy storage systems could have a major impact on United States energy consumption patterns with continued research, nontechnical factors such as regulatory strategies and pricing policies are likely to become decisive to their large-scale use. See ENERGY SOURCES. Fritz Kalhammer; Thomas R. Schneider

Bibliography. G. Beckmann and P. V. Gilli, *Thermal Energy Storage*, 1984; F. Dinter, M. A. Geyer, and R. Tamme, *Thermal Energy Storage for Commercial Applications*, 1993; J. Jensen *Fundamentals of Energy Storage*, 1985; F. R. Kalhammer, Energy-storage systems, *Sci. Amer.*, 241(6):56-65, December 1979; S. Karkainen and Y. El-Maghary (eds.), *Energy Storage Systems in Developing Countries*, 1988; D. Linden (ed.), *Handbook of Batteries*, 2d ed., 1995; B. Kilkis and S. Kakac (eds.), *Energy Storage Systems*, 1989; M. N. Wilson, *Superconducting Magnets*, 1987; F. De Winter, *Solar Collectors, Energy Storage, and Materials*, 1991.

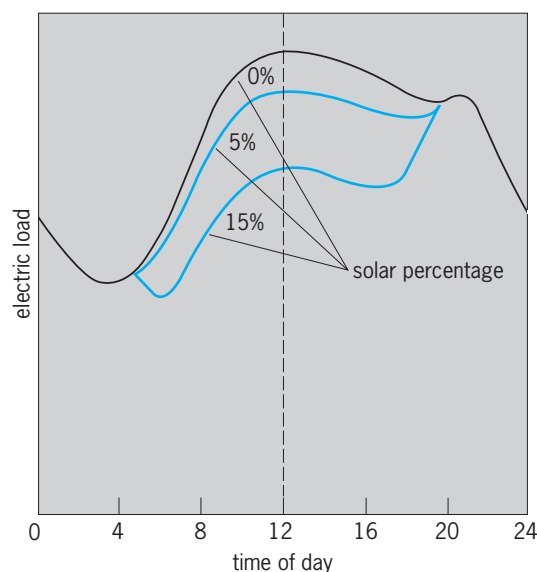


Fig. 6. Impact of solar-electric generation on utility load curve during a typical 24-h period.

Engine

A machine designed for the conversion of energy into useful mechanical motion. The principal characteristic of an engine is its capacity to deliver appreciable mechanical power, as contrasted to a mechanism such as a clock, whose significant output is motion. By usage an engine is usually a machine that burns or otherwise consumes a fuel, as differentiated from an electric machine that produces mechanical power without altering the composition of matter. Similarly, a spring-driven mechanism is said to be powered by a spring motor; a flywheel acts as an inertia motor. By definition a hydraulic turbine is not an engine, although it competes with the engine as a prime source of mechanical power. *See* ENERGY CONVERSION; HYDRAULIC TURBINE; MOTOR; PRIME MOVER; WATERPOWER.

Applications. A fuel-burning engine may be stationary, as a donkey engine used to lift cargo between wharf and ship, or it may be mobile, like the engine in an aircraft or automobile. Such an engine may be used for both fixed service and mobile operation, although accessory modifications that adapt the engine to its particular purpose are preferable. For example, the fan that draws air through the radiator of a water-cooled fixed engine is large and fitted in a baffle, whereas the fan of a similar but mobile engine can be small and unbaffled because considerable air is driven through the radiator by means of ram action as the engine propels itself along. *See* AIRCRAFT ENGINE; AUTOMOTIVE ENGINE; MARINE ENGINE.

Some types of engine can be designed for economic efficiencies in fixed service but not in mobile operation. Thus, the steam engine is widely used in central electric generator stations but is obsolete in mobile service. This is chiefly because, in a large ground installation, the furnace and boiler can be fitted with means for using most of the available heat. The engine proper can be a reciprocating (piston) or a rotating (turbine) type. Because shaft rotation is by far the most used form of mechanical motion, the turbine is the more common form of modern steam engine. For railroad service the steam engine has given place to diesel and gasoline internal combustion engines and to electric motors. *See* BOILER; LOCOMOTIVE; POWER PLANT; STEAM ENGINE; STEAM-GENERATING UNIT; STEAM TURBINE.

Types. Traditionally, engines are classed as external or internal combustion. External combustion engines consume their fuel or other energy source in a separate furnace or reactor. *See* NUCLEAR REACTOR; STEAM-GENERATING FURNACE.

Strictly, the furnace or reactor releases chemical or nuclear energy into thermal energy, and the engine proper converts the heat into mechanical work. The principal means for the conversion of heat to work is a gas or vapor, termed the working fluid. By extension, an engine which derives its heat energy from the Sun by solar radiation to working fluid in a boiler can be considered an external combustion type. To avoid loss of or contamination from nuclear fuel, the reactor and boiler are separated from (and

may also be shielded from) the engine. *See* SHIP NUCLEAR PROPULSION.

The working fluid takes on energy in the form of heat in the boiler and gives up energy in the engine, the engine proper being a thermodynamic device. The device may be a turbine for stationary power generation or a nozzle for long-range vehicular propulsion. *See* NUCLEAR POWER.

In an engine used for propulsion, the rearward velocity with which the working fluid is ejected and, thus, the forward acceleration imparted to the vehicle depend on the temperature of the fluid. For practical purposes, temperature is limited by the engine materials that serve to contain the chemical combustion or nuclear reaction. To achieve higher exhaust velocity, the working fluid may be contained by non-material means such as electric and magnetic fields, in which case the fluid must be electrically conductive. *See* ION PROPULSION.

The engine proper is then a magnetohydrodynamic device receiving electric energy from a separate fuel-consuming source such as a gas turbine and electric generator or a nuclear reactor and electric generator or possibly from direct electric conversion of nuclear or solar radiation. *See* INTERPLANETARY PROPULSION; SPACECRAFT PROPULSION; THERMOELECTRICITY.

A further basis of classification concerns the working fluid. If the working fluid is recirculated, the engine operates on a closed cycle. If the working fluid is discharged after one pass through boiler and engine, the engine operates on an open cycle. Closed-cycle operation assures the purity of the working fluid and avoids the discharge of harmful wastes. The open cycle is simpler. Thus the commonest types of engine use atmospheric air in open cycles both as the principal constituent of their working fluids and as oxidizer for their fuels.

If open-cycle operation is used, the next modification is to heat the working fluid directly by burning fuel in the fluid; the engine becomes its own furnace. Because this internal combustion type engine uses the products of combustion as part of the working fluid, the fuel must be capable of combustion under the operating conditions in the engine and must produce a noncorrosive and nonerosive working fluid. Such engines are the common reciprocating gasoline and diesel units. *See* DIESEL ENGINE; INTERNAL COMBUSTION ENGINE; ROTARY ENGINE; STIRLING ENGINE.

At low speeds the combustion process is carried out intermittently in a cylinder to drive a reciprocating piston. At high speed, however, friction between piston and cylinder walls and between other moving parts dissipates an appreciable portion of the developed power. Thus, where high power is developed at high speed, performance is improved by continuous combustion to drive a turbine wheel. *See* BRAYTON CYCLE; CARNOT CYCLE; DIESEL CYCLE; GAS TURBINE; OTTO CYCLE; TURBINE PROPULSION.

Engine shaft rotation may be used in the same way as in a reciprocating engine. However, for high-velocity vehicular propulsion, the energy of the

working fluid may be converted into thrust more directly by expulsion through a nozzle. Once the vehicle is in motion, the turbine can be omitted. Alternatively, instead of drawing atmospheric oxygen into the combustion chamber, the engine may draw both oxidizer and fuel from storage tanks within the vehicle, or the combustion chamber may contain the full supply of fuel and oxidizer. See JET PROPULSION; RAMJET; ROCKET PROPULSION; TURBOJET; TURBOPROP; TURBORAMJET.

Despite all the variation in structure, mode of operation, and working fluid—whether of moving parts, moving fields, or only moving working fluid—these machines are basically means for converting heat energy to mechanical energy. See THERMODYNAMIC PROCESSES.

Frank H. Rockett; D. L. Anglin

Bibliography. W. H. Crouse and D. L. Anglin, *Automotive Engines*, 8th ed., 1994; L. Guzzella and C. H. Onder, *Introduction to Modeling and Control of Internal Combustion Engine Systems*, 2004; R. Stone, *Introduction to Internal Combustion Engines*, 3d ed., 1999.

Engine cooling

A cooling system in an internal combustion engine that is used to maintain the various engine components at temperatures conducive to long life and proper functioning. Gas temperatures in the cylinders may reach 4500°F (2500°C). This is well above

the melting point of the engine parts in contact with the gases; therefore it is necessary to control the temperature of the parts, or they will become too weak to carry the stresses resulting from gas pressure. The lubricating oil film on the cylinder wall can fail because of chemical changes at wall temperatures above about 400°F (200°C). Complete loss of power may take place if some spot in the combustion space becomes sufficiently heated to ignite the charge prematurely on the compression stroke. See INTERNAL COMBUSTION ENGINE.

Fortunately, a thin protective boundary of relatively stagnant gas of poor heat conductivity exists on the inner surfaces of the combustion space. If the outer cylinder surface is placed in contact with a cool fluid such as air or water and there is sufficient contact area to cause a rapid heat flow, the resulting drop in temperature produced by the heat flow in the inside boundary layer keeps the temperature of the cylinder wall much closer to the temperature of the coolant than to the temperature of the combustion gas. The quantity of heat that crosses the stagnant boundary layer and must be carried away by the coolant is a function of the Reynolds number of the gas existing in the cylinder. In terms of practical engine quantities, the heat flow to the coolant varies approximately as in the relationship: (charge density \times piston speed)^{0.8}. At full throttle and normal piston speed, this heat flow amounts to about 15% of the energy of the incoming fuel. See REYNOLDS NUMBER.

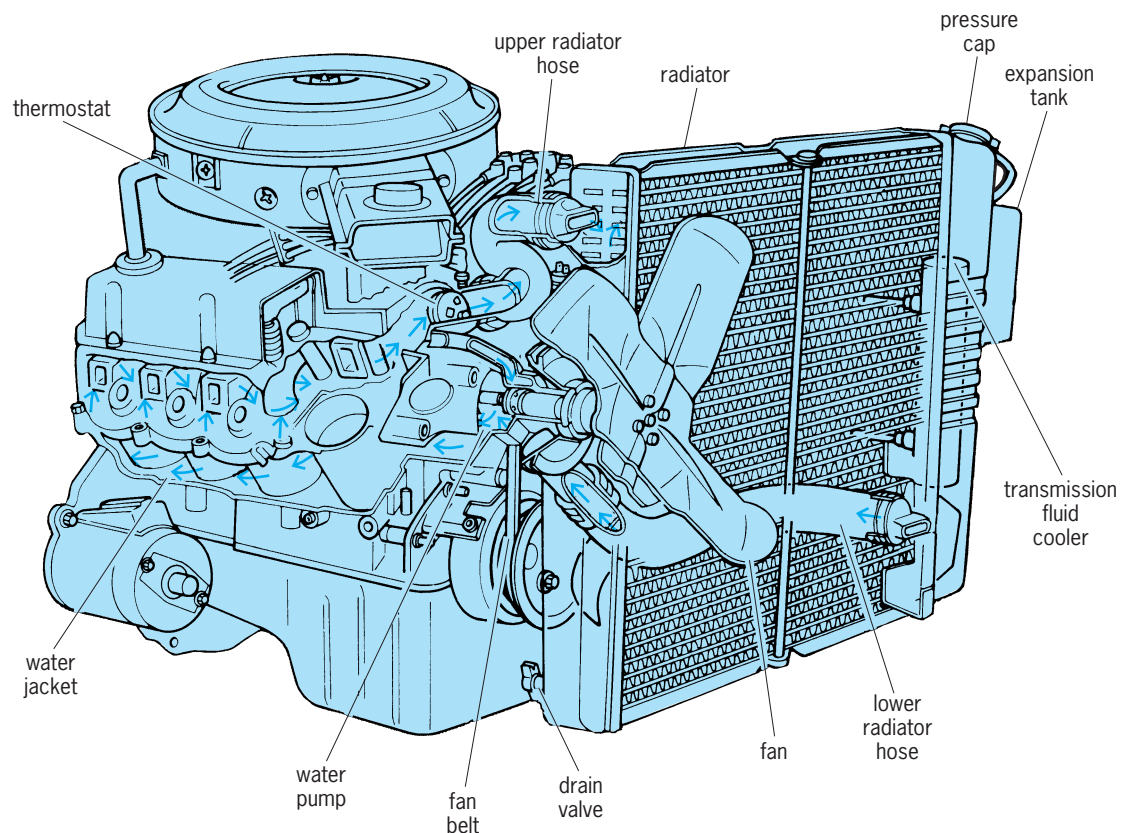


Fig. 1. Cooling system of a V-8 automotive spark-ignition engine. The arrows show the direction of coolant flow through the engine water jackets and cooling system. (Ford Motor Co.)

Liquid cooling. If the coolant is water, it is usually circulated by a pump through jackets surrounding the cylinders and cylinder heads. The water is circulated fast enough to remove steam bubbles that may form over local hot spots and to limit the water's temperature rise through the engine to about 15°F (8°C). In most engines in automotive and industrial service, the warmed coolant is piped to an air-cooled heat exchanger called a radiator (Fig. 1). The airflow required to remove the heat from the radiator is supplied by an electric or engine-driven fan; in automotive applications the airflow is also supplied by the forward motion of the vehicle. The engine and radiator may be separated and each placed in the optimum location, being connected through piping. To prevent freezing, the water coolant is usually mixed with ethylene glycol. See ANTIFREEZE MIXTURE; ETHYLENE GLYCOL; HEAT EXCHANGER.

Low water-jacket temperature contributes to corrosive wear of engine parts and increases piston friction losses. High water-jacket temperature increases the coolant loss by evaporation or by actual boiling. In automotive engines, temperature of the water jacket is often automatically maintained near 195°F (91°C) by a thermostat placed in the line from the engine to the radiator. When the engine outlet water is too cool, as when the engine is first started, the water is prevented from entering the radiator, and is usually recirculated in the engine block until warm enough to open the thermostat.

Many marine diesel engines have a closed cooling system in which fresh water is circulated through the engine. Instead of flowing through a radiator, the heated water flows through a heat exchanger, or cooler, immersed in seawater so that the heat is carried away by it instead of by air. In some engines, the heated water is used in a jacket to cool the exhaust manifold.

Air cooling. Engines are often cooled directly by a stream of air without the interposition of a liquid medium. The heat-transfer coefficient between the cylinder and airstream is much less than with a liquid coolant, so that the cylinder temperatures must be much greater than the air temperature to transfer to the cooling air the heat flowing from the cylinder gases. To remedy this situation and to reduce the cylinder wall temperature, the outside area of the cylinder, which is in contact with the cooling air, is increased by finning. The heat flows easily from the cylinder metal into the base of the fins, and the great area of finned surface permits heat to be transferred to the cooling air (Fig. 2). The ideal fin shape depends upon the conductivity of the fin material. In general, the fin is thickest at the base to permit heat flow from the cylinder. The fin should taper to a thin edge to give a good temperature gradient along its length. For reasons of mechanical strength, fins are usually made thicker than necessary for heat-transfer considerations. See HEAT TRANSFER.

High-output cylinders require many closely spaced fins. In these engines the area between adjacent cylinders is blocked off with sheet-metal baffles, which are also shaped to follow the fin tips part

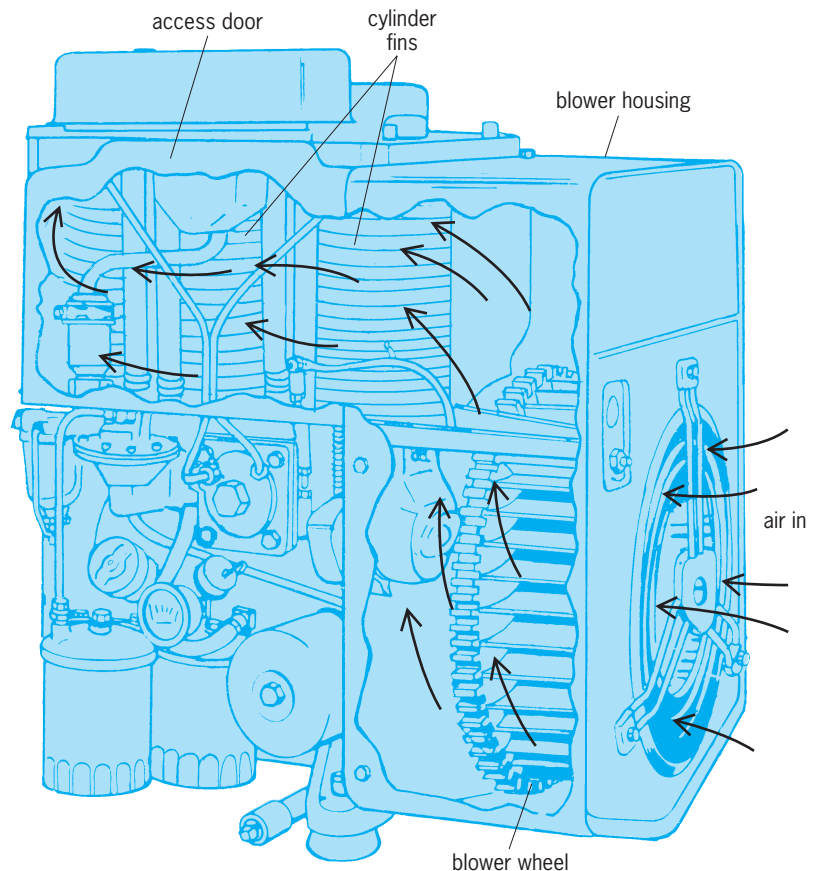


Fig. 2. Airflow around the cylinders and fins of an air-cooled engine. (Onan Corp.)

way around the cylinder. A pressure is built up in front of the baffles by means of a fan or because of the forward motion of the vehicle (ram effect). The pressure differential between front and rear of the engine forces the cooling air through the spaces between the fins.

The power required to cool depends upon the quantity of cooling air used and the velocity at which it passes the fins. To keep the power required for cooling to a minimum, the fins should be long and close together so that a large heat-transfer area is served by a small coolant flow area. The temperature difference between fins and cooling air should be kept as high as possible so that less air velocity will be required. Cylinder temperatures of air-cooled engines are sometimes controlled by louvers or flaps, which may be set to restrict the cooling airflow until the engine becomes warm.

Adiabatic engine. This type of diesel engine has no cooling system through which heat is lost. Ceramic materials, which have higher melting points than metals, retain the heat generated in the combustion chamber and cylinder. This has the potential to make the adiabatic diesel engine a more efficient and economical power plant. See ADIABATIC PROCESS; CERAMICS.

Augustus R. Rogowski; Donald L. Anglin

Bibliography. W. H. Crouse and D. L. Anglin, *Automotive Fuel, Lubricating, and Cooling Systems*, 1981; Society of Automotive Engineers, *Engine Coolants, Cooling System Materials, and Components*, 1993.

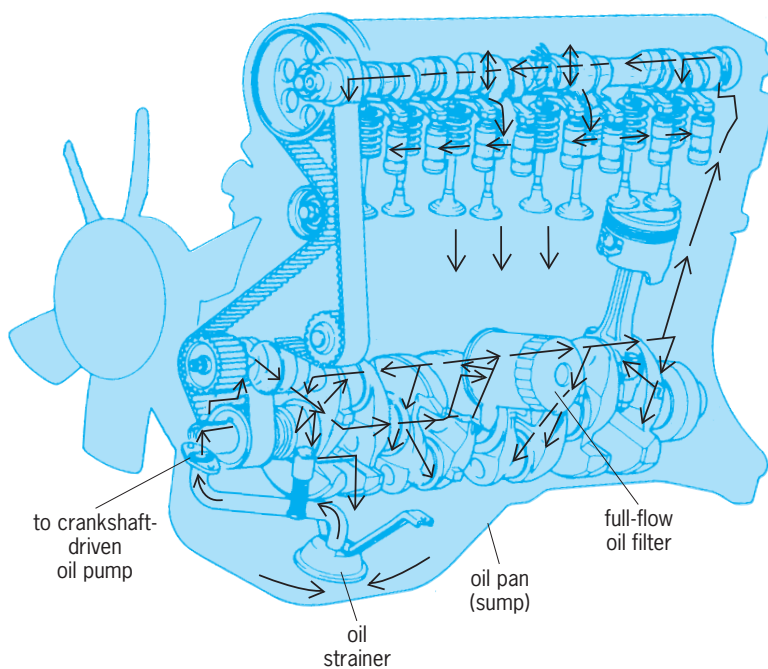
Engine lubrication

In an internal combustion engine, the system for providing a continuous supply of oil between moving surfaces during engine operation. This viscous film, known as the lubricant, lubricates and cools the power transmission components while removing impurities, neutralizing chemically active products of combustion, transmitting forces, and damping vibrations. See INTERNAL COMBUSTION ENGINE; LUBRICANT.

Lubricant. Automotive and other four-stroke Otto-cycle engines are generally lubricated with petroleum-base oils that contain chemical additives to improve their natural properties. Synthetic oils are used in gas turbines and may be used in other engines. See OTTO CYCLE.

Probably the most important property of oil is the absolute viscosity, which is a measure of the force required to move one layer of the oil film over the other. If the viscosity is too low, a protecting oil film is not formed between the parts. With high viscosity too much power is required to shear the oil film, and the flow of oil through the engine is retarded. Viscosity tends to decrease as temperature increases. Viscosity index (VI) is a number that indicates the resistance of an oil to changes in viscosity with temperature. The smaller the change in viscosity with temperature, the higher the viscosity index of the oil. See VISCOSITY.

Lubricating system. Small two-stroke cycle engines may require a premix of the lubricating oil with the fuel going into the engine, or the oil may be injected into the ingoing air-fuel mixture. This is known as a total-loss lubricating system because the oil is consumed during engine operation.



Pressurized lubricating system for an automobile engine. Arrows show the flow of oil through the engine. (Toyota Motor Sales U.S.A. Inc.)

Most automotive engines have a pressurized or force-feed lubricating system in combination with splash and oil mist lubrication (see *illus.*). The lubricating system supplies clean oil cooled to the proper viscosity to the critical points in the engine, where the motion of the parts produces hydrodynamic oil films to separate and support the various rubbing surfaces. The oil is pumped under pressure to the bearing points, while sliding parts are lubricated by splash and oil mist. After flowing through the engine, the oil collects in the oil pan or sump, which cools the oil and acts as a reservoir while the foam settles out. Some engines have an oil cooler to remove additional heat from the oil. See WEAR.

Oil filter. This filter removes impurities from the oil passing through it and helps maintain the lubricating ability of the oil during the maintenance intervals. The impurities removed include solid particles such as combustion residue, metal particles, and dust.

Two types of filters have been used on internal combustion engines, full-flow and bypass. The full-flow filter protects the engine by filtering the entire oil-pump output and traps any particles during the first pass of the oil. However, a bypass valve is required to ensure continued supply of oil to the engine should the filter become clogged. Most automotive engines have pleated-paper filters of this type. Bypass filters remove about 10% or less of the oil passing through the lubricating system, and return this oil to the sump after cleaning. Some engines have both full-flow and bypass filters for the lubricating oil.

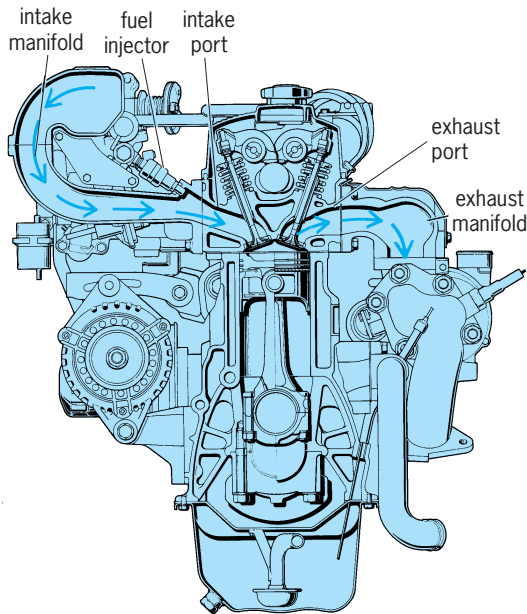
Donald L. Anglin

Bibliography. *Bosch Automotive Handbook*, 1986; W. H. Crouse and D. L. Anglin, *Automotive Fuel, Lubricating, and Cooling Systems*, 1981; Deere and Co. Staff (eds.), *Engines*, 7th ed., 1991.

Engine manifold

An arrangement or collection of pipes or tubing with several inlet or outlet passages through which a gas or liquid is gathered or distributed. The manifold may be a casting or fabricated of relatively light material. Manifolds are usually identified by the service provided, such as the intake manifold and exhaust manifold on the internal combustion engine (see *illus.*). Some types of manifolds for handling oil, water, and other fluids such as engine exhaust gas are often called headers. In the internal combustion engine, the intake and exhaust manifolds are an integral part of multicylinder engine construction and essential to its operation. See INTERNAL COMBUSTION ENGINE.

Intake manifold. The engine intake manifold is a casting or assembly of passages through which air or an air-fuel mixture flows from the air-intake or throttle valves to the intake valve ports in the cylinder head or cylinder block. In a spark-ignition engine with a carburetor or throttle-body fuel injection, the intake manifold carries an air-fuel mixture. In an engine with port fuel injection or in a diesel engine, the intake manifold carries only air. For the



Cross-sectional view of a fuel-injected spark-ignition engine showing the intake and exhaust manifolds. Arrows show the flow of the air through the intake manifold and of the exhaust gases through the exhaust manifold. (Ford Motor Co.)

diesel engine, the air should be inducted with a minimum of pressure drop.

The purpose of the intake manifold is to distribute the air or air-fuel mixture uniformly to each of the cylinders and to assist in the vaporization of the fuel. The problems of manifold design are most severe when a large number of cylinders is charged through a single carburetor. Careful design is required to prevent unequal flow, localized condensation of fuel, and power losses due to overheating of the charge. The downdraft manifold is most frequently employed with automotive engines. To maximize volumetric efficiency, some engines have a tuned intake manifold, in which the port cross-sectional area and length are adjusted to a size that fills the cylinders most efficiently at a given engine speed, thereby causing engine torque to peak at that speed.

A leaking intake system may allow unfiltered air to reach the cylinders. In a spark-ignition engine, a leak may cause detonation, misfire, and driveability and exhaust-emission problems during engine operation, and evaporative hydrocarbon emissions when the engine is not operating. Leakage of air into the air-fuel mixture manifold unbalances the engine by producing lean mixtures at the leak and upsetting the calibration of the fuel-metering system. See AIR POLLUTION; FUEL INJECTION.

A leaking intake system on a naturally aspirated diesel engine may not be extremely critical. However, an explosive combustible vapor may be aspirated. Scavenging air manifolds or pressure-charged intakes lose valuable air, both quantitatively and in terms of pressure, if leaks occur. See SUPERCHARGER.

Exhaust manifold. The engine exhaust manifold is a casting or assembly of passages through which

the products of combustion leave the exhaust-valve ports in the cylinder head or cylinder block and enter the exhaust piping system. The purpose of the exhaust manifold is to collect and carry these exhaust gases away from the cylinders with a minimum of back pressure. The pulsating flows in the exhaust manifold results in pressure variations that may improve scavenging or may result in increased back pressure for one or more cylinders. Since the frequency of the pressure variations for a given engine and exhaust system varies with engine speed, the exhaust manifold should be designed to operate efficiently over the expected speed range of the engine. The entire exhaust system, including the exhaust manifold, catalytic converter (in automotive engines), muffler, and piping, affects the efficiency of combustion-gas evacuation from the cylinders.

The volume of exhaust gas is greater than the volume of air handled by the intake system, and the exhaust gases are at much higher temperature than the intake. Some small pressure drop is tolerated in the exhaust system, but exhaust back pressure represents a direct loss of engine power. Stationary engines intended to operate indoors have water-jacketed exhaust manifolds which provide cooling while limiting the heat dissipated into the engine room.

Exhaust manifolds operate at high temperatures and may be subject to erosive or corrosive attack while pressures are low or moderate. Engines with turbochargers may have multiple exhaust manifolds to prevent pulsations or interferences in the gas flow to the power turbine. An engine with an exhaust-driven turbocharger will lose power and efficiency if leakage exists ahead of the power turbine. Exhaust manifolds may also be tuned to improve cylinder scavenging and reduce back pressure in the exhaust system.

Air manifold. Some automotive spark-ignition engines have an air-distribution manifold as part of the exhaust-emission control system. This manifold distributes and proportions air to the individual exhaust ports through external tubing or integral passageways.

Donald L. Anglin

Bibliography. W. H. Crouse, *Automotive Engines*, 8th ed., 1994; W. H. Crouse and D. L. Anglin, *Automotive Fuel, Lubricating, and Cooling Systems* 1981; C. F. Taylor, *The Internal-Combustion Engine in Theory and Practice*, vol. 2, 1985.

Engineering

Most simply, the art of directing the great sources of power in nature for the use and the convenience of humans. In its modern form engineering involves people, money, materials, machines, and energy. It is differentiated from science because it is primarily concerned with how to direct to useful and economical ends the natural phenomena which scientists discover and formulate into acceptable theories. Engineering therefore requires above all the creative imagination to innovate useful applications of natural

phenomena. It is always dissatisfied with present methods and equipment. It seeks newer, cheaper, better means of using natural sources of energy and materials to improve the standard of living and to diminish toil.

Types of engineering. Traditionally there were two divisions or disciplines, military engineering and civil engineering. As knowledge of natural phenomena grew and the potential civil applications became more complex, the civil engineering discipline tended to become more and more specialized. The practicing engineer began to restrict operations to narrower channels. For instance, civil engineering came to be concerned primarily with static structures, such as dams, bridges, and buildings, whereas mechanical engineering split off to concentrate on dynamic structures, such as machinery and engines. Similarly, mining engineering became concerned with the discovery of, and removal from, geological structures of metalliferous ore bodies, whereas metallurgical engineering involved extraction and refinement of the metals from the ores. From the practical applications of electricity and chemistry, electrical and chemical engineering arose.

This splintering process continued as narrower specialization became more prevalent. Civil engineers had more specialized training as structural engineers, dam engineers, water-power engineers, bridge engineers; mechanical engineers as machine-design engineers, industrial engineers, motive-power engineers; electrical engineers as power and communication engineers (and the latter divided eventually into telegraph, telephone, radio, television, and radar engineers, whereas the power engineers divided into fossil-fuel and nuclear engineers); mining engineers as metallic-ore mining engineers and fossil-fuel mining engineers (the latter divided into coal and petroleum engineers).

As a result of this ever-increasing utilization of technology, people and their environments have been affected in various ways—some good, some bad. Sanitary engineering has been expanded from treating the waste products of humans to also treating the effluents from technological processes. The increasing complexity of specialized machines and their integrated utilization in automated processes has resulted in physical and mental problems for the operating personnel. This has led to the development of bioengineering, concerned with the physical effects upon humans, and management engineering, concerned with the mental effects.

Integrating influences. While the specialization was taking place, there were also integrating influences in the engineering field. The growing complexity of modern technology called for many specialists to cooperate in the design of industrial processes and even in the design of individual machines. Interdisciplinary activity then developed to coordinate the specialists. For instance, the design of a modern structure involves not only the static structural members but a vast complex including moving parts (elevators, for example); electrical machinery and power distribution; communication systems; heat-

ing, ventilating, and air conditioning; and fire protection. Even the structural members must be designed not only for static loading but for dynamic loading, such as for wind pressures and earthquakes. Because people and money are as much involved in engineering as materials, machines, and energy sources, the management engineer arose as another integrating factor.

Typical modern engineers go through several phases of activity during their careers. Formal education must be broad and deep in the sciences and humanities underlying the particular field. Then comes an increasing degree of specialization in the intricacies of the discipline, also involving continued postscholastic education. Normal promotion thus brings interdisciplinary activity as the engineer supervises various specialists. Finally, the engineer enters into the management function by interweaving workers, money, materials, machines, and energy sources into completed processes for the use of humankind.

For specific articles on various engineering disciplines *see* CHEMICAL ENGINEERING; CIVIL ENGINEERING; ELECTRICAL ENGINEERING; INDUSTRIAL ENGINEERING; MANUFACTURING ENGINEERING; MARINE ENGINEERING; MECHANICAL ENGINEERING; METHODS ENGINEERING; MINING; NUCLEAR ENGINEERING; SCIENCE; TECHNOLOGY.

Joseph W. Barker

Engineering, social implications of

The rapid development of human ability to bring about drastic alterations of the environment has added a new element to the responsibilities of the engineer. Traditionally, the ingredients for sound engineering have been sound science and sound economics. Today, sound sociology must be added if engineering is to meet the challenge of continued improvement in the standard of living without degradation of the quality of the environment.

Despite the fact that present and evolving engineering practices must meet the criteria of scientific and economic validity, these same practices generally cause societal problems of new dimensions. Consider, for example, exhaust gases emitted from tens of millions of internal combustion engines, both stationary and moving; stack gases from fossil-fuel-burning plants generating steam or electric power; gaseous and liquid effluents and solid waste from incinerators and waste-treatment systems; strip mining of coal and mineral ores; noise issuing from automotive vehicles, aircraft, and factory and field operations; toxic, nondegradable or long-lived chemical and particulate residues from ore reduction, chemical processing, and a broad spectrum of factory and mill operations; dust storms, soil erosion, and disruption of ground-water quality and quantity accompanying intensified mechanized farming in conjunction with massive irrigation and fresh-water diversion. *See* ECOLOGY, APPLIED.

Progress often results in the substitution of one set of problems for another. For example, in nuclear

electric power generating plants, replacement of fossil fuels by nuclear fuels relieves the burden of atmospheric pollution from stack gas emissions. Lower thermal efficiency of a nuclear plant, however, results in higher heat rejection rates and increased thermal pollution of sources of cooling water or air. The attendant consequences on atmospheric conditions or on the viability of aquatic life in the affected water are of great concern in the short and long terms. Ultimately, the cost and benefit considerations of nuclear power must be all-inclusive; in addition to usual considerations of economic length of plant life and so forth, one must account for all the economic and societal costs of the entire fuel cycle, from mining and refinement through use and ultimate recycling or safe disposal. The long-term effects of very low levels of radiation exposure (as such studies become available) will be an additional factor to consider.

The modern engineer must be increasingly conscious of the societal consequences of technological innovation. See ENGINEERING. Eugene A. Avallone

Engineering and architectural contracts

The legal documents pursuant to which most professional engineering and architectural services are rendered. A contract is the statement of promises agreed to by two or more parties in order to effectuate a specific goal. For the law to enforce these promises, certain elements must be present: mutuality of obligation, that is, each of the parties must have an obligation to the others; a meeting of the minds, that is, it must appear that the parties agree; and all material elements of the proposed transaction must be provided for.

Types. Engineers and architects, (that is, design professionals,) typically draw up three types of contract. (1) The design contract is an instrument under which the design professionals provide services to their client, who is usually the owner of the project. (2) The design subcontract is an instrument under which the original group of design professionals may retain other design professionals to assist them. (3) The construction contract is an agreement between the owner and construction contractors. It consists of the plans and specifications produced by the engineer, and such other provisions as may be needed to define the relative rights and responsibilities of owner and contractor in implementing the plans and specifications. Although usually not a party to the construction contract, the design professional normally prepares this document and is usually involved in carrying out the activities detailed in the design contract. Indeed, the purpose of most design contracts is the preparation of the construction contract.

Need for written contract. Except where the law requires certain types of agreements to be in writing, an oral contract is as enforceable as a written one. Within the United States, the law in most states requiring that certain types of contracts be in writing (generally called the Statute of Frauds) mandates written contracts in certain circumstances such as

the purchase of real property. The purpose of requiring a written contract is to lessen disputes as to what was agreed upon.

Unlike the services of virtually all other learned professionals (for example, doctors, dentists, lawyers, and accountants), the services of design professionals are almost always rendered according to a written contract. The principal reason is that in the construction process many parties are involved (for example, the owner; architect; surveyor; soils, structural, electrical, heating, ventilating, and air-conditioning engineers; the various prime contractors and subcontractors; and manufacturers and suppliers). Without written agreements to define the rights and obligations of each of these parties, there would be hopeless confusion.

Contract enforcement. Contracts are enforced by the legal system. Contracts are interpreted and the rights of the parties are determined either by courts or by arbitration (most often under the Construction Industry Rules as promulgated and administered by the American Arbitration Association). Arbitrators have no power to enforce their decisions, but the laws in virtually all states provide that arbitrators' decisions will be enforced by the legal system.

Where a matter is submitted to the court, it attempts to apply established legal principles to enforce the terms of the contract and to give each party what was agreed to, even though the results might appear to be inequitable. Arbitrators, however, are free to ignore the law and make rulings based upon their personal views of what is equitable. Although it might appear that an equitable result having no legal foundation is better than an unfair legal result, disregard of legal principles by arbitrators destroys one of the functions of a legal system—to provide predictability so that contracts can be entered into with an awareness on the part of the signers of how their contracts will be interpreted.

Since both litigation in court and arbitration are time-consuming and costly processes, more and more construction disputes are resolved by alternate dispute resolution techniques such as mediation. Although the suggestions of a mediator are wholly unenforceable (unlike the rulings of an arbitrator), mediation has proved to be a very efficient approach to the resolution of construction industry disputes.

The most effective tool for avoiding disputes remains the carefully drawn contract clearly setting forth the rights and obligations of the parties and anticipating most of the situations that might arise during the course of a construction project.

Tort liability. While contract law enforces the rights of the parties to the contract, tort law enforces the rights and duties of those who are not parties to the contract (that is, in privity). A design professional has an obligation to persons not related by contract, such as the general public, construction workers, or construction contractors (third parties). In general, the design professional has the duty to be free from negligence in the performance of services. Where it is reasonably foreseeable that the negligent performance of services will cause harm to such third

parties, and such parties do suffer harm, the design professional will be liable for the damages suffered. This rule exists in the United States wherever the damage suffered by the third party is personal injury, loss of life, or physical damage to property. In most states the design professional is also liable to third parties if the professional's negligence causes only an economic loss; however, a significant minority of states do not impose liability for economic loss.

It is less clear to which third parties is owed the duty to exercise reasonable care. Courts usually look to the design professional's contracts to determine the class of persons to whom the design professional owes duties.

Standard forms of agreement. Design professionals frequently utilize standard contract forms that are intended to be used (with certain modifications) on a wide variety of projects. The American Institute of Architects (A.I.A.) publishes a wide range of such standard documents, as does the Engineering Joint Contract Documents Committee (EJCDC). The EJCDC, which is composed of representatives from the American Consulting Engineers Council, the American Society of Civil Engineers, the Construction Specifications Institute, and the National Society of Professional Engineers, publishes sample documents covering every phase of the design and construction process. The A.I.A. documents focus on building projects involving architects.

A standard contract is a widely applicable document that, while it covers the subject, does not fit the endeavor precisely, as a custom-tailored document would. Although the majority of the A.I.A. and EJCDC documents provide for additions of project-specific clauses and selecting alternative provisions to accommodate alternative methods of payment, what emerges is still a basic document that has been altered—it is still not custom tailored for its intended use. However, not all consumers can afford custom tailoring, and not all projects warrant such cost. Furthermore, many attorneys are not sufficiently versed in the construction process to produce a good custom-tailored contract, and most design professionals are not sufficiently versed in the law to produce a legally effective contract. Accordingly, the use of a standard document may be a better choice than an inadequate custom-tailored document. The greatest value of the standard documents may be to serve as a check list of what properly belongs in the contract, and to provide certain general clauses that serve to limit the design professionals' liability.

There are certain key provisions of the standard documents, provisions which in some form or other should appear in the design professionals' contracts. These include description and limitation of services, limitations of parties to whom duty is owed, and specific details of the construction contract.

Description and limitation of services. The key provisions in the contract are the establishment of the duties of each party. In the case of the owner, the obvious duties relate to the amount and timing of payment to the design professional (design contract) or to the con-

struction contractor (construction contract). In the design contract, a detailed description of the design professional's duties is essential, especially where there may be others rendering services to the owner in connection with the design (for example, surveyors or soils engineers). During construction there are usually even more parties involved in providing goods and services, and there may be considerable overlap. A number of important concerns must be addressed, such as (1) whether the design professional who has performed the design is responsible for seeing that the construction contractor follows the plans, and if so, whether this duty applies only to the owner or also to the construction contractor; (2) whether the design professional is responsible if the contractor's operations result in injury to workers or others; and (3) whether the design professional is responsible if the method of construction chosen by the contractor is not feasible or unnecessarily costly.

A listing of specific services that the design professional is to perform is often insufficient to clearly define his or her responsibilities. The reason is that an obligation to perform certain other related services might be implicit in performing the listed services. For example, a requirement that the design professional design a building to be used as a diner might be interpreted to also require him or her to perform the property survey, collect and analyze soil samples, design the foundation, specify food-handling equipment, or perform environmental analyses. If it is not the intent of the parties that the design professional perform certain of these services, the contract should explicitly specify what is not to be done as well what is to be done.

Limiting parties to whom duty is owed. The law often looks to the design professional's contracts to determine whether a duty was owed to a third party. If the design contract requires the design professional to see that the construction contractor employs adequate safety precautions, the law will say that the design professional owes a duty to workers and the general public to see that the workplace is safe. If the design contract requires the design professional to inspect the work to guarantee to the owner that the construction conforms to the plans and specifications, the design professional has guaranteed to the owner that all construction will be in accordance with the contract and the design professional will be liable to the owner if the construction contractor has not followed the plans and specifications.

If it is not the intention that the design professional be responsible to third parties for the construction contractor's safety precautions, the design contract should explicitly so state and the design professional should not exercise any control over the contractor's operations. If it is not the intention that the design professional be the guarantor of the construction contractor's performance, the design contract should explicitly say this. Both the A.I.A. and the EJCDC standard documents provide language that clearly limits the engineer's responsibilities in these two areas.

Construction contract provisions. Although the design professional is not a party to the construction contract, he or she is nevertheless affected by that document. First, it is usually the design professional who prepares the construction contract documents, including the plans and specifications that define the construction contractor's duties. Second, the legal system looks to the construction contract should there be a question as to whether a duty is that of the design professional or the construction contractor, or both. For example, if the construction contract requires the contractor to prepare supplementary plans or details and submit them to the design professional for review and approval, it is reasonably clear the design professional has a duty to check the construction contractor's work. If the construction contract makes the construction contractor exclusively responsible for safety precautions and the means and methods of construction, in the absence of specific language also assigning this responsibility to the design professional the courts will not find that the design professional had any responsibilities in these areas. Again, the A.I.A. and EJCDC standard construction contract documents clearly define the relative responsibilities between the design professional and construction contractor.

Subcontracts for professional services. Many design professionals may be involved in a single project. While some owners contract separately with the various design professionals, most prefer to deal with a single professional who will be responsible for the performance of all design services. In the case of a building project the owner contracts with an architect, who in turn is expected to retain the services of structural and mechanical engineers and any other professional services required. The considerations necessary in preparing these subcontracts are essentially identical to the key provisions for the design contract. Both the EJCDC and the A.I.A. publish standard forms of subcontract agreements (but subject to the limitations of any standard contract as discussed earlier). See ARCHITECTURAL ENGINEERING; ENGINEERING.

Donald A. Ostrower

Engineering design

Engineering is concerned with the creation of systems, devices, and processes useful to, and sought by, society. The process by which these goals are achieved is engineering design.

The process can be characterized as a sequence of events as suggested in Fig. 1, with the recognition that no final, universally accepted description of so complex an intellectual and physical exercise, applicable to an enormously broad spectrum of products and processes, is possible. The process may be said to commence upon the recognition of, or the expression of, the need to satisfy some human want or desire, the "goal," which might range from the detection and destruction of incoming ballistic missiles to a minor kitchen appliance or fastener.

Concept formulation. Since the human aspiration is usually couched in nonspecific terms as a sought-for goal, the first obligation of the engineer is to develop more detailed quantitative information which defines the task to be accomplished in order to satisfy the goal, labeled on Fig. 1 as task specification. At this juncture the scope of the problem is defined, and the need for pertinent information is established. The source of the original request is questioned to establish the correspondence between the developing specifications and the initial definition of the goal.

But to know that a need exists and to have started on the task of qualitatively and quantitatively defining its substance and bounds should not be confused with the generation of ideas for possible solutions to the problem. This creative stage is called the concept formulation. When great strides in engineering are made, this represents ingenious, innovative, inventive activity; but even in more pedestrian situations where rational and orderly approaches are possible, the conceptual stage is always present.

The concept does not represent a solution, but only an idea for a solution. Initially evisceral and ephemeral, it can only be described in broad, qualitative, frequently graphical terms. Concepts for possible solutions to engineering challenges arise initially as mental images which are recorded first as sketches or notes and then successively tested, refined, organized, and ultimately documented by using standardized formats.

At this point it is important to note that the necessarily two-dimensional description of Fig. 1 and this sequential textual description should not be construed as an implication that the goal, task specification, and concept always appear in a simple, sequential temporal order. In fact, a central characteristic of the design process is the unpredictable emergence of, and iterations between, the various steps. This stochastic character can be suggested by drawing two-way connections between the various steps. For example, the definition of a task and emergence of a concept might precipitate restudy and possible alteration of the original goal. Consider the not infrequent experience of the serendipitous emergence of a new and interesting concept, and then subsequent

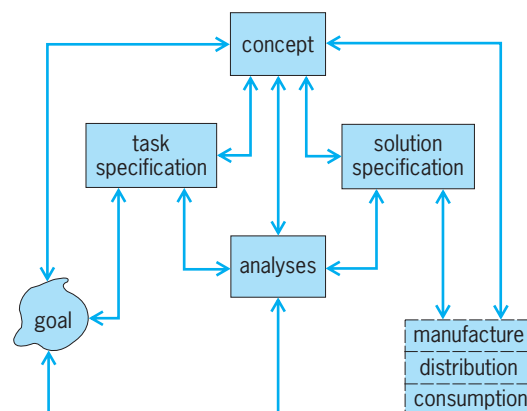


Fig. 1. Engineering design process.

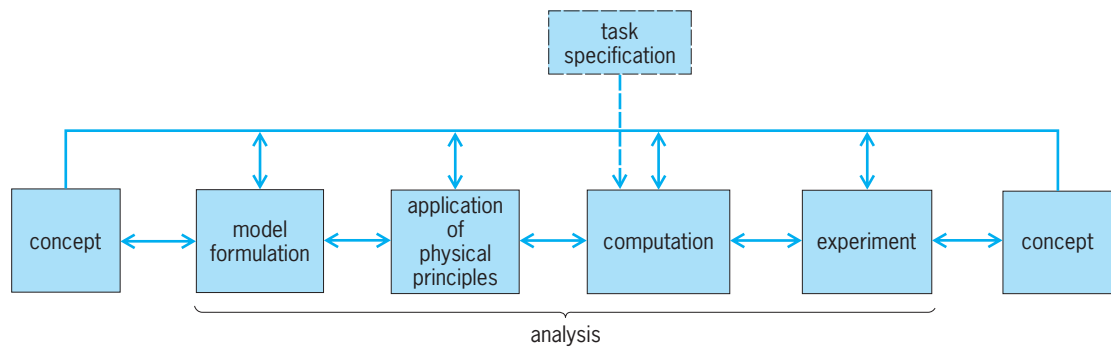


Fig. 2. Role of analysis in concept formulation.

search for a goal (application or market) to which to apply the idea.

Concepts are accompanied and followed by, sometimes preceded by, acts of evaluation, judgment, and decision. It is in fact this testing of ideas against physical, economic, and functional reality that epitomizes engineering's bridge between the art of innovation and science. The process of analysis is sometimes intuitive and qualitative, but it is often mathematical, quantitative, careful, and precise. The iteration of concept and analysis invariably gives rise to a focusing and sharpening, possibly a complete change, of the concept. Frequently, in light of the analyses conducted, what can or might be accomplished is illuminated and task specifications changed. Even goals may be altered. That which was sought may be beyond accomplishment, or perhaps solutions more sophisticated or more useful than initially undertaken may prove achievable.

Production considerations can have a profound influence on the design process, especially when high-volume manufacture is anticipated. Evolutionary products manufactured in large numbers, such as the automobile, are tailored to conform with existing production equipment and techniques such as assembly procedures, interchangeability, scheduling, and quality control. New techniques such as those associated with space exploration, where volume production is not a central concern, factor into the engineering design process in a very different fashion.

Similarly, the design process must anticipate and integrate provisions for distribution, maintenance, and ultimate replacement of products. Well-conceived and executed engineering design will encompass the entire product cycle from definition and conception through realization and demise and will give due consideration to all aspects.

The iterations of analysis and concept and considerations of manufacture develop information which defines a sequence of progressively specific solutions. In final form the solution specification consists of all drawings, materials and parts lists, manufacturing information, and so on necessary for construction of the device, system, or process.

It is important to note that, while the culmination of a particular engineering design process defines but one solution, inspection of the process during

its evolution indicates a complex series of discrete partial, temporary, interim solutions which are compared one with another, and out of this comparison emerges the final compromise solution. The limitations on this scrutiny of alternative approaches are the practical contingencies of time and personnel, or other ways of measuring resources. When these assets reach their budgetary limits, the most satisfactory solution is accepted; however, in as dynamic a field as engineering the final design is not necessarily an ultimate or optimal best.

Figure 1 provides a broad overview of the engineering design process. Some expansion of the interaction between conceptualization and analysis is warranted, as shown in Fig. 2. On the left- and right-hand sides is concept; everything between is characterized as analysis.

Any physical entity, existing or hypothesized, of any degree of complexity cannot be analyzed in its entirety because of inadequate knowledge of the relevant physical laws, or inadequate time or facilities for the required computation, or a combination of these shortcomings. For these reasons, plus its inherent initial vagueness, the concept cannot be analyzed completely. Instead, simplified models are deliberately and precisely defined by applying established physical principles and laws to describe the model via mathematical equations. By using numerical values from the task specifications for parameters, the requisite computational tasks are performed.

In many engineering situations, particularly those where there is no body of experience with similar geometries, materials, and so on, the model from which the analysis is derived cannot be confidently assumed to characterize completely all significant attributes of the ultimate physical system. In some cases certain relevant processes are not completely understood; other times adequate resources to perform all pertinent analyses are not available. Thus, recourse to experiment is necessary. Since nature is the final arbiter of all physical proposals, however analyzed, experiment or test always precedes final acceptance of any proposal.

As previously suggested for Fig. 1, where all possible interconnections occurred between various stages and the process underwent a dynamic evolution, complex interactions occur also in the concept-analysis loop of Fig. 2. For example, the model is

defined, based on the concept recognizing existing physical knowledge about included processes and phenomena, and the effort that will be involved in the reduction of the resulting mathematical equations into quantitative results is anticipated. For example, an initial model might be a greatly simplified, or perhaps oversimplified, characterization of the concept in order that the computational results might quickly identify the utility of the idea or its ranges of applicability. Subsequently, more refined models of aspects of the concept of greater subtlety might be subject to scrutiny.

Where the function of the part is critical or the physical knowledge is inadequate, the engineer might, in fact, eliminate computation and move rapidly to experiment, perhaps on an analog or scaled-down version of important attributes of the concept.

Hierarchy of design. An adequate description of the engineering design process must have both general validity and applicability to a wide variety of engineering situations: tasks simple or complex, small- or large-scale, short-range or far-reaching. That is to say, there is a hierarchy of engineering design situations.

Systems engineering occupies one end of the spectrum. The typical goal is very broad, general, and ambitious, and concepts are concerned with the interrelationships of a variety of subsystems or components which, when taken together, make up the system to accomplish the desired goal. See SYSTEMS ENGINEERING.

At a subsidiary level of the design problem hierarchy, the same engineering design process applies to creation of a device which might be one component of the overall system. And at the most detailed end of this hierarchy the same process diagrams the engineering design of a single element of a component. Obviously, as the engineering design process is applied to create one of these several elements, components, or systems, different phases of the process come into play in different ways and to different degrees, depending upon the particular problem.

As an illustration of this hierarchy of design problems, all of which exhibit the schema of the engineering design process, consider a particular goal, the interplanetary inspection of the atmosphere of Venus by the *Mariner 2* space probe. At the system level this study encompasses considerations of possible Earth-Venus trajectories based on astronomical events, launching times, booster-spacecraft weight-thrust combinations, and so on. Projections of foreseeable booster capabilities balanced against estimates of the spacecraft weight, combined with the astronomical "windows" to Venus which result from Earth-Venus juxtaposition, give rise to detailed analyses of conceivable space paths which pose the need for midcourse maneuver capabilities.

The scientific measurements desired and the assumed Venus atmosphere postulate the functional characteristics of instruments: weight, volume, and power consumption. Control and telemetry requirements to and from the spaceship likewise augment



Fig. 3. Flight configuration of *Mariner 2* spacecraft.

projected weight, volume, and power consumption needs. Combinations of these requirements and others give rise to concepts for the spaceship design, the flight configuration of which is shown in Fig. 3. *Mariner 2* is an interesting example of new design, there being no established precedent as to what a Venus probe should look like. See SPACE PROBE.

Each major component of the spacecraft—structure, internal power supply, solar collecting vanes, radio telemetry antenna, and so on—undergoes a similar evolution from requirements and concepts through analysis. Each major subsystem is composed of subcomponents, each of which is composed of elements; Fig. 4 shows the detailed configuration of the planetary boom. The same process by which the major spacecraft arrangement decisions were made applies again; individual concepts and analyses are permuted in order to arrive at the design configuration of elements, such as this planetary boom. Thus, the same general overall diagrams of Figs. 1 and 2 characterizing the engineering design process can be applied in an ever broadening and encompassing fashion from the most minor



Fig. 4. Planetary boom of the *Mariner 2* spacecraft.

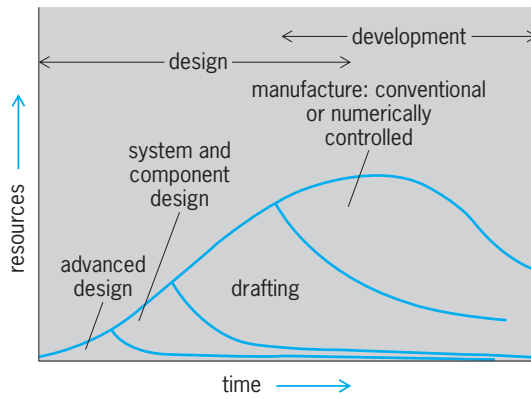


Fig. 5. Elapse of time and resources in an engineering design project showing various stages in sequence.

components to the most major system aspects of an overall design. In point of fact, in an actual situation the order is usually reversed, the process being applied first to the most sweeping question, and progressively to the more detailed aspects, although it is important to understand that the overall design is the sum of all its parts, that each part makes a necessary

contribution to the whole, and that the performance of the whole is dependent on each part.

The project is interesting as an example, not only because of its novel characteristics, adventuresome quest, and dramatic success, but because it also demonstrates the dynamics of the design process. The spaceship weight, of course, depended upon the particular booster rocket designated to bear it into interplanetary space. As the booster itself was under development during the spaceship development, the permissible weight of the Mariner vehicle varied in accordance with the predicted capability of the booster. Shortly before one of the scheduled departure times, major reductions in the vehicle were made mandatory by a change in the booster configuration. But this impermanence and the variability of the “ground rules” and “boundary conditions” of a design situation are always present, although perhaps not as dramatically as in this case.

Another aspect of the Mariner project particularly appropriate to this discussion is the importance of time and schedule to the success of the mission. Trajectory “windows” to Venus are only available at certain discrete, foreseeable times. Either the vehicle-booster combination would be ready to launch at the designated time, or else there would be no possibility of success. Again, most design situations are characterized by this urgency of time and schedule.

Finally, while the government-funded NASA program is perhaps not the best example of cost as the vital common denominator of all industrial and commercial engineering ventures, even in this case there was a budget, not a matter of profit or return on investment, but rather the limited skilled worker-hours available to meet the commitment.

Time-source dynamics. Another dimension of the dynamics of the engineering design process is the elapse of time and expenditure of worker effort in the evolution of an engineering design project. **Figure 5** plots time as the abscissa and resources (worker-hour or dollars) as the ordinate. The various stages of the engineering design process are set out in time sequence from left to right.

Goal refinement, task specification, and first-order concept and analyses iterations are conducted by one to a few engineers in the early stages to establish the feasibility of the idea and to block out possible approaches. This is called the advanced design stage.

As the design concept becomes more specific and substantive, more and more engineers, technicians, and draftsmen become involved in the project. For example, in the case of a modern aircraft probably fewer than a dozen extremely talented engineers carry through the early feasibility and configuration studies, while at a later stage hundreds of engineers, drafters, expeditors, and coordinating personnel are involved. This deployment aspect of the design process cannot be overemphasized. In projects of significant size, the problem of coordinating and integrating the efforts of the many participants of different talents and skills becomes itself a major consideration. While some of the coordination involves judgment and decisions, much of this coordination is

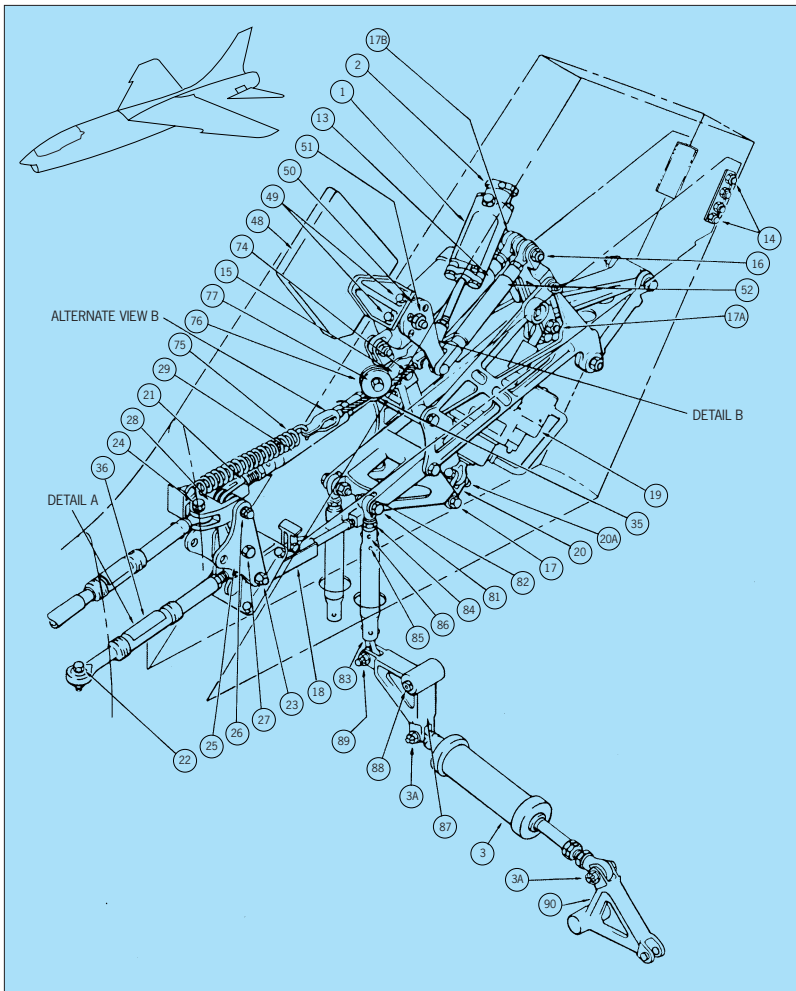


Fig. 6. One of the 10,000 sheets involved in the design engineering of a military aircraft. Numbers refer to the various components and are an indication of the tremendous bookkeeping involved in a project of this type.

purely clerical and some involves prosaic application of standard reference material. See PERT.

Ultimately the process culminates crescendolike in the solution specification—manufacturing drawings and specifications, parts lists, and so on. Where automated manufacturing techniques are warranted, the efforts of computer programmers are devoted to the transformation of graphic and numeric manufacturing information into a form decipherable by computers and machine-tool directors. Ultimately, the physical parts are realized through the manufacturing process.

The efforts of this large group of people demand a high level of coordination and integration if each of the thousands of separate parts of the aircraft are to satisfy its function, be compatible with one another, and be adequately strong and yet of minimum weight and volume. **Figure 6** suggests a few of the individual elements which must undergo design, one sheet of a 2-ft-high (0.67-m) pile of similar diagrams for a typical jet aircraft. The parts shown represent only a few of the better than 10,000 parts that make up the typical jet aircraft.

Use of the computer in design. Reference to the role of analysis in Fig. 1 and to computation in Fig. 2 highlights the ever increasing use of the computer, both analog and digital, in the engineering design process. As a high-powered successor to the slide rule and desk calculator, the computer is used routinely for much of the calculation, computation, and data reduction which constitute a major activity in design. Where repetitive series of calculations are carried out, programs are prepared to delegate to the computer more and more of the responsibility for analysis. Where economically justified, the overall engineering design process for a product is mechanized via computer programming. For example, **Fig. 7** describes the flow chart of a branching computer program which, given vital data on the requirements of an electrical power transformer, carries out the design automatically. In the process the program, having made assumptions on core size and copper windings, calculates temperature gradient, impedance, and so on; checks these against pre-programmed constraints; and optimizes its choices considering, among other things, the current cost of copper. **Figure 8** illustrates typical iterations of the computer design as it “homes” in on the lowest weight and cost design satisfying the initial data and the internal constraints. See COMPUTER; COMPUTER PROGRAMMING.

The speed, memory, and accuracy of the computer to iteratively calculate, store, sort, collate, and tabulate have greatly enhanced its use in design and encouraged the study, on their own merits, of the processes and subprocesses in the design process. These include optimization or sensitivity analysis, reliability analysis, and simulation as well as design theory. See DIGITAL COMPUTER.

Optimization analysis, given a model of the design and using linear and nonlinear programming, determines the best values of the parameters consistent with stated criteria and further studies the effects of

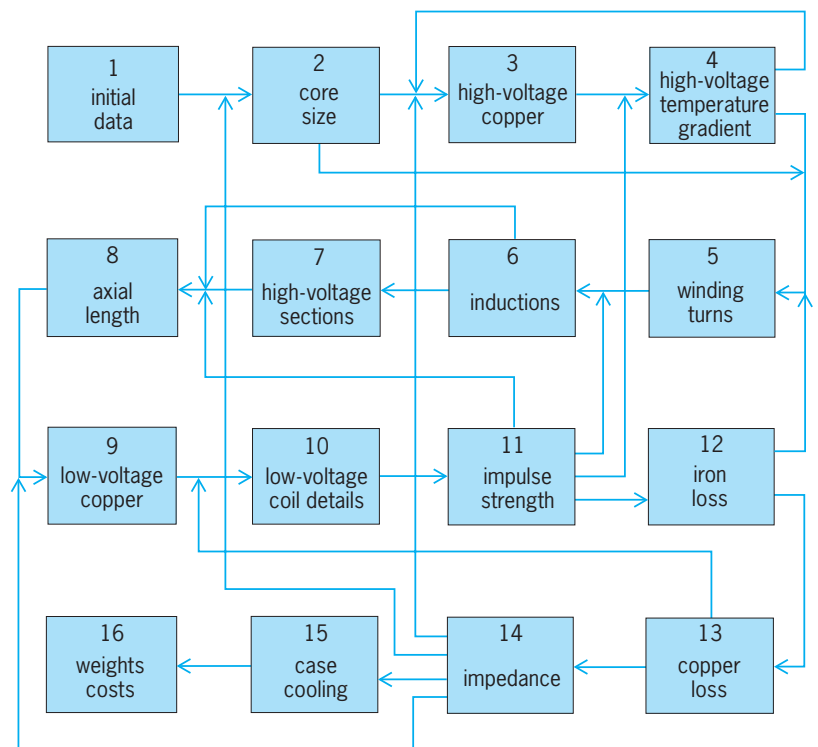


Fig. 7. Flow chart of a computer program.

variations in the values of the parameters. See OPTIMIZATION.

Reliability is a special case of optimization where the emphasis is to choose or evaluate a system so as to maximize its probability of successful operation, for example, the reliability of electronics. See CIRCUIT (ELECTRONICS); RELIABILITY, AVAILABILITY, AND MAINTAINABILITY.

Simulation as of dynamic systems, is mathematical modeling to study the response of a design to various inputs and disturbances. The analog computer was formerly widely used for simulation through its physical modeling of the mathematic analytical relationships of the proposed design. The digital computer's use of numerical data adapts it more readily to nonlinear or probabilistic situations by using random or Monte Carlo techniques, as well as to those situations requiring higher accuracy. See ANALOG COMPUTER; MONTE CARLO METHOD; SIMULATION.

Decision theory deals with the general question of how to choose between a great number of alternatives according to established criteria. It proposes models of the decision process as well as defining techniques, that is, programs or algorithms, of calculation by which to make choices. See DECISION THEORY.

Artificial intelligence or its less pretentious subdiscipline, expert systems, strives to identify and codify human thought processes and produce computer algorithms and programs for automatic problem solving, including facets of the design process. See ALGORITHM; ARTIFICIAL INTELLIGENCE; EXPERT SYSTEMS.

Although each of these techniques represents either a gross oversimplification or but a part of the

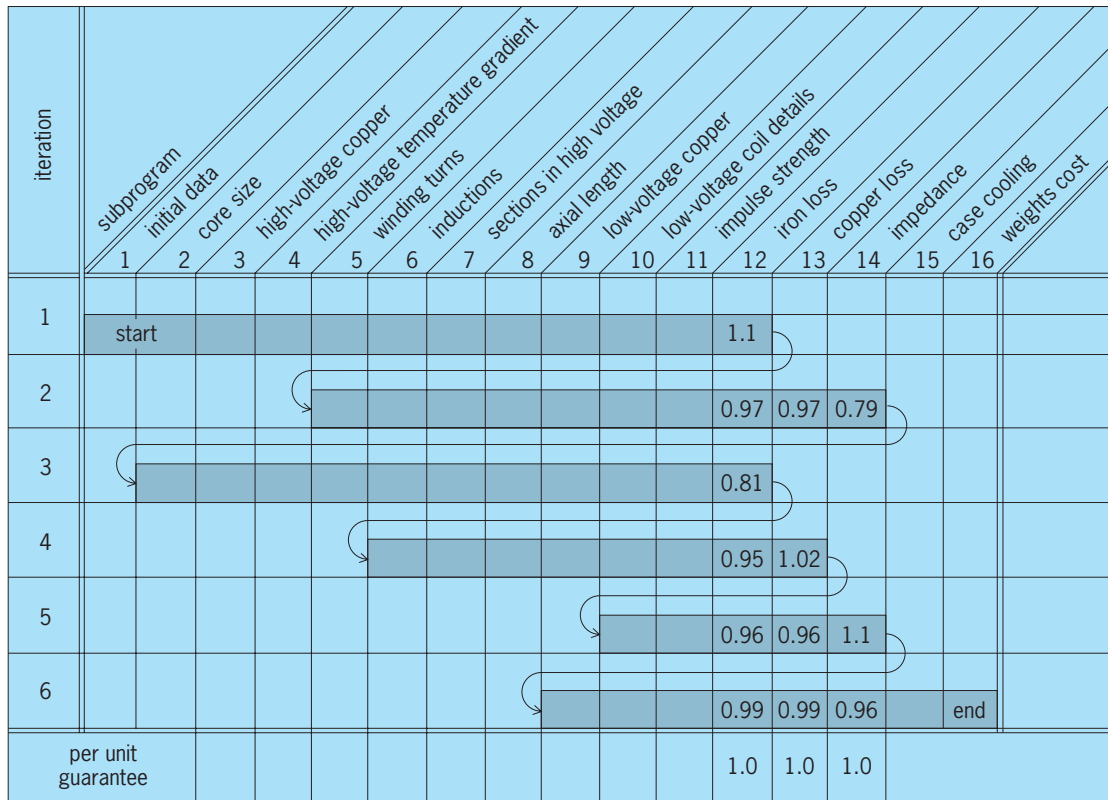


Fig. 8. Iterations of the computer program shown by the flow chart in Fig. 7.

overall design process of Fig. 1, their study and application are useful. Even where not strictly applicable, the awareness of a model of a process enhances its qualitative evaluation by providing a guide to action. The utility and comprehensiveness of these techniques are already advancing and will mature rapidly, largely because of the ability of the high-speed digital computer to carry out evaluations which heretofore were impractical. Finally, the record-keeping capacity of the computer makes possible the recording and reutilization of past action. This integration of prior results into new endeavors represents a “mechanization of experience.”

Computer utilization. Notwithstanding the present contribution and unlimited promise of the computer to the design process, there are certain limitations. In reflecting on the model of the engineering design process and especially on the interaction between concept and analysis, it is useful to consider the nature of the activities involved. Concepts, or really new ideas, usually do not come about as a result of any ponderous, systematic, organized effort; rather they appear in an undisciplined, creative, spectacular fashion. The short-duration exultation of concept creation is then followed by usually very much longer, more regular, systematic periods of analysis. The typical use of a computer does not lend itself to the kind of person-analysis interaction implicit in the engineering design process and essential to its successful negotiation. Usually, the solution process for a piece of work must first be well thought out and known in the most intricate detail; only then can the user write a program to describe to the com-

puter the sequence of steps necessary. But when the nature of the problem is only vaguely grasped and much learning has to be done, the computer may not be as helpful as pencil and paper.

The reason is that in the thinking process one needs to advance in steps and to test these steps frequently. While the steps may be large, as in the first gross examination of an idea, or small, as in an exercise of refinement, these tests need not be elaborate or even precise. What is wanted first is only a qualitative result or perhaps a quantitative result of only moderate precision, a confirmation or a denial of tentative guess work.

By conventional computer programming, it is virtually impossible to obtain quick answers to small discrete problems, even though the computer can work very rapidly. To solve a problem, the user must first prepare a program, a detailed ritual of calculations and comparisons of quantities for the computer to carry out. It is hoped that data so processed will yield the solution.

If there is a mistake in the program, the computer either will detect the mistake and refuse to waste time by trying to run the program, or if the mistake is more subtle, perhaps an error in principle, will go ahead and compute a mass of meaningless nonsense. This dependence on programming, waiting, and error finding and correcting frustrates the formulation and testing of concepts essential to a creative effort.

Another point is that the creative process is, virtually by definition, unpredictable. The sequence of the steps is never known at the beginning. If

it were, the whole process could be accomplished by the computer since the information prerequisite to the computer program would be available. Indeed, the creative process is the process of learning how to accomplish the desired result.

Clearly what is needed if the computer is to be of greater use in the creative process is a more intimate and continuous interchange between engineer and machine. This interchange must be of such a nature that all forms of thought congenial to humans, whether verbal, symbolic, numerical, or graphical, are also understood and acted upon by the machine in ways appropriate to the user's purpose. In reflecting on this human-machine symbiosis with the goal of designing a system with which to design, it is useful to identify the special attributes of each partner. See HUMAN-MACHINE SYSTEMS.

The speed of the computer is prodigious when measured against any term of reference and especially so when measured against the appropriate one, the human mind. A computer performs millions of arithmetic or logical operations per second. This rate outstrips by many orders of magnitude the human neurological responses to premeditated cerebrations.

The machine's memory is likewise extraordinary, although restricted to an extremely narrow but nonetheless useful class of memory impression. It can store binary numbers of 32 digits and recall them unerringly at microsecond speeds, while the human memory system is very much slower and not nearly so reliable. Of course, in contrast, the human memory encompasses an indescribable universe of types and kinds, far beyond the range of a computer. The reliability routinely exercised by the computer is unattainable by people. Once properly instructed, the computer executes its routines faithfully, untiringly, yielding to no boredom or carelessness, and introducing no humanlike errors. See COMPUTER STORAGE TECHNOLOGY.

Counterposed against the computer's assets stands an open-ended list of human attributes. In terms of the engineering design process, these human characteristics would include reflection on, and evaluation of, the social, esthetic, and economic aspects of the original goals; the formulation of previously unforeseen and unanticipated questions at many points in the investigation; an unflagging curiosity about the way things are done and about the way they are proposed to be done; the flexibility of mind to shift from one approach to another, to sort out the significant from a great mass of information and misinformation on the basis of very few discernible criteria; powers of mental association which can detect, correlate, and amplify useful relevancies between bodies of normally disparate information; the devising and structuring of new and unusual approaches; the willingness, ability, and intellectual integrity to make decisions, frequently on the basis of inadequate data; the ability to recognize that the job must get on and some decisions must be made rapidly, though the right to reverse decisions whenever facts so indicate is reserved; exer-

cise of judgment based on prior experience in making decisions; and finally the desire and willingness to develop such adjudicative ability through experience. While the computer demonstrates the capability of encroaching upon some of these attributes, but only to the extent that humans thoroughly understand them and can program them, clearly others are forever denied it. See COMPUTER-AIDED DESIGN AND MANUFACTURING; COMPUTER-AIDED ENGINEERING.

Time-sharing. Communicating directly with a computer by keyboard with the computer replying immediately on a cathode-ray-tube monitor removes the obstacle to thought caused by waiting hours for a reply. The user has no compunctions about posing smaller problem fragments to the computer, in a way more consonant with the small steps of the creative process. The results are immediately available, and the course of further action is guided more precisely.

Such intimate "conversations" with a computer are scarcely economically feasible when only one person uses the computer at a time because, measured against the speed of a computer, human beings are intolerably slow. A person likes to ponder a problem and often needs time to decide what to try next. Even when a decision is reached, the instruction to the computer takes several seconds to type on the keyboard, and during even this short time the computer could perform hundreds of thousands to millions of calculations were it free to do so. It is logical to arrange to have large main-frame computers in conversation with a large number of people, so that when one user is idle the computer can turn to another and answer that question. Should all users be idle at once, the computer can work on some large problem left as a backlog. This type of effective computer utilization is quite routine, and is called time-sharing. On a rough average, if a user is intensely busy for 1 h, the computer can discharge its responsibilities to this particular user in short bursts that total about 3 min. The rest of the time it serves the other users. See MULTIACCESS COMPUTER.

The dramatic reductions in cost and increases in power of contemporary minicomputers have expanded their use in time-share environments. The advent of the microprocessor, or computer on a chip, made economic the personal computer where only one keyboard operator can afford to monopolize a whole computer system. The power of such stand-alone systems is as yet still very limited for design application; therefore, microcomputers are currently used, together with their keyboards and graphic displays, as "intelligent" terminals on larger minicomputers or mainframe computers, providing the operator access to the computing power and vast memory of the larger machine. See MICROCOMPUTER; MICROPROCESSOR.

Graphical input/output. In many fields of design—notably architecture; design of airplanes, automobiles, and ships; consumer products; in almost all mechanical design and in electronic computer circuit design—the designer works largely in visual terms.

Initially limited to interpreting only binary digits, the digital computer was subsequently organized to understand ordinary numbers and words and combinations thereof via the various programming languages. Then, in the early 1960s Ivan Sutherland established bilateral graphical communication between the human and the digital computer. The essential link between the pictorial cognitive style of the design engineer and the speed, memory, and reliability of the digital computer had been forged, but there was a lag of more than a decade before the availability and economy of adequate main-frames and improved graphic interfaces made computer-aided graphics a practical reality. The computer has become an active partner in the act of drawing, so that it can provide a certain superskill in preparing the drawing once the intentions of the human operator have been made clear. *See* COMPUTER GRAPHICS.

Computer-aided design (CAD). The engineering design process, as aided by CAD capability, can be suggested via several examples.

Selecting from a task menu presented on a graphics console, a designer, using keyboard and cursor or light-pen, can then draw on the console and into the computer the elemental parts of a machine; have the computer embody them as solids; assemble, section, and orient them under operator control; and finally document the parts following traditional drafting and specification standards, all electronically. The production engineer, with access to the common data base, can then invoke computer-aided manufacturing (CAM), optimize the machine tool cutter paths, with the coded information transferred to the actual numerically controlled (N/C) machine tool for part fabrication.

Comprehensive CAD systems provide for human-computer interaction in engineering analysis of the evolving design. A recurrent question in the concept-analysis iteration of Fig. 2 is: Will the strength be adequate? This is addressed in CAD, with a computer graphics portrayal of a hypothetical wrench turning a bolt. Using a stress-analysis technique known as the finite element method, the computer automatically superimposes a mesh into the region the engineer identifies as a potentially highly stressed region in the wrench. Then, having been provided the wrench material strength properties and the applied torque, the computer's finite element method program calculates, and the CAD graphics presents, coded in color, the resulting stress distribution. The engineer must then evaluate these data, accept the results, or intervene and alter the geometry, material, acceptable torque, and so on. The essence of the engineering design process has been retained in CAD, but the process has been vastly facilitated and accelerated. More iterations of more alternatives can be explored more rigorously and more rapidly when compared with traditional engineering design. The results are faster and better design and improved productivity. *See* FINITE ELEMENT METHOD.

Use of CAD (and CAM) will relegate increasingly to the computer tasks it can be programmed to perform automatically, leaving to the human the exercise of

goal definition and setting and judgmental evaluation and decision. In one integrated system, progressive solutions of an automatic part design and analysis are produced by a computer program which, given functional geometrical constraints on a part and its structural strength requirements, automatically iterates geometrical designs and finite-element-method analyses to produce the minimum-mass part which satisfies the constraints.

Computer effectiveness. The effect of using computers in the design and development stages of a manufacturing process is twofold. First, the costs are greatly reduced because fewer people are involved in the stages, and less time is required to accomplish the tasks. For example, the translation of a designer's sketch of the shape of an airplane fuselage into a full-size, precise, geometrical shape description suitable for manufacturing purposes using traditional methods used to take about 6 months. Conventional, that is, non-CAD, computer methods reduce this time to a few weeks. But the emergent computer techniques will cause a further reduction of this time, perhaps to seconds. The results will be far more accurate than have previously been obtained and will be tailored much more closely to manufacturing requirements.

Second, an important hidden benefit is the greatly increased lead time which a CAD/CAM manufacturer gains over a noncomputerized (or less computerized) competitor; the former can be ready much sooner to begin production. The national military and the commercial aspects of this acceleration are obvious.

Robert W. Mann

Bibliography. J. L. Adams, *Conceptual Blockbusting*, 3d ed., 1990; G. L. Glegg, *The Design of Design*, 1969; M. P. Groover and E. W. Zimmers, Jr., *CAD/CAM: Computer-Aided Design and Manufacturing*, 1984; F. H. Jones, *Computer Aided Architecture and Design*, 1994; J. E. Shigley and L. D. Mitchell, *Mechanical Engineering Design*, 5th ed., 1989; M. F. Spotts, *Design of Machine Elements*, 7th ed., 1997.

Engineering drawing

A graphical language used by engineers and other technical personnel associated with the engineering profession. The purpose of engineering drawing is to convey graphically the ideas and information necessary for the construction or analysis of machines, structures, or systems.

In colleges and universities, engineering drawing is usually treated in courses with titles like Engineering Graphics. Sometimes these courses include other topics, such as computer graphics and nomography. *See* COMPUTER GRAPHICS.

The basis for much engineering drawing is orthographic representation (projection). Objects are depicted by front, top, side, auxiliary, or oblique views, or combinations of these. The complexity of an object determines the number of views shown. At times, pictorial views are shown. *See* DESCRIPTIVE GEOMETRY; PICTORIAL DRAWING.

Engineering drawings often include such features as various types of lines, dimensions, lettered notes, sectional views, and symbols. They may be in the form of carefully planned and checked mechanical drawings, or they may be freehand sketches. Usually a sketch precedes the mechanical drawing. Final drawings are usually made on tracing paper, cloth, or Mylar film, so that many copies can be made quickly and cheaply by such processes as blueprinting, ammonia-developed (diaz) printing, or lithography. See PHOTOCOPYING PROCESSES.

Section drawings. Many objects have complicated interior details which cannot be clearly shown by means of front, top, side, or pictorial views. Section views enable the engineer or detailer to show the interior detail in such cases. Features of section drawings are cutting-plane symbols, which show where imaginary cutting planes are passed to produce the sections, and section-lining (sometimes called cross-hatching), which appears in the section view on all portions that have been in contact with the cutting plane. When only a part of the object is to be shown in section, conventional representation such as revolved, rotated, or broken-out section is used

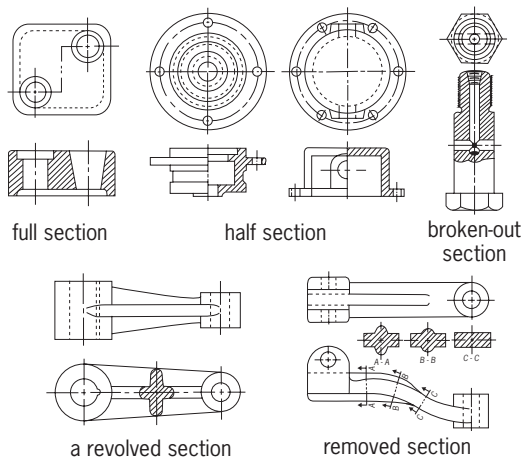


Fig. 1. Section views. (After T. E. French and C. J. Vierck, Engineering Drawing, McGraw-Hill, 4th ed., 1978)

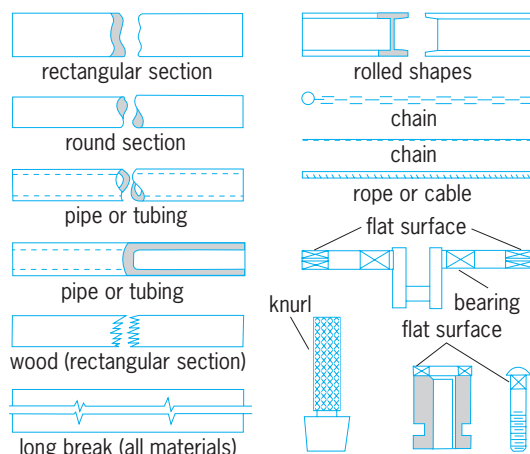


Fig. 2. Conventional breaks and other symbols to indicate details. (After T. E. French and C. J. Vierck, Engineering Drawing, McGraw-Hill, 4th ed., 1978)

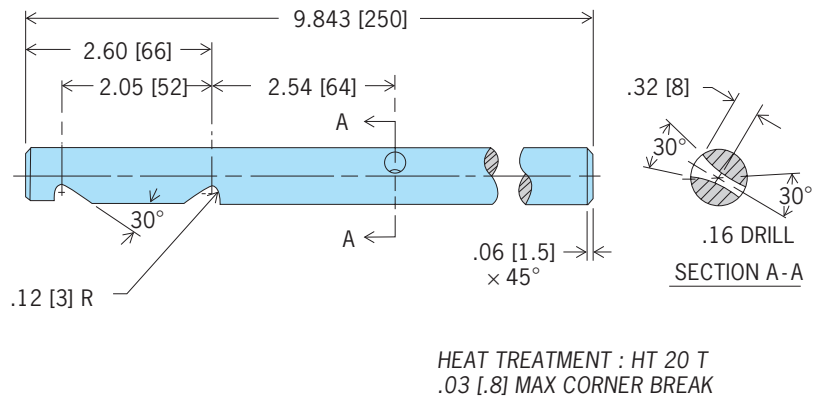


Fig. 3. Machine part which has been dimensioned in inches and millimeters (brackets).

(Fig. 1). Details such as flat surfaces, knurls, and threads are treated conventionally, which facilitates the making and reading of engineering drawings by experienced personnel (Fig. 2). Thus, certain engineering drawings will be combinations of top and front views, section and rotated views, and partial or pictorial views.

Dimensioning. In addition to describing the shape of objects, many drawings must show dimensions, so that workers can build the structure or fabricate parts that will fit together. This is accomplished by placing the required values (measurements) along dimension lines (usually outside the outlines of the object) and by giving additional information in the form of notes which are referenced to the parts in question by angled lines called leaders. In drawings of large structures the major unit is the foot, and in drawings of small objects the unit is the inch. A drawing containing dual dimensioning, inches and millimeters (in brackets), is shown in Fig. 3. In metric dimensioning, the basis unit may be the meter, the centimeter, or the millimeter, depending upon the size of the object or structure.

Working types of drawings may differ in styles of dimensioning, lettering (inclined lowercase, vertical uppercase, and so on), positioning of the numbers (aligned, or unidirectional—a style in which all numbers are lettered horizontally), and in the type of fraction used (common fractions or decimal fractions). If special precision is required, an upper and lower allowable limit are shown. Such tolerance, or limit, dimensioning is necessary for the manufacture of interchangeable mating parts, but unnecessarily close tolerances are very expensive.

Layout drawing. Layout drawings of different types are used in different manufacturing fields for various purposes. One is the plant layout drawing, in which the outline of the building, work areas, aisles, and individual items of equipment are all drawn to scale (Fig. 4). Another type is the aircraft, or master, layout, which is drawn on glass cloth or on steel or aluminum sheets. The object is drawn to full size with extreme accuracy. The completed drawing is photographed with great precision, and a glass negative made. From this negative, photo templates are made on photosensitized metal in various sizes and for

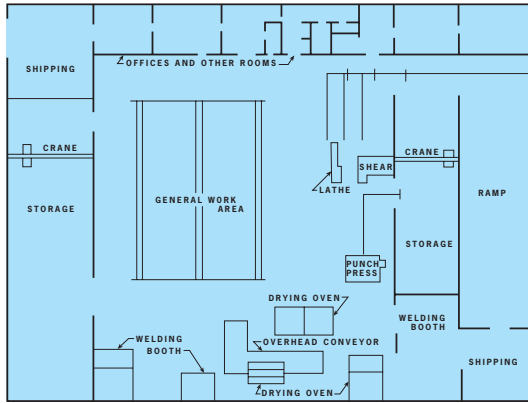


Fig. 4. A plant-layout drawing. (After F. Zozzora, *Engineering Drawing, McGraw-Hill, 2d ed., 1958*)

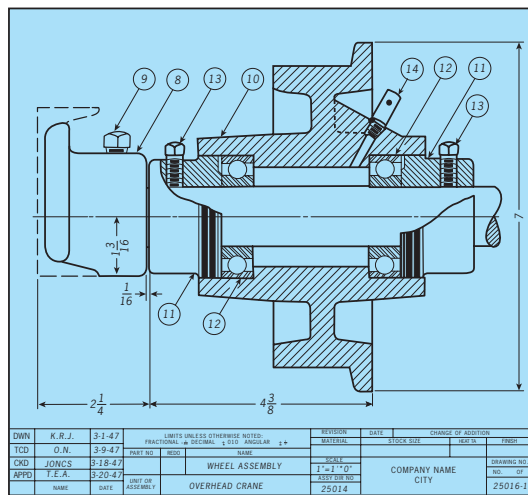


Fig. 5. A unit, or general, assembly. (After T. E. French and C. J. Vierck, *Engineering Drawing, McGraw-Hill, 4th ed., 1978*)

different purposes, thereby eliminating the need for many conventional detail drawings. Another type of layout, or preliminary assembly, drawing is the design layout, which establishes the position and clearance of parts of an assembly.

Assembly drawings. A set of working drawings usually includes detail drawings of all parts and an assembly drawing of the complete unit. Assembly drawings vary somewhat in character according to their use, as design assemblies or layouts; working drawing assemblies; general assemblies; installation assemblies; and check assemblies. A typical general assembly may include judicious use of sectioning and identification of each part with a numbered balloon (Fig. 5). Accompanying such a drawing is a parts list, in which each part is listed by number and briefly described; the number of pieces required is stated and other pertinent information given. Parts lists are best placed on separate sheets and typewritten to avoid time-consuming and costly hand lettering.

Schematic drawings. Schematic or diagrammatic drawings make use of standard symbols which indicate the direction of flow. In piping and electrical schematic diagrams, symbols recommended by the

American National Standards Institute (ANSI), other agencies, or the Department of Defense (DOD) are used. In contrast to Fig. 6, the fixtures or components are not labeled in most schematics because the readers usually know what the symbols represent.

Additional information is often lettered on schematic drawings, for example, the identification of each replaceable electrical component. Etched-circuit drawing has revolutionized the wiring of electronic components. By means of such drawing, the wiring of an electronic circuit is photographed on a copper-clad board, and unwanted areas are etched away. On electrical and other types of flow diagrams,

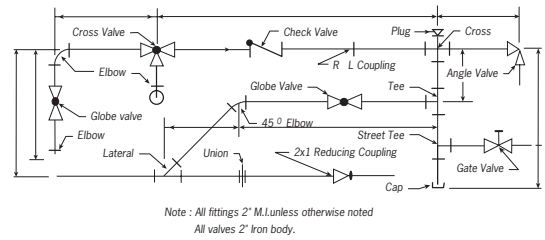


Fig. 6. Piping diagrammatic drawing. (After T. E. French and C. J. Vierck, *Engineering Drawing, McGraw-Hill, 1953*)

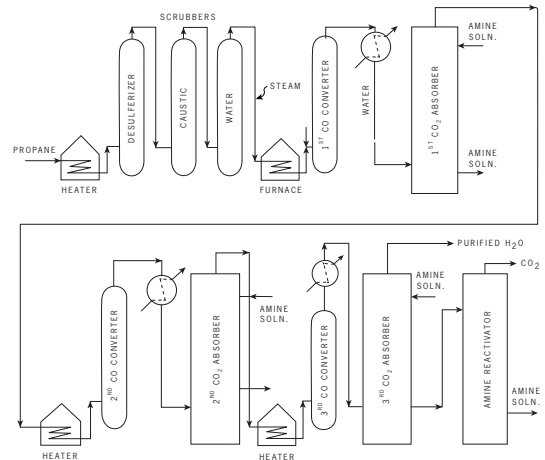


Fig. 7. Chemical engineering flow diagram. (After P. H. Groggins, ed., *Unit Processes in Organic Synthesis, 5th ed., McGraw-Hill, 1958*)

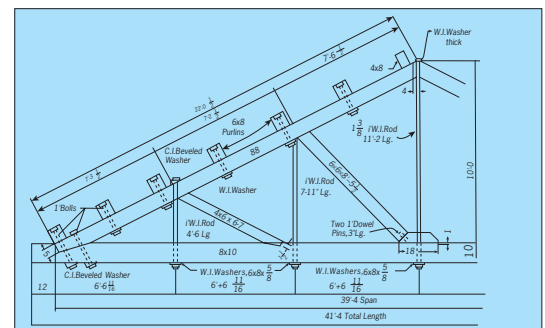


Fig. 8. Structural drawing of wood truss. A drawing of a steel truss would look similar, but would include symbols for such structural shapes as angles and channels. (After T. E. French and C. J. Vierck, *Engineering Drawing, McGraw-Hill, 4th ed., 1978*)

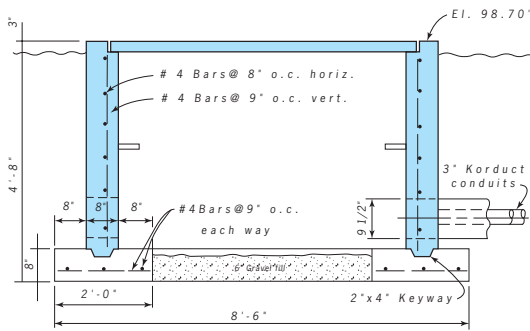


Fig. 9. Drawing of a reinforced concrete structure.

all single lines (often with arrows showing direction of flow) are drawn horizontally or vertically (Fig. 7); there are few exceptions. In some flow diagrams, rectangular enclosures are used for all items. Lettering is usually placed within the enclosures. See SCHEMATIC DRAWING; WIRING DIAGRAM.

Structural drawings. Structural drawings include design and working drawings for structures such as building, bridges, dams, tanks, and highways (Figs. 8 and 9). Such drawings form the basis of legal contracts. Structural drawings embody the same principles as do other engineering drawings, but use terminology and dimensioning techniques different from those shown in previous illustrations. See DRAFTING; GRAPHIC METHODS; NOMOGRAPH.

Charles J. Baer

Bibliography. C. J. Baer and J. R. Ottaway, *Electrical and Electronics Drawing*, 5th ed., 1986; T. E. French et al., *Mechanical Drawing*, 14th ed., 1993; T. E. French and C. J. Vierck, *Engineering Drawing and Graphic Technology*, 14th ed., 1993; C. Jensen and J. Helsel, *Engineering Drawing and Design*, 5th ed., 1997; J. R. Wirshing and R. H. Wirshing, *Civil Engineering Drafting*, 1983.

Engineering geology

The application of education and experience in geology and other geosciences to solve geological problems posed by civil engineering works. The branches of the geosciences most applicable are surficial geology, petrofabrics, rock and soil mechanics, geohydrology, and geophysics, particularly exploration geophysics and earthquake seismology. This article discusses some of the practical aspects of engineering geology.

Terminology. The terms engineering geology and environmental geology often seem to be used interchangeably. Specifically, environmental geology is the application of engineering geology in the solution of urban problems; in the prediction and mitigation of natural hazards such as earthquakes, landslides, and subsidence; and in solving problems inherent in disposal of dangerous wastes and in reclaiming mined lands.

Another relevant term is geotechnics, the combination of pertinent geoscience elements with civil engineering elements to formulate the civil engineer-

ing system that has the optimal interaction with the natural environment.

Engineering properties of rock. The civil engineer and the engineering geologist consider most hard and compact natural materials of the earth crust as rock, and their derivatives, formed mostly by weathering processes, as soil. A number of useful soil classification systems exist. Because of the lack of a rock classification system suitable for civil engineering purposes, most engineering geology reports use generic classification systems modified by appropriate rock-property adjectives. See ROCK; ROCK MECHANICS; SOIL MECHANICS.

Rock sampling. The properties of a rock element can be determined by tests on cores obtained from boreholes. These holes are made by one or a combination of the following basic types of drills: the rotary or core drill, the cable-tool or churn drill, and the auger. The rotary type generally is used to obtain rock cores. The rotary rig has a motor or engine (gasoline, diesel, electric, or compressed air) that drives a drill head that rotates a drill rod (a thick-walled hollow pipe) fastened to a core barrel with a bit at its end. Downward pressure on the bit is created by hydraulic pressure in the drill head. Water or air is used to remove the rock that is comminuted (chipped or ground) by the diamonds or hard-metal alloy used to face the bit. The core barrel may be in one piece or have one or two inner metal tubes to facilitate recovery of soft or badly broken rock (double-tube and triple-tube core barrels). The churn-type drill may be used to extend the hole through the soil overlying the rock, to chop through boulders, occasionally to deepen a hole in rock when core is not required or to obtain drive samples of the overburden soils. When the rock is too broken to support itself, casing (steel pipe) is driven or drilled through the broken zone. Drill rigs range in size from those mounted on the rear of large multiwheel trucks to small, portable ones that can be packed to the investigation site on a person's back or parachuted from a small plane. See DRILLING, GEOTECHNICAL.

The rock properties most useful to the engineering geologist are compressive and triaxial shear strengths, permeability, Young's modulus of elasticity, erodability under water action, and density (in pounds per cubic foot, or pcf).

Compressive strength. The compressive (crushing) strength of rock generally is measured in pounds per square inch or kilograms per square centimeter. It is the amount of stress required to fracture a sample unconfined on the sides and loaded on the ends (Fig. 1). If the load P of 40,000 lb is applied to a sample with a diameter of 2 in. (3.14 in.^2), the compressive stress is $40,000 \div 3.14 = 12,738 \text{ lb/in.}^2$ ($177,920 \text{ N} \div 0.00203 \text{ m}^2 = 87,645 \text{ kN/m}^2$). If this load breaks the sample, the ultimate compressive strength equals the compressive stress acting at the moment of failure, in this case $12,738 \text{ lb/in.}^2$. The test samples generally are cylindrical rock cores that have a length-to-diameter ratio (L/D) of about 2. The wide variety of classification systems used for rock results in a wide variation in compressive strengths for rocks

Compressive strength of rocks*

Type of rock	No. of tests	Weighted mean, kN/m ² †	Minimum and maximum values, kN/m ² †
Igneous			
Basalt	195	112,390	4,140–383,360
Diabase	6	150,655	4,140–275,100
Diorite	26	167,550	78,600–310,275
Gabbro	10	245,875	110,000–399,220
Granite	140	160,515	17,925–332,340
Rhyolite	44	162,100	11,720–453,690
Sedimentary			
Dolomite	62	87,840	6,205–358,540
Limestone	211	74,190	1,380–260,630
Sandstone	257	63,430	2,070–328,200
Shale	67	66,605	690–230,980
Siltstone	14	108,525	3,445–315,790
Metamorphic			
Amphibolite	14	152,240	24,820–280,625
Gneiss	103	133,760	35,850–292,350
Marble	31	100,600	30,340–273,730
Quartzite	25	292,350	25,510–628,825
Schist	16	50,265	6,895–162,030

*Based on data in W. R. Judd, *Statistical Relationships for Certain Rock Properties*, DDC AD735376, 1971, and *Statistical Methods to Compile and Correlate Rock Properties*, Purdue Res. Tech. Rep. 2, 1969.

†1 kN/m² = 0.1450 lb/in.²

having the same geologic name. The **table** gives a statistical evaluation of the compressive strengths of several rocks commonly encountered in engineering geology.

Most laboratory tests show that an increase in moisture in rock causes a decrease in its compressive strength and elastic modulus; what is not generally known, however, is that the reverse situation shown in **Fig. 2** has been encountered in certain types of volcanic rocks. In sedimentary rock the compressive strength is strongly dependent upon the quality of the cement that bonds the mineral grains together (for example, clay cement gives low strength) and upon the quantity of cement (a rock may have only a small amount of cement, and despite a strong bond between the grains, the strength is directly related to the inherent strength of the grains). Strength test

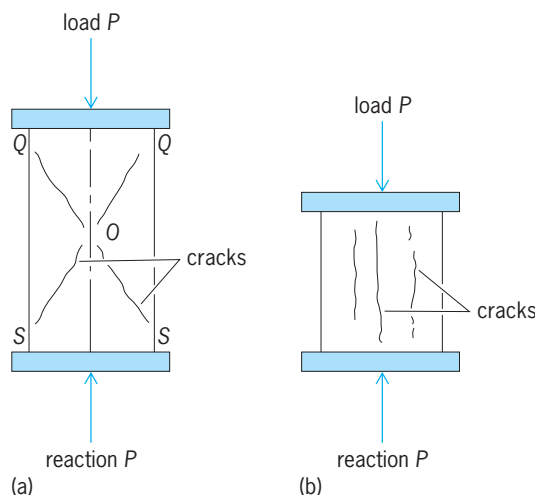


Fig. 1. Unconfined compression test. (a) Shear failure, showing failure planes QS. (b) Tension failure. (After D. P. Krynine and W. R. Judd, *Principles of Engineering Geology and Geotechnics*, McGraw-Hill, 1957)

results are adversely affected by microfractures that may be present in the sample prior to testing, particularly if the microfractures are oriented parallel to the potential failure planes.

The value of compressive strength to be used in an engineering design must be related to the direction of the structure's load and the orientation of the bedding, discontinuities, and structural weaknesses in the foundation rock. This relationship is important because the highest compressive strength usually is obtained when the compressive stress is normal to the bedding. Conversely, the highest Young's modulus of elasticity (E) usually results when the compressive stress parallels the bedding. When these strength and elastic properties apparently are not affected by the direction of applied load, the rock is described as isotropic; if load applied parallel to the bedding provides physical property data that are significantly different than those obtained when the load is applied normal (perpendicular) to the bedding, the rock is anisotropic or aeolotropic. If the physical components of the rock element or rock system have equal dimensions and equal fabric relationships, the rock is homogeneous; significant variance in these relationships results in a heterogeneous rock. Most rocks encountered in foundations and underground works are anisotropic and heterogeneous.

Shear in rocks. Shearing stresses tend to separate portions of the rock (or soil) mass. Faults and folds are examples of shear failures in nature. In engineering structures, every compression is accompanied by shear stresses. For example, an arch dam compresses the abutment rock and, if the latter is intersected by fissures or weak zones, it may fail in shear with a resulting tensile stress in the dam concrete that may rupture the concrete. The application of loads over long periods of time on most rocks will cause them to creep or even to flow like a dense fluid (plastic flow). See STRUCTURAL GEOLOGY.

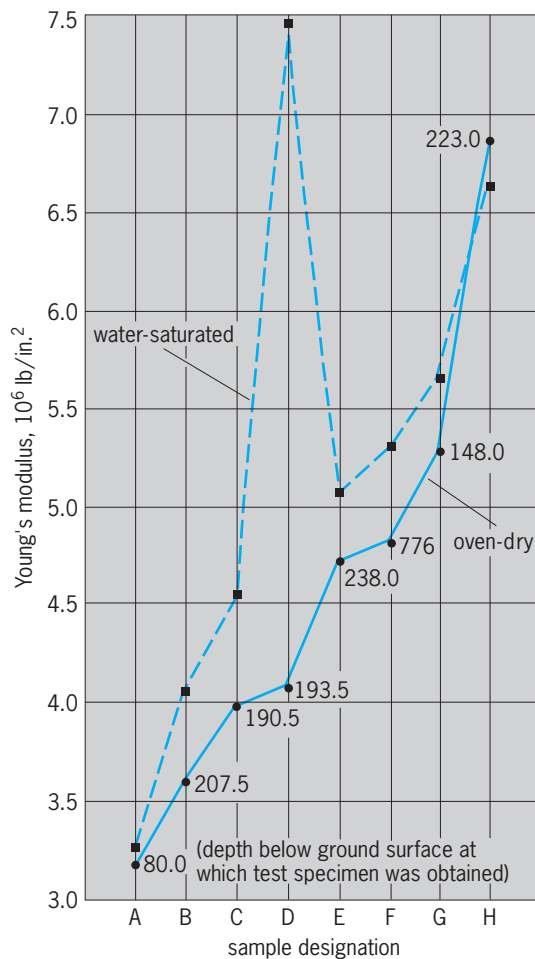


Fig. 2. Increase in Young's modulus caused by saturation of dacite porphyry. (After J. R. Ege and R. B. Johnson, *Consolidated tables of physical properties of rock samples from area 401, Nevada Test Site, U.S. Geol. Surv. Tech. Letter Pluto-21, 1962*)

Ambient stress. This type of stress in a rock system is actually potential energy, probably created by ancient natural forces, recent seismic activity, or nearby human-caused disturbances. Ambient (residual, stored, or primary) stress may remain in rock long after the disturbance is removed. An excavation, such as a tunnel or quarry, will relieve the ambient stress by providing room for displacement of the rock, and thus the potential energy is converted to kinetic energy. In tunnels and quarries, the release of this energy can cause spalling, the slow outward separation of rock slabs from the rock mass; when this movement is rapid or explosive, a rock burst occurs. The latter is a different phenomenon from a rock bump, which is a rapid upward movement of a large portion of a rock system and, in a tunnel, can have sufficient force to flatten a steel mine car against the roof or break the legs of a person standing on the floor when the bump occurs. See ROCK BURST; UNDERGROUND MINING.

One of two fundamental principles generally is used to predict the possible rock load on a tunnel roof, steel, or timber supports, or a concrete or steel lining: (1) The weight of the burden (the rock and

soil mass between the roof and the ground surface) and its shear strength control the load, and therefore the resultant stresses are depth-dependent, or (2) the shear strength of the rock system and the ambient stresses control the stress distribution, so the resultant loading is only indirectly dependent upon depth. The excavation process can cause rapid redistribution of these stresses to produce high loads upon supports some distance from the newly excavated face in the tunnel. The geometry and span of the opening also influence the stress distribution. Lined tunnels can be designed so that the reinforced concrete or steel lining will have to carry only a portion of the ambient or burden stresses. See TUNNEL.

Construction material. Rock as a construction material is used in the form of dimension, crushed, or broken stone. Broken stone is placed as riprap on slopes of earth dams, canals, and river banks to protect them against water action. Also, it is used as the core and armor stone for breakwater structures. For all such uses, the stone should have high density (± 165 lb/ft³ or 2650 kg/m³), be insoluble in water, and be relatively nonporous to resist cavitation. Dimension stone (granite, limestone, sandstone, and some basalts) is quarried and sawed into blocks of a shape and size suitable for facing buildings or for interior decorative panels. For exterior use, dimension stone preferably should be isotropic (in physical properties), have a low coefficient of expansion when subjected to temperature changes, and be resistant to deleterious chemicals in the atmosphere (such as sulfuric acid). Crushed stone (primarily limestone but also some basalt, granite, sandstone, and quartzite) is used as aggregate in concrete and in bituminous surfaces for highways, as a base course or embankment material for highways, and for railroad ballast (to support the ties). When used in highway construction, the crushed stone should be resistant to abrasion as fine stone dust reduces the permeability of the stone layer; the roadway then is more susceptible to settling and heaving caused by freezing and thawing of water in the embankment. Concrete aggregate must be free of deleterious material such as opal and chalcedony; volcanic rocks containing glass, devitrified glass, and tridymite; quartz with microfractures; phyllites containing hydromica (illites); and other rocks containing free silica (SiO₂). These materials will react chemically with the cement in concrete and release sodium and potassium oxides (alkalies) or silica gels. Preliminary petrographic analyses of the aggregate and chemical analysis of the cement can indicate the possibility of alkali reactions and thus prevent construction difficulties such as expansion, cracking, or a strength decrease of the concrete. See CAVITATION; CONCRETE; PETROGRAPHY; STONE AND STONE PRODUCTS.

Geotechnical significance of soils. Glacial and alluvial deposits contain heterogeneous mixtures of pervious (sand and gravel) and impervious (clay, silt, and rock flour) soil materials. The pervious materials can be used for highway subgrade, concrete

aggregate, and filters and pervious zones in earth embankments. Dam reservoirs may be endangered by the presence of stratified or lenticular bodies of pervious materials or ancient buried river channels filled with pervious material. Deep alluvial deposits in or close to river deltas may contain very soft materials such as organic silt or mud. An unsuitable soil that has been found in dam foundations is open-work gravel. This material may have a good bearing strength because of a natural cement bond between grains, but it is highly pervious because of the almost complete lack of fine soil to fill the voids between the gravel pebbles. *See* DELTA; FLOODPLAIN.

Concrete or earth dams can be built safely on sand foundations if the latter receive special treatment. One requirement is to minimize seepage losses by the construction of cutoff walls (of concrete, compacted clay, or interlocking-steel-sheet piling) or by use of mixed-in-place piles 3 ft (0.9 m) or more in diameter. The latter are constructed by augering to the required depth but not removing any of the sand. At the desired depth, cement grout is pumped through the hollow stem of the auger, which is slowly withdrawn while still rotating; this mixes the grout and the sand into a relatively impervious concrete pile. The cutoff is created by overlapping these augered holes. Some sand foundations may incur excessive consolidation when loaded and then saturated, particularly if there is a vibratory load from heavy machinery or high-velocity water in a spillway. This problem is minimized prior to loading by using a vibrating probe inserted into the sand or vibratory rollers on the sand surface or by removing the sand and then replacing it under vibratory compaction and water sluicing.

Aeolian (windblown) deposits. Loess is a relatively low-density (± 0.044 ton/ft³ or 1.4 metric tons/m³) soil composed primarily of silt grains cemented by clay or calcium carbonate. It has a vertical permeability considerably greater than the horizontal. When a loaded loess deposit is wetted, it rapidly consolidates, and the overlying structure settles. When permanent open excavations ("cuts") are required for highways or canals through loess, the sides of the cut should be as near vertical as possible: Sloping cuts in loesses will rapidly erode and slide because of the high vertical permeability. To avoid undesirable settlement of earth embankments, the loess is "prewettered" prior to construction by building ponds on the foundation surface. Permanently dry loess is a relatively strong bearing material. Aeolian sand deposits present the problem of stabilization for the continually moving sand. This can be done by planting such vegetation as heather or young pine or by treating it with crude oil. Cuts are traps for moving sand and should be avoided. The failure of Teton Dam in 1976 indicated, among other factors, that when loessial or silty soils are used for core materials in dam embankments, it is important to take special measures to prevent piping of the silts by carefully controlled compaction of the core, by using up- and downstream filters, and by extraordinary treatment of the foundation rock. *See* LOESS.

Organic deposits. Excessive settlement will occur in structures founded on muskeg terrain. Embankments can be stabilized by good drainage, the avoidance of cuts, and the removal of the organic soil and replacement by sand and gravel or, when removal is uneconomical, displacement of it by the continuous dumping of embankment material upon it. Structures imposing concentrated loads are supported by piling driven through the soft layers into layers with sufficient bearing power. *See* MUSKEG; TUNDRA.

Residual soils. These soils are derived from the in-place deterioration of the underlying bedrock. The granular material caused by the in-place disintegration of granite generally is sufficiently thin to cause only nominal problems. However, there are regions (such as California, Australia, and Brazil) where the disintegrated granite (locally termed DG) may be hundreds of feet thick; although it may be competent to support moderate loads, it is unstable in open excavations and is pervious. A thickness of about 200 ft (60 m) of DG and weathered gneiss on the sides of a narrow canyon was a major cause for construction of the Tumut-1 Power Plant (New South Wales) in hard rock some 1200 ft (365 m) underground. Laterite (a red clayey soil) derived from the in-place disintegration of limestone in tropical to semitropical climates is another critical residual soil. It is unstable in open cuts on moderately steep slopes, is compressible under load, and when wet produces a slick surface that is unsatisfactory for vehicular traffic. This soil frequently is encountered in the southeastern United States, southeastern Asia, and South America.

Clays supporting structures may consolidate slowly over a long period of time and cause structural damage. When clay containing montmorillonite is constantly dried and rewetted by climatic or drainage processes, it alternately contracts and expands. During the drying cycle, extensive networks of fissures are formed that facilitate the rapid introduction of water during a rainfall. This cyclic volume change of the clay can produce uplift forces on structures placed upon the clay or compressive and uplift forces on walls of structures placed within the clay. These forces have been known to rupture concrete walls containing 10.75-in.-diameter (19-mm) steel reinforcement bars. A thixotropic or "quick" clay has a unique lattice structure that causes the clay to become fluid when subjected to vibratory forces. Various techniques are used to improve the foundation characteristics of critical types of clay: (1) electroosmosis that uses electricity to force redistribution of water molecules and subsequent hardening of the clay around the anodes inserted in the foundation; (2) provision of adequate space beneath a foundation slab or beam so the clay can expand upward and not lift the structure; (3) bellling, or increasing in size, of the diameter of the lower end of concrete piling so the pile will withstand uplift forces imposed by clay layers around the upper part of the pile; (4) treatment of the pile surface with a frictionless coating (such as poly(tetrafluoroethylene) or a loose wrapping of asphalt-impregnated paper) so the upward-moving

clay cannot adhere to the pile; (5) sufficient drainage around the structure to prevent moisture from contacting the clay; and (6) replacement of the clay by a satisfactory foundation material. Where none of these solutions is feasible, the structure then must be relocated to a satisfactory site or designed so it can withstand uplift or compressive forces without extensive damage. *See* CLAY.

Silt may settle rapidly under a load or offer a “quick” condition when saturated. For supporting some structures (such as residences), the bearing capacity of silts and fine sands can be improved by intermixing them with certain chemicals that will cause the mixture to “set” or harden when exposed to air or moisture; some of the chemicals used are sodium silicate with the later application of calcium chloride, bituminous compounds, phenolic resins, or special cements (to form “soil cement”). The last mixture has been used for surfacing secondary roads, for jungle runways in Vietnam, and as a substitute for riprap of earth dams. Some types of silt foundations can be improved by pumping into them soil-cement or clay mixtures under sufficient pressures to create large bulbs of compacted silt around the pumped area.

Geotechnical investigation. For engineering projects, these investigations may include preliminary studies, preconstruction or design investigations, consultation during construction, and the maintenance of the completed structure.

Preliminary studies. These are made to select the best location for a project and to aid in formulating the preliminary designs for the structures. The first step in the study is a search for pertinent published material in libraries, state and federal agencies, private companies, and university theses. Regional, and occasionally detailed reports on local geology, including geologic maps, are available in publications of the U.S. Geological Survey; topographic maps are available from that agency and from the U.S. Army Map Service. Oil companies occasionally will release the geologic logs of any drilling they may have done in a project area. Air photos and other remote sensing techniques such as pulsed or side-looking radar or false color can be used to supplement map information (or may be the only surficial information readily available). The U.S. Geological Survey maintains a current index map of the air-photo coverage of the United States. The photos are available from that agency, the U.S. Forest Service, the Natural Resources Conservation Service, and commercial air-photo companies; for some projects, the military agencies will provide air-photo coverage. The topographic maps and air photos can be used to study rock outcrop and drainage patterns, landforms, geologic structures, the nature of soil and vegetation, moisture conditions, and land use by humans (cultural features). Airborne geophysical techniques using magnetometers or gravimeters also may be useful to delineate surface and subsurface geologic conditions. *See* AERIAL PHOTOGRAPHY; LITERATURE OF SCIENCE AND TECHNOLOGY; REMOTE SENSING; TOPOGRAPHIC SURVEYING AND MAPPING.

Field reconnaissance may include the collection of rock and soil specimens; inspection of road cuts and other excavations; inspection of the condition of nearby engineering structures such as bridges, pavements, and buildings; and location of sources of construction material. Aerial reconnaissance is essential at this stage and can be performed best in helicopters and second-best in slow-flying small planes.

Preconstruction. Surface and subsurface investigations are required prior to design and construction. Surface studies include the preparation of a detailed map of surficial geology, hydrologic features, and well-defined landforms. For dam projects, a small-scale geologic map (for example, 1:5000) is made of the reservoir area and any adjacent areas that may be directly influenced by the project; in addition, a large-scale geologic map (for example, 1:500) is required of the specific sites of the main structures (the dam, spillway, power plant, tunnels, and so on). [This preferred means of designating map scales can be used for either customary or metric units. It means 1 unit of measurement on the map is equal to 5000 similar units on the ground; for example, 1 cm measured on the map is equal to 5000 cm measured on the ground.] These maps can be compiled by a combination of field survey methods and aerial mapping procedures. They should have a grid system (coordinates) and show the proposed locations for subsurface investigations.

Subsurface investigations are required to confirm and amplify the surficial geologic data. These may include test pits, trenches, short tunnels (drifts or adits), and the drilling of vertical, horizontal, or oblique (angle) boreholes. Geologic data obtained by these direct methods can be supplemented by indirect or interpreted data obtained by geophysical methods on the surface or in subsurface holes and by installation of special instruments to measure strain or deformation in a borehole or tunnel.

The geology disclosed by subsurface investigations is “logged” on appropriate forms. Tunnel logs display visual measurements of features and joint orientations (strike and dip); rock names and a description of their estimated engineering properties; alteration, layering, and other geologic defects; the location and amount of water or gas inflow; the size and shape of blocks caused by fracturing or jointing and the width of separation or the filling material between blocks; and the irregularities in the shape of the tunnel caused by the displacement of blocks during or after excavation (rock falls, rock bursts, chimneying, and overbreakage). Geophysical seismic methods may be used to define the thickness of loosened rock around the tunnel; geoaoustical techniques that detect increases in microseismic noise during tunneling may be used to determine if the excavation is causing excessive loosening in the tunnel rock. This detection of “subaudible rock noise” occasionally is used to detect the potential movement of rock slopes in open excavations.

The borehole data can be logged on a form such as shown in **Fig. 3**. These data can be obtained by direct examination of the core, by visual inspection of

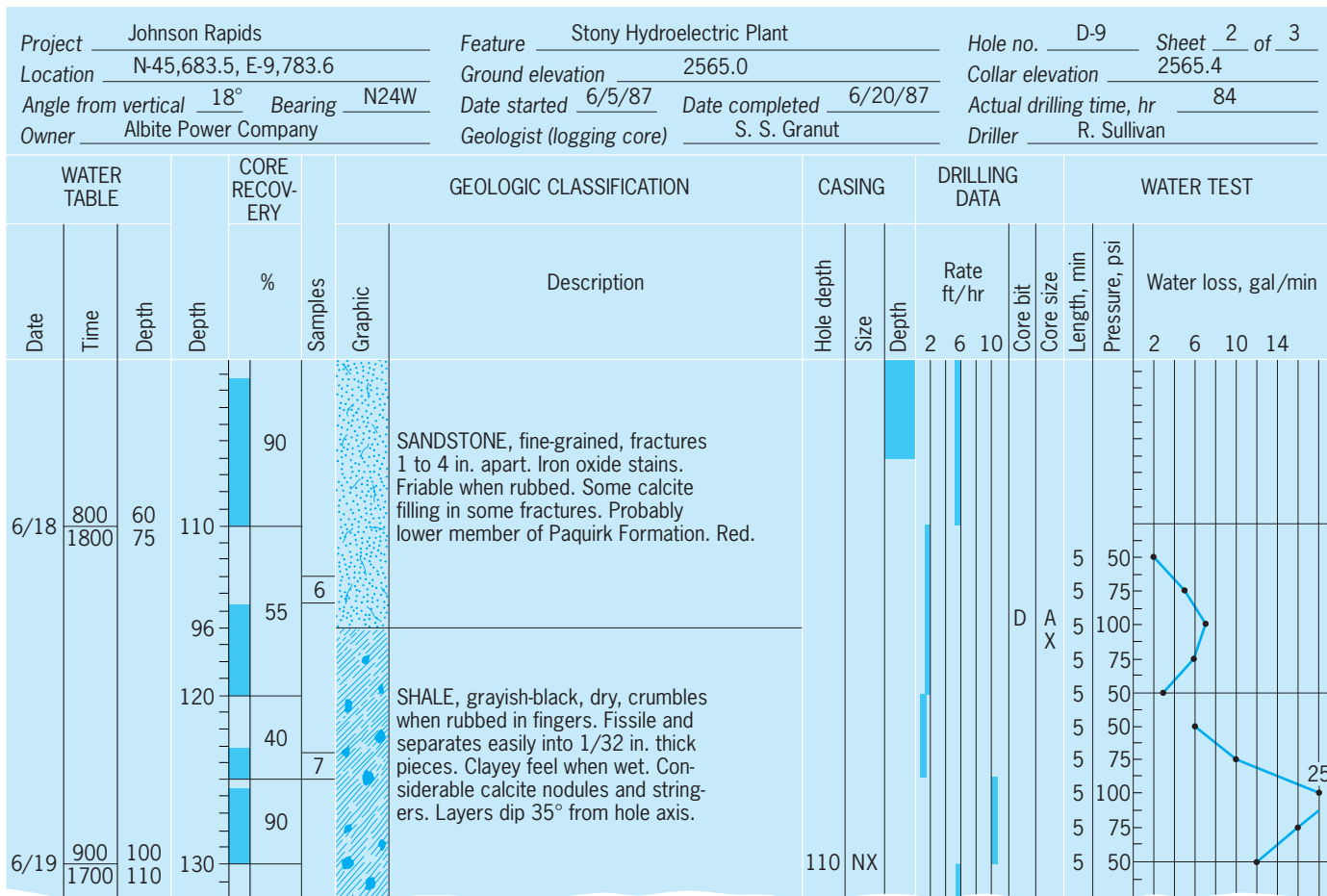


Fig. 3. Log recording of geological information from a borehole.

the interior of the borehole using a borehole camera (a specially made television camera) or a stratoscope (a periscopelike device), or by geophysical techniques. Direct viewing of the interior of the hole is the only positive method of determining the in-place orientation and characteristics of separations and of layering in the rock system. The geophysical techniques include use of gamma-gamma logging that evaluates the density of the rock surrounding the borehole or at depths as great as 150 ft (45 m) beneath the gamma probe; neutron logging to determine the moisture content of the rock system by measuring the depth of penetration of the neutrons; traversing the borehole with a sonic logger that, by calibration, measures differences in the velocity of wave propagation in different strata (and thus can determine in place Young's modulus of elasticity and the thickness of each stratum encountered by the borehole); and electric logging that uses differences in the electrical resistivity of different strata to define their porosity, moisture content, and thickness.

Occasionally a hole is drilled through a talus deposit containing the same type of rock as the underlying rock in place (bedrock). Because of the similarity in rock types, the talus-bedrock contact sometimes is best identified by determining the orientation of the remnant magnetism in the core: the magnetic lines in the core will have a regular orientation, but the

talus magnetism will have random directions. This method is useful only in rocks that contain appreciable remnant magnetism such as some basalts.

Geophysical seismic or electrical resistivity methods also can be used on the ground surface to define the approximate depth of bedrock or various rock layers. The results require verification by occasional boreholes, but this is an inexpensive and satisfactory technique for planning and design investigations. The seismic methods are not useful when it is necessary to locate soft strata (wherein the seismic waves travel at relatively low velocity) that are overlain by hard strata (that have higher wave velocity); the latter conceal and block the signal from the soft strata. Also, difficulties may occur when the strata to be located are overlain by soil containing numerous large boulders composed of rock having higher velocities than the surrounding soil, or when the soil is very compact (such as glacial till) because its velocity characteristics may resemble those of the underlying bedrock. Another problem is that the seismic method seldom can identify narrow and steep declivities in the underlying hard rock (because of improper reflection of the waves).

Construction. Geotechnical supervision is desirable during construction in or on earth media. The engineering geologist must give advice and keep a record of all geotechnical difficulties encountered

during the construction and of all geological features disclosed by excavations. During the operation and maintenance of a completed project the services of the engineering geologist often are required to determine causes and assist in the preparation of corrective measures for cracks in linings of water tunnels, excessive settlement of structures, undesirable seepage in the foundations of dams, slides in canal and other open excavations, overturning of steel transmission-line towers owing to a foundation failure, and rock falls onto hydroelectric power plants at the base of steep canyon walls. The engineering geologist also is considered an important member of the team assigned to the task of assessing the safety of existing dams as now required by federal legislation.

Legal aspects. An important consideration for the engineering geologist is the possibility of a contractor making legal claims for damages, purportedly because of unforeseen geologic conditions (generally referred to legally as charged conditions) encountered during construction. Legal support for such claims can be diminished if the engineering geologist supplies accurate and detailed geologic information in the specification and drawings used for bidding purposes. These documents should not contain assumptions about the geological conditions (for a proposed structure), but they should show all tangible geologic data obtained during the investigation for the project: for example, an accurate log of all boreholes and drifts and a drawing showing the boundaries of the outcrops of all geological formations in the project area. The engineering geologist should have sufficient experience with design and construction procedures to formulate an investigation program that results in a minimum of subsequent uncertainties by a contractor. Numerous uncertainties about the geologic conditions not only can result in increased claims but also may cause a contractor to submit a higher bid (in order to minimize risks) than if detailed geologic information were available.

Special geotechnical problems. In arctic zones, structures built on permafrost may be heaved or may cause thawing and subsequent disastrous settlement. The growth of permafrost upward into earth dams seriously affects their stability and permeability characteristics. Obtaining natural construction materials in permafrosted areas requires thawing of the borrow area to permit efficient excavation; once excavated, the material must be protected against refreezing prior to placement in the structures. Permafrost in rock seldom will cause foundation difficulties. In planning reservoirs, it is essential to evaluate their watertightness, particularly in areas containing carbonate or sulfate rock formation or lava flows. These formations frequently contain extensive systems of caverns and channels that may or may not be filled with claylike material or water. Where extensive openings occur, grouting with cement slurry or chemicals can be used as a sealant; however, as demonstrated by the 1976 failure of the Teton Dam, such measures are not always successful. Sedimentation studies are required for the design of efficient harbors or reservoirs because soil carried by moving

water will settle and block or fill these structures. In areas with known earthquake activity, structural design requires knowledge of the intensity and magnitude of earthquake forces. The prevention and rehabilitation of slides (landslides) in steep natural slopes and in excavations are important considerations in many construction projects and are particularly important in planning reservoirs, as was disastrously proved by the Vaiont Dam catastrophe in 1963. See EARTHQUAKE ENGINEERING; PERMAFROST.

Geohydrologic problems. In the foundation material under a structure, water can occur in the form of pore water locked into the interstices or pores of the soil or rock, as free water that is moving through openings in the earth media, or as included water that is a constituent or chemically bound part of the soil or rock. When the structure load compresses the foundation material, the resulting compressive forces on the pore water can produce undesirable uplift pressures on the base of the structure. Free water is indicative of the permeability of the foundation material and possible excessive water loss (from a reservoir, canal, or tunnel); uplift on the structure because of an increase in hydrostatic head (caused by a reservoir or the like); or piping, which is the removal of particles of the natural material by flowing water with a consequent unfilled opening that weakens the foundation and increases seepage losses.

The possibility of excessive seepage or piping can be learned by appropriate tests during the boring program. For example, water pressure can be placed on each 5-ft (2-m) section of a borehole, after the core is removed, and any resulting water loss can be measured. The water pressure is maintained within the 5-ft section by placing an expandable rubber ring (packer) around the drill pipe at the top of the test section and then sealing off the section by using mechanical or hydraulic pressure on the pipe to force expansion of the packer. When only one packer is used, because it is desired to test only the section of hole beneath it, it is a "single-packer" test. In a double-packer test, a segment of hole is isolated for pressure testing by placing packers at the top and the bottom of the test section. The best information on the permeability characteristics of the rock can be obtained by the use of three or more increments of increasing and then decreasing water pressure for each tested length of hole. If the water loss continues to increase when the pressure is decreased, piping of the rock or filling material in fractures may be occurring or fractures are widening or forming. The water-pressure test can be supplemented by a groutability test in the same borehole. This test is performed in the same way as the water test except, instead of water, a mixture of cement, sand, and water (cement grout) or a phenolic resin (chemical grout) is pumped under pressure into the test section. The resulting information is used to design cutoff walls and grout curtains for dams. The pressures used in water-pressure or grouting tests should not exceed the pressure exerted by the weight of the burden between the ground surface and the top of the test section. Excessive test pressure can cause

uplift in the rock, and the resulting test data will be misleading.

Included or pore water generally is determined by laboratory tests on cores; these are shipped from the borehole to the laboratory in relatively impervious containers that resist loss of moisture from the core. The cores with their natural moisture content are weighed when received and then dried in a vacuum oven at about 110°F (45°C) until their dry weight stabilizes. The percentage of pore water (by dry weight) is $(\text{wet weight} - \text{dry weight}) \times 100 \div \text{dry weight}$.

Temperatures up to 200°F (about 90°C) can be used for more rapid drying, provided the dried specimens are not to be used for strength or elastic property determinations. (High temperatures can significantly affect the strength because the heat apparently causes internal stresses that disturb the rock fabric or change the chemical composition of the rock by evaporation of the included moisture.)

Protective construction. Civilian and military structures may be designed to minimize the effects of nuclear explosions. The most effective protection is to place the facility in a hardened underground excavation. A hardened facility, including the excavation and its contents, is able to withstand the effects generated by a specified size of nuclear weapon. These effects include the amount of displacement, acceleration, and particle velocity that occurs in the earth media and the adjacent structure. Desirable depths and configurations for hardened facilities are highly dependent upon the shock-wave characteristics of the surrounding earth media, for example, the type of rock, discontinuities in the rock system, free water, and geologic structure. Therefore, prior to the design and construction of such facilities, extensive geotechnical field and laboratory tests are performed, including an accurate geologic map of the surface and of the underground environment that will be affected by the explosion. The map should show the precise location and orientation of all geologic defects that would influence the wave path, such as joints, fractures, and layers of alternately hard and soft rock. See EXPLOSIVE; NUCLEAR EXPLOSION.

Application of nuclear energy. The use of nuclear energy for the efficient construction of civil engineering projects has been investigated in the Plowshare Program. Examples include rapid excavation, increasing production of natural gas by opening fractures in the reservoir rock, expediting production of low-grade copper ore by causing extensive fracturing and possible concentration of the ore, and by the underground "cracking" of oil shale. The production of electrical energy by nuclear fission requires engineering geology inputs during the planning and design of the power plant; for example, a major question to be answered is the presence or absence of faults and an estimate of when the last movement on the fault occurred. This question of "active" faults also is of increasing concern in the siting of dams.

Waste disposal. Another geotechnical problem occurs in the use of nuclear energy for generation of power or radioisotopes: safe disposal of the radioactive waste products. These products can be

mixed with concrete and buried in the ground or ocean, but geohydrologic or oceanographic conditions must not be conducive to the deterioration of the concrete. One proposed solution is to excavate large caverns in rock or salt a thousand or more meters deep; however, such a solution must consider possible contamination of ground water in the event that the waste products' containers leak. The disposal of toxic chemical or biological products in deep wells no longer is considered safe. See RADIOACTIVE WASTE MANAGEMENT. William R. Judd

Bibliography. F. G. Bell, *Engineering Properties of Soil and Rocks*, 4th ed., 2000; B. M. Das, *Principles of Geotechnical Engineering*, 5th ed., 2001; R. E. Goodman, *Engineering Geology: Rock in Engineering Construction*, 1993; M. E. Harr, *Groundwater and Seepage*, 1992; R. B. Johnson and J. V. DeGraff, *Principles of Engineering Geology*, 1989; R. F. Leggett and A. W. Hatheway, *Geology and Engineering*, 3d ed., 1988; R. F. Leggett and P. K. Karrow, *Handbook of Geology in Civil Engineering*, 3d ed., 1982; P. H. Rahn, *Engineering Geology: An Environmental Approach*, 2d ed., 1996; Q. Zaruba, *Landslides and Their Control*, 2d ed., 1982.

Enopla

A class of the phylum Rhynchocoela which is divided into the orders Hoplonemertini (free-living) and Bdellonemertini or Bdellomorpha (symbiotic). The proboscis in the Hoplonemertini is armed with stylets, and in both orders the mouth is anterior to the brain, the nervous system is inside the body wall musculature, and the vascular system is well developed.

The Bdellonemertini contain the single genus *Malacobdella* with four species, all of which are symbiotic with mollusks. See ANOPLA; BDELLONEMERTINI; HOPLONEMERTINI; NEMERTEA. J. B. Jennings

Enoplida

An order of nematodes characterized by having the cuticle of the cephalic region doubled or formed into a cap or helmet. The helmet results from the infolding of the cuticle over the extreme anterior end of the esophagus. The stoma may or may not possess armature in the form of teeth or movable jaws. Even though the stoma is surrounded by esophageal tissue, the anterior portion of the movable jaws is of cheilostomal origin. The well-developed esophagus is nearly cylindrical, and the esophagointestinal valve is prominent. The five esophageal glands are located in the posterior region of the esophagus; however, the orifices for these glands are located anterior to the nerve ring. The dorsal and two anterior subventral glands open at the base of the stoma or through the stomatal teeth. The amphid apertures may be transverse slits or elongate ovals. The cephalic sensilla are most often in two whorls: an anterior

circumoral whorl of 6 and a posterior whorl of 10 (6 + 4); in some taxa this whorl is clearly separated into the ancestral state of a cirlet of 6 and one of 4. A single medioventral excretory cell is usually present; the orifice is generally at the level of or anterior to the nerve ring. Medioventral supplementary organs may or may not be present on males. The male spicules are paired and accompanied by a gubernaculum. Caudal glands and a cuticular spinneret are found in males and females.

There are two enoplid superfamilies, Enploidea and Oxystominoidea. The Enploidea are free-living marine nematodes comprising small to very large species in five well-defined families: Leptosomatidae, Lauratonematidae, Phanodermatidae, Thoracostomopsidae, and Enoplidae. While the feeding habits are largely unknown, it has been postulated that the group is predaceous or omnivorous. The group is almost exclusively marine in habitat, but a few species have been reported from brackish water or soils. Oxystominoidea contains species that maintain the most ancestral of characters known in the phylum. They are separated from other Enoplida by not having the lip region (=cephalic region) set off by a groove. See NEMATATA (NEMATODA).

Armand R. Maggenti

Enstatite

A pyroxene close in composition to magnesium metasilicate (MgSiO_3). Like other pyroxenes, enstatite is typified by a pair of nearly perpendicular cleavages, both of which are parallel to the elongate dimension of the characteristically prismatic crystals. These cleavages reflect an atomic structure dominated by unbroken, parallel, linear chains of SiO_4 tetrahedra (silicon atoms surrounded by four oxygens). See PYROXENE.

The most common form of enstatite in nature is orthorhombic; it is part of a continuous solid-solution series with ferrosilite (FeSiO_3). By convention, the name enstatite is restricted to members of this solution series that are optically positive, which includes only compositions with less than about 12 mol % of the iron component. Pure magnesian enstatite is colorless, with refractive indices $n\alpha = 1.650$, $n\beta = 1.653$, $n\gamma = 1.658$; the substitution of iron increases the refractive indices and produces gray to pale-green coloration. Solid solution of calcic and aluminous components in orthoenstatite is common, but is limited to a few mole percent. See SOLID SOLUTION.

Compositions very close to MgSiO_3 may crystallize in a variety of different structures. Upon heating at atmospheric pressure, pure orthoenstatite transforms at 985°C (1805°F) to another orthorhombic mineral, protoenstatite, which is then stable up to its melting point near 1557°C (2835°F). A monoclinic form, clinoenstatite, is stable over a wide range of temperatures at pressures in excess of about 7.5–10 gigapascals ($1.088\text{--}1.450 \times 10^6$ lb/in.²). Clinoenstatite may also originate outside its field of stability, either by the inversion of protoenstatite on cooling

or by intense shearing of orthoenstatite. Crystals of clinoenstatite produced by these mechanisms are repetitively twinned at very fine scale on the crystallographic plane (100). See CRYSTAL STRUCTURE.

Magnesium-rich members of the enstatite-ferrosilite series are commonly found in mafic and ultramafic rocks of the Earth's crust, such as basalts, gabbros, peridotites, dunites, serpentinites, and some mafic granulites. In these occurrences, associated minerals are typically olivine, calcic clinopyroxene, and calcic plagioclase. Alteration of enstatite to amphibole or to serpentine is common in such rocks. Enstatite is also an important constituent of the Earth's mantle, as indicated by its presence (along with olivine, diopside, spinel, or garnet) in ultramafic nodules included in lavas that originate at great depths, and in larger bodies of rock, with similar mineralogy, exposed in some deeply rooted mountain belts. Orthoenstatite and clinoenstatite are also found in both iron and stony meteorites.

Because the amounts of calcium and aluminum incorporated in orthopyroxene depend strongly on temperature and pressure, enstatite is an important geologic thermometer and barometer. Estimates of temperatures and pressures of mineral formation, based upon enstatite compositions in natural occurrences, permit geologists to determine environments of origin for many rocks of the Earth's lower crust and upper mantle. See GEOLOGIC THERMOMETRY; SILICATE MINERALS.

William D. Carlson

Bibliography. W. D. Carlson, Subsolidus phase equilibrium on the forsterite-saturated join $\text{Mg}_2\text{Si}_2\text{O}_6\text{--CaMgSi}_2\text{O}_2$ at atmospheric pressure, *Amer. Mineralog.*, 73:232–241, 1988; D. H. Lindsley, Pyroxene thermometry, *Amer. Mineralog.*, 68:477–493, 1983; C. T. Prewitt, *Pyroxenes*, vol. 7 of *Reviews in Mineralogy*, 1980.

Enteropneusta

A class of Hemichordata with approximately 70 species commonly known as acorn worms. They are free-living, solitary animals with a very soft cylindrical body lacking external appendages. The size is highly variable, from 1 in. to 8 ft (2.5 cm to 2.5 m), and the color ranges from white to shades of violet. The body is covered with cilia and mucus and is always divided into proboscis, collar, and trunk (Fig. 1). Sometimes the animal smells like iodoform, and sometimes it shows luminescence.

The acorn worms usually live in shallow waters, buried in the sandy or muddy bottoms, but some species occur at depths of more than 9900 ft (3000 m). *Saxipendium coronatum* lives on hydrothermal vents. The species that live along the shore, sometimes in U-shaped or spiral burrows, can be detected at low tide by the presence of a spiral casting. The proboscis and collar (the acorn) are used in excavating burrows. During this activity, water, slime, sediment, and organic particles enter the mouth and pass into the gut. Seawater is filtered

and expelled through the gill slits, while the organic matter and sediment are retained to be digested.

The proboscis, or protosome, is generally conical and is continuous with the dorsal edge of the collar. The protoceol is more or less filled by muscle and connective tissue, confined to the posterior part, and lined by peritoneum. The protoceol opens to the dorsal surface by means of a proboscis pore.

The collar, or mesosome, is a short cylinder with a highly glandular epithelium usually differentiated into zones. There are two mesocoels, which extend anteriorly into the proboscis stalk. Each mesocoel opens to the exterior by way of a collar canal and pore ending into the first gill slit.

The trunk is differentiated into four regions: the branchial region, externally recognized by two longitudinal rows of dorsal gill pores; the genital region, characterized by the gonads, which occur dorsolaterally (in some genera, laterally from the genital wings); the hepatic region, distinguishable by a darker color and sometimes by the presence of external sacculations. The final region, the abdominal, is simply tubular with a terminal anus. The metacoels are closed.

The nervous system consists of an intraepidermal nervous plexus that becomes thickened to form the middorsal and midventral nerve cords. The dorsal trunk cord continues into the collar through the collar coelom to form the neurochord, or collar cord, which is isolated from the epidermis. A ciliated groove in the proboscis may be a chemoreceptor. Other sense organs are lacking.

The gut is a straight epithelial tube. The mouth opens at the anterior margin of the collar, and a buccal tube whose roof forms the buccal diverticulum by evagination occupies the center of the collar. The pharynx usually is differentiated into a ventral digestive part and a dorsal branchial part that is perforated by gill slits, which are primarily oval and elongated dorsoventrally. In all enteropneusts, however,

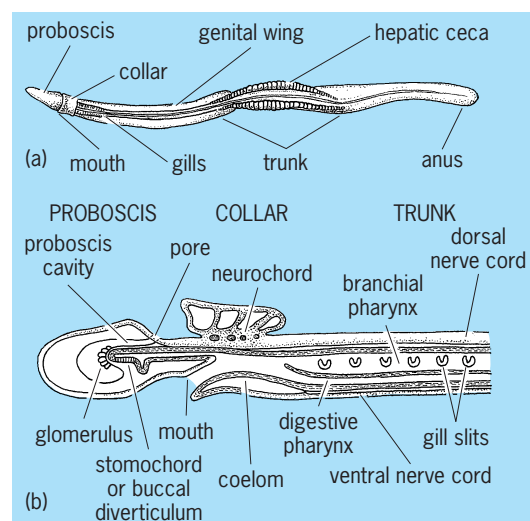


Fig. 1. *Saccoglossus*, the tongue worm or acorn worm. (a) Dorsal view. (b) Median section of anterior portion. (After T. I. Storer et al., *General Zoology*, 6th ed., McGraw-Hill, 1979)

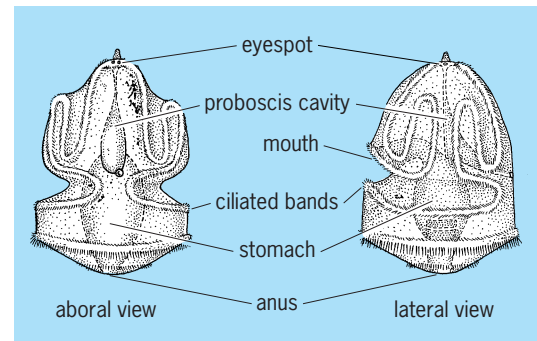


Fig. 2. *Tornaria* larva. (After T. I. Storer et al., *General Zoology*, 6th ed., McGraw-Hill, 1979)

the gill slits are transformed into U-shaped apertures because of a downgrowth from the dorsal end of the slit, the tongue bar. The gill slits lead to branchial sacs that open on the body surface by the gill pores. A muscular esophagus that molds the food particles follows the pharynx and leads to the hepatic region, where two rows of dorsal sacculations are distinguishable in some species. The intestine continues as a tube to the terminal anus.

The blood vascular system, which is open, consists of two main longitudinal vessels lined by peritoneal cells that lack an endothelium. Also present are a venous sinus, a central sinus, and a heart vesicle. The glomerulus is lined by podocytes.

The sexes are distinct. The sacciform gonads occur in from one to several longitudinal rows at the sides of the digestive tube in the branchial and genital regions. Each gonad opens externally in a pore through which the gametes are shed. Either the eggs are relatively yolky and develop directly, with a free-swimming larval stage that hatches into a miniature acorn worm, or they lack yolk and develop through a planktonic tornaria larval stage (Fig. 2). Only one species (*Xenopleura vivipara*) is viviparous, and only one (*Balanoglossus capensis*) reproduces asexually, but all enteropneusts have strong regeneration power. See HEMICHORDATA. Jesús Benito

Enterovirus

A genus of the family Picornaviridae. Enteroviruses include the polioviruses (three types), the coxsackieviruses (more than 29 types), and the echoviruses (more than 34 types). They are small (17–28 nanometers in diameter), contain ribonucleic acid (RNA), and are resistant to ether. In contrast to the rhinoviruses, the enteroviruses multiply chiefly in the alimentary tract and are stable under acid conditions (pH 3–5) for 1–3 h. They are protected by magnesium chloride against inactivation by heat. The polioviruses, most echoviruses, and a number of the coxsackieviruses can be grown in cell cultures of monkey origin, as well as in human cells. Certain of the coxsackieviruses will not grow in cultures; they are usually studied by infecting newborn mice. Those strains that can be grown in tissue cultures usually grow best if the cultures are kept

stationary and incubated at 98.5°F (37°C). See RIBONUCLEIC ACID (RNA).

Enteroviruses are widespread during summer and fall in temperate climates, but may circulate throughout the year in tropical areas. Different enteroviruses may produce the same symptoms; on the other hand, the same enterovirus may cause more than a single set of symptoms. The majority of enterovirus infections are benign and inapparent. However, when these viruses invade tissues other than the enteric tract, serious disease may result, as when poliovirus invades the spinal cord or when some of the coxsackievirus types invade the heart muscle. See ANIMAL VIRUS; COXSACKIEVIRUS; ECHOVIRUS; PICORNAVIRIDAE; POLIOMYELITIS; RHINOVIRUS; VIRUS CLASSIFICATION.

Joseph L. Melnick

Bibliography. H. Fraenkel-Conrat, *The Viruses: Catalogue, Characterization and Classification*, 1985; R. Rott and W. Goebel, *The Molecular Basis of Viral and Microbial Pathogenesis*, 1987; A. Scott, *Pirates of the Cell: The Story of Viruses from Molecule to Microbe*, 1987.

Enthalpy

For any system, that is, the volume of substance under discussion, enthalpy is the sum of the internal energy of the system plus the system's volume multiplied by the pressure exerted by the system on its surroundings. This may be expressed as $U + PV = H$, where U is the system's internal energy, P the pressure of the system, V the system's volume, and H the enthalpy of the system. The sum of $U + PV$ is given the special symbol H primarily as a matter of convenience because this sum appears repeatedly in thermodynamic discussion. Consistent units must, of course, be used in expressing the terms in the above equation. Previously, enthalpy was referred to as total heat or heat content, but these terms are misleading and should be avoided. Enthalpy is, from the viewpoint of mathematics, a point function, as contrasted with heat and work, which are path functions. Point functions depend only on the initial and final states of the system undergoing a change; they are independent of the paths or character of the change. Mathematically, the differential of a point function is a complete or perfect differential. See CALCULUS; MATHEMATICS; MAXWELL'S EQUATIONS.

Because the absolute value of internal energy of even a simple system is usually unknown, recorded values of enthalpy are relative values measured above some convenient but arbitrarily chosen datum. Thus in the steam tables of Keenan and Keyes, the datum is liquid water at 32°F (0°C) and under its own vapor pressure. At this state water is assumed to have an enthalpy equal to zero. Under this assumption the internal energy of water in this state is a negative quantity equal to PV . No complication is introduced by this fact, although visualization of negative energies of this kind may be disturbing to some. There is limited utility for absolute enthalpies because only the changes in enthalpy are measurable. It is instruc-

tive to examine the utility of the enthalpy function in terms of some simple but important thermodynamic processes.

The first law of thermodynamics is merely a statement of the law of conservation of energy. The first law alone indicates that:

1. For a chemical reaction carried out at constant pressure and temperature with no work performed except that resulting from keeping the internal and external pressure equal to each other as the volume changes, the change in enthalpy of the system (the material taking part in the chemical reaction) is numerically equal to the heat that must be transferred to maintain the above-mentioned conditions. This heat is often loosely referred to as the heat of reaction. More properly, it is the enthalpy change for the reaction.

2. So-called heat balances on heat exchangers, furnaces, and similar industrial equipment that operate under steady flow conditions are really enthalpy balances.

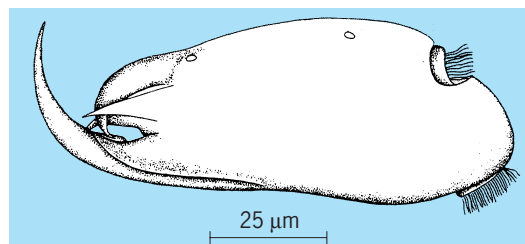
3. The work developed in a steadily running adiabatic engine or turbine is equivalent to the enthalpy change of the fluid passing through the engine.

4. The adiabatic, irreversible, steady flow of a stream of materials through a porous plug or a partially opened valve under circumstances where the change in kinetic energy of flow is negligible (a Joule-Thomson process) results in no change in enthalpy of the flowing stream. Although no change in enthalpy results from this process, there is a loss in the energy available for doing work as a result of the pressure drop across the plug or valve. See ENTROPY; THERMODYNAMIC PRINCIPLES; THERMODYNAMIC PROCESSES.

Harold C. Weber; William A. Steele

Entodiniomorphida

An order of the Spirotrichia. These are strikingly different-looking ciliates, covered with a smooth, firm pellicle. They are devoid of external ciliature except for the adoral zone of membranelles and, occasionally, one or two other tufts or zones of other specialized cilia. Internal organization of the body is very specialized and complex. These organisms are considered to be highly evolved. Entodiniomorphids occur exclusively as endocommensals of herbivorous mammals, either in the rumen and reticulum of ruminants or in the colon of certain higher mammals. The ophryoscolecids, which comprise the



Ophryoscolex, an entodiniomorphid.

majority of species in this order, are found in fantastic abundance in their hosts; for example, there have been estimations of as many as 10 billion per cow. Herbivores can survive free of their harmless protozoan guests, but the reverse is not true. *Epidinium*, *Opbryoscolex* (see **illus.**), and *Entodinium* occur in ruminants. *Troglodytella* is found in the colon of anthropoid apes. See CILIOPHORA; PROTOZOA; SPIROTRICHIA.

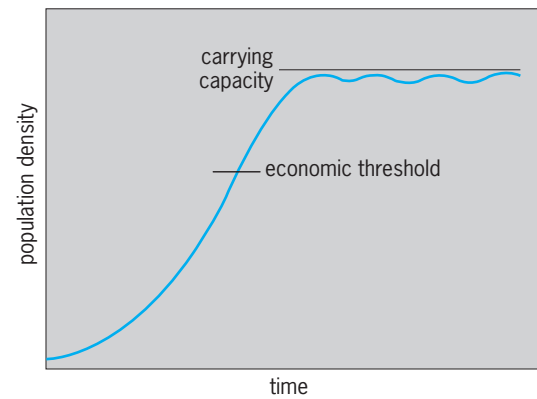
John O. Corliss

Entomology, economic

The study of insects that have a direct influence on humanity. Though this includes beneficial as well as harmful species, most attention is devoted to the latter and how they become pests and are controlled. The emphasis on managing harmful insects reflects the immediacy and seriousness of pest problems, particularly the destruction of agricultural products and the transmission of disease. These are highly visible problems, whereas the benefits gained from useful insects are, in most cases, not so clearly understood or so well documented economically.

Economic thresholds. Central to the definition of a pest is determination of the economic threshold. Any insect population, when introduced into a favorable environment, increases numbers until reaching an environmental carrying capacity (see **illus.**). In pest insects, there exists a density above which the insect population interferes with human health, comfort, convenience, or profits. When this economic threshold is reached, a decision must be made to utilize some control measure to prevent further increase in numbers. Precise estimation of the economic threshold is difficult, since it depends on a myriad of factors such as weather, crop conditions, stage in the life cycle of the insect, market value of the crop, and cost of control, each of which varies. Often, the presence of even a single insect is sufficient to warrant control measures, for instance, when that insect is a flea harboring the plague bacillus, or a mosquito transmitting malaria. Also, consumer expectations in most markets of the United States, Canada, and western Europe are for insect-free produce, so the economic threshold is very low on items that people eat. However, some crops, such as field corn or forest trees, have a higher threshold since they can tolerate a low amount of insect damage before yield is reduced. Economic thresholds may also be higher for insects that damage only the inedible portions of crop plants such as the leaves of beans, tomatoes, and apple trees. In any case, knowledge of the amount of injury which is due to different densities of insects is an important prerequisite for efficient management.

Harmful insects. Insects cause damage in various ways. Direct damage occurs when insects eat foods destined for human consumption or otherwise decrease resource yields. Direct damage to crops by locusts has been a major factor limiting grain production in some semiarid countries of Asia and Africa. Damage to structures by termites and other wood-chewing insects runs to billions of dollars an-



Growth of an insect population continues until reaching the environmental carrying capacity.

nally. Indirect damage may occur when an insect, feeding on nonedible portions of a crop plant, reduces yield of the edible portion. Red mites and aphids respectively feed on leaves and sap of apple trees, which results in lowered yield of fruit. Many of the most serious insect pests are those that affect humans indirectly by transmitting pathogens. Mosquitoes transmit malaria, yellow fever, viral encephalitis, and filarial roundworms. Such vectors (transmitters) are number one among insect pests. Fleas, ticks, lice, houseflies, cockroaches, and many others are capable of transmitting diseases to humans, while many aphids and leafhoppers are major vectors of plant virus diseases. Indirect damage may also occur from the mere inclusion of insects or their fragments in fresh or processed foods, especially when prohibited by law. Some persons exhibit extreme psychological reaction to a real or imagined insect infestation, and this, too, is damaging.

Pest management. Management of insect pests begins with prevention. Many of the United States' most noxious insects have been imported from overseas: most domestic cockroaches, the gypsy moth, Japanese beetle, corn borer, housefly, cabbageworm, and codling moth are just a few. Most major pests of foreign origin were introduced accidentally during the nineteenth century. Some North American insects have spread elsewhere, for example—the Colorado potato beetle to Europe and the fall webworm to Japan. To stem the flow of insect invasions, the federal government's Animal and Plant Health Inspection Service maintains inspection facilities for the examination of all incoming shipments of plant or other material that may harbor pests in order to intercept infested items.

Once a pest is established, its spread can sometimes be slowed by an efficient system of local quarantines, early detection, and local eradication. The gypsy moth, long an established defoliator of forests in New England, is slowly extending its range into the midwestern United States, with isolated individuals as distant as Oregon and California. States outside the infested region maintain a vigorous program of detection: sticky traps baited with sex attractant lure male moths, and when several are located at one site, a thorough inspection is followed by application of

insecticide. Gypsy moths lay their eggs on the undersides of vehicles, so that inspection and treatment of autos and trucks that have traveled in the infested area slows the westward spread of this pest.

Insecticides. Once noxious insects are firmly established in a region, there are many techniques to reduce their numbers below the economic threshold. The most widely used method is the application of synthetic chemical insecticides. Insecticides are probably the only efficient control technique for insects with an exceedingly low economic threshold such as occurs on commercial floral crops, fresh and processed fruits and vegetables, or nursery stock. There are scores of available insecticides which vary greatly in their characteristics; extension specialists at the county level or state university should be consulted for the latest information on available compounds. Most insecticides are poisonous to other animals, and handling requires appropriate caution. Insecticides were once regarded as a panacea for pest problems, but the development of resistant strains of major insect pests, together with the rising cost of materials and application and legal restrictions, has led to recognition that insecticides are more efficiently utilized in a program that integrates them with other techniques in a framework of total crop management. Research has been directed toward development of insecticides that minimize hazards to people and their environment. *See* INSECTICIDE.

Biological control. For insects whose economic threshold lies somewhat higher than the artificially low demands of "clean" produce, there is a wider choice of control techniques. Insects introduced into a new homeland where they become numerous sometimes have had their numbers checked in their old homeland by natural enemies. Economic entomologists have effectively reduced densities of several pests by releasing parasites or predators. Imported natural enemies have had major impact in controlling scale insects, aphids, whiteflies, alfalfa weevils, and many others. Natural enemies that are mobile and relatively restricted in diet are most effective in biological control. Diseases of insects often are important in halting outbreaks by killing large numbers in a short time. A few pathogens of caterpillars, mosquito larvae, Japanese beetle larvae, and others are formulated as biological rather than chemical insecticides and are applied commercially. *See* INSECT CONTROL, BIOLOGICAL.

Cultural practices. All populations are limited by their environment, and much can be done to augment the activity of natural control by manipulation of cultural practices. Crop rotation is a standard agronomic practice that often reduces damage due to insects. Rotation of alfalfa or soybeans with corn reduces populations of corn rootworms, wireworms, and white grubs. The physical disruption of autumnal plowing and disking destroys many insects, such as the corn borer, wheat stem sawfly, and cereal leaf beetle, that could overwinter in stubble or on the soil surface. Plant breeders and entomologists have cooperated to produce varieties of corn, wheat, soybeans, alfalfa, and other crops which retain high yields but

which resist attacks from some of their more serious pests. The action of natural enemies can be enhanced by cultural practices such as leaving fencerows or other preserves for refuges at harvest. Strip cutting of alfalfa in the southwestern states has been particularly valuable in preserving predators of alfalfa aphids and caterpillars. *See* AGRICULTURAL SOIL AND CROP PRACTICES; BREEDING (PLANT).

Sanitation. The cleanup of breeding and gathering sites is useful, especially in management of medically important insects, many of which have evolved resistance to the commoner insecticides. Housefly larvae feed in decaying vegetable compost and dung; efficient disposal of garbage, manure, and sewage brings enormous relief from fly problems and associated diseases. Similarly, draining standing water with its crop of mosquito larvae reduces their numbers near human habitations. Personal hygiene and proper care of pets virtually eliminate problems with lice and fleas.

Special programs. Rarely, a unique program is developed to control a single pest insect. The screwworm fly, whose larvae infest open wounds of livestock, has been controlled in the United States and much of Mexico by weekly release of 180,000,000 sterile males. These mate with wild females (who only mate once), and the ensuing eggs do not hatch. Obviously, this technique is limited to species in which males are undamaging, though despite technical problems it is being developed for use against codling moths. The method is also being developed to control some species of mosquitoes, in which only females bite.

Integrated control. Most successful programs of insect pest management rely on integrated control or the use of several methods in concert to control a complex of pests. Apples, for example, are attacked by a host of insect pests, led by species such as the codling moth and apple maggot that attack the fruit directly and consequently have a low economic threshold (1/100 apples). Other pests, notably aphids and European red mites, feed on leaves and sap and are only pests when their numbers become enough to reduce yield of fruit.

Application of chemical insecticides to control codling moths killed predators of red mites and caused these creatures to become damaging secondary pests. Eventually, both the red mite and its major predators become partially resistant to most conventional insecticides. In the integrated control program, reduced dosages of insecticides conserve predators of the red mite while still yielding acceptable control of the fruit eaters.

The alfalfa weevil, whose larvae eat alfalfa tips in spring, is usually managed by use of resistant plant varieties and judiciously planned cutting times, with some additional control from six species of introduced parasites. Occasionally, the economic threshold is exceeded enough for an insecticide to be applied. Farmers must watch carefully as populations increase, or the crop may also be inspected by an agricultural consultant. The consultant's pest management scouts observe the fields regularly, and the farmers receive weekly (or more frequent) reports

on the status of their crops. Pest management programs involving crop scouting are now operational for corn, cotton, soybeans, fruits, tomatoes, tobacco, alfalfa, and other crops.

Beneficial insects. It has been estimated that the dollar value gained from a single insect, the honeybee, equals the loss from damage plus cost of control for all pests combined. Honeybees are managed for their honey and beeswax, but their most valued service is pollination of crop plants. Nearly all fruits and many vegetables, ornamental plants, and seed crops require pollination by honeybees or other insects. Bees of many species are the chief pollinators, though wasps, flies, moths, butterflies, and beetles pollinate as well.

Silk is produced by larvae of the silkworm, an insect so thoroughly domesticated that it cannot climb its food plant, mulberry, with its degenerated legs. The silkworm apparently no longer survives in the wild. In many uses, silk has been more recently replaced by less expensive synthetic materials.

The economic value of silkworms and honeybees is rather easily estimated from the cash value of their products. Other insects may be equally beneficial, but their value is not so easily calculated. Foremost among these are predatory insects of several orders. Some, chiefly parasitic wasps and predatory beetles, have been imported specifically for control of noxious pests. Others, such as dragonflies, mantises, lace wings, and such, are voracious predators and doubtless are a factor determining the environmental carrying capacity for many pests. These predators may prevent other insects from ever reaching an economic threshold and thus from becoming pests, as was demonstrated when predators of the European red mite were destroyed by chemical insecticides, allowing the mite population to reach damaging numbers on apple crops.

Innumerable insect species are scavengers, quietly but efficiently breaking down the remains of dead plants and animals. Their economic activity goes unnoticed, save when scavengers such as termites forsake logs on the forest floor for sills and siding on a summerhouse, or when flies invade a home. A lack of scavenging insects would, however, result in a great increase of decomposing organic material lying about.

Plant-eating insects have been set to beneficial use when their diets consist mainly of unwanted weeds. Alligatorweed in the southeastern United States and klamath weed in the northwestern states have been controlled by imported beetles, while prickly-pear cactus was similarly eliminated from Australian rangeland by larvae of moths from South America. Pasture thistles may soon be controlled by beetles in the eastern United States and adjacent Canada.

In many parts of the world, particularly the tropics and subtropics, termites, grubs, locusts, and other insects are routinely eaten by people, and for some they are a major source of dietary protein. In places like the United States, human consumption of insects is limited to expensive novelty items such as

chocolate-covered ants, fried grasshoppers, and the like.

Certain rare and showy butterflies and beetles are sought after so that they have considerable economic worth. Conservation of rare and endangered insects incurs some expense as well. Habitat management to conserve rare insects is a valid and growing concern of economic entomologists.

Finally, insects have rendered invaluable service to science, and thus to humanity, as easily reared experimental animals for investigation of basic principles of genetics, biochemistry, development, and behavior. Pomace flies (genus *Drosophila*) have been extremely useful in this regard, and even cockroaches and houseflies are helpful in testing the effectiveness of new chemical insecticides and other insect control methods. See INSECTA. David J. Horn

Bibliography. R. H. Davidson and W. F. Lyon, *Insect Pests of Farm, Garden, and Orchard*, 8th ed., 1987; D. J. Horn, *Ecological Approach to Pest Management*, 1988; R. L. Metcalf and W. H. Luckmann, *Introduction to Insect Pest Management*, 3d ed., 1994; R. E. Pfadt, *Fundamentals of Applied Entomology*, 4th ed., 1985.

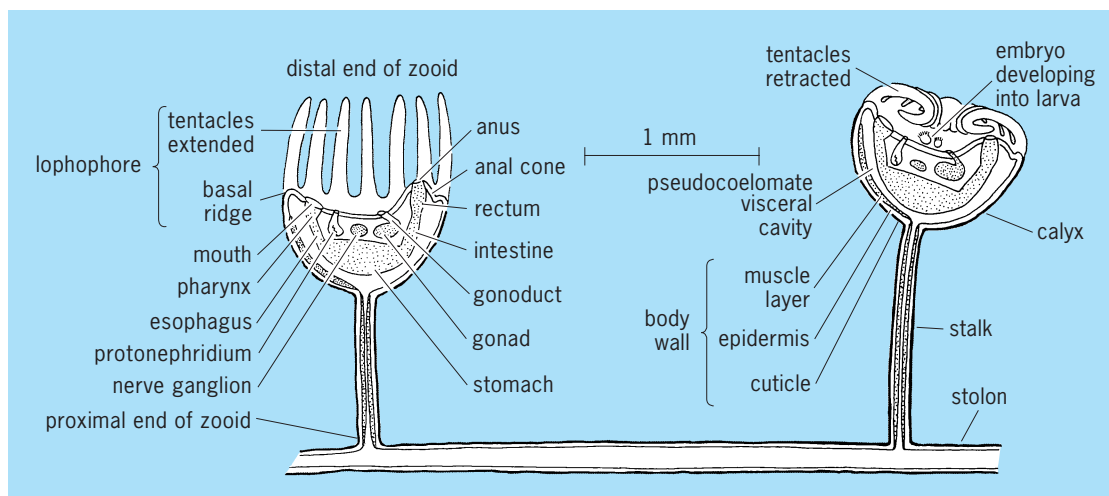
Entoprocta

A phylum of sessile, aquatic, often colonial invertebrates having a looped gut with both mouth and anus situated inside a cirlet of tentacles, a pseudocoelomate body cavity, and no hardened skeleton.

Form and function. Entoprocts, for example, *Pedicecellina*, are mainly marine, and may be encountered as small colonies on seaweeds and bryozoans near low tidemark. The colonies consist of threadlike, creeping stolons and erect zooids arising at intervals from them (see *illus.*). Noncolonial forms have similar zooids but lack the stolons. Each zooid is goblet-shaped, consisting of a basal stalk supporting a rounded or pyriform viscera-containing dilatation, the calyx. The zooids (or individuals) are usually in the range 0.016–0.08 in. (0.4–2.0 mm), but may reach 0.2 in. (5 mm) in species with long stalks. Usually undifferentiated, the stalk has alternating flexible and rigid sections in *Barentsia*.

The body wall comprises a cuticle-covered epidermis above a layer of muscle fibers. These fibers are particularly developed in the stalk. The body cavity, morphologically a pseudocoel, contains connective tissue and fluid. In the colonial forms, cells at the junction of stalk and calyx contract rhythmically, perhaps pumping or compressing fluid.

The U-shaped alimentary canal occupies most of the calyx, below a cirlet of 6–36 ciliated tentacles. The entoprocts feed on particles filtered from the surrounding water. Long lateral cilia drive water inward between the tentacles, while shorter frontal cilia (that is, on the inner face) convey particles trapped in mucus down the tentacles. These frontal tracts join an annular ciliated groove, which dips into the mouth. The stomach is globular; the anus is raised on a papilla, leaving a space or atrium between



Morphological features of entoproct zooids, illustrated schematically in vertical cross section.

the rectum and the mouth. Two protonephridia discharge into the atrium, except in the fresh-water *Urnatella*, which has many nephridia of known osmoregulatory function. Adverse stimuli cause the tentacles to fold inward; the calyx is then closed by a muscular membrane.

Ecology and life history. In one order, Pedicellinida, zooids are colonial and free-living; in the other, Loxosomatida, they are noncolonial and live in association with various other invertebrates, for example, tubicolous polychaetes, sipunculids, sponges, and bryozoans. The loxosomatids are energy commensals, collecting food from the respiratory or feeding current of the host animal.

Entoprocts reproduce asexually and sexually. Loxosomatids bud off a succession of daughter zooids from the paired blastogenic zones on the calyx. These disperse only a short way from the parent, giving rise to characteristically dense aggregations on the host. Pedicellinids bud from the stolon: daughter zooids do not break free, so that a colony develops.

Dispersal is by means of sexually produced larvae. Zooids may be unisexual or hermaphroditic. Gonads are paired, with the ducts uniting before opening into the coelom. Fertilization is internal. Young embryos become attached to a part of the atrial wall, via which they obtain nourishment. The larva is possibly of the trochophore type, but the ciliary girdle is basal rather than equatorial. The *Pedicellina* larva has obvious similarity to an inverted and nontentaculate adult calyx. It attaches with mouth and anus toward the substrate, and metamorphosis involves reorientation of the gut to the adult position. Loxosomatid larvae may metamorphose or give rise directly to external or internal zooid buds, the latter type being released by rupture of the larva.

History and classification. The name Entoprocta was introduced in 1869 as a subdivision of the phylum Bryozoa, but was recognized as an independent phylum in 1888. Despite some possible similarities in larval morphology and metamorphosis, the differences between Entoprocta and Ectoprocta are numerous and fundamental. See BRYOZOA.

Alternative names are Calysozoa and Kamptozoa. In view of the confusing similarity between the words Entoprocta and Ectoprocta, preference should be given to the distinctive and apposite name Calysozoa. The phylum is divided into two orders, Loxosomatida (or Solitaria), with about 100 species, and Pedicellinida (or Coloniales), with about 30. No fossils are known.

John S. Ryland

Entropy

A function first introduced in classical thermodynamics to provide a quantitative basis for the common observation that naturally occurring processes have a particular direction. Subsequently, in statistical thermodynamics, entropy was shown to be a measure of the number of microstates a system could assume. Finally, in communication theory, entropy is a measure of information. Each of these aspects will be considered in turn. Before the entropy function is introduced, it is necessary to discuss reversible processes.

Reversible processes. Any system under constant external conditions is observed to change in such a way as to approach a particularly simple final state called an equilibrium state. For example, two bodies initially at different temperatures are connected by a metal wire. Heat flows from the hot to the cold body until the temperatures of both bodies are the same. As another example, a vessel containing a gas is connected through a stopcock to an evacuated vessel. When the stopcock is opened, the gas expands to fill the whole of the available space uniformly. It is common experience that the reverse processes never occur if the systems are left to themselves; that is, heat is never observed to flow from the cold to the hot body, nor will the gas compress itself into one of the vessels. Max Planck classified all elementary processes into three categories: natural, unnatural, and reversible.

Natural processes do occur, and proceed in a direction toward equilibrium. Unnatural processes move

away with equilibrium and never occur. If $A \rightarrow B$ is a natural process between states A and B, then $B \rightarrow A$ is an unnatural process. A reversible process is an idealized natural process that passes through a continuous sequence of equilibrium states. Consider the evaporation of a liquid in the presence of its vapor at a pressure P . Let the equilibrium vapor pressure of the liquid be p . If $P < p$, liquid evaporates as a natural process. If $P > p$, evaporation is an unnatural process and will not occur; indeed, the opposite process—condensation—will take place. Finally, if $P = p$, both processes of condensation and evaporation are reversible and can be initiated by a very slight increase or decrease in the external pressure P .

A useful idea is that a reversible process may be exactly reversed by an infinitesimal change in the external conditions. If a hot object is placed adjacent to a much colder object, the heat-flow direction cannot be reversed by small changes in the temperature of either object. In reversible processes, work is accomplished through small pressure differences, and heat transfer occurs through small temperature differences.

Entropy function. The state function entropy S uses the foregoing discussion on a quantitative basis. The function is not derived in this article; but, rather, some of its properties are stated, and its implications are discussed mainly by example. Entropy is related to q , the heat flowing into the system from its surroundings and to T , the absolute temperature of the system. The important properties for this discussion are:

1. $dS > q/T$ for a natural change. $dS = q/T$ for a reversible change. It is necessary to introduce both S and T together. A formal derivation would show T^{-1} as an integrating factor leading to the complete differential dS .

2. The entropy of the system S is made up of the sum of all the parts of the system so that $S = S_1 + S_2 + S_3 \dots$. See HEAT; TEMPERATURE; THERMODYNAMIC PRINCIPLES.

Heat flow. Consider two bodies, α and β , at different temperatures separated by an adiabatic (no heat transfer) wall. If the two bodies are connected by a fine wire that allows a small heat flow q from α to β , then $dS_\alpha = -q/T_\alpha$ and $dS_\beta = q/T_\beta$.

For the whole system, Eq. (1) holds. If $T_\alpha > T_\beta$,

$$dS = dS_\alpha + dS_\beta = q \left(\frac{1}{T_\beta} - \frac{1}{T_\alpha} \right) \quad (1)$$

$dS > 0$, and heat flows from α to β as a natural process. The process could be continued until $T_\alpha = T_\beta$ and $dS = 0$.

Once the constraint of the adiabatic wall is abrogated, the entropy increases to a maximum value, and T_α becomes equal to T_β . This is a special case of the most important notion in thermodynamics; that is, the system will assume that equilibrium state which maximizes the entropy at constant energy, consistent with the constraints. See HEAT TRANSFER.

Nonconservation of entropy. In his study of the first law of thermodynamics, J. P. Joule caused work to

be expended by rubbing metal blocks together in a large mass of water. By this and similar experiments, he established numerical relationships between heat and work. When the experiment was completed, the apparatus remained unchanged except for a slight increase in the water temperature. Work (W) had been converted into heat (Q) with 100% efficiency. Provided the process was carried out slowly, the temperature difference between the blocks and the water would be small, and heat transfer could be considered a reversible process. The entropy increase of the water at its temperature T is $\delta S = Q/T = W/T$.

Since everything but the water is unchanged, this equation also represents the total entropy increase. The entropy has been created from the work input, and this process could be continued indefinitely, creating more and more entropy. Unlike energy, entropy is not conserved. See CONSERVATION OF ENERGY.

Although the heat transfer is considered to be reversible in order to calculate the entropy increase, the overall process of converting work into heat is irreversible. The frictional process that converts kinetic energy into the heat of the metal blocks is a natural process. In fact, the impossibility of the reverse process is Lord Kelvin's statement of the second law of thermodynamics. Heat cannot be completely converted into work without other changes occurring in the surroundings. For example, a gas in a cylinder can be expanded reversibly by extracting heat from a large constant-temperature bath. All of the heat extracted from the bath is converted into work, but eventually the pressure of the gas system would be reduced to an unusable level. The system has changed, and the process cannot continue indefinitely. If one tries to convert heat into work through a system undergoing a cycle so that the system will return to its initial state, one finds that only a portion of the heat input does work and that the remainder must be rejected to a lower temperature; this is just the process which takes place in a heat engine. See THERMODYNAMIC CYCLE; THERMODYNAMIC PROCESSES.

Degradation of energy. Energy is never destroyed. But in the Joule friction experiment and in heat transfer between bodies, as in any natural process, something is lost. In the Joule experiment, the energy expended in work now resides in the water bath. But if this energy is reused, less useful work is obtained than was originally put in. The original energy input has been degraded to a less useful form. The energy transferred from a high-temperature body to a lower-temperature body is also in a less useful form. If another system is used to restore this degraded energy to its original form, it is found that the restoring system has degraded the energy even more than the original system had. Thus, every process occurring in the world results in an overall increase in entropy and a corresponding degradation in energy. R. Clausius stated the first two laws of thermodynamics as: "The energy of the world is constant. The entropy of the world tends toward a maximum."

Increasing entropy and mixing. Once the atomic theory of matter is accepted, the entropy concept

can be made much clearer. It is then found through statistical thermodynamics that the increase of entropy toward its maximum value at equilibrium corresponds to the change of the system toward its most probable state consistent with the constraints. The most probable state represents the most mixed or most random state. Mixing must be given a broad interpretation which includes particle or configurational mixing, and spreading of energy over the particles or thermal mixing. Diffusion of one gas into another represents obvious configurational mixing and increased entropy. Irreversible expansion of a gas represents configurational mixing of the molecules over the available space. Heat flow represents spreading of the kinetic energy between the particles. Friction spreads the kinetic energy of the body over the constituent particles. Sometimes the energy-spread entropy increase and the configurational entropy increase are not compatible, and a compromise is struck. A subcooled liquid adiabatically crystallizes to a lower configurational entropy but gains even more entropy through the additional energy levels made available. The same sort of behavior occurs in partially miscible liquids—some configurational entropy is sacrificed in order to gain a large amount of energy-spread entropy. See STATISTICAL MECHANICS.

Absolute entropy. The third law of thermodynamics (Nernst's heat theorem) refers to the vanishing of entropy at zero temperature. In 1912 Planck proposed that the theorem applied to pure crystalline solids. However, the theorem is now known to be applicable to gases and, by all reasonable expectation, is applicable to any system. Thus, any substance at finite temperatures has an absolute entropy, the value of which can be determined from either calorimetric or spectroscopic data. Absolute entropies, together with thermochemical data, are very useful in the calculation of equilibrium compositions of reaction systems.

The statistical viewpoint is that a thermodynamic state at finite temperatures corresponds to many microstates. During an observation the microstates of a system undergo continuous rapid transitions. Since entropy is proportional to the logarithm of the number of available microstates, the Nernst theorem implies that the thermodynamic state at zero temperature corresponds to a single microstate. At zero temperature, even a ferromagnetic material should exist in a single state, fully magnetized in a direction determined by its inevitable interactions with the environment.

William F. Jaep

Measure of information. The probability characteristic of entropy leads to its use in communication theory as a measure of information. The absence of information about a situation is equivalent to an uncertainty associated with the nature of the situation. This uncertainty, designated H , is the entropy of the information about the particular situation, Eq. (2),

$$H(p_1, p_2, \dots, p_n) = - \sum_{k=1}^n p_k \log p_k \quad (2)$$

where p_1, p_2, \dots, p_n are the probabilities of mutu-

ally exclusive events, the logarithms are taken to an arbitrary but fixed base, and $p_k \log p_k$ always equals zero if $p_k = 0$. For example, if $p_1 = 1$ and all other p 's are zero, the situation is completely predictable beforehand; there is no uncertainty and so the entropy is zero. In all other cases the entropy is positive. See INFORMATION THEORY.

In introducing entropy of an information space, C. E. Shannon described a source of information by its entropy H in bits per symbol. The ratio of the entropy of a source to the maximum rate of signaling that it could achieve with the same symbols is its relative entropy. One minus relative entropy is the redundancy of the source.

Frank H. Rockett

Bibliography. H. B. Callen, *Thermodynamics and an Introduction to Thermostatistics*, 2d ed., 1985; K. G. Denbigh, *Principles of Chemical Equilibrium*, 4th ed., 1981; R. M. Gray, *Entropy and Information Theory*, 1991; J. N. Kapur, *Measures of Information and Their Application in Science and Engineering*, 1994; K. S. Pitzer, *Thermodynamics*, 3d ed., 1995; R. C. Tolman, *The Principles of Statistical Mechanics*, 1980; K. Wark, *Thermodynamics*, 6th ed., 2000.

Environment

The sum of all external factors, both biotic (living) and abiotic (nonliving), to which an organism is exposed. Biotic factors include influences by members of the same and other species on the development and survival of the individual. Primary abiotic factors are light, temperature, water, atmospheric gases, and ionizing radiation, influencing the form and function of the individual.

For each environmental factor, an organism has a tolerance range, in which it is able to survive. The intercept of these ranges constitutes the ecological niche of the organism. Different individuals or species have different tolerance ranges for particular environmental factors—this variation represents the adaptation of the organism to its environment. The ability of an organism to modify its tolerance of certain environmental factors in response to a change in them represents the plasticity of that organism. Alterations in environmental tolerance are termed acclimation. Exposure to environmental conditions at the limit of an individual's tolerance range represents environmental stress. See ADAPTATION (BIOLOGY); ECOLOGY; PHYSIOLOGICAL ECOLOGY (ANIMAL); PHYSIOLOGICAL ECOLOGY (PLANT).

Abiotic factors. All the physical factors which affect an organism constitute the abiotic environment.

Light radiation. The spectrum of electromagnetic radiation reaching the Earth's surface is determined by the absorptive properties of the atmosphere (Fig. 1). Biologically, the most important spectral range is 300–800 nanometers, incorporating ultraviolet, visible, and infrared radiation. Visible light provides the energy source for most forms of life. Light absorbed by pigment molecules (chlorophylls, carotenoids, and phycobilins) is converted into chemical energy through photosynthesis. Photosynthetic bacteria

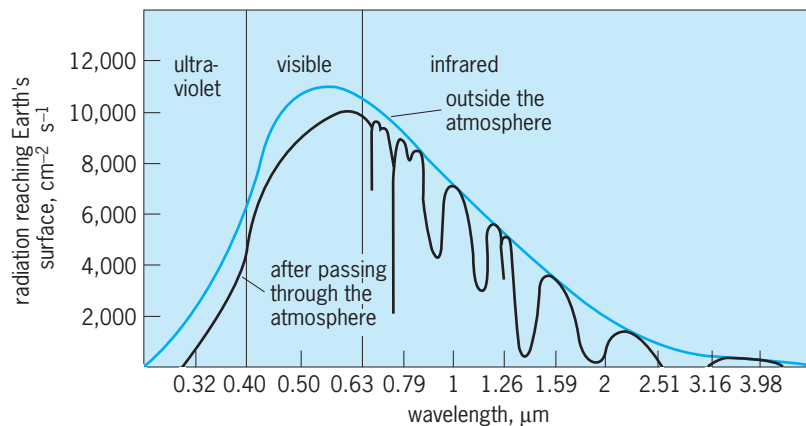


Fig. 1. Spectrum of electromagnetic radiation reaching the Earth's surface. Light quality is determined by the spectrum of light available from the Sun and by the absorptive properties of the atmosphere.

containing bacteriochlorophyll can absorb light into the far-red portion of the spectrum, up to 1000 nm, allowing them to survive in locations where the availability of visible light is low. Light availability is especially important in determining the distribution of plants. Plants receiving too little light (low light stress) may be unable to maintain an average photosynthetic rate high enough to support net growth. Excess light (high light stress) causes damage to leaf tissue, which in extreme cases may kill the plant. Photosynthetic organisms can exist within a wide range of light intensities. Full sunlight in the tropics is around $2000 \mu\text{mol photons} \cdot \text{m}^{-2} \cdot \text{s}^{-1}$. Photosynthetic organisms have survived in locations where the mean light is as low as 0.005% of this value. *See* INSOLATION; PHOTOSYNTHESIS; SOLAR RADIATION.

In addition to providing energy, light is important in providing an organism with information about its surroundings. The human eye, for example, is able to respond to wavelengths of light between 400 and 700 nm—the visible range. Within this range, sensitivity is greatest in the green part of the spectrum. This is the portion of the spectrum that plants absorb least, and so is the principal part of the spectrum to be reflected. The detection of light is not confined to animals with recognizable eyes. Many photosynthetic bacteria exhibit phototaxis, movement toward or away from light. Plants use light cues to determine the direction of growth and the extent of stem elongation. In such responses, light quality sometimes may be more important than the intensity of light. Plants possess photoreceptors that respond to ultraviolet, blue, red, and far-red light. Of particular importance in plant development is the ratio of red to far-red light. Red light is absorbed efficiently by leaves, whereas they are relatively transparent to far-red light. Hence, a low red to far-red ratio indicates that the leaf is being shaded by another leaf on the same or another plant. Such signals allow plants to optimize their growth form to maximize light capture and minimize shading from neighbors.

Temporal variation in light also provides an important stimulus. Life forms from bacteria upward

are able to detect and respond to daily light fluctuations. Such a response may be directly controlled by the presence or absence of light (diurnal rhythms) or may persist when the variation in light is removed (circadian rhythms). In the latter case, regulation is through an internal molecular clock, which is able to predict the daily cycle. Such circadian clocks are normally reset by light on a daily basis. Processes controlled by circadian clocks range from the molecular (gene expression) to the behavioral (for example, sleep patterns in animals or leaf movements in plants). The ability to detect changes in day length is also essential to many biological processes. Day length changes are used by organisms to predict seasonal changes. For example, shortening of days, indicating the onset of winter, is used as a cue to trigger hibernation or migration in animals, and flowering, seed set, and leaf loss in plants. *See* PHOTOPERIODISM.

Ultraviolet radiation represents a form of stress. It has the ability to break chemical bonds and so may lead to damage to proteins, lipids, and nucleic acids. Damage to deoxyribonucleic acid (DNA) may result in genetic mutations. The amount of ultraviolet radiation reaching the Earth's surface is thought to be increasing due to depletion of the ozone layer in the stratosphere. Ozone is responsible for absorbing a large proportion of ultraviolet radiation reaching the outer atmosphere. As ozone is destroyed by the action of pollutants such as chlorofluorocarbons, the proportion of ultraviolet radiation reaching the surface of the Earth rises.

Water. Water is ubiquitous in living systems, as the universal solvent for life. A sufficient supply of water is essential for biological activity. Many organisms have evolved the ability to survive prolonged periods in the total absence of water, but this is achieved only through the maintenance of an inactive state. Life first evolved in water and organisms later migrated onto land. Water availability remains a primary environmental factor limiting survival on land. Primitive land organisms possess little or no ability to conserve water within their cells and are termed poikilohydric. Examples include amphibians and primitive plants such as mosses and liverworts. These are confined to places where water is in plentiful supply, for example close to ponds, or they must be able to tolerate periods of desiccation. Lichens can survive total water loss and rapidly regain activity upon rewetting. Such organisms must be able to minimize the damage caused to cellular structures when water is lost. Dehydration causes irreversible damage to membranes and proteins. This damage can be prevented by the accumulation of protective molecules termed compatible solutes.

Homeohydric organisms possess a waterproof layer that restricts the loss of water from the cells. Such waterproofing is never absolute, as there is still a requirement to exchange gas molecules and to absorb organic or mineral nutrients through a water phase. Hence, lungs and guts in mammals or stomata and roots in plants represent potential sites of water loss. Nevertheless, water conservation allows

organisms to live in environments in which the water supply is extremely low. In extremely arid environments, behavioral adaptations may allow the water loss to be minimized. Animals may be nocturnal, emerging when temperatures are lower and hence evaporation minimized. Cacti possess a form of photosynthesis, crassulacean acid metabolism (CAM), that allows them to separate gas exchange and light capture. Carbon dioxide is taken up and stored during the night, allowing the stomata to remain closed during the daytime. At the other extreme, inundation with water may be detrimental to organisms not adapted to wet environments. Asphyxiation through drowning in animals is paralleled in plants, where waterlogging of soils may lead to oxygen stress in the roots. Adaptations to survival in waterlogged soils includes air spaces in stems, aerenchyma, that allow stems to act like breathing tubes passing oxygen down to the roots and into the soil. *See* GROUND-WATER HYDROLOGY; OSMOREGULATORY MECHANISMS; PLANT-WATER RELATIONS.

Temperature. Temperature is a determinant of survival in two ways: (1) as temperatures decrease, the movement of molecules slows and the rate of chemical reactions declines; (2) temperature determines the physical state of water.

The slowing of metabolic activity at low temperatures is illustrated in reptiles. Such poikilothermic animals, unable to maintain their internal temperature, are typically inactive in the cold of morning. They bask in the sun to increase their body temperature and so become active. High temperatures will cause the three-dimensional structure of proteins to break down, preventing the organisms from functioning. Organisms adapted to extremely high temperatures need more rigid proteins that maintain their structure. Temperature also affects the behavior of cell membranes, made up of lipids and proteins in a liquid crystalline state. At low temperatures, the membrane structure becomes rigid and liable to break. At high temperatures, it becomes too fluid and again liable to disintegrate. In adapting to different temperatures, organisms alter the composition of the lipids in their membranes, whose melting temperature is thereby changed. This outcome also applies to storage lipids. Hence, cold-water fish are a useful source of oils, whereas mammals, with their higher body temperature, contain fats. The effect of temperature on membranes is thought to be a key factor determining the temperature range that an organism is able to survive.

The effect of temperature on the physical state of water is essential to determining the availability of that water to organisms. Poikilothermic organisms may find that the water in their cells begins to freeze at low temperatures. The temperature at which this occurs depends on a number of factors. For example, the presence of solutes in the cells depresses the freezing temperature. Water in cells can also remain liquid at temperatures below the freezing point (a phenomenon termed supercooling) depending on the absence of sites where ice crystals can form. Finally, certain species can survive total freezing, down

to 77 K (-196°C or -321°F), through the prevention of ice crystal formation altogether. Water is frozen in an amorphous state, without forming crystals, which would otherwise damage cellular structures. To survive low temperatures, cells must be able to survive desiccation, and so low-temperature tolerance involves the formation of compatible solutes. High temperatures increase the rate of evaporation of water. Hence, where water supply is limiting, an organism's ability to survive high temperatures is impaired.

Mammals and birds, homeothermic organisms, are able to regulate their internal temperature, limiting the effects of external temperature variations. Temperature still acts as an environmental constraint in such organisms, however. Cooling is achieved through sweating and hence loss of water. Heat is produced through the metabolism of food, and hence survival in cold climates requires a high metabolic rate. *See* CRYPTOBIOSIS; HIBERNATION AND ESTIVATION; THERMOREGULATION.

Atmospheric gases. The atmosphere on Earth is thought to be determined to a large extent by the presence of life. At the same time, organisms have evolved to survive in the atmosphere as it is. The atmospheric constituents with the most direct biological importance are oxygen (O_2) and carbon dioxide (CO_2). Oxygen makes up approximately 20% of the atmosphere and is due to the occurrence of oxygenic photosynthesis. This process involves the simultaneous uptake of CO_2 to make sugars. Aerobic respiration involves the reverse of this process, the release of CO_2 and the uptake of O_2 to form water. Hence, the current atmosphere represents the balance of previous biological activity. For most terrestrial organisms, neither CO_2 nor O_2 is limiting in the atmosphere; however, the need to get either or both of these gases to cells may represent a limitation on size or on the ability to tolerate water stress. Limitation of either gas may be important in aquatic environments, where the concentration of each is significantly lower. Conversely, oxygen may be toxic to certain anaerobic bacteria, which lack mechanisms to prevent oxidative damage. Such organisms are confined to anaerobic environments, such as waterlogged soils.

Nitrogen is also required by all organisms but cannot be used by most in the gaseous form. Nitrogen fixation, the conversion of N_2 gas into a biologically useful form, occurs in some species of bacteria and cyanobacteria or may be caused by lightning.

Atmospheric gases are important in determining the climate and the light environment. Absorption of electromagnetic radiation by the atmosphere determines the spectrum of light reaching the Earth's surface. Absorption and reflectance of infrared radiation by greenhouse gases such as CO_2 and water vapor regulate temperature. *See* PHOTORESPIRATION; RESPIRATION.

Ionizing radiation. Ionizing radiation entering the atmosphere as cosmic rays or generated through the decay of radioactive material may have a significant impact on the survival of organisms. High-level

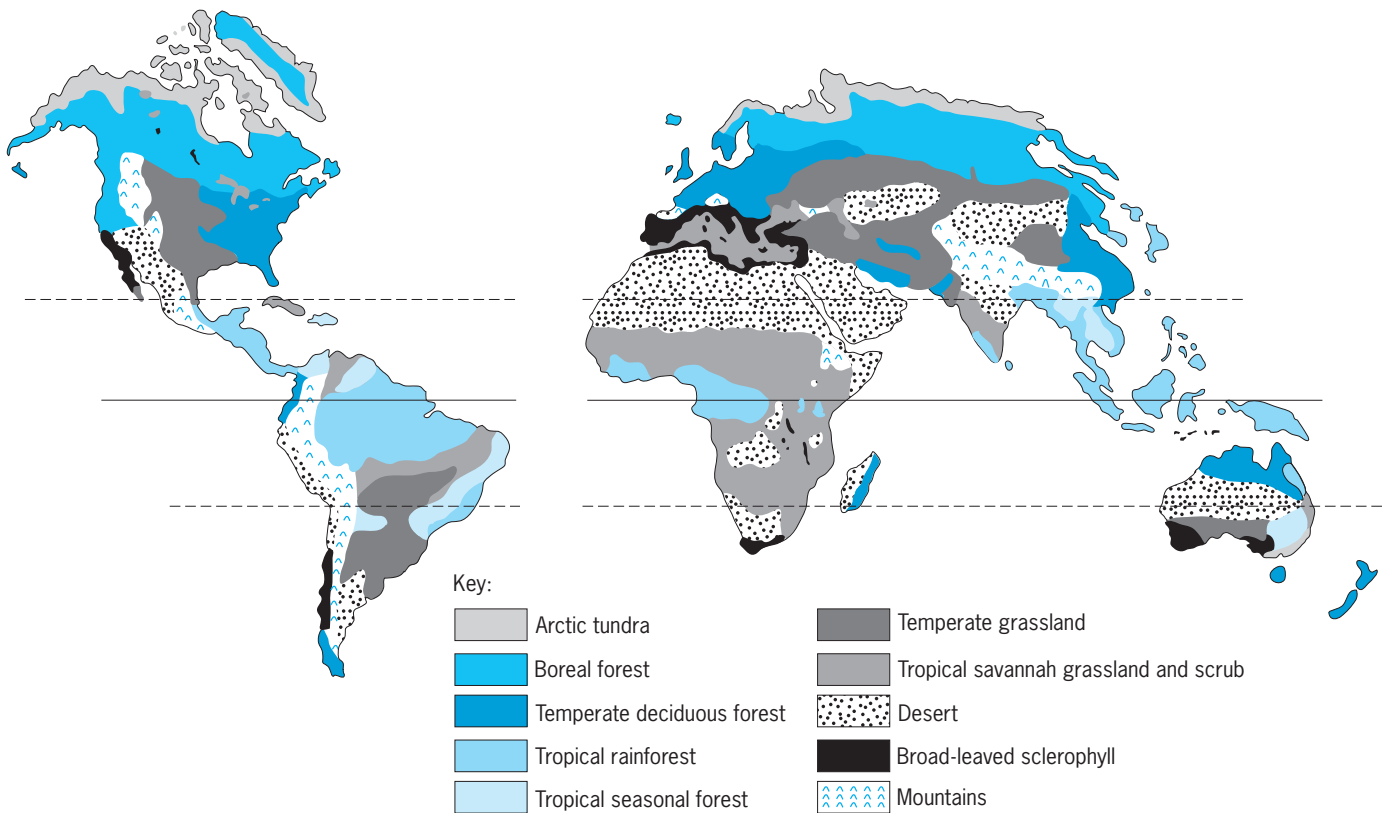


Fig. 2. Global distribution of major vegetation types (biomes).

exposure to radiation may directly damage tissues, ultimately killing the individual. Lower levels of radiation may cause damage to DNA, leading to the accumulation of genetic mutations. These mutations are frequently deleterious, but they also give rise to beneficial genetic variations. See PHOTORESPIRATION; RESPIRATION.

Other factors. Among other environmental factors determining the range and distribution and form of organisms are mechanical stimuli such as wind or water movement, and the presence of metals, inorganic nutrients, and toxins in the air, soil, or food. Wind, for example, may be a key factor in determining vegetation height, with plants tending to be shorter but having stronger stems when growing in exposed locations. Metals and other ions in the soil or food may be essential to the survival of an organism, as is true of magnesium, iron, and calcium. The concentration of these ions may represent a limit on the organisms in a particular location. Alternatively, ions may be nonessential or toxic (such as lead or aluminum). All inorganic ions may be toxic if present in sufficient quantities. Environmental factors particularly relevant in aquatic environments include water pressure and salinity.

Biotic factors. The biotic environment of an individual is made up of members of the same or other species. Intraspecific interactions involve the need to breed with other individuals, to gain protection through living in a group, and to compete for resources such as food, light, nutrients, and space. Low population densities exert a negative effect on

an individual due to failure to breed or exposure to predation. High population densities are stressful as a result of increased competition. The optimal population density depends on the availability of resources and on the behavior, size, and structure of the organism. Interspecific interactions may also be positive or negative. For example, symbiotic relationships involve the mutual benefit of the individuals involved, whereas competition for resources is deleterious to both. Symbiotic relationships may be positive to one species and negative to the other. For example, a parasite gains benefit from its host but exerts a negative influence on that individual. Such relationships are not simple. Although predation exerts a negative influence on the population as a whole, the success of an individual may be enhanced if a predator removes one of its conspecific competitors.

The impact of humans on natural environments may be considered a special case among biotic factors. Humans alter their environment in ways that exceed the impact of all other organisms. For example, living in heated houses creates a hot, relatively dry environment that can be exploited by other organisms. Hence, cockroaches have exploited humans to extend their ecological range. Heating involves the release of greenhouse gases into the atmosphere, contributing to climate alterations over the entire planet. This in turn has impacts on the distribution of all other species. The release of pollutants into the environment brings organisms into contact with stresses to which they were not previously

exposed. This causes the evolution of new varieties, eventually perhaps new species, adapted to the polluted environments. See AIR POLLUTION; BIOSPHERE; HUMAN ECOLOGY; WATER POLLUTION.

Limiting factors and environmental stress. For any given organism, it is often possible to identify a factor in the environment that limits survival and growth. The limiting factor may change through time. Such a change may cause the organism to be at the limit of or outside its tolerance range for that or another environmental factor. In such cases, the organism is said to suffer stress. The relationship between the limitation of one factor and the stress observed may not be simple. For example, a plant suffering from water limitation closes its stomata, restricting the entry of CO₂ into the leaf and so reducing photosynthesis. The plant may then be vulnerable to light stress, due to light absorption exceeding the rate at which the energy is used.

If the stress to which an individual is exposed is extreme, it may result in irreversible damage and death. Exposure to moderate stress, however, results in a period of acclimation within the organism that allows it to adjust to the new conditions. Organisms exposed gradually to new conditions usually have a higher chance of survival than those exposed suddenly. This concept is familiar to gardeners who gradually expose plants to cold, hardening them off over several weeks before leaving them outdoors totally unprotected. See POPULATION ECOLOGY.

Environment and biomes. Where a particular environmental factor (or combination of factors) dominates the growth and development of organisms, it is often found that the adaptations and gross features of the landscape will be the same, even when the actual species are different. Thus, mediterranean vegetation is found not only around the Mediterranean Sea but also in California and South Africa, where the conditions of hot dry summers and warm wet winters occur. Regions with similar environmental conditions are classed as biomes (Fig. 2). The occurrence of such global vegetation types clearly illustrates the role played by the environment in determining the form and function of individual species.

Gites Johnson

Bibliography. R. Brewer, *The Science of Ecology*, Saunders College Publishing, 1994; P. Colinvaux, *Ecology 2*, Wiley, 1993; R. M. M. Crawford, *Studies in Plant Survival*, Blackwell Scientific, 1989; B. Freedman, *Environmental Ecology*, Academic Press, 1995; E. P. Odum, *Ecology: A Bridge Between Science and Society*, Sinauer Associates, 1997.

Environmental engineering

The division of engineering concerned with the environment and management of natural resources. The environmental engineer places special attention on the biological, chemical, and physical reactions in the air, land, and water environments and on improved technology for integrated management sys-

tems, including reuse, recycling, and recovery measures.

Scope. Environmental engineering began with the need for acceptable drinking water and for management of liquid and solid wastes. In fact, the terms public health and sanitary were used to identify the engineers seeking solutions for eliminating waterborne disease in the 1800s. Later, urbanization and industrialization contributed significantly to the formation of unsanitary conditions in many areas. Abatement of air and land contamination became new challenges for the environmental engineer, followed by toxic-waste and hazardous-waste concerns. The environmental engineer is also instrumental in the mitigation and protection of wildlife habitat, preservation of species, and the overall well-being of ecosystems.

During project development, environmental engineers conduct a life-cycle analysis (cost-effectiveness analysis) of project alternatives that consider the capital and the operation and maintenance costs over the life of the project. The National Environmental Policy Act and other legislation, regulations, and policies contain provisions to expand the alternative analysis to reduce ecological, health, and safety risks during project development. This process embraces the concepts often called green engineering.

Green engineering is the early application of environmentally favorable engineering project concepts in the design and development of processes, products, and facilities. Embracing the basic concepts of pollution prevention and waste minimization, green engineering considers innovative and alternative technologies. The project should be feasible and economical, while minimizing the generation of air, water, and solid-waste pollution at the source, as well as the risk to public health and the environment.

The U.S. Environmental Protection Agency has a Design for the Environment (DFE) program that includes a life-cycle assessment. This process provides a tool to examine the environmental impacts of products over their entire life cycle from materials acquisition to manufacturing, use, and disposition. Design for the Environment provides decision-makers with information, tools, and incentives to make informed decisions that integrate risk, performance, and cost issues. The life-cycle analysis evaluates impacts on human health, atmospheric resources/air quality, water quality, ecological health, and natural-resource consumption.

Education of professionals. Traditionally, environmental engineers drew their basic education and training from civil engineering programs. In order to broaden their perspective and capabilities, contemporary environmental engineers undertake course work and postgraduate training in a number of professional areas, including biology, epidemiology, chemical engineering, chemistry, and hydrology. Since environmental engineers deal with sensitive public issues, training in public education, risk assessment, public policy, and social sciences is desirable.

After satisfying experience and testing requirements, the environmental engineer obtains

professional engineering registration. Professional associations of interest to the environmental engineer include the American Academy of Environmental Engineers, the Air and Waste Management Association, the American Water Works Association, the Water Environment Federation, and the Solid Waste Association of North America.

Professional activities. Principal areas of employment for practicing environmental engineers include consulting, industry, and government. Other environmental engineers are academic or research faculty, or they direct the development and production of pollution-control equipment.

Consulting environmental engineers apply their knowledge and experience under contract to public and private clients. Their services may range from studies to preliminary design, final design, construction management, and operation and maintenance services. Specific areas of service include site evaluations, environmental impact studies and assessments, assistance in obtaining permits, and expert witness testimony.

Within the industrial sector, there is considerable demand for environmental professionals to achieve and maintain compliance with environmental statutes and regulations. Increasing public sensitivity to the environment and demands for a clean environment have expanded industry's need for environmental engineering technology.

Environmental engineers in public service provide technical expertise in all levels of government. Though most publicized for their role in the development of regulations and in enforcement, they also are involved in numerous research and development programs and technology transfer activities.

The principal environmental engineering specialties are air-quality control, water supply, wastewater disposal, stormwater management, solid-waste management, and hazardous-waste management. Other specialties include industrial hygiene, noise control, oceanography, and radiology.

Air-quality control. The term air pollution is used to describe the presence in the atmosphere of one or more contaminants in quantities, and with characteristics, that will be injurious to, or unreasonably interfere with, public health and welfare or other natural environmental processes. The extent of air-pollution problems ranges from relatively small areas, such as an industrial park impacted by one or more emission sources, to urban areas impacted by a number of contaminant sources.

Contaminants are categorized as particulate matter and gases and their associated forms, including dust, smoke, fumes, mist, and vapor. The primary gaseous air contaminants are carbon monoxide, hydrocarbons, nitrogen oxides, and sulfur oxides.

Meteorological and topographical factors contribute to the creation and continuation of air pollution under specific site conditions. Temperature inversions prevent upward diffusion, and very low wind speeds allow emissions to remain near their source. Some terrains cause emissions to follow specific patterns from one area to another.

Generally, sources of air contaminants may be classified as stationary, mobile, or fugitive. Respectively, they are attributed to point sources such as industrial stack emissions, transportation activities such as automobile emissions, and uncontrolled (fugitive) sources such as windblown dusts from stockpiles.

The environmental engineer is instrumental in developing particulate and gas controls for all sources of air contaminants. Source control is the first abatement method considered. Controls for particulates include settling chambers, inertial separators, wet scrubbers, and fabric filters. Gas controls include absorption, adsorption, condensation, flaring, and incineration. Other areas of practice address acid rain issues, fugitive emissions, odor control, indoor air quality, and noise abatement.

Indoor air quality is of increasing concern for health officials and the public. Industrial hygienists and various engineering disciplines are instrumental in identifying and controlling indoor air pollution. Environmental engineering principles can be used in technical assessments for identifying air pollutants and selecting source controls.

Indoor air quality control measures are very source-specific. The basic abatement strategies are source control, ventilation improvements, and air cleaners. Typically, source control is the most cost-effective abatement measure. *See* AIR POLLUTION; AIR POLLUTION, INDOOR.

Water supply. Historically, the environmental engineer has found methods to provide ample quantities of quality drinking water for domestic use as well as quality water for commercial and industrial uses. Water-supply issues include demand projections, quality requirements, surface-water and ground-water source evaluations, ground-water production, surface-water collection and storage, surface-water treatment, saline-water treatment, nonconventional water production, and treated-water distribution. *See* WATER RESOURCES; WATER SUPPLY ENGINEERING.

The typical surface-water treatment plant uses chemicals for enhancing removal of suspended solids and for disinfection. Physical treatment processes include simple settling and filtration. In sequence, the unit processes are rapid mix, coagulation, flocculation, sedimentation, filtration, and disinfection. Auxiliary systems are needed for chemical feed facilities and for sludge handling. Processes that are more specialized include carbon adsorption, ion exchange, and softening. *See* WATER SOFTENING; WATER TREATMENT.

Treatment of other sources generally requires site-specific determination of raw-water quality. Often, the quantity of demand will influence the cost-effective selection of treatment processes. Ground-water supplies require well development and treatment, such as aeration, softening, and disinfection. Likewise, brackish and saline waters require site-specific determination of treatment processes. Typical processes include membrane technology, such as reverse osmosis and electrodialysis. *See* RAW WATER; WATER DESALINATION.

The water distribution system includes service and distribution lines, transmission mains, and storage facilities. Elevated storage tanks with gravity distribution and ground storage tanks with distribution pumping are designed to provide the quantity and pressure required to satisfy system demands.

Wastewater disposal. Wastewater is the combination of liquid- and water-transported wastes from homes, commercial buildings, industrial facilities, and institutions, along with any ground-water infiltration and surface-water and stormwater flow that may enter the sewer system. As a minimum measure, treatment is required for suspended solids and dissolved organics. Special processes may be necessary to achieve removal of specific pollutants, such as phosphorus from a municipal source or heavy metals from an industrial plating facility. *See* SEWAGE COLLECTION SYSTEMS; SEWAGE DISPOSAL; WASTEWATER REUSE.

Minimum levels of treatment have been established by regulation. For example, in the United States 85% removal of oxygen-demanding organics and suspended solids followed by disinfection is the minimum level of treatment for domestic wastewaters. Additional treatment is dictated by the assimilative capacity of the receiving stream and by downstream water uses. *See* SEWAGE TREATMENT.

Physical processes are used to remove suspended solids that may damage or interfere with subsequent pumping and treatment units. Screens remove debris and other large solids, and gravity or aerated grit chambers capture sandy matter. Normally, gravity sedimentation is used to remove finer (organic) suspended solids. For special applications, centrifugation, dissolved air flotation, and filtration are used to remove suspended solids. *See* CENTRIFUGATION; FILTRATION; FLOTATION.

Dissolved organics generally are treated with biological processes. The more common systems are aerobic (with oxygen) and include aerobic-pond, trickling-filter, and activated-sludge processes. Concentrated wastes, such as primary sludges, or high-strength industrial wastewaters, such as meat-processing or brewery wastes, are considered for anaerobic (without oxygen) treatment processes. Sludges, principally from biological processes, require special handling. The sequence of processes includes stabilization, conditioning, dewatering, drying, and residual disposal. Land application and land-filling are the most practiced means of final disposal. Special concerns for land-applied and composted sludges arise because of the concentrating of contaminants, such as heavy metals, and presence of pathogens in these sludges. *See* SEWAGE SOLIDS.

Of the many other types of treatment process, one of the most important is natural systems, which historically included pond systems and, later, various modes of land application. Through technology, natural and constructed wetlands can provide high-quality effluents. *See* WETLANDS.

Stormwater management. While the study of stormwater management includes all elements of the hydrologic cycle, it focuses on how humans affect

the production, movement, and control of surface runoff. In a natural system, the rate of surface runoff is controlled by the rainfall rate, soil conditions, vegetation, and subsurface geology. *See* HYDROLOGY.

Urbanization creates large impervious areas that increase the quantity and peak rate of runoff. Rainfall then washes deposited materials directly into surface waters, causing stream pollution. Organics create oxygen demands, nutrients accelerate lake eutrophication, and heavy metals accumulate in bottom sediments. *See* EUTROPHICATION; STREAM POLLUTION; WATER POLLUTION.

Environmental engineers apply modern stormwater management practices to use natural and engineered systems to minimize environmental damage. A complete stormwater management program contains many elements, including on-site infiltration and detention, collection and transport systems, regional flood control, and major stream-channel improvements.

In the past, sanitary wastes and stormwater were collected in the same (combined) sewer. During heavy rainfalls, the sewers would overflow, creating water pollution problems. Many of these combined sewers have been separated, while the remaining overflows are treated as point sources. Stormwater management practices are used to lessen the requirements and costs for combined sewer overflow treatment facilities.

Solid-waste management. The collection, transport, processing, and disposal of solid wastes is an important area of environmental engineering. Solid wastes are those materials that are deemed by their owner to possess no value and are discarded. They are generated by almost every activity, and the amount varies by source, season, geography, and time. As land becomes more limited and regulation increases, the solid-waste generator and handler must employ new and improved technology to reduce the quantities of materials requiring disposal. Furthermore, recovery and reuse are important elements of solid-waste management.

Historically, solid-waste disposal consisted of open dumping. However, the modern method of disposal uses double-lined landfills with collection and controls for gases and leachate. Other disposal means include composting and various incineration processes, which also may be used for co-disposal of wastewater treatment sludges. These disposal means typically require controls for created pollutants, such as leachate and odor from compost operations, and chemical and particulate emission from incinerator combustion.

Recovery and reuse are practiced widely. Source or central-facility separation is used for a variety of products, including paper, glass, plastics, ferrous metals, and nonferrous metals. Also, refuse-derived fuels may be used for energy production, and yard wastes may be composted to produce a humus soil conditioner.

Integrated solid-waste management plans include recovery and reuse as essential elements. Source reduction techniques and cost-effective recovery of solid wastes for reuse in industrial production

improve energy recovery and production, provide direct economic benefits, and lessen the overall requirements for solid-waste processing and disposal.

Source reduction is a priority pollution prevention measure and includes design, manufacture, or use of materials to reduce the amount and/or toxicity of materials before they enter the solid-waste management process. Examples include managing nonproduct organic wastes such as food scraps and yard trimmings, extending the useful life of products to postpone disposal, designing packaging to reduce the quantity for waste disposal, reusing existing products or packaging, and segregating and collecting waste materials onsite for recycle at residential, commercial, and manufacturing generator sites.

Principal energy recovery technologies are incineration furnace walls with integrated boiler tubes (waterwalls); modular incineration with heat recovery boilers or heat exchangers; refuse-derived fuels for co-fired or dedicated boilers; pyrolysis into gaseous, liquid, or solid fuels; and anaerobic digestion or landfill gas recovery for methane fuel development.

Hazardous-waste management. Whether they are pesticides from agricultural lands, gasoline leaks from service stations, heavy metals from plating solutions, medical wastes from hospitals, or radioactive wastes from nuclear power plants, hazardous wastes are present throughout the world as by-products of growth in developing nations. They pose unreasonable risks to human health and safety, property value and use, and all other components of the environment.

Hazardous-waste treatment involves liquid-waste treatment, solid-waste treatment, solidification and stabilization, thermal destruction, and land disposal. Remedial action is characterized as surface-water control, air-pollution control, or in-place treatment.

A significant portion of the technology for the treatment and disposal of hazardous wastes is refinement or adaptation of proven practices in air-quality control, wastewater treatment, and solid-waste management. Also, the environmental engineer must understand hydrogeology to assess the subsurface disposition of hazardous wastes. *See* AQUIFER.

At a hazardous-waste treatment or disposal site, the routes of the waste into the environment must be obstructed. Principal routes include surface-water contamination from runoff or overflows; ground-water contamination from leaks or leachate; air contamination from open burning, evaporation, or fugitive dusts; fire and explosion; and health risks from human contact.

U.S. Environmental Protection Agency regulations define hazardous wastes in specific lists based on characteristics of ignitability, corrosivity, reactivity, and an analytic procedure known as EP toxicity. Household wastes, domestic sewage, and certain other wastes are excluded from these regulations. In the United States, regulations for hazardous wastes include requirements of a manifest system for tracking wastes from generators through storage, transport, treatment, and disposal. In addition to treat-

ment and disposal, it is necessary to provide the means to satisfy these regulatory controls as well as means of pollution prevention using waste reduction through process changes, including recovery and reuse. *See* CONSERVATION OF RESOURCES; HAZARDOUS WASTE; RADIOACTIVE WASTE MANAGEMENT.

Robert A. Corbitt

Bibliography. American Water Works Association et al., *Water Treatment Plant Design*, 4th ed., 2004; R. A. Corbitt (ed.), *Standard Handbook of Environmental Engineering*, 2d ed., 1999; H. M. Freeman (ed.), *Standard Handbook of Hazardous Waste Treatment and Disposal*, 2d ed., 1997; Metcalf and Eddy, Inc., *Wastewater Engineering: Treatment, Disposal, and Reuse*, 4th ed., 2002; J. Pichtel, *Waste Management Practices: Municipal, Hazardous, and Industrial*, 2005; M. P. Wanielista and Y. A. Yousef, *Stormwater Management*, 1993.

Environmental fluid mechanics

The study of the flows of air and water, of the species carried by them, and of their interactions with geological, biological, social, and engineering systems in the vicinity of a planet's surface. The environment on the Earth is intimately tied to the fluid motion of air (atmosphere), water (oceans), and species concentrations (air quality). In fact, the very existence of the human race depends upon its abilities to cope within the Earth's environmental fluid systems.

Meteorologists, oceanologists, geologists, and engineers study environmental fluid motion. Weather and ocean-current forecasts are of major concern, and fluid motion within the environment is the main carrier of pollutants. Biologists and engineers examine the effects of pollutants on humans and the environment, and the means for environmental restoration. Air quality in cities is directly related to the airborne spread of dust particles and of exhaust gases from automobiles. The impact of pollutants on drinking-water quality is especially important in the study of ground-water flow. Likewise, flows in porous media are important in oil recovery and cleanup. Lake levels are significantly influenced by climatic change, a relationship that has become of some concern in view of the global climatic changes that may result from the greenhouse effect (whereby the Earth's average temperature increases because of increasing concentrations of carbon dioxide in the atmosphere). *See* AIR POLLUTION; GREENHOUSE EFFECT; WEATHER FORECASTING AND PREDICTION.

Scales of motion. Environmental fluid mechanics deals with the study of the atmosphere, the oceans, lakes, streams, surface and subsurface water flows (hydrology), building exterior and interior airflows, and pollution transport within all these categories. Such motions occur over a wide range of scales, ranging from eddies on the order of centimeters to large recirculation zones the size of continents. This range accounts in large part for the difficulties associated with understanding fluid motion within the environment. In order to impart motion (or inertia)

to the atmosphere and oceans, internal and external forces must develop. Global external forces consist of gravity, Coriolis, and centrifugal forces, and electric and magnetic fields (to a lesser extent). The internal forces of pressure and friction are created at the local level, that is, on a much smaller spatial scale; likewise, these influences have different time scales. The winds and currents arise as a result of the sum of all these external and internal forces.

Global motion is the largest scale (greater than 5000 km or 3000 mi); synoptic scale motion is the next largest (100–1000 km or 60–600 mi). Mesoscale motion occurs over regional areas (10–100 km or 6–60 mi); local motion is commonly referred to as microscale motion (less than 100 m or 300 ft). Humans live in the microscale associated with atmospheric motion, that is, the boundary layer of air that extends about 1 km (0.6 mi) above the Earth's surface. The terrain of the Earth as well as ocean surface conditions significantly affects both the microscale and mesoscale motion of the atmosphere, and hence weather conditions. The scales of motion range from eddies on the order of centimeters to huge masses extending over thousands of kilometers. See MESOMETEOROLOGY; MICROMETEOROLOGY.

El Niño, the periodic flow of warm waters along the western coast of South America, disrupts the coastal ocean and the upwelling of cold waters, producing large amounts of precipitation, along with widespread destruction of plankton, fish, and sea birds (which prey on the fish). Major El Niño events occurred in 1925, 1941, 1957–1958, 1972–1973, 1982–1983, and 1992. It has been determined that the events are caused by changes in the surface winds over the western tropical Pacific, which periodically release and drive warm waters eastward to the South American continent. See EL NIÑO.

The study of air pollution falls within the category of environmental fluid mechanics, because the air within the lower atmosphere steers (or advects) and diffuses pollutants. Atmospheric winds near the Earth's surface are generally turbulent and gusty, which helps to clear polluted areas; the velocity varies with altitude, local stability (level of turbulence), and roughness of the terrain. However, when the winds are calm, stagnant conditions can occur which subsequently prevent pollutants from being cleared from a city, resulting in high levels of bad air quality and smog. Of particular importance on the mesoscale level is acid rain (whereby rainfall removes sulfates and nitrates within the atmosphere), which has resulted in serious environmental damage. Likewise, mixing of pollutants into the upper atmosphere can cause long-term changes in the ozone layer (even though the causes, such as propellants within spray cans, may have been generated within the microscale layer). The 1991 explosion of Mount Pinatubo in the Philippines resulted in the discharge of many tons of particulates into the Earth's atmosphere; these particulates in turn acted as seed nuclei for precipitation, and were the cause of much of the flooding and climatic changes over the following few years. See WEATHER MODIFICATION.

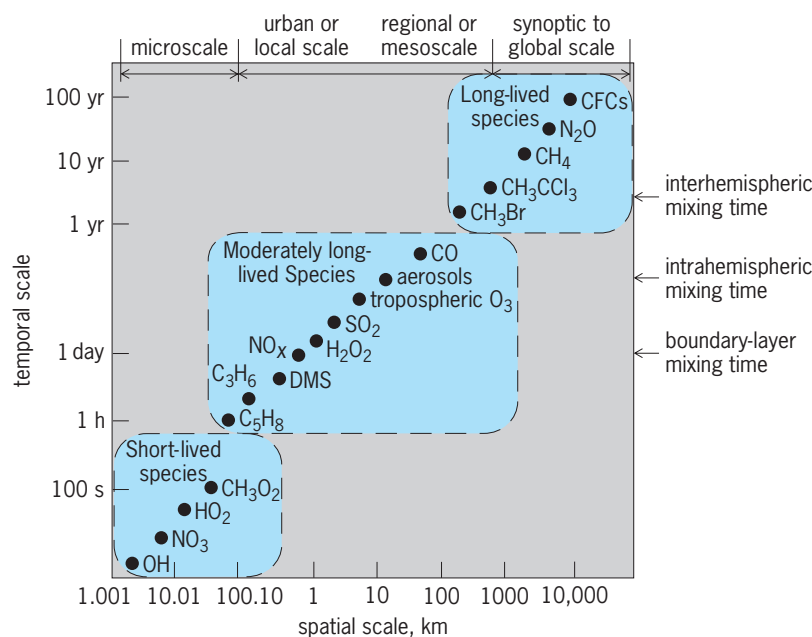


Fig. 1. Spatial and temporal scales of variability for atmospheric constituents. 1 m = 3.3 ft; 1 km = 0.6 mi. (After J. H. Seinfeld and S. N. Pandis, *Atmospheric Chemistry and Physics*, Wiley, 1998)

The various species found in the atmosphere have a wide range of lifetimes (residence time in the atmosphere). Species with short lifetimes have small spatial scales; those with lifetimes of years have spatial scales comparable to the entire atmosphere (Fig. 1). For example, the hydroxyl radical (OH) has a very short lifetime and small scale; methane (CH_4) has a lifetime of nearly 10 years and can become mixed over the entire Earth.

Governing equations. The foundations of environmental fluid mechanics lie in the same conservation principles as those for fluid mechanics, that is, the conservation of mass, momentum (velocity), energy (heat), and species concentration (for example, water, humidity, other gases, and aerosols). The differences lie principally in the formulations of the source and sink terms within the governing equations, and the scales of motion. These conservation principles form a coupled set of relations, or governing equations, which must be satisfied simultaneously. The governing equations consist of nonlinear, independent partial differential equations that describe the advection and diffusion of velocity, temperature, and species concentration, plus one scalar equation for the conservation of mass. In general, environmental fluids are approximately newtonian, and the momentum equation takes the form of the Navier-Stokes equation. An important added term, neglected in small-scale flow analysis, is the Coriolis acceleration, $2\Omega \times V$, where Ω is the angular velocity of the Earth and V is the flow velocity. See CONSERVATION LAWS (PHYSICS); CONSERVATION OF ENERGY; CONSERVATION OF MASS; CONSERVATION OF MOMENTUM; CORIOLIS ACCELERATION; DIFFERENTIAL EQUATION; DIFFUSION; FLUID-FLOW PRINCIPLES; NAVIER-STOKES EQUATION; NEWTONIAN FLUID.

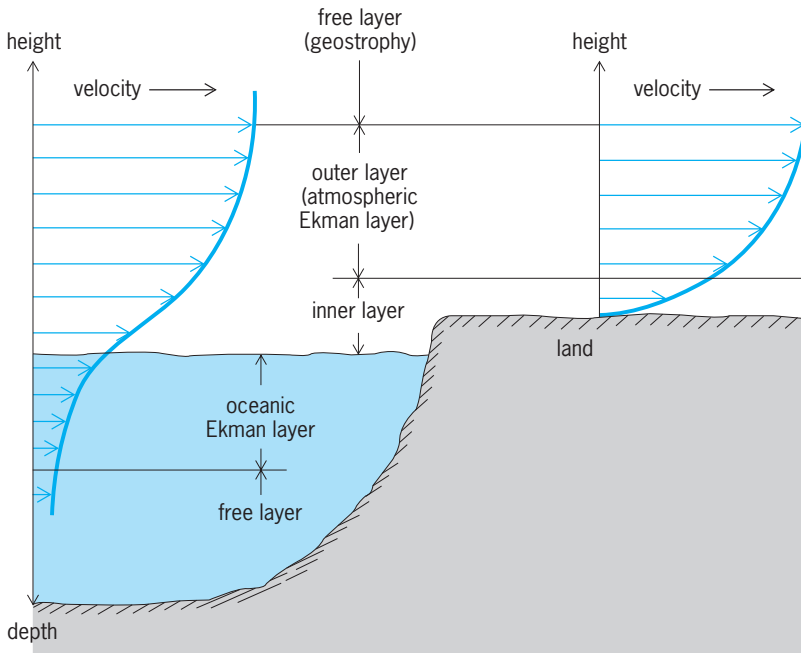


Fig. 2. Velocity distribution in the atmosphere and ocean. (After S. Eskinazi, *Fluid Mechanics and Thermodynamics of Our Environment*, Academic Press, 1975)

Driving mechanisms of flow. The mechanisms which drive the flow patterns in the atmosphere and oceans are vastly different. The atmosphere is thermodynamically driven, with the major source of energy coming from solar radiation. Short-wave radiation traverses the air and becomes partially absorbed by the land and oceans, which reemit the radiation at longer wavelengths. Long-wave radiation heats the atmosphere from below, creating convection currents in the atmosphere. See ATMOSPHERIC GENERAL CIRCULATION.

In the oceans, periodic gravitational forces of the Sun and Moon generate tides; in addition, the ocean surface is affected by wind stress that drives most of

the ocean currents. Local differences between the air and sea temperatures generate heat fluxes, evaporation, and precipitation, which ultimately act as thermodynamical forces that create or modify wind-driven currents. See OCEAN CIRCULATION.

Environmental layers. Fortunately, not every term in the Navier-Stokes equation is important in all layers of the environment. The horizontal component of the motion is usually the most significant and is subjected to maximum frictional forces at atmosphere-ocean interfaces. This frictional force causes the formation of a boundary layer in which the velocity of air at the surface of the Earth is zero (relative to the Earth), and the velocity at the surface of the ocean is a minimum equal to the surface velocity of the water. The ocean current is primarily generated by the wind; hence, the water velocity at the surface is a maximum and decreases in depth, again as a result of frictional forces. In both instances, frictional forces cause strong velocity gradients and vorticity (rotation) within the boundary layer. Figure 2 shows the velocity distribution in the atmosphere and ocean. See BOUNDARY-LAYER FLOW.

The thickness of the atmospheric boundary layer varies with the wind speed, degree of turbulence, and type of surface. For atmospheric flows, the layer is on the order of 1 km (0.6 mi) thick; within the ocean, it may be 30 m (100 ft) thick. Beyond this layer, the environmental flow is typically considered to be viscous-free (without turbulent shear), or inviscid. The rougher the terrain, or the larger the surface obstructions, the thicker the atmospheric boundary layer becomes, and the more gradual the increase of velocity with height (Fig. 3). The influence of the ground on the wind profile extends from a few hundred meters to over 500 m (1640 ft), depending on the roughness of the surface. Above this height, velocity is established from upper level meteorology. The wind speed is proportional to some power of height (empirically determined from experiments) above the surface.

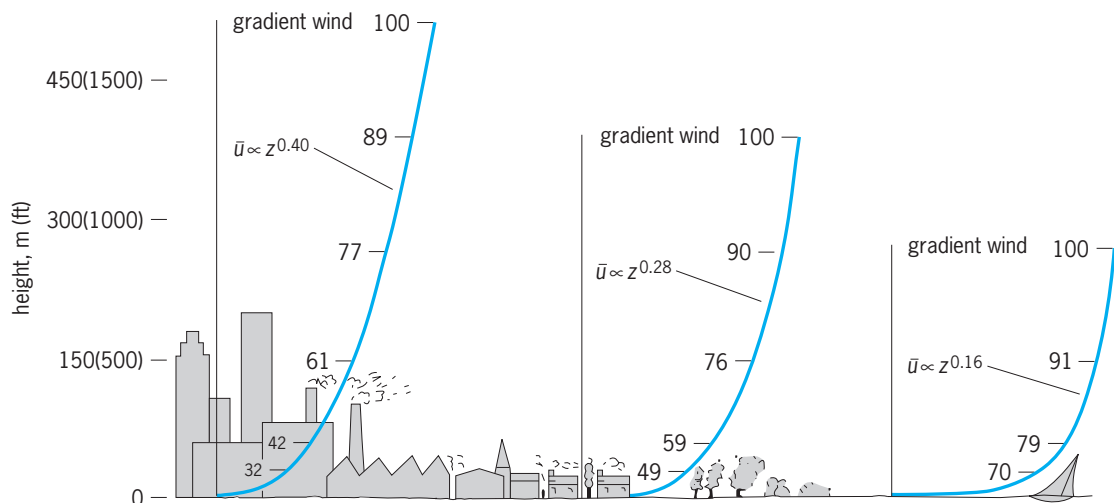


Fig. 3. Atmospheric boundary-layer profiles (plots of average wind speed \bar{u} versus height z) over different terrains. Wind speeds are expressed as percentages of the upper level wind (referred to as the gradient wind) above the boundary (or surface) layer. (After E. J. Plate, *Aerodynamic Characteristics of Atmospheric Boundary Layers*, AEC Critical Review Series, U.S. Department of Energy, 1971)

Length, velocity, and time scales in the Earth's atmosphere and oceans*			
Phenomenon	Length scale, km (mi)	Velocity scale, m/s (mi/h)	Time scale
Atmosphere			
Sea breeze	5–50 (3–30)	1–10 (2–20)	12 h
Mountain waves	10–100 (6–60)	1–20 (2–40)	Days
Weather patterns	100–500 (60–300)	1–50 (2–100)	Days to weeks
Prevailing winds	Global	5–50 (10–100)	Seasons to years
Climatic variations	Global	1–50 (2–100)	Decades
Ocean			
Internal waves	1–20 (0.6–12)	0.05–0.5 (0.1–1)	Minutes to hours
Coastal upwelling	1–10 (0.6–6)	0.1–1 (0.2–2)	Several days
Large eddies, fronts	10–200 (6–120)	0.1–1 (0.2–2)	Days to weeks
Major currents	50–500 (30–300)	0.5–2 (1–4)	Weeks to seasons
Large-scale gyres	Basin scale	0.01–0.1 (0.02–0.2)	Decades

*From B. Cushman-Roisin, *Introduction to Geophysical Fluid Dynamics*, Prentice Hall, 1994.

Because there are no shear stresses, the motion of the inviscid layer is governed only by the advection, pressure, and body-force terms. In atmospheric flows, the rotation of the Earth strongly influences this layer of flow, generally referred to as the geostrophic layer. Just above the surface, the mean velocities are small; the advection terms and the Coriolis force (which depends on the velocity) are negligible compared to the shear forces (viscous terms) which appear to be constant in this inner layer. However, within the outer, or Ekman, layer advection is still negligible and the viscous forces are small; this part of the boundary layer is in equilibrium with the Coriolis, pressure, and Reynolds stresses (turbulence). The **table** shows typical scales of length, velocity, and time for both atmospheric and oceanic motions. Oceanic motions are slower and more confined, and tend to evolve more slowly, than atmospheric motions.

Relative importance of terms. The key to being able to obtain solutions to the Navier-Stokes equation lies in determining which terms can be neglected in specific applications. For convenience, problems can be classified on the basis of the order of importance of the terms in the equations utilizing nondimensional numbers based on various ratios of values. The Rossby number (Ro) is the ratio of the advection (or inertia) forces to the Coriolis force, $Ro = V/L\Omega$, where V is velocity, Ω is the Earth's angular velocity, and L is a specified reference length. When the Rossby number is much less than 1, the inertia forces become insignificant, implying that these types of flows are more geostrophic. The ratio of the viscous to Coriolis forces is defined by the Ekman number, $Ek = \mu/\rho\Omega H^2$, where ρ is density, μ is viscosity, and H is a vertical reference height (or thickness). The ratio of inertia to viscous forces is referred to as the Reynolds number, $Re = \rho VL/\mu$. The Rossby number divided by the Ekman number yields the Reynolds number, that is, $Re = \rho VL/\mu = (V/\Omega L)(\Omega \rho H^2/\mu)(L^2/H^2) = (Ro/Ek)(L/H)^2$. When the Rossby number is large and the Ekman number is small, the motion is geostrophic; when the Rossby number is small and the Ekman number large, an Ekman-type boundary layer develops. As the Reynolds number increases, the ratio of the flow velocity to viscos-

ity increases (that is, the advection terms become more important than the viscous terms), with the flow eventually becoming turbulent. Since the Ekman number is generally small and the geometric ratio (L/H) is large (Ro is on the order of unity), the Reynolds number for geophysical flow is generally large and the flow turbulent. See DIMENSIONLESS GROUPS; GEOSTROPHIC WIND; REYNOLDS NUMBER; TURBULENT FLOW; VISCOSITY.

Measurements. Because of the scales of motion and time associated with the environment, and the somewhat random nature of the fluid motion, it is difficult to conduct full-scale, extensive experimentation. Likewise, some quantities (such as vorticity or vertical velocity) resist direct observations. It is necessary to rely on the availability of past measurements and reports (as sparse as they may be) to establish patterns, especially for climate studies. However, some properties can be measured with confidence.

Both pressure and temperature can be measured directly in the atmosphere and ocean with conventional instruments. In the ocean, depth is typically calculated from measured pressures obtained from instruments lowered into the sea. In the atmosphere, ground precipitation, radiative heat fluxes, and moisture content can be accurately measured. Likewise, the salinity of the ocean can be determined from electrical conductivity, and the levels monitored at shore stations. Concentration samples, collected at receptor sites over long periods of time, are examined to determine specific concentration levels and particulate sizes. These data are used to determine isopleth (concentration) levels and exposures over various atmospheric and oceanic conditions. Occasionally, inert tracer gases are released into the atmosphere to determine wind directions as well as atmospheric diffusion (turbulence levels) and plume trajectories.

Vector quantities such as horizontal winds and currents are typically measured by using anemometers and current meters. Anemometers atop buildings and towers, and current meters attached to mooring lines at fixed depths, offer fine temporal readings but are too expensive to adequately cover large areas. Instruments are routinely deployed on drifting platforms in the ocean, and balloons

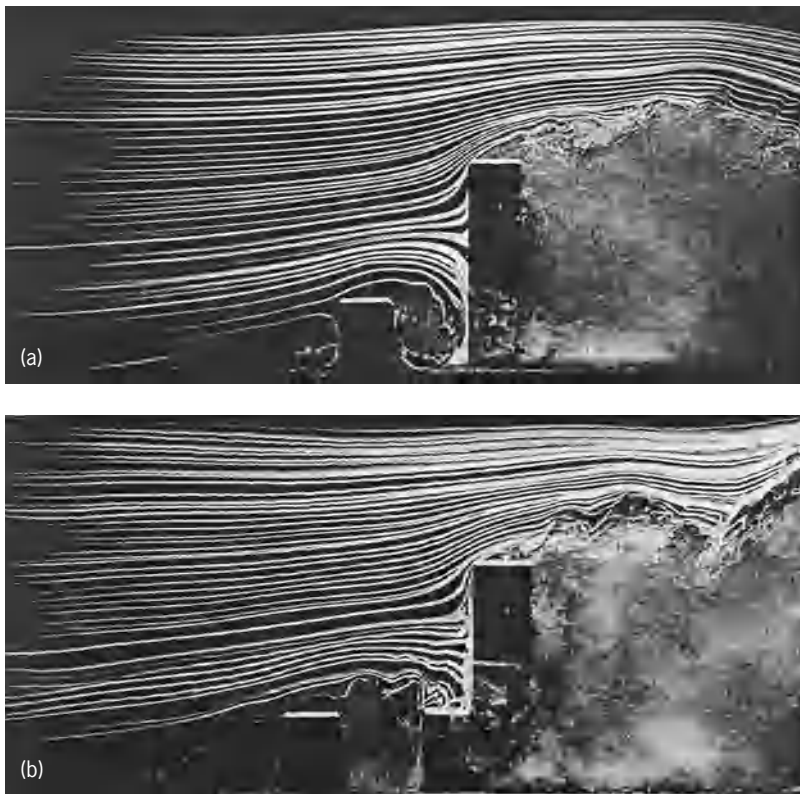


Fig. 4. Flow around two models of a tall building showing how minor design modifications can make a large difference in wind velocity at the pedestrian level. (a) Flow with vortex between the two buildings. (b) Flow with vortex removed by a slight change in the shape of the tall building. (From H. Thomann, *Wind effects on buildings and structures*, *Amer. Sci.*, 63:278–287, 1975)

are released in the atmosphere. (However, such measurements are mixed in time and space.) Measuring the three-dimensional velocity components simultaneously and obtaining meaningful three-dimensional heat fluxes is difficult, and essentially relegated to small-scale laboratory experiments.

Advances utilizing satellite imagery, Doppler radar, acoustic sounding, and lidar (laser) have made it possible to obtain highly detailed data, including turbulence information, over much broader spatial distances. Doppler radar has yielded three-dimensional velocity data and rotational characteristics within thunderstorms that can be used reliably to predict the onset of tornadoes. See DOPPLER RADAR; LIDAR; METEOROLOGICAL INSTRUMENTATION; METEOROLOGICAL SATELLITES.

Modeling. There are two types of modeling strategies: physical models and mathematical models. Physical models are small-scale (laboratory) mock-ups that can be measured under variable conditions with precise instrumentation. Such modeling techniques are effective in examining wind effects on buildings and species concentrations within city canyons (flow over buildings; Fig. 4). Generally, a large wind tunnel is needed to produce correct atmospheric parameters (such as Reynolds number) and velocity profiles. Mathematical models (algebra- and calculus-based) can be broken down further into either analytical models, in which an exact solution exists, or numerical models, whereby approximate

numerical solutions are obtained using computers. See WIND TUNNEL.

By far the most interesting and widely used models are the numerical models. The reason for their popularity is that it is possible to model more of the actual physics of the flow, that is, solve the Navier-Stokes equation, rather than make assumptions and eliminate key components of the physics just to obtain a solution. Although the Navier-Stokes equation is nonlinear, the partial differential equation can now be solved with some measure of confidence and reliability. In many instances dealing with environmental flows, the use of supercomputers is required. See SUPERCOMPUTER.

Numerical methods. Several broad classes of solution techniques are employed to solve the various derivatives and terms of the Navier-Stokes equation. The most common and widely used numerical methods are finite difference schemes (which are based on the use of truncated Taylor-series expansions); finite element schemes (which use an integral approach with local weighting and basis, or shape, functions); spectral methods (in which dependent variables are transformed to wave-number space by using a global basis function, such as the Fourier transform); pseudospectral methods (which use truncated spectral series to approximate derivatives); interpolation techniques (whereby polynomials are used to approximate the dependent variables in one or more spatial directions); and particle methods (which use lagrangian particles whose trajectories are calculated within a conventional eulerian grid). Such numerical models depend strongly on boundary and initial conditions; care must be exercised to correctly initialize and specify all variables at the boundaries of the computational model. All these schemes except the particle methods require knowledge of properties such as viscosity, dispersion coefficients, and thermal conductivity; particle methods require no constitutive models for particle viscosity or thermal conductivity, but do require a large number of particles for an accurate description of the flow field. The most popular modeling approaches are the finite difference, finite element, and interpolation schemes, especially for mesoscale and synoptic-scale simulations. See COMPUTATIONAL FLUID DYNAMICS; FINITE ELEMENT METHOD; INTERPOLATION; NUMERICAL ANALYSIS; SIMULATION.

Capabilities. The continuing rapid improvement in computational hardware has made it possible to model more complicated problems and include more physics (or mathematical terms) in the governing equations. Simulations of environmental fluid flow over microscale and mesoscale regions without simplifications of the equations of motion are now fairly common. Arrays consisting of millions of nodes can be calculated within a few hours on supercomputers, and three-dimensional graphical displays can be generated on work stations. By using satellite, radar, and conventional surface observations as input data to meteorological models, reasonably accurate local forecasts can be made for up to several days. Advances in numerical techniques as well as

computer hardware will continue, making it possible to perform more detailed calculations over broader expanses with improved accuracy over longer forecast periods.

Examples. An example of the simulation of fluid flow over a building complex is shown in Fig. 5. The three-dimensional equations of motion and species transport were solved using an adaptive finite element method. The flow field at the x - y midplane of the three-dimensional model domain (Fig. 5a) shows the development of a series of eddies as the flow moves downstream of the buildings. The flow field at the x - z vertical midplane (Fig. 5b) shows the formation of a large eddy between the two sets of buildings. Lagrangian particles, introduced upstream of the building array, clearly show the dispersion of concentration over and around the buildings (Fig. 5c).

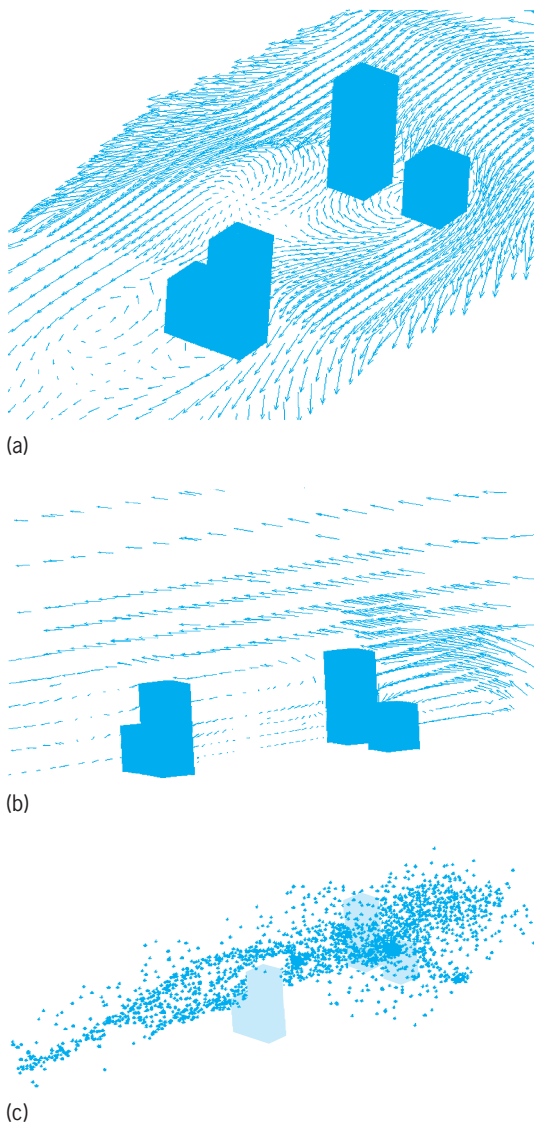


Fig. 5. Numerical simulation of flow over a building complex, carried out on a SGI-Cray Origin 2000 parallel computer. (a) Wind field in the horizontal (x - y) midplane. (b) Wind field in the vertical (x - z) midplane. (c) Lagrangian particles depicting species transport. (After D. W. Pepper et al., eds., *Development and Application of Computer Techniques to Environmental Studies*, WITPress, 1998)

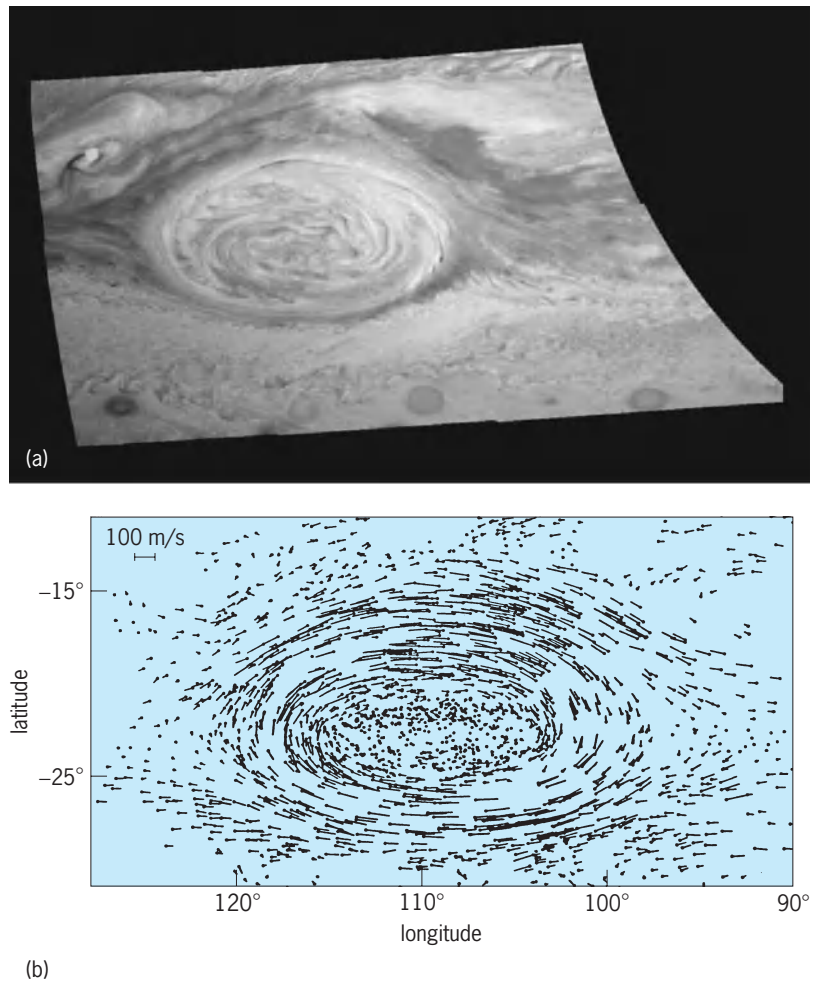


Fig. 6. Great Red Spot of Jupiter. (a) Image from *Galileo* spacecraft (NASA). (b) Velocity field in and around the Great Red Spot obtained by tracking small cloud formations in sequential *Voyager 1* images. (From T. E. Dowling and A. P. Ingersoll, *Potential vorticity and layer thickness variations in the flow around Jupiter's Great Red Spot and White Oval BC*, *J. Atmos. Sci.*, 45:1380-1396, 1988)

Particles are being pulled into the large eddy between the two sets of buildings.

The flow circulation around the Great Red Spot of Jupiter (Fig. 6a) is an example of the turbulent nature of fluid flow on a large scale. The velocity field in and around the Great Red Spot was obtained by tracking small cloud features over time (Fig. 6b). See FLUID MECHANICS; JUPITER.

Darrell W. Pepper

Bibliography. B. Cushman-Roisin, *Introduction to Geophysical Fluid Dynamics*, 1994; T. N. Krishnamurti and L. Bounoua, *An Introduction to Numerical Weather Prediction Techniques*, 1996; D. W. Pepper and J. C. Heinrich, *The Finite Element Method: Basic Concepts and Applications*, 1992; D. W. Pepper, P. Zannetti, and C. Brebbia (eds.), *Development and Application of Computer Techniques to Environmental Studies*, WITPress, 1998; M. L. Salby, *Fundamentals of Atmospheric Physics*, 1996; J. H. Seinfeld and S. N. Pandis, *Atmospheric Chemistry and Physics: From Air Pollution to Climate Change*, 1998; P. Zannetti, *Air Pollution Modeling*, 1990.

Environmental geology

The branch of geology that deals with the ways in which geology affects people. Examples of the effect of geology on human civilizations include (1) the ways that fertile soils develop from rocks and how these soils can become polluted by human activities; (2) how rocks and soils move down-slope to destroy roads, houses, and other human constructions; (3) sources of surface and subsurface water supplies and how they become polluted; (4) why floods occur where they do and how human activities affect floods; (5) locations of earthquakes and volcanic eruptions and the dangers they pose; (6) location of mineral resources such as copper, oil and gas, and uranium, and how mining these resources can pollute the environment; (7) how human activities can pollute the atmosphere and cause global warming, sea-level rise, and ozone depletion.

Soil. Most environmental concerns center on three materials essential to human life: soil, water, and air. Soils develop from chemical reactions between rocks, water, and atmospheric gases such as oxygen and carbon dioxide. The most abundant solid product of these reactions is a variety of clay minerals, which are major contributors to agricultural fertility. Ions on clay surfaces and decaying organic matter are the sources of nutrient elements for growing plants. Clay minerals are also responsible for destructive events, such as when they absorb water and cause soil to swell, resulting in foundation cracking and other structural problems in buildings. Changes in clay minerals as they lie in the soil are also the cause of many types of landslides. *See* CLAY MINERALS; SOIL.

The major environmental problems affecting soil are erosion and pollution. Almost half of American farmland is losing fertile topsoil faster than it can be formed. Most of the erosion results from poor farming practices and overzealous logging, which deplete the soil of plants and root systems that hold the soil in place. Most soil pollution is caused by synthetic chemical pesticides which not only kill the pests but also are absorbed by plants, which are then eaten by people and other animals. *See* EROSION; SOIL CONSERVATION.

Water. The amount of water on the Earth's surface is constant and more than enough to supply human needs. Water problems are generally concerned with local supply and its pollution. Serious water shortages in the United States occur in many western states because of increasing populations, as in Phoenix, Las Vegas, and Los Angeles.

In areas where surface water supplies are inadequate, fresh water can commonly be obtained from underground sources. The water occurs in the open spaces between the grains in sedimentary rocks and has accumulated in the rocks over many thousands of years. When this water is withdrawn faster than it is naturally replenished by rain- and snowfall, the underground supply can become seriously depleted, as is occurring today in many western cities. Such depletion is also a serious problem in the heavily irrigated agricultural areas of the central United States,

because most irrigation water is drawn from subsurface supplies. Removal of subsurface water can also result in ground subsidence, a serious problem in central California, Houston, and New Orleans. *See* GROUND-WATER HYDROLOGY.

Pollution. Most water pollution problems in the United States are caused by synthetic organic compounds and heavy elements released in wastewater discharges, many by the military as well as by chemical plants. Oil refineries are also major polluters. Other sources of pollution include leaking sanitary landfills, pesticides in runoff from farms, and mixtures of sulfuric acid and cyanide from mining operations. *See* HAZARDOUS WASTE; WATER POLLUTION.

Hazardous waste water can be disposed of by underground injection, incineration, evaporation from open pits, or export to other areas. The most common domestic method is underground injection, but because the underground geology is inadequately known, liquid hazardous waste injected deep underground has seeped into underground water supplies in about half the states. About 20% of hazardous waste is incinerated, but the Environmental Protection Agency (EPA) estimates that 60% could be disposed of by this safer method. Incineration is more expensive than subsurface disposal, however. Disposal by evaporation from surface impoundments is very unsafe and has caused ground-water contamination in about three-quarters of the 50 states.

Floods. Floods are the most commonly experienced natural calamities and are the cause of about 60% of declared disasters. Flooding may be caused by hurricanes from the ocean or excessive water from landward areas. Hurricanes cannot be controlled, and the value of property destroyed by them has increased considerably over time because of people's desire to live near the seashore. Loss of life has decreased, however, because of improved forecasting and early warning. Flooding resulting from excessive rain on the land surface has been exacerbated by the human desire to live near river channels, to pave land surfaces and increase the rate of runoff, and to denude hill slopes by excessive logging. Without vegetation and porous soil, the rate of runoff of rainwater is increased and river channels cannot hold all the water they receive. The massive flooding in central China in 1998 was attributed to the stripping of vegetation over large areas.

Atmosphere. Concerns about air quality revolve around increases in the amount of carbon dioxide, decreases in the amount of ozone 10–20 mi above the surface, and pollution of the air that sustains life in the lower atmosphere. Increasing carbon dioxide levels result largely from the burning of coal and oil, the main sources of energy in industrial societies. The amount of carbon dioxide in the air has increased by nearly 30% during the past 150 years. Carbon dioxide absorbs radiation from the Earth's surface and has caused a measurable increase in global temperatures (greenhouse effect); this global warming may in turn result in rising sea levels (from melting glaciers), endangering coastal communities over the long term. *See* GREENHOUSE EFFECT.

The concentration of ozone in the stratosphere

absorbs most of the Sun's ultraviolet (UV) radiation, thus protecting people from excessive skin cancers, damaged immune systems, and other maladies. The thickness of the ozone layer has decreased significantly during the past 50 years because of the release to the atmosphere of chlorofluorocarbons (CFCs) which destroy ozone molecules. Production of CFCs has now nearly stopped worldwide, but the ozone layer will not restore itself completely for about another 50 years.

Noxious gases and particulates in the air result largely from combustion of coal and oil in automobiles and factories. It is the cause of many lung problems and premature deaths in industrial countries. See AIR POLLUTION.

Earthquakes and volcanoes. Earthquakes and volcanoes tend to be associated geographically and are prominent around the Pacific rim. Japan and the American west coast are particularly vulnerable, and both have suffered major quakes in recent years. In the United States, California is the most seriously affected by earthquakes, with major property damage and loss of life. Attempts to predict the timing and location of earthquakes have not yet been successful despite very large research efforts in many countries over several decades.

Dangers from volcanic eruptions include mud and rock flows, lavas, and acid gases. Most dangerous to people are the debris flows, which can move downslope at speeds as high as 50 mi/h for distances of tens of miles. Entire towns have been buried by such flow with much loss of life. Lava flows move slowly and hence cause few deaths, but their power is irresistible and property damage is extensive. Eruptions can sometimes be predicted, but surprises are common. See EARTHQUAKE; VOLCANO. Harvey Blatt

Bibliography. H. Blatt, *Our Geologic Environment*, Prentice Hall, Upper Saddle River, NJ, 1997; N. K. Coch, *Geobazards*, Prentice Hall, Upper Saddle River, NJ, 1995; T. E. Graedel and P. J. Crutzen, *Atmospheric Change*, W. H. Freeman, New York, 1993; P. Gralla, *How the Environment Works*, Ziff-David, Emeryville, CA, 1994.

Environmental management

The development of strategies to allocate and conserve resources, with the ultimate goal of regulating the impact of human activities on the surrounding environment. "Environment" here usually means the natural surroundings, both living and inanimate, of human lives and activities. However, it can also mean the artificial landscape of cities, or occasionally even the conceptual field of the noosphere, the realm of communicating human minds.

Approaches. Environmental management is a mixture of science, policy, and socioeconomic applications. It focuses on the solution of the practical problems that humans encounter in cohabitation with nature, exploitation of resources, and production of waste. In a purely anthropocentric sense, the central problem is how to permit technology to evolve continuously while limiting the degree to which this pro-

cess alters natural ecosystems. Environmental management is thus intimately intertwined with questions regarding limiting economic growth, ensuring an equitable distribution of consumable goods, and conserving resources for future generations. Environmental management is a response to the increasing seriousness of the human impact on natural ecosystems. With a smaller global population base and a less pervasive use of technology, the environment might be able to recuperate on its own from human misuse, but it is now widely recognized that in many cases positive intervention is necessary if the environment is to recover.

There is, however, considerable disagreement about the course that such intervention should take, which has created a plurality of approaches to managing the environment. "Deep ecology" was born in the 1960s with the rise of movements that renounced technological development and decried the political basis of power and autocracy. However, "shallow" ecologists sought a compromise with those who argued that the solution to the world's environmental problems can come only through the generation of more technology. Environmental managers therefore fall within a broad spectrum that extends from conservationists to technocrats, from those who would limit human interference in nature to those who would increase it in order to guide natural processes along benign paths. Hence both conservationists and developers are represented. It is hoped that they will come together over the need to make economic development sustainable, without it being undermined by long-term damage to resources and habitats. This is the intention of the United Nations Convention on Environment and Development (the process that began at the Earth Summit in Rio de Janeiro in 1992), though underfunding and lack of commitment at the national level have severely limited the extent to which it has changed the global course of environmental management.

Participants. Participants in the process of environmental management fall into seven main groups: (1) governmental organizations at the local, regional, national, and international levels, including world bodies such as the United Nations Environment Programme and the U.N. Conference on Environment and Development; (2) research institutions, such as universities, academies, and national laboratories; (3) bodies charged with the enforcement of regulations, such as the U.S. Environmental Protection Agency; (4) businesses of all sizes and multinational corporations; (5) international financial institutions, such as the World Bank and International Monetary Fund; (6) environmental nongovernmental organizations, such as the World Wildlife Fund for Nature; and (7) representatives of the users of the environment, including tribes, fishermen, and hunters. The agents of environmental management include foresters, soil conservationists, policy-makers, engineers, and resource planners. The main link between these diverse groups of people is the need for accountability in the use of nature's riches. However, though there is much collaboration, relationships are often adversarial as objectives differ among the groups.

Intellectually, environmental management has assumed the status of a multidisciplinary academic field dedicated to furthering the human stewardship of natural resources. Biologists participate by virtue of their links with ecology and their interest in flora, fauna, and habitats. Geographers have a long-standing involvement in the ecology of people's relationships with their surroundings—in other words, in the two-way interaction between human communities and the landscapes that offer opportunities to prosper but that limit activities. Economists have become involved through the need to value environmental goods and services, assess the costs of pollution, and calculate materials balances. Engineers have created the field of ecotechnology, which is dedicated to the practical restoration of degraded or polluted environments. Finally, political scientists are the source of much environmental policy and regulatory know-how. Though these and other disciplines have a hand in environmental management, the links between them are complex, and as a result the field is remarkably diffuse. Indeed, it tends to lack a core of integrated concepts that might give it a clear popular identity. Some common themes of environmental management are as follows:

1. Bilateral and multilateral environmental treaties (transboundary ecological management)
2. Design and use of decision-support systems (practical utilization of environmental data; expert systems for environmental management)
3. Environmental policy formulation and enactment (participatory planning and public consultation regarding environmental programs)
4. Estimation, analysis, and management of environmental risk (risk perception and communication studies)
5. Formulation of environmental regulations (for dumping of wastes, emission of pollutants, and extraction of resources; monitoring and policing compliance)
6. Impacts and management of recreation and tourism (design and implementation of environmentally friendly "ecotourism" programs)
7. Natural resource conservation (designation and management of parks, preserves, and other protected areas; designation and protection of wilderness areas)
8. Positive environmental economics (economic justifications for investment in environmental protection)
9. Promotion of positive environmental values by education, debate, and information dissemination
10. Reduction of adverse environmental impacts
11. Resource evaluation and management
12. "Scoping" and investigation of environmental impacts (design of policies, norms, and procedures to limit impacts)
13. Strategies, methods, and programs for the rehabilitation of damaged environments (postpollution clean-up processes)

Improvements. The need to improve management of the environment has given rise to several new techniques. There is environmental impact analysis,

which was first formulated in California and is codified in the U.S. National Environmental Policy Act (NEPA). Through the environmental impact statement, it prescribes the investigatory and remedial measures that must be taken in order to mitigate the adverse effects of new development. In this sense it is intended to act in favor of both prudent conservation and participatory democracy.

Another technique is environmental auditing, which uses the model of the financial audit to examine the processes and outcomes of environmental impacts. It requires value judgments, which are usually set by public preference, ideology, and policy, to define what are regarded as acceptable outcomes. Audits use techniques such as life-cycle analysis and environmental burden analysis to assess the impact of, for example, manufacturing processes that consume resources and create waste.

The atmosphere, surface and subsurface waters, growing plants, minerals, and so on, are sometimes considered to be beneficial resources. The process of exploitation often involves risks to the user, and if these are magnified the resource may turn into a hazard—for example, when excesses of water generate destructive floods. Similarly, pollution has been defined as "a resource in the wrong place at the wrong time." Thus the process of environmental management can be considered one of limiting resource usage to its more benign forms, thus reducing risks and hazards, and using human ingenuity to transform pollutants into recycled resources. This involves some thorny problems. For example, no human activity is completely devoid of risk, and few environmental pathologies have ever been eliminated. Moreover, even full-scale recycling is not without costs; for example, the energy required to reclaim waste materials will usually be generated at the expense of at least some pollution. In this sense, one of the most salutary lessons of recent decades has been that the so-called benign generators of electrical power—the renewable sources based on winds, tides, solar radiation, or waves—involve potentially large costs in terms of how they modify landscapes. This has limited their attractiveness in relation to nonrenewable sources, and has demonstrated the need to broaden the analyses that feed into decision-making about the environment so that hidden and unexpected costs are given their full weight. Thus, as the field has evolved, it has become correspondingly more sophisticated in its treatment of the variables that are considered when formulating policy.

New challenges. All of the main environmental problems of the late twentieth and early twenty-first centuries fall under the environmental management field. Most problems are controversial. Tropical deforestation, ozone depletion, and global warming have fueled debate over strategies for the management of the global environment. Transboundary pollution and the international exploitation of resources (for example, the appropriation of raw materials in one country and the patenting of their genetic derivatives in another) have underlined the need for bilateral, and often multilateral, agreements about sharing responsibilities. Radiation

emissions, toxic waste issues, and hazardous material spills have emphasized the need for secure and standardized methods of treating pollutants. The production of organic chemicals, for example, grew by more than two orders of magnitude during the second half of the twentieth century, and there was a corresponding growth in the number of catastrophic pollution episodes.

Environmental management has risen to meet many of these challenges. The field has expanded from a purely governmental preserve to one that encompasses the private sector as well. Indeed, the manufacture of pollution control equipment and the institutional management of environmental hazards have turned into growth areas. Yet the successes must be seen against a backdrop of deepening environmental crisis. Relentless population pressure, the unfettered nature of international capital, and the exposure of a record of significant environmental mismanagement in eastern Europe are examples of remaining problems.

David Alexander

Bibliography. M. Atchia and S. Tropp (eds.), *Environmental Management: Issues and Solutions*, 1995; R. S. Dorney, *The Professional Practice of Environmental Management*, 1987; T. O'Riordan (ed.), *Environmental Science for Environmental Management*, 1995; R. Welford (ed.), *Corporate Environmental Management: Systems and Strategies*, 1996; R. R. White, *Urban Environmental Management*, 1994; G. A. Wilson and R. L. Bryant, *Environmental Management: New Directions for the 21st Century*, 1997.

Environmental radioactivity

Radioactivity that originates from natural and anthropogenic sources, including radioactive materials in food, housing, and air, radioactive materials used in medicine, nuclear weapon tests in the open atmosphere, and radioactive materials used in industry and power generation.

Natural radioactivity, which is by far the largest component to which humans are exposed, is of both terrestrial and extraterrestrial (cosmic) origin. About 340 nuclides are known in nature, of which 70 are radioactive and are found mainly among the heavy elements. Three nuclides which are responsible for most of the terrestrial component are potassium-40, uranium-238, and thorium-232.

Terrestrial nuclides. Potassium-40 has a half-life of 1.3×10^9 years and decays, through calcium-40, to stable argon-40. This nuclide is present to the extent of about 0.01% in natural potassium, thereby imparting a specific activity of approximately 800 picocuries per gram (3.0×10^4 becquerels per kilogram) of potassium. Potassium is of course an essential element of the human body, which contains about 0.1 microcurie (3.7×10^3 Bq) of the radioactive form of the element. See POTASSIUM.

Uranium-238, the most abundant of the uranium isotopes, decays through a chain of 18 nuclides to stable lead-206. This chain includes radionuclides of 10 heavy elements, the half-lives of which range

from 4.5×10^9 years (uranium-238) to 1.64×10^{-4} s (polonium-214). See URANIUM.

An important characteristic of the uranium-238 chain is that it includes a radioactive noble gas, radon-222, with a half-life of 3.8 days. Radon emanates from the Earth's crust, diffuses into the atmosphere, and eventually decays to radionuclides of lead and polonium that attach themselves electrically to the dust normally present in the atmosphere. The dose of the lung from radon daughter products attached to atmospheric dust is much higher than the dose received from natural sources by any other part of the body. See RADON.

The radon daughters are also a source of natural "fallout" which, when released into the soil, results in broad-leaved plants having relatively high concentrations of both lead-210 and polonium-210. In the case of tobacco, the polonium is volatilized in the course of smoking and delivers a higher-than-normal dose to the lung. It has been suggested that this phenomenon may play a role in the production of lung cancer from tobacco smoking, but this suggestion must be regarded as highly speculative.

Thorium-232 has a half-life of 1.4×10^9 years and decays through a series of 12 nuclides to a stable isotope of lead (lead-208). Human exposure from this series results mainly from the energetic gamma radiation from thallium-208. Internal irradiation results from absorption of radium-228 (6.7-year half-life) which finds its way into food. See THORIUM.

Dose rates. The average person in the United States receives 80–180 mrem/year (0.8–1.8 millisieverts/year) from natural sources of ionizing radiation, depending on the organ considered (see **table**). Most of this dose originates from radioactive materials in the Earth's crust. The external dose due to cosmic rays is an average of about 28 mrem/year (0.28 mSv/year), a value that increases with altitude due to reduced shielding of cosmic radiation by the atmosphere. At Denver, which is located 1 mi (1.6 km) above sea level, the annual dose from cosmic radiation is more than 50 mrem/year (0.5 mSv/year), or about twice the average dose. See COSMIC RAYS.

The human body is also exposed to radionuclides in food and water. Potassium-40 is the most important of these, with radium-226 and radium-228 of perhaps less importance from the point of view of the dose delivered.

A number of radionuclides are produced in nature by the interaction of cosmic rays with atmospheric gases, and these add slightly to the exposure received from natural sources. Carbon-14 is the most important of these nuclides which, like potassium-40, is present in the body to the extent of about 0.1 microcurie (3.7×10^3 Bq), but which delivers a dose of only about 1 mrem/year (10 μ Sv/year) because of the unusual softness, that is, low energy, of its beta emission. See COSMOGENIC NUCLIDE.

There are wide deviations from the average doses shown in the table. Thus, at one extreme, miners working underground in the presence of radioactive ore can be exposed to such high levels of atmospheric radon that they develop lung cancer. This

Average dose equivalent rates from various sources of natural background radiation in the United States

Source	Dose, mrem/year*				
	Gonads	Lung	Bone surfaces	Bone marrow	Gastrointestinal tract
Cosmic radiation	28	28	28	28	28
Cosmogenic radionuclides	0.7	0.7	0.8	0.7	0.7
External terrestrial	26	26	26	26	26
Inhaled radionuclides	—	100	—	—	—
Radionuclides in the body	27	24	60	24	24
Rounded totals	80	180	120	80	80

*1 mrem/year = 10^{-2} mSv/year.

was observed in mines in central Europe that were worked for precious metals long before radioactivity was discovered. Uranium miners in the southwestern United States have also been subject to a high incidence of lung cancer because of lack of control of radon.

Radioactive anomalies. There are also geographical areas where the levels of natural radioactivity are unusually high. Six types of anomalies that can be important from the point of view of population exposure are: monazite sands and other placers, alkaline intrusives and granites of the Conway type in New Hampshire, bauxites and intensely weathered soils, uraniferous phosphate rock (and soils), ground waters enriched in radium and radon, and black shales and related organic accumulations.

There are two locations, in India and in Brazil, where the external radiation levels are greatly elevated by the presence of thorium-bearing monazite sands on or near ocean beaches. The local inhabitants are thus exposed to doses several times those received under normal circumstances. In Kerala, India, more than 25% of the 70,000 people living in the monazite area receive a dose greater than 500 mrem/year (5 mSv/year), and a small percentage receive between 1000 and 2000 mrem/year (10 and 20 mSv/year). In Guarapari, Brazil, the 6000 residents receive a mean dose of about 600 mrem/year (6 mSv/year).

Another type of anomaly, of which the outstanding example is the Morro do Ferro near Pocos de Caldas, Brazil, is found in regions of alkaline intrusives. The Morro do Ferro is a hill located in the state of Minas Gerais. The plants absorb so much radioactivity that they autoradiograph themselves when placed on radiation-sensitive film (see *illus.*). It has been estimated that the lungs of rodents living underground on the Morro do Ferro receive as much as 3000 mrem/year (30 mSv/year) from inhalation of a short-lived (53 s) radon isotope found in the thorium chain. The source of the radioactivity on the Morro do Ferro is a deposit of about 13,000 tons (12,000 metric tons) of thorium located near the summit of the hill. The deposit has been in place for perhaps as long as 80,000,000 years, and is in an advanced state of weathering. Rainfall averages about 67 in./year (170 cm/year).

Because of marked chemical similarities between thorium and plutonium, a study of the history of the deposit and the manner in which it has been mobilized by hydraulic and other forces will provide useful information for predicting the behavior over geologic time of plutonium and other actinides in a radioactive waste repository that has been breached.

It is known that water from deep wells may have high concentrations of radium. About 1,000,000



Autoradiograph of plant from the Morro do Ferro, Brazil, made by the plant's natural radioactivity.

people in northern Illinois and southern Iowa drink water with abnormally high radium concentrations.

Technological modification. The natural radioactive environment can also be altered by human activities, as described below.

Building construction. The radon concentration within buildings is influenced by the materials of construction. Construction materials that often have high concentrations of radium (the parent of radon) include pumice, concrete containing alum shale, and granite. Tailings from uranium and phosphate mining have been particularly troublesome in this regard, and it has been necessary to place restrictions on the use of these materials in construction of buildings. In Grand Junction, Colorado, the use of uranium mine tailings in home construction resulted in such high indoor radon concentrations that remedial action was necessary.

It is ironic that the need to conserve energy is resulting in higher exposures to radon. Weather stripping and other measures that reduce the rate of building ventilation cause less dilution of the radon that diffuses from the materials of construction, and hence tend to increase the radon concentration of the air within such buildings.

Combustion of fossil fuels. All fossil fuels contain naturally occurring radionuclides. Natural gas contains radon, which is sometimes present in detectable concentrations at the place of combustion.

Both oil and coal contain uranium and thorium and their daughter products. The radionuclides are discharged into the atmosphere in measurable amounts that are probably not significant to public health. It has been found that the radioactivity disseminated by burning coal is greater than that discharged into the atmosphere during normal operation of nuclear power reactors.

Aircraft travel. Flying at high altitudes increases the exposure from cosmic rays, and it is estimated that the dose received during a transcontinental trip across the United States is about 2 mrem (20 μ Sv).

Medical procedures. Procedures in which x-rays or radionuclides are used for diagnostic or therapeutic purposes are also a major source of radiation exposure, second in importance only to natural radioactivity. Doses of iodine-131 and iodine-125 are administered the most frequently. Other important nuclides include phosphorus-32, cobalt-60, gold-198, iron-59, and technetium-99. These nuclides are usually administered internally. See ISOTOPIC IRRADIATION; RADIOACTIVE TRACER.

The use of diagnostic x-rays subjects the general population to far more exposure than radiopharmaceuticals. Estimates of the per capita dose from diagnostic x-rays range from 75 mrem/year (0.75 mSv/year) in New Orleans and 50 mrem/year (0.5 mSv/year) in New York to as little as 1.25 mrem/year (12.5 μ Sv/year) in Thailand. The per capita dose tends to increase with the general level of economic development of the country and the availability of medical services. However, in all countries, developed and developing, the per capita dose can

be minimized by improving the techniques used. See RADIOGRAPHY; RADIOLOGY.

Nuclear weapons testing. Environmental radioactivity from the testing of nuclear weapons diminished after the major nuclear powers declared a moratorium on atmospheric testing in 1963. Relatively small-scale tests have been conducted by France, China, and India, but these have not added significantly to the levels of radioactivity previously disseminated from tests by the United States, Soviet Union, and United Kingdom prior to the moratorium.

The residual radioactivity from weapons tests is largely from cesium-137 and strontium-90. Cesium-137 is chemically similar to potassium and tends to be distributed more or less uniformly throughout the soft tissues. The dose commitment (the mean lifetime dose from cesium-137 already deposited) was about 27 mrem (0.27 mSv) in the United States in 1976. Cesium-137 has a half-life of 30 years.

Strontium-90, which has a half-life of 38 years, is chemically similar to calcium and therefore tends to deposit in the skeleton. The dose from strontium-90 is highest in individuals whose skeletons were formed during the period of heavy weapons testing, from about 1954 to 1962. Among this segment of the population, the bone dose was annually about 2 mrem (20 μ Sv) and the 50-year dose in the temperate zone was estimated to have been about 87 mrem (0.87 mSv) in 1976. See RADIOACTIVE FALLOUT.

Nuclear power plants. Although various national and international regulatory organizations have proposed guidelines that limit the per capita dose received by individuals in the general population to 170 mrem/year (1.7 mSv/year), it has become evident that nuclear power plants can be routinely operated so that the general population will not be exposed to more than 1% of this limit. The U.S. Nuclear Regulatory Commission requires that the reactors be operated so that the dose to the maximum exposed individual, usually a hypothetical person located at the fence line, is not more than 10 mrem/year (0.1 mSv/year).

The radioactive wastes discharged into the environment by a normally operated nuclear power reactor are minimal. As noted above, it has been shown that, at least so far as airborne radioactivity is concerned, coal- and oil-fired power plants discharge higher quantities of radioactivity into the environment.

Two of the oldest privately operated power reactors in the United States are Dresden I, a boiling water reactor operated by Commonwealth Edison Company near Chicago, and Yankee Rowe, a pressurized water reactor operated by the Yankee Atomic Electric Company in Rowe, Massachusetts. Surveys made by the U.S. Public Health Service after the reactors were in operation for about 10 years showed that there was no detectable residual radioactivity in the environs. Other reactors that have been studied have been found to expose people in the vicinity to no more than a few mrem per year, and these small levels should be reduced even further as the technology of reactor operation is improved.

A serious accident at the Three Mile Island Unit Two nuclear reactor early in 1979 resulted in worldwide concern about the danger from exposure of nearby Pennsylvania populations. Fortunately, despite extensive damage to the reactor core brought about by a series of mechanical and human errors, the radioactivity was well contained and exposures to the general population were minimal. The average dose, due to the accident, to the population in a 50-mi (80-km) radius was less than 2 mrem ($20 \mu\text{Sv}$).

The principal source of this low-level radiation exposure in the past was gaseous plumes from boiling water reactors. It was at one time permissible to discharge the waste gases after only a 30-min holdup, and this permitted the discharge of short-lived noble gases in copious amounts, thereby accounting for the dose from the plume. However, in order to meet the more stringent standards set by the Nuclear Regulatory Commission in 1975, it is now necessary to provide longer holdup times to permit essentially complete decay of the shorter-lived noble gases.

The most significant nuclides in the liquid wastes are cesium-137 and cesium-134. The former is produced in both weapons tests and nuclear reactors. However, cesium-134 is produced only in nuclear reactors and provides a label by which it is possible to differentiate environmental cesium due to weapons testing from that due to reactors.

The standard methods of treating the low-level radioactive wastes are sufficient to limit the maximum dose to humans who consume fish or shellfish to a few mrem per year. The Environmental Protection Agency has forecast that by the year 2000 the per capita dose from the operation of the several hundred reactors (and associated fuel-reprocessing services) will be 0.4 mrem/year ($4 \mu\text{Sv}/\text{year}$), which is less than 1% of the average dose people receive from nature and well within the range of variability in the natural radiation background. See NUCLEAR POWER; NUCLEAR REACTOR; RADIOACTIVE WASTE MANAGEMENT; RADIOACTIVITY.

Merril Eisenbud

Bibliography. M. Eisenbud et al., *Environmental Radioactivity*, 4th ed., 1997; R. L. Kathren, *Radioactivity in the Environment*, 1984; *Report of the President's Commission on the Accident at Three Mile Island*, 1979; V. Schultz and F. W. Whicker (eds.), *Radioecology: Nuclear Energy and the Environment*, 2 vols., 1982; L. Tommasino et al., *Radon Monitoring in Radioprotection, Environmental Radioactivity, and Earth Science*, 1990.

Environmental test

The evaluation of a physical system (engineering product) in conditions which simulate one or more of the environments that may harm the system or adversely affect its performance. Such testing is typically undertaken during the development of new products, with the aim of establishing or improving their reliability and durability.

Role of market forces. Almost every engineering product, from familiar household devices such as refrigerators to military hardware such as high-

performance aircraft, undergoes some type of environmental testing before being offered to the public or government for sale. The increasing competitiveness of the world marketplace is a major impetus behind the drive for improved product reliability and durability, and has led to the expanded use of environmental tests by manufacturers worldwide.

In addition to the evaluation of a final, finished product, environmental testing can play an important role throughout a product's design/development cycle to ensure that the materials and manufacturing processes employed can meet the stresses imposed by the environment in which the product is likely to operate. By not waiting until a finished product is evaluated, manufacturers can use environmental testing to eliminate costly redesigns late in the design/development cycle.

Facilities. Because it is necessary to precisely control the environmental factors which define the test (for example, temperature, vibration level, or altitude), environmental testing is typically conducted in specially designed facilities, or environmental chambers. Some environmental chambers can generate extremely high and low temperatures and humidity levels. Others can simulate corrosive environments such as salt sprays. The expense of developing and maintaining facilities for conducting the wide variety of environmental tests that new products must undergo is often beyond the resources of the product manufacturers themselves. This has led to the formation of companies specializing in environmental tests. The increasing use of environmental test facilities is mirrored in the expenditures for the test equipment itself.

Standards and specifications. Products and equipment intended for military use are often subjected to the harshest and most variable of environmental conditions. The military has pioneered the creation of well-documented standards and specifications for the evaluation and testing of any products or equipment which it will purchase. In the United States, these standards and specifications are often cited in Department of Defense solicitations as requirements. Examples of standards documents (typically identified as MIL STD- x with x being a three- or four-digit number) published by the Department of Defense are MIL STD 167-1 *Mechanical Vibration of Shipboard Equipment*; MIL STD-202 *Test Methods for Electronic and Electrical Component Parts*; MIL STD-750 *Test Methods for Semiconductor Devices*; and MIL STD-810 *Environmental Test Methods*. Although a number of military standards and specification documents are being reviewed and canceled for the sake of efficiency, these documents still constitute an important source of information for environmental testing.

Civilian organizations, such as the Society of Automotive Engineers, publish standards for automotive and aerospace equipment. In addition, nearly 100 countries have adopted the International Organization for Standardization (ISO) 9000 series of standards for quality management and quality assurance. These standards are being implemented by thousands of manufacturing and service organizations,

both public and private. The ISO 9000 family of standards represents an international consensus of good management practice and is intended to give organizations guidelines as to what constitutes an effective quality management system. Of the ISO 9000 family standards, ISO 9003 covers quality assurance obligations of the manufacturer in the areas of final inspection and testing. Among the guidelines provided by ISO 9003 are those dedicated to (1) developing procedures to inspect, test, and verify that final products meet all specified requirements before they are sold; (2) developing procedures to control and calibrate the testing equipment; (3) ensuring that every product is identified as having passed or failed the required tests. ISO 9003 will have a significant impact upon the standardization of environmental testing in both civilian and military endeavors.

Types of tests. An incomplete but representative list of environmental tests requiring dedicated test chambers includes tests for altitude, dust, explosiveness, flammability, fungus, humidity, icing, acoustic vibration, overpressure, rain, salt, fog, sand, temperature, and wind. Tests typically not requiring chambers but still utilizing dedicated mechanical testing equipment include acceleration, fatigue cycling, transportation simulation, shock, and vibration. The size of environmental test chambers can range from less than a cubic meter for testing component parts (for example, exposing a sample of an automobile tire to controlled levels of ozone) to 1,000,000 m³ (for example, testing entire aircraft in “cold-soak” conditions).

A representative example of environmental testing is provided by an ISO standard, also adopted by the Society for Automotive Engineers, that specifies environmental tests for aircraft equipment whose operating environment includes exposure to acoustic vibration (or noise). The standard specifies laboratory conditions for testing the ability of equipment in subsonic and supersonic civilian aircraft to withstand the stresses of acoustic vibration. Two types of tests are required to be performed: equipment functioning and acoustic endurance.

Equipment functioning refers to the ability of the equipment in question, for example, an avionics subsystem, to perform as desired in the presence of acoustic vibration representative of operational use. Tests for functioning are typically conducted both during equipment development and for certification of the equipment, once produced. In aircraft applications, acoustic vibration typically finds its source in external aerodynamics (that created by the aircraft passing through the atmosphere) or internal aerodynamics and combustion (that created by the jet engines). Noise or acoustic vibration in these environmental tests is quantified in decibels (dB). For example, sound pressure level (SPL) is defined as $SPL = 10 \cdot \log_{10}(p_{rms}^2/p_{ref}^2)$ decibels, where p_{rms} refers to the root-mean-square (rms) sound pressure (in pascals) in the acoustic source (for example, a jet engine) and p_{ref} is a reference sound pressure level, taken to be 0.0002 pascal (Fig. 1). Typically, sound pressure levels in excess of 125 dB are employed in the acoustic environmental tests for aircraft equip-

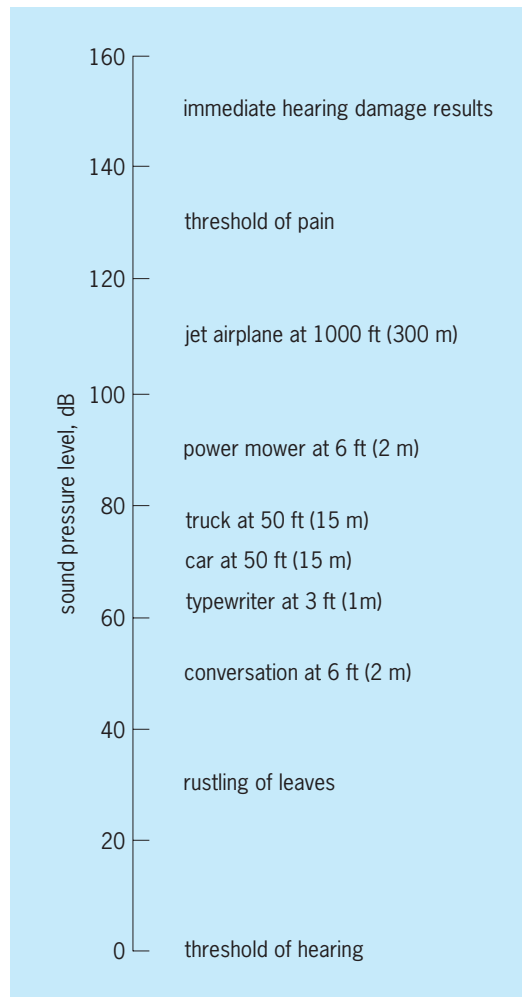


Fig. 1. Sound pressure levels of familiar sources.

ment. In addition to decibel rating, the frequency of acoustic vibration is used as an experimental variable.

Acoustic endurance refers to the length of time that equipment can perform to specification in the presence of noise of specified power. In environmental tests of acoustic endurance, so called wide-band noise is used, that is, noise with power distributed evenly across a broad range of frequencies. The level of power (test severity grade) is determined by the user according to either of two procedures: (1) If measurements of in-flight sound pressure levels at the location where the equipment is to be installed are available or can be calculated, they can be used to determine test sound pressure levels. (2) If such measurements are unavailable, the test levels may be obtained based on the type of aircraft in which the equipment is to be installed and the location of the equipment relative to the noise sources. See ACOUSTIC NOISE; AIRCRAFT TESTING; MECHANICAL VIBRATION; SOUND; SOUND PRESSURE; VIBRATION.

An example of an environmental test chamber is the large “sunspot chamber” for Earth orbital and deep-space simulations of spacecraft, operated by the National Aeronautics and Space Administration (Fig. 2). The chamber is equipped with a shroud that is liquid-nitrogen-cooled. Heat lamps are used to



Fig. 2. Large "sunspot chamber" for Earth orbital and deep-space simulations of spacecraft, located at the Marshall Space Flight Center. (NASA)

simulate the radiance of the Sun. The chamber itself consists of a vertical cylinder 3.2 m (10.5 ft) in diameter and 3.7 m (12 ft) in height. To simulate the vacuum in deep space, the chamber pressure can be lowered to 1.3×10^{-4} pascal.

Environmental engineering specialists. The importance and ubiquity of environmental testing has led to the formation of an engineering subspecialty. An environmental engineering specialist is defined to be a person who is skilled in one or more environmental engineering areas, which include but are not limited to natural and induced environments and their effects on material; expertise in measuring and analyzing field environmental conditions; formulating environmental test criteria; determining when environmental laboratory tests are appropriate and valid substitutes for natural field/fleet environmental tests; and evaluating the effects of specific environments on material.

Challenger accident. The accident involving the NASA space shuttle *Challenger* illustrates the vital importance of environmental testing and the results of either inadequate testing or use of systems in environments which have not been adequately represented in tests. On January 28, 1986, the *Challenger* was launched from Cape Canaveral for a 7-day mission. Approximately 73 seconds into the flight, the vehicle exploded and the crew perished. After a lengthy investigation, the cause of the disaster was traced to the failure of an O-ring seal in one of the solid rocket boosters, which allowed hot exhaust gases to leak. The O-ring failure was attributed to the freezing temperatures of the launch environment, which severely degraded the flexibility of the O-ring, a flexibility that was required in the vibration and aerodynamic load environment of the ascent. See SPACE FLIGHT; SPACE SHUTTLE; SYSTEMS ENGINEERING.

Ronald A. Hess

Environmental toxicology

A broad field of study encompassing the production, fate, and effects of natural and synthetic pollutants in the environment. The breadth of this field depends

on the definition of environment. It can be defined as narrowly as the home and workplace or as broadly as the entire Earth and its biosphere. Environmental toxicology is truly an interdisciplinary science. The effects of a pollutant on the environment depend on the amount released (the dose) and its chemical and physical properties. Pollutants can be grouped according to their origin and effects.

Nutrients. Pollution from nutrients is generally a problem of aquatic systems. Carbon, nitrogen, and phosphorus are essential nutrients and, when present in excess, can result in an overstimulation of microbial and plant growth. Nutrients enter the environment in runoff from fertilized agricultural areas, in effluents from municipal and industrial wastewater treatment facilities, and in dead and decaying plant material. Excess microbial growth can result in deoxygenation of a body of water. The overstimulation of algae and plant growth in aquatic systems due to excess nutrients of human origin is called cultural eutrophication. Decreased clarity of water, overgrowth of emergent vegetation, increased rate of sedimentation and filling-in of lake basins, and greatly shortened lifespan of an aquatic system are all characteristics of cultural eutrophication. See EUTROPHICATION.

Microbial pollution. Pathogenic bacteria and protozoa can be a major source of pollution in areas that receive untreated sewage, items from ocean dumping, and improperly discarded hospital waste. Toxic metabolites of fungal origin (mycotoxins) are also potential pollutants. Most important of the broad range of mycotoxins are those that contaminate human food and the feed of domestic animals. These contaminants include the ergot alkaloids, tricothecins, and aflatoxins. Ergot alkaloids are produced by *Claviceps* spp. and are strong neurotoxic and vasoconstrictive agents. Tricothecins are produced by a variety of species of the genus *Fusarium* and are potent inflammatory agents and tissue irritants. Aflatoxins are produced by species of the genus *Aspergillus* and contaminate improperly stored grains and nuts. Aflatoxins have been implicated in a variety of diseases in poultry and humans, have caused cancer in experimental animals, and have been linked epidemiologically to cancer in humans. Aflatoxins are the only mycotoxins to be regulated by the U.S. Food and Drug Administration. See AFLATOXIN; ERGOT AND ERGOTISM; TOXIN.

Suspended solids. Forest fires, volcanic eruptions, and dust storms can be major sources of suspended materials. These materials can also originate in runoff from agricultural areas, construction and mining sites, and roads and other paved areas. Truck and automobile exhaust and industrial discharge to the atmosphere are also sources of suspended solids. These materials reduce the amount of light in the atmosphere and bodies of water, interfering with photosynthesis. Physical irritation and abrasion, and obliteration of habitat by silting are the primary toxic effects of suspended solids.

Atmospheric pollutants. Metabolic processes and natural combustion and thermal activity (such as forest fires and volcanoes) can release large amounts of

gaseous by-products to the atmosphere. However, natural inputs are minor compared to atmospheric pollutants due to human activity. Although most anthropogenic air pollution is produced by the various forms of transportation, emission from stationary sources of fuel combustion (for example, factories and power plants) are responsible for the greatest amount of hazardous materials released.

Sulfur oxides, a primary combustion product of coal, are converted to acids in the atmosphere and precipitated in rain. In areas of limited buffering capacity such as northeastern North America and Scandinavia, pH levels are known to be greatly reduced, with loss of forests and aquatic life. *See* ACID RAIN.

Nitrogen oxides are formed during high-temperature combustion of petroleum products in refining operations and internal combustion engines. As with sulfur oxides, nitrogen oxides can form acid precipitation. They can also react with hydrocarbons and sunlight to form ozone and peroxyacetyl nitrates, the major constituents of photochemical smog.

Chlorofluorocarbons are important refrigerants and aerosol propellants. These gases accumulate in the upper atmosphere and have been implicated in a decrease in the protective ozone layer surrounding the Earth. Because the ozone layer absorbs harmful ultraviolet radiation from the Sun, it has been predicted that a decrease in this layer will result in an increase in skin cancer. In addition, the biological structures of aquatic and terrestrial ecosystems are strongly influenced by solar ultraviolet radiation, and most organisms exist very near their threshold of tolerance.

Carbon monoxide, formed by the incomplete combustion of fossil fuels, is primarily derived from transportation sources. It interferes with the oxygen-carrying function of hemoglobin in the bloodstream. The effects of carbon monoxide are acute but reversible, and thus pose a relatively small health hazard compared to other forms of atmospheric pollution. However, repeated anoxic episodes due to carbon monoxide can permanently damage the blood-brain barrier and the white matter of the brain.

Carbon dioxide is a by-product of respiration and fossil fuel combustion. Increasing amounts of carbon dioxide in the atmosphere have caused concern over a global warming trend, the so-called greenhouse effect. *See* AIR POLLUTION; GREENHOUSE EFFECT.

Metals. All living organisms require certain metals for physiological processes. These elements, when present at concentrations above the level of homeostatic regulation, can be toxic. In addition, there are metals that are chemically similar to, but higher in molecular weight than, the essential metals (heavy metals). Metals can exert toxic effects by direct irritant activity, blocking functional groups in enzymes, altering the conformation of biomolecules, or displacement of essential metals in metallo-enzymes.

Solvents. Organic solvents are used widely and in large amounts in industries, laboratories, and homes

(*see* **table**). They are released to the atmosphere as vapor and can pose a significant inhalation hazard. Improper storage, use, and disposal have resulted in the contamination of surface and ground waters and drinking water. *See* WATER POLLUTION.

Pesticides. The pesticides represent an important group of materials that can enter the environment as pollutants. They are highly toxic, and many nontarget organisms can suffer harmful effects if misuse or unintended release occurs. *See* PESTICIDE.

Insecticides. The most diverse group of pesticides is the insecticides and acaricides. Insecticides have been used for hundreds of years, but it has only been since World War II that synthetic insecticides have been widely used.

The chlorinated hydrocarbons include dichlorodiphenyl trichloroethane (DDT) and its analogs, the cyclodienes and related compounds, and the benzene hexachlorides (BHCs). Beginning in the 1940s, chlorinated hydrocarbons were used extensively in mosquito control and agriculture. These compounds exert their toxic effect by interfering with the transmission of nerve impulses, inhibiting the transport of ions across the axonal membrane. Many are very persistent in the environment and undergo biomagnification in food chains. Toxic effects on top predators such as birds and the contamination of human food supplies raised concern over their use. After the 1960s, most were replaced with other, less persistent compounds. *See* FOOD WEB.

Organophosphates inhibit the enzyme acetylcholinesterase, resulting in an overstimulation and excitation of nerves. Because of their mode of action, organophosphates are toxic not only to target insects but also to nontarget insects (bees and aquatic insects), birds, wildlife, and fish and other aquatic life. With few exceptions (malathion), organophosphates are toxic to mammals and humans. They are easily degraded and do not persist in the environment. However, because of their high nontarget organism toxicity, their use is limited to areas where undesired exposure can be minimized. *See* ACETYLCHOLINE.

Carbamate insecticides exhibit neurotoxic action similar to the organophosphates, by inhibiting acetylcholinesterase. However, this action is more readily reversible than that of the organophosphates.

Botanical insecticides are natural and synthetic derivatives of toxic plant materials and have been used for hundreds of years. Nicotine and its analogs are neuroactive agents. Rotenones are electron-transport inhibitors and are used as insecticides and pesticides. Pyrethroid insecticides were originally derived from chrysanthemum plants; however, most pyrethroids in use today are synthetic derivatives of natural plant toxins. Pyrethroids are nerve poisons, acting through interference of ion transport along the axonal membrane. The pyrethroids are the most widely used insecticides, primarily because of their low mammalian toxicity. Although not very toxic to mammals, pyrethroids are extremely toxic to nontarget arthropods, such as bees, aquatic insects, and crustaceans, and to nonmammalian vertebrates such as fish.

Uses and effects of selected organic solvents		
Solvent	Uses	Effects
Aliphatic hydrocarbons Pentane, hexanes, heptanes, octanes	Commercial products and solvents	Depression of central nervous system and liver pathology
Halogenated aliphatic hydrocarbons Methylene chloride	Paint remover, aerosol solvent	Depression of central nervous system, respiratory poison
Chloroform	Chemical intermediate, solvent	Liver pathology, carcinogen
Carbon tetrachloride	Industrial and laboratory solvent	Lipid peroxidation, liver and renal pathology
Aromatic hydrocarbons Benzene	Organic synthesis, solvent, printing	Hematopoietic toxicity, immunosuppression, leukemia
Toluene	Solvent, chemical intermediate, paints, rubber	Narcotic, central nervous system pathology
Xylene	Solvent, chemical intermediate, pesticides, adhesives	Central nervous system pathology
Alcohols Methanol	Commercial and laboratory solvent, paints, fuel additive	Toxic metabolites (formaldehyde), optic nerve damage
Isopropyl alcohol	Cosmetics, glass cleaning solutions	Central nervous system depressant
Glycols Ethylene glycol	Heat exchangers, antifreeze, hydraulic fluids	Toxic metabolites

Fumigants are volatile substances used as soil pesticides and to control insects in stored products and scale insects on citrus. Common fumigants include ethylene dibromide, ethylene dichloride, ethylene oxide, carbon disulfide, methyl bromide, hydrogen cyanide, phosphine, and chloropicrin. Because they are nonselective and toxic to humans, many fumigants are now restricted or banned.

Organic thiocyanate insecticides are mild general poisons, and have been used as fly sprays and fumigants that are also toxic to nontarget animals and plants. Dinitrophenols were important as early insecticides and acaricides and still have limited use as dormant acaricides. They are toxic to mammals and plants. Fluoroacetate derivatives are general toxicants that form lethal metabolic products.

Inorganic insecticides are general toxicants and have been largely replaced with synthetic organic insecticides. However, two classes of inorganic insecticides are still being used, arsenicals and fluorides. Common arsenicals include lead and calcium arsenates and sodium arsenite. Sodium fluoride, sodium fluoraluminate, and sodium fluorosilicate are common fluorides. Both groups are stomach poisons, increasing in toxicity with increasing metal content. See INSECTICIDE.

Herbicides. These chemicals are selectively toxic to plants. Examples include the chlorophenoxy growth stimulators 2,4-dichlorophenoxyacetic acid (2,4-D) and 2,4,5-trichlorophenoxyacetic acid (2,4,5-T), the protein synthesis-inhibiting alachlor, the defoliant paraquat, and the chlorophenolic contact herbicides. Although herbicides have in large part not been an environmental problem because of their selectivity and low persistence, some can be very toxic to nontarget organisms. Paraquat, for example, can cause severe pulmonary symptoms. See HERBICIDE.

Fungicides. Fungicidal compounds are used widely to treat seed grains and wood. Common ones are pentachlorophenol and the organomercurials. See FUNGICIDE AND FUNGICIDE.

Fossil fuels. Coal and petroleum-derived materials and by-products are major environmental pollutants. The world's economy is highly dependent on fossil fuels in energy production, industry, and transportation. Widespread use has led to enormous releases to the environment of distillate fuels, crude oils, runoff from coal piles, exhaust from internal combustion-fired power plants, industrial emissions, and emissions from municipal incinerators. Point-source leaks and spills, and non-point-source emissions have resulted in environmental contamination with millions of tons of petroleum hydrocarbons each year. Spills of crude oil and fuels have caused wide-ranging damage in the marine and fresh-water environments. Oil slicks and tar in shore areas and beaches can ruin the esthetic value of entire regions.

The toxicity of polycyclic aromatic hydrocarbons is perhaps one of the most serious long-term problems associated with the use of petroleum. They comprise a large class of petroleum compounds containing two or more benzene rings. They are present in fossil fuels and are formed in the incomplete combustion of organic materials. Polycyclic aromatic hydrocarbons are formed in nature by long-term, low-temperature chemical reactions in sedimentary deposits of organic materials and in high-temperature events such as volcanoes and forest fires. The major source of this pollution is, however, human activity. Polycyclic aromatic hydrocarbons accumulate in soil, sediment, and biota. At high concentrations, they can be acutely toxic by disrupting membrane function. Many cause sunlight-induced toxicity in humans and fish and other aquatic organisms. In addition, long-term, chronic toxicity has been demonstrated in a wide variety of organisms. Through metabolic activation, some polycyclic aromatic hydrocarbons form reactive intermediates that bind to deoxyribonucleic acid (DNA). For this reason many of these hydrocarbons are mutagenic, teratogenic, or carcinogenic. They are also suspected of interfering in the reproduction of aquatic life. Low

rates of reproduction and high rates of larval deformities and mortality have been observed in fish exposed to polycyclic aromatic hydrocarbons. See FOSIL FUEL; MUTAGENS AND CARCINOGENS.

Other synthetic organic compounds. Polychlorinated biphenyls (PCBs) are produced by the chlorination of biphenyl, giving rise to mixtures of up to 210 possible products. They have been used worldwide in electrical equipment, vacuum pumps, hydraulic fluids, heat-transfer systems, lubricants, and inks. The related polybrominated biphenyls (PBBs) have been used as fire retardants. Major sources of polychlorinated biphenyls have included leaks from waste disposal facilities, vaporization during combustion, and disposal of industrial fluids. They have been identified in environments and organisms worldwide. Being lipid-soluble, they accumulate in the biota, and biomagnification in food chains can be demonstrated. They induce activity in hepatic enzymes and cause long-term liver damage, including cancer. Their use has been largely restricted or eliminated. Environmental concentrations are decreasing, but with their persistence they remain significant pollutants. See POLYCHLORINATED BIPHENYLS.

Chlorinated dibenzo-*p*-dioxins and dibenzofurans are formed during the heating of chlorophenols, and have been identified as potential contaminants in the herbicide 2,4,5-T. They can be formed during the incineration of municipal wastes, polychlorinated biphenyls, or plant materials treated with chlorophenols. In certain organisms (guinea pig) they are some of the most toxic known substances, exhibiting hepatotoxic, carcinogenic, and immunosuppressive activity. The dioxins and dibenzofurans have been implicated in several spectacular incidents, including the spraying of Agent Orange during the Vietnam War, the Yusho disease in Japan, the Seveso herbicide plant explosion in Italy, and the discovery of contamination at Love Canal, New York, at Times Beach, Missouri, and in the Great Lakes region of North America. See ENVIRONMENTAL ENGINEERING; TOXICOLOGY.

James T. Oris

Bibliography. M. O. Amdur, *Casarett and Doull's Toxicology: The Basic Science of Poisons*, 5th ed., 1996; D. W. Connell and G. J. Miller, *Chemistry and Ecotoxicology of Pollution*, 1984; E. Hodgson and P. E. Levi, *A Textbook of Modern Toxicology*, 2d ed., 1997; T. N. Veziroglu (ed.), *The Biosphere: Problems and Solutions*, 1984.

Enzyme

A catalytic protein produced by living cells. The chemical reactions involved in the digestion of foods, the biosynthesis of macromolecules, the controlled release and utilization of chemical energy, and other processes characteristic of life are all catalyzed by enzymes. In the absence of enzymes, these reactions would not take place at a significant rate. Several hundred different reactions can proceed simultaneously within a living cell, and the cell contains a comparable number of individual enzymes, each of which controls the rate of one or more of these reactions.

The potentiality of a cell for growing, dividing, and performing specialized functions, such as contraction or transmission of nerve impulses, is determined by the complement of enzymes it possesses. Some representative enzymes, their sources, and reaction specificities are shown in the **table**.

Characteristics

Enzymes can be isolated and are active outside of the living cell. They are such efficient catalysts that they accelerate chemical reactions measurably, even at concentrations so low that they cannot be detected by most chemical tests for protein. Like other chemical reactions, enzyme-catalyzed reactions proceed only when accompanied by a decrease in free energy; at equilibrium the concentrations of reactants and products are the same in the presence of an enzyme as in its absence. An enzyme can catalyze an indefinite amount of chemical change without itself being diminished or altered by the reaction. However, because most isolated enzymes are relatively unstable, they often gradually lose activity under the conditions employed for their study.

Chemical nature. All enzymes are proteins. Their molecular weights range from about 10,000 to more than 1,000,000. Like other proteins, enzymes consist of chains of amino acids linked together by peptide bonds. An enzyme molecule may contain one or more of these polypeptide chains. The sequence of amino acids within the polypeptide chains is characteristic for each enzyme and is believed to determine the unique three-dimensional conformation in which the chains are folded. This conformation, which is necessary for the activity of the enzyme, is stabilized by interactions of amino acids in different parts of the peptide chains with each other and with the surrounding medium. These interactions are relatively weak and may be disrupted readily by high temperatures, acid or alkaline conditions, or changes in the polarity of the medium. Such changes lead to an unfolding of the peptide chains (denaturation) and a concomitant loss of enzymatic activity, solubility, and other properties characteristic of the native enzyme. Enzyme denaturation is sometimes reversible. It is possible, in some cases, to obtain conditions under which the denatured enzyme refolds into its native conformation and regains its original catalytic activity. See AMINO ACIDS; PROTEIN.

Many enzymes contain an additional, nonprotein component, termed a coenzyme or prosthetic group. This may be an organic molecule, often a vitamin derivative, or a metal ion. The coenzyme, in most instances, participates directly in the catalytic reaction. For example, it may serve as an intermediate carrier of a group being transferred from one substrate to another. Some enzymes have coenzymes that are tightly bound to the protein and difficult to remove, while others have coenzymes that dissociate readily. When the protein moiety (the apoenzyme) and the coenzyme are separated from each other, neither possesses the catalytic properties of the original conjugated protein (the holoenzyme). By simply mixing the apoenzyme and the coenzyme together, the fully active holoenzyme can often be reconstituted. The

Some representative enzymes, their sources, and reaction specificities

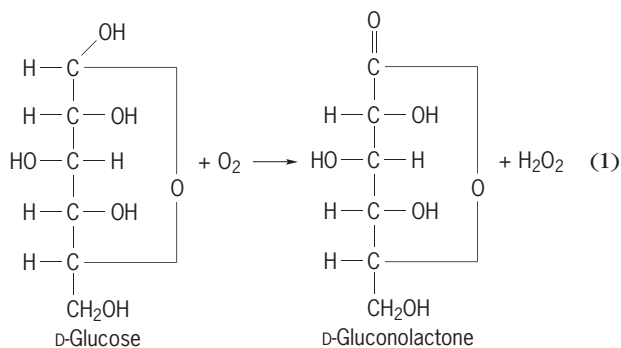
Enzyme	Some sources	Reaction catalyzed
Pepsin	Gastric juice	Hydrolysis of proteins to peptides and amino acids
Urease	Jack bean, bacteria	Hydrolysis of urea to ammonia and carbon dioxide
Amylase	Saliva, pancreatic juice	Hydrolysis of starch to maltose
Phosphorylase	Muscle, liver, plants	Reversible phosphorylation of starch or glycogen to glucose-1-phosphate
Transaminases	Many animal and plant tissues	Transfer of an amino group from an amino acid to a keto acid
Phosphohexose isomerase	Muscle, yeast	Interconversion of glucose-6-phosphate and fructose-6-phosphate
Pyruvic carboxylase	Yeast, bacteria, plants	Decarboxylation of pyruvate to acetaldehyde and carbon dioxide
Catalase	Erythrocytes, liver	Decomposition of hydrogen peroxide oxygen and water
Alcohol dehydrogenase	Liver	Oxidation of ethanol to acetaldehyde
Xanthine oxidase	Milk, liver	Oxidation of xanthine and hypoxanthine to uric acid

same coenzyme may be associated with many enzymes which catalyze different reactions. It is thus primarily the nature of the apoenzyme rather than that of the coenzyme which determines the specificity of the reaction. See BIOLOGICAL OXIDATION; COENZYME.

The complete amino acid sequence of several enzymes has been determined by chemical methods. By x-ray crystallographic methods even the exact three-dimensional molecular structure of a few enzymes has been deduced. See X-RAY CRYSTALLOGRAPHY.

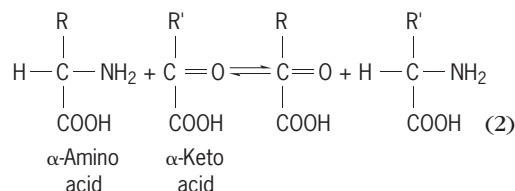
Classification and nomenclature. Enzymes are usually classified and named according to the reaction they catalyze. The principal classes are as follows.

Oxidoreductases. These enzymes, which catalyze reactions involving electron transfer, play an important role in cellular respiration and energy production. Some of them participate in the process of oxidative phosphorylation, whereby the energy released by the oxidation of carbohydrates and fats is utilized for the synthesis of adenosine triphosphate (ATP) and thus made directly available for energy-requiring reactions. Most oxidoreductases require prosthetic groups or coenzymes, which participate in the catalytic reaction by serving as intermediate electron carriers. These coenzymes include the flavin nucleotides, the pyridine nucleotides, heme, and various metal ions, such as iron, copper, and molybdenum. Some oxidoreductases, the oxygenases, catalyze the direct incorporation of oxygen into their substrate. Oxidases are enzymes which utilize molecular oxygen as an electron acceptor, while dehydrogenases remove hydrogen atoms from their substrates and transfer them to an acceptor other than oxygen. For example, glucose oxidase catalyzes the reaction shown in reaction (1). This enzyme utilizes



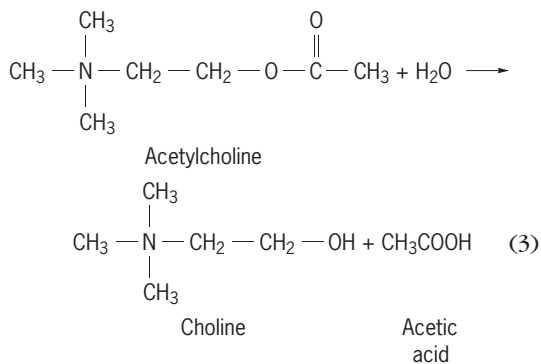
flavin adenine dinucleotide, a derivative of the vitamin riboflavin, as coenzyme. See RIBOFLAVIN.

Transferases. These enzymes catalyze the transfer of a particular chemical group from one substance to another. Thus, transaminases transfer amino groups, transmethylases transfer methyl groups, and so on. An important subclass of this group are the kinases, which catalyze the phosphorylation of their substrates by transferring a phosphate group, usually from ATP, thereby activating an otherwise metabolically inert compound for further transformations. The reaction catalyzed by transaminases may be represented by reaction (2). These enzymes require



pyridoxal phosphate or pyridoxamine phosphate, derivatives of vitamin B₆, as coenzymes. See VITAMIN B₆.

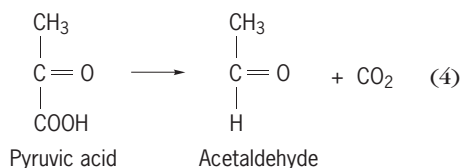
Hydrolases. These are enzymes which catalyze the hydrolysis of proteins (proteinases and peptidases), of nucleic acids (nucleases), of starch (amylases), of fats (lipases), of phosphate esters (phosphatases), and of other substances. Many hydrolases are secreted by the stomach, pancreas, and intestine and are responsible for the digestion of foods. Others participate in more specialized cellular functions. For example, cholinesterase, which catalyzes the hydrolysis of acetylcholine, as shown in reaction (3), plays



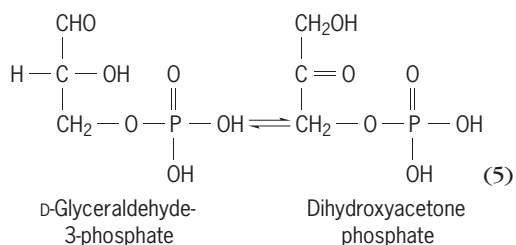
an important role in the transmission of nervous impulses. See ACETYLCHOLINE.

Lyases. These enzymes catalyze the nonhydrolytic cleavage of their substrate with the formation of a double bond. Examples are decarboxylases, which remove carboxyl groups as carbon dioxide, and dehydrases, which remove a molecule of water. The reverse reactions are catalyzed by the same enzymes, but are difficult to demonstrate in some cases because of an unfavorable equilibrium.

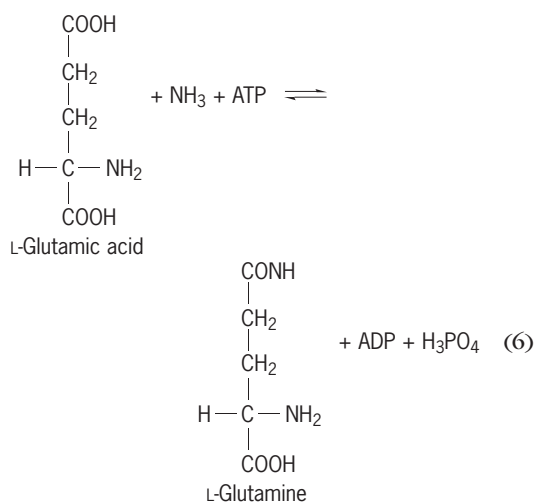
Thus pyruvic decarboxylase, which requires thiamine pyrophosphate and magnesium ions for activity, catalyzes the reaction that is shown in reaction (4).



Isomerases. These are enzymes which catalyze the interconversion of isomeric compounds. For example, triose phosphate isomerase catalyzes the reaction shown in reaction (5).



Ligases or synthetases. These enzymes catalyze endergonic syntheses coupled with the exergonic hydrolysis of ATP. They allow the chemical energy stored in ATP to be utilized for driving reactions uphill. An example is glutamine synthetase, which catalyzes the reaction shown in reaction (6). This enzyme, which



is abundant in brain tissue, requires magnesium or manganese ions for activity.

Active site. Since enzymes are large molecules and their substrates are often of low molecular weight, it appears probable that only a small portion of the enzyme protein, the active site, comes into contact with the substrate and is directly involved in catalyzing the reaction. The active site consists of a

few amino acid residues, not necessarily adjacent to each other on the polypeptide chain but brought into proximity by the folding of the chain or chains. The active site also includes any coenzyme which may be required for activity. The function of the remainder of the protein molecule may be primarily to hold the components of the active site in the proper relative position and orientation. Part of the active site is involved in binding the substrate, while another part is responsible for the making or breaking of chemical bonds. There is much evidence to indicate that, for some enzymes, the binding of the substrates produces a conformational change which brings some reactive group of the enzyme in the proper position for the reaction to take place.

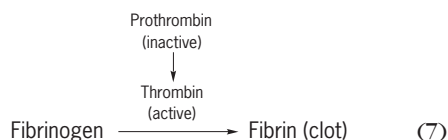
Because of their special environment at the active site of enzymes, coenzymes and certain amino acid residues acquire unique properties and reactivities which they do not possess in the free state. For example, when the proteinase chymotrypsin is treated with diisopropylfluorophosphate, one serine residue at the active site reacts rapidly with this reagent, resulting in complete inactivation of the enzyme. The other 27 serine residues of the protein do not react, and no reaction takes place when the protein molecule is first unfolded, that is, when it is denatured. On the other hand, amino acid residues remote from the active site of an enzyme may often be removed or chemically modified without affecting the catalytic activity.

Intracellular organization. Enzymes are not distributed uniformly throughout the cell but are localized within various subcellular structures or compartments. For example, the enzymes responsible for glycolysis, the degradation of glucose to pyruvate, are located in the soluble, nonparticulate portion of the cell, whereas the enzymes that catalyze the further oxidation of pyruvate to carbon dioxide and water and couple this oxidation to the production of ATP are localized in specialized particles, the mitochondria. Similarly, while the activation of amino acids for protein synthesis takes place in free solution, the assembly of the amino acids into proteins takes place on a particle, the ribosome. Other enzymes are localized on the cell membrane, in the nucleus, and in other subcellular particles. See CELL (BIOLOGY); MITOCHONDRIA; RIBOSOMES.

Control of enzyme activity. In order to preserve the stability of its internal chemical environment, a living organism must possess the means of regulating the rate of its metabolic reactions. This is accomplished by regulating both the amount and the activity of the enzymes catalyzing these reactions. Some of the control mechanisms involved are very complex and not yet completely understood, although much progress has been made in this field. A few examples are mentioned below.

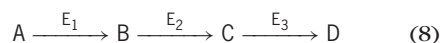
Proenzyme activation. One way by which the activity of enzymes is controlled is as follows. The enzyme is produced in an inactive form, termed a zymogen or proenzyme, and converted into the active form only when needed. This is the case with thrombin, an enzyme which causes blood to clot by catalyzing the conversion of the soluble protein, fibrinogen, to

the insoluble fibrin. Blood does not normally contain thrombin but does contain its inactive precursor, prothrombin. It is only when prothrombin is converted to thrombin (by another enzyme) that clotting results. The mechanism is shown in reaction (7). An-



other example of enzyme activation occurs with the proteinase trypsin, which is stored and secreted by the pancreas as the inactive trypsinogen. In the intestine, trypsinogen is converted to trypsin by the enzyme enterokinase. The conversion of trypsinogen to trypsin is also catalyzed by trypsin itself, the activation being therefore autocatalytic. The activation involves the splitting of a peptide from one end of the trypsinogen molecule.

Feedback inhibition. Other mechanisms of enzyme control are shown by reaction sequence (8), which



depicts a hypothetical metabolic pathway in which the starting material A is converted to the metabolite D required by the cell. One type of control, termed feedback inhibition, results from a specific, reversible, inhibition of enzyme E_1 by the end product D. This phenomenon is obviously of great advantage to the cell, since it allows the concentration of the product D to be accurately regulated: When D reaches a certain concentration, it decreases the rate of its own synthesis by preventing the conversion of A to B; the intermediates, B and C, do not accumulate; and the starting material, A, is spared for other reactions for which it may be needed. As soon as the concentration of D decreases, the inhibition of enzyme E_1 is relieved and the synthesis of D resumes. See ENZYME INHIBITION.

This type of control is also termed allosteric control, and the enzyme involved, an allosteric enzyme. The allosteric enzyme contains a specific binding site for the controlling metabolite, D, such that when D is bound to this site the enzyme is inhibited. Positive allosteric control, or allosteric activation, has also been observed. In this case, the binding metabolite, which may be the product of another pathway, causes activation of the enzyme. Some enzymes are susceptible to allosteric control, both positive and negative, by a number of different metabolites. See ALLOSTERIC ENZYME.

Induction and repression. Controls also exist at the level of enzyme synthesis. One of these controls is enzyme induction. For example, if the inducible enzymes E_1 , E_2 , and E_3 are required to utilize substance A, the enzymes are not produced by the cell when A is not available. Only in the presence of an inducer, which may be A or a related compound, are enzymes E_1 , E_2 , and E_3 synthesized. Another, related, type of control is enzyme repression. In this case the control is exerted by the end product of a metabolic pathway, for example, compound D. Thus, if enzymes

E_1 , E_2 , and E_3 are repressible enzymes, the accumulation of D prevents their synthesis by the cell. Enzyme repression and induction involve complex regulatory mechanisms which are under direct genetic control.

Purification. One of the prerequisites for a study of the structure and mechanism of action of an enzyme is its availability in pure form. The enzyme must therefore first be extracted from the cells or tissue in which it is found and separated from many proteins and other substances present in the crude extract. Because of the instability of most enzymes, however, their isolation requires special techniques, some of which are described below.

Cell destruction. If the enzyme is present in cells, the cells must first be broken. This can be accomplished by grinding with abrasives, by freezing and thawing, or by exposure to ultrasonic vibrations. Although many enzymes are readily released in soluble form by these methods, some remain associated with membranes or subcellular particles and may be difficult to solubilize.

Heating with substrate. Some enzymes can be purified by taking advantage of the fact that they are stabilized in the presence of their substrate and can then be heated for a short time to a temperature at which less stable proteins are denatured and precipitated.

Selective precipitation. Enzymes contain many ionizable groups, such as carboxyl groups, which may be negatively charged, and amino groups, which may be positively charged. Since the degree of ionization of each group depends on the pH of the solution, the total charge on the protein varies greatly with pH. The pH at which the protein has no net charge, its isoelectric point, is usually the pH of minimal solubility. Thus an enzyme can often be selectively precipitated from a solution by adjusting the pH to its isoelectric point.

Electrophoresis. An enzyme can also be separated from proteins with different electric charges by electrophoresis. This method consists of applying an electric field to the solution and allowing each protein to migrate in a direction and at a rate which depend on its charge. See ELECTROPHORESIS.

Ammonium sulfate fractionation. Another useful method of separating proteins from each other is by ammonium sulfate fractionation. The solubility of proteins decreases as the concentration of ammonium sulfate is raised, and when the salt concentration is sufficiently high, the protein precipitates. Since different proteins precipitate at different ammonium sulfate concentrations, it is possible, by increasing the salt concentration in steps and collecting the precipitate that appears at each stage, to obtain considerable purification.

Dialysis. Salts and organic impurities of low molecular weight can be easily removed from enzymes by dialysis. This consists of placing the enzyme solution in a bag made from a semipermeable membrane, such as cellophane, and immersing it in water or in a buffer solution. The protein remains in the bag, while solutes of low molecular weight diffuse out. An enzyme may be separated from smaller or larger proteins by gel filtration chromatography or by

preparative ultracentrifugation. See CHROMATOGRAPHY; ULTRACENTRIFUGE.

Other techniques. Other methods frequently used for enzyme purification include precipitation at low temperatures with organic solvents, such as acetone or alcohol; selective adsorption on gels, such as calcium phosphate or alumina; and chromatography on ion-exchange resins or modified cellulose. For the complete purification of an enzyme, a combination of these methods must be employed. No generally applicable enzyme purification procedure exists. The most satisfactory method must be determined empirically for each enzyme.

The final step in the purification of many enzymes is crystallization. The first crystalline enzyme, urease, was prepared by J. B. Sumner in 1926, and since then many other enzymes have been crystallized.

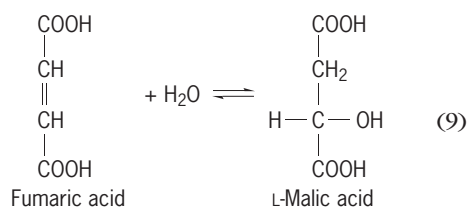
Isozymes. It has been found that many enzyme preparations, even after crystallization, still consist of a number of distinct proteins, each capable of catalyzing the same reaction. These different enzymes, originating from the same organism or even from the same cells, have been named isozymes. They can usually be separated from each other by electrophoretic or chromatographic techniques. For example, in many tissues the oxidation of lactic acid to pyruvic acid is found to be catalyzed by five different isozymes. These isozymes differ not only in electrophoretic mobility but also in some catalytic properties. Each lactic dehydrogenase molecule is composed of four subunits, of which the two main types are the H and M subunits. The existence of five isozymes is accounted for by the fact that any combination of H and M subunits may assemble in a group of four to form an enzyme molecule. Thus there exist some lactic dehydrogenase molecules containing four H subunits (H_4), some with four M subunits (M_4), and some with intermediate properties having the compositions H_3M , H_2M_2 . It may be that one of the physiological functions of isozymes is to give the cell more flexibility in metabolic regulation.

Specificity. The high degree of specificity exhibited by enzymes is one of their most characteristic properties.

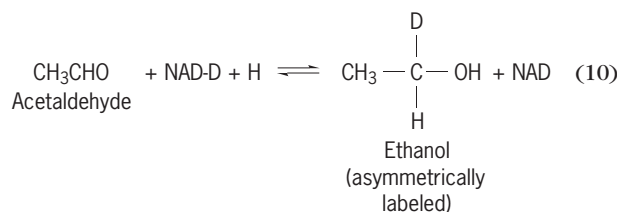
Group specificity. The majority of enzymes catalyze only one type of reaction and act on only one compound or on a group of closely related compounds. There must exist between an enzyme and its substrate a close fit, or complementarity. This relationship has been compared to that between a lock and key. In many cases, a small structural change, even in a part of the molecule remote from that altered by the enzymatic reaction, abolishes the ability of a compound to serve as a substrate. An example of an enzyme highly specific for a single substrate is urease, which catalyzes the hydrolysis of urea to carbon dioxide and ammonia. A variety of urea derivatives which have been tested are unaffected by this enzyme. On the other hand, some enzymes exhibit a less restricted specificity and act on a number of different compounds that possess a particular chemical group. This is termed group specificity.

Stereospecificity. A remarkable property of many enzymes is their high degree of stereospecificity, that

is, their ability to discriminate between asymmetric molecules of the right-handed and left-handed configurations. An example of a stereospecific, as well as a group-specific, enzyme is L-amino acid oxidase. This enzyme catalyzes the oxidation of a variety of amino acids of the type $R-CH(NH_2)COOH$. The rate of oxidation varies greatly, depending on the nature of the R group, but only amino acids of the L configuration react. Another enzyme, D-amino acid oxidase, specifically attacks amino acids of the D configuration. Conversely, when an enzyme catalyzes the formation of an asymmetric compound from symmetric substrates, the product is almost always optically active, that is, only one of the possible stereoisomers of the product is formed. For example, only the L isomer of malic acid is formed by the enzyme fumarase from fumaric acid and water, as shown in reaction (9).

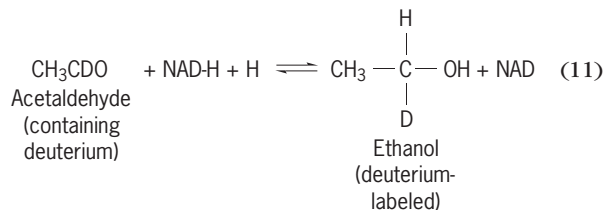


An enzyme may also act asymmetrically on a symmetrical substrate to yield a symmetrical product. Such specificity can be detected by the use of isotopes. Thus alcohol dehydrogenase, which catalyzes the conversion of ethanol to acetaldehyde, removes a specific hydrogen atom from carbon atom 1 of ethanol. This hydrogen atom is transferred to the nicotinamide ring nicotinamide adenine dinucleotide (NAD). If reduced NAD containing deuterium (D) is used in the reversal of this reaction, asymmetrically labeled ethanol is produced, as shown in reaction (10). No deuterium is found in



the acetaldehyde arising from enzymatic dehydrogenation of this ethanol. See NICOTINAMIDE ADENINE DINUCLEOTIDE (NAD).

The enzyme, however, cannot distinguish between hydrogen and deuterium. This is shown by the fact that when acetaldehyde containing deuterium on carbon atom 1 is hydrogenated with unlabeled NAD, the other isomer of deuterium-labeled ethanol is formed, as shown in reaction (11). On subsequent



enzymatic dehydrogenation of this ethanol, all of the

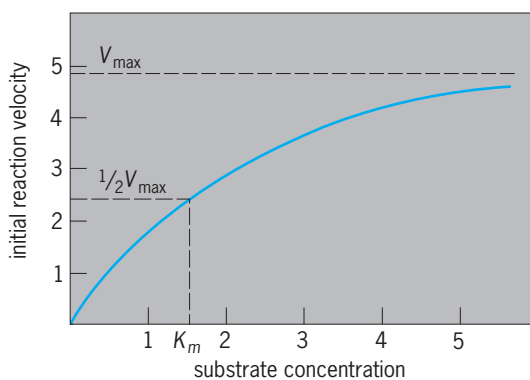
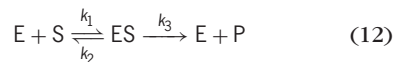


Fig. 1. Initial velocity of an enzyme-catalyzed reaction as a function of substrate concentration.

deuterium is found in the acetaldehyde. See PROCHIRALITY; STEREOCHEMISTRY.

Kinetics. Much information can be obtained concerning the mechanism of an enzyme-catalyzed reaction by studying the effect of pH, temperature, and the concentration of substrates, coenzymes, inhibitors, and activators on the velocity of the reaction. Such studies were already being carried out in the nineteenth century, long before anything was known about the chemical nature of enzymes.

Michaelis-Menten equation. It has been known for a long time that, unlike the velocity of many uncatalyzed chemical reactions, which increases linearly with the concentration of one of the reactants, the velocity of enzymatic reactions levels off and approaches a maximal value as the substrate concentration increases (**Fig. 1**). A mathematical theory accounting for this behavior was first proposed by V. Henri in 1903 and extended by L. Michaelis and L. M. Menten in 1913. The theory was based on the idea that the reaction velocity is proportional to the concentration of an enzyme-substrate complex ES formed by a specific combination of the enzyme E with its substrate S . This complex was assumed to react further, yielding a product P and regenerating a free enzyme. These reactions can be represented by reaction (12), where k_1 , k_2 , and k_3 are rate con-



stants for the three reactions, whose velocities are v_1 , v_2 , and v_3 , respectively.

The existence of an enzyme-substrate complex has been amply confirmed experimentally, and in some cases such a complex has been directly observed.

Under the conditions usually employed in kinetic studies, the concentration of enzyme is very small compared to that of substrate. When enzyme and substrate are mixed, the concentration of ES increases rapidly until a steady state is reached in which the rate of decomposition of the complex equals its rate of formation. Such a steady state is usually reached in less than a second and may last for several minutes. During this time, the velocity of the overall reaction is essentially constant. This velocity, which may be measured as the rate of product for-

mation, is termed the initial velocity of the reaction. An expression for the initial velocity v_0 in terms of the substrate concentration $[S]$ can now be derived simply.

According to the law of mass action, $v_1 = k_1[S][E]$, $v_2 = k_2[ES]$, and $v_3 = k_3[ES]$; by definition, $v_0 = v_3$. The total enzyme concentration $[E_t]$ is $[E] + [ES]$, and the steady-state conditions state that $v_1 = v_2 + v_3$. These relations can be solved for v_0 to give Eq. (13).

$$v_0 = \frac{k_3[E_t]}{1 + \frac{k_2 + k_3}{k_1[S]}} \quad (13)$$

From Eq. (13), it is seen that, for any given substrate concentration, the initial velocity is directly proportional to the enzyme concentration. As the substrate concentration tends to infinity, the initial velocity approaches its maximal value (v_{\max}), as expressed by notation (14).

$$v_{\max} = k_3[E_t] \quad (14)$$

If K_m is defined as $k_2 + k_3/k_1$, the expression for v_0 becomes Eq. (15). This is the Michaelis-Menten equa-

$$v_0 = \frac{v_{\max}}{1 + \frac{K_m}{[S]}} \quad (15)$$

tion, and the Michaelis constant K_m represents the substrate concentration required to obtain one-half the maximum velocity with a given concentration of enzyme. When k_3 is small compared to k_2 , K_m approximates the dissociation constant of the enzyme-substrate complex and is therefore a measure of the affinity of an enzyme for its substrate. The equation predicts that, when v_0 is plotted as a function of $[S]$, a rectangular hyperbola should be obtained, as shown in **Fig. 1**. This closely approximates the results found with a number of enzymes. Most enzymatic reactions, however, involve more than one substrate, and a number of different enzyme-substrate complexes may be formed before the appearance of the product and the regeneration of free enzyme. It may also be necessary in some cases to take into account the effect of product, inhibitors, or activators on the velocity of the reaction. With new instrumentation, it has been possible to study the kinetics of the transient, pre-steady-state portion of the reaction and to determine the value of the individual rate constants for some enzymatic reactions.

Temperature optima. Chemical reactions are usually accelerated by an increase in temperature, because molecules require a minimum quantity of energy to react and at higher temperatures a greater proportion of the molecules possess this energy. When the rate of an enzymatic reaction is determined as a function of temperature, it is found that the rate increases with temperature up to a point and then abruptly decreases. This is because most enzymes are unstable and undergo rapid denaturation above a critical temperature. The temperature optimum for many enzymatic reactions is about 37°C (99°F).

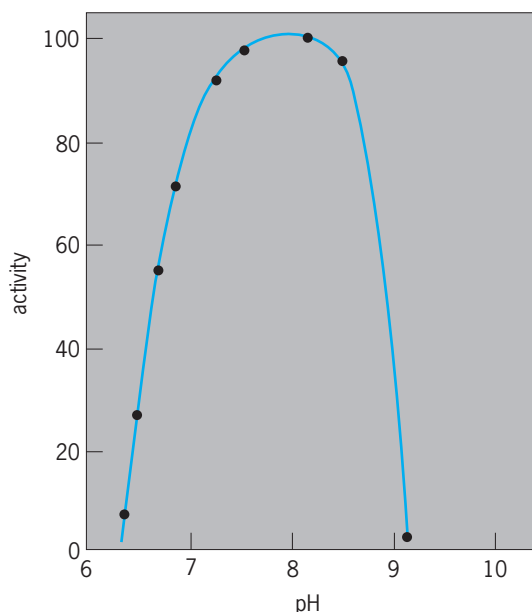


Fig. 2. Effect of pH on the velocity of the reaction catalyzed by the enzyme carbamyl phosphate synthetase.

pH optima. Enzyme activity is also influenced considerably by the hydrogen-ion concentration of the reaction medium (Fig. 2). An enzyme may exist in several different ionic forms, and the ionic species that predominates is a function of pH. Only one of these forms can be catalytically active, or one form may be more active than the others. The pH optimum for different enzymes varies over a broad range, from about pH 2 for pepsin to about pH 10 for arginase. Some enzymes exhibit a broad optimal range extending over several pH units, while others have a very sharp pH optimum. The pH optimum for an enzyme, like the temperature optimum, may vary, depending on the substrate being used and on other experimental conditions.

Daniel Wellner

Immobilization. In many instances, the use of soluble enzymes is not feasible because of high cost. Immobilization of the enzymes offers the benefit of reuse of the catalyst, thus reducing expenses, and has a second possible advantage of enhanced reaction stability.

The principal techniques for enzyme immobilization are adsorption, covalent attachment, and entrapment within polymer gels. Many enzymes have a preponderance of either positive or negative amino acid side chains on the surface, and electrostatic attraction at multiple points between a solid support and an enzyme leads to a stable complex. Hydrophobic forces may also contribute to stability of noncovalent complexes of enzyme with water-insoluble materials. Covalent attachment requires reagents that create chemical linkages (between enzyme and support). For entrapment, the enzyme is mixed with polymerizable elements, such as vinyl monomers or "prepolymers." The polymerization yields a gel containing the enzyme entrapped within it.

Properties of immobilized enzymes. In most cases, enzymatic activity is reduced as a result of attachment of

enzymes to solid supports. Loss of activity may occur either because of chemical reaction of a functional group at or near the enzyme active center or because of physical denaturation. When charged supports are used, microenvironmental effects result in the perturbation of the pH-rate profile and the apparent reactant-binding capability of the enzyme. Both effects are reduced by increased ionic strength of the bulk solution. The transport of reactant to the surface of an enzyme particle can be limited by the unstirred (Nernst) layer surrounding the particle. The supply of reactant to enzyme molecules in the interior of the supporting particle may also be rate-limiting.

Immobilized enzyme reactors. The attachment of enzymes to insoluble supports permits their use in reactors which resemble those employed in the chemical process industry. Tank reactors are often preferred for intermittent, small-scale production. Viscous reactant solutions and enzyme particles with low activity might also dictate the use of tank reactors. For large-scale processes, continuous packed-bed reactors are frequently selected. Upward-flow reactors are used when the enzyme particles are compressible and when gaseous products are produced. The fluidized-bed reactor, where the enzyme particles are kept in suspension by the upward flow of reactant, can be used with viscous reactant solutions or when gaseous products are produced.

The development of ultrafiltration membranes and hollow-fiber devices has led to the construction of reactors which depend on the fact that the enzyme (of high molecular weight) can be retained by the membrane, while the product (of low molecular weight) is allowed to pass. Such reactors are not appropriate for processes in which proteases are present by choice or as contaminants; proteases are often subject to autolysis, which means that they digest themselves as well as the reactant. Contaminating proteases can severely limit the lifetime of an enzyme in this type of reactor. In contrast, when the enzyme preparation is covalently immobilized on particles, the contaminants are also immobilized (or washed through), and proteolytic degradation is eliminated.

Application of immobilized enzymes. A large-scale process based on immobilized enzyme technology is the production of high-fructose corn syrup, much of which is used as a sweetener for soft drinks. In this process, glucose syrup (made enzymatically from corn-starch) is passed over a bed of immobilized glucose isomerase; about 50% of the glucose is isomerized to fructose. Because fructose is sweeter than glucose, high-fructose corn syrup is sweeter than the original glucose solution; its sweetness is equivalent to a solution of sucrose (table sugar) of comparable concentration.

In the production of 6-amino penicillanic acid from penicillin G, the phenylacetyl moiety is released from the amide linkage, but the lactam amide is left intact. The transformation occurs quantitatively under very mild conditions. This application is an excellent example of the exploitation of enzymatic specificity.

In the production of optically pure amino acids

from a racemic mixture with immobilized acylase, the enzyme exhibits preference for one optical isomer; only the L enantiomer is deacylated. The resulting L-amino acid crystallizes, and the acyl-D-amino acid remains. The process can be run continuously by rracemization of the D isomer to the D,L mixture.

Chemical analysis is a very important area of application for immobilized enzymes. Here again, enzyme specificity is exploited. Enzymes can selectively transform a single compound in a complex mixture. The immobilized enzyme may be coupled with a thermal detector or an ion-selective electrode. Coating the latter with an enzyme-polymer layer results in an "enzyme electrode." An even simpler approach is to cover the tip of a conventional electrode with an ultrafiltration membrane which contains the desired enzyme solution. A cation-selective electrode immersed in a solution of glutamate dehydrogenase has been successfully used for the detection of glutamic acid. *See* ION-SELECTIVE MEMBRANES AND ELECTRODES.

Medical applications also appear promising. For example, a deheparinization method for blood has been proposed. The use of extracorporeal devices for blood oxygenation during bypass surgery or kidney dialysis requires heparinization of the blood to prevent clotting. However, heparin treatment can lead to hemorrhagic complications. To prevent these problems, an immobilized heparinase reactor is placed in the external circuit, prior to the point of blood reentry. Heparinase in the column completely eliminates the heparin anticoagulant activity, and heparinization is restricted to the external circuit. Immobilized heparinase may find applications in other extracorporeal devices used in enzyme therapy.

Certain types of leukemia are asparagine-dependent; when the tumor cells are deprived of asparagine, proliferation slows dramatically. The enzyme L-asparaginase degrades asparagine. Depletion of asparagine by direct injection of asparaginase into a patient is accompanied by two problems which relate to the fact that the asparaginase most commonly used is of bacterial origin. The reticuloendothelial system recognizes the enzyme as alien and clears it from the bloodstream. Also, life-threatening allergic reactions may result following administration of the foreign protein asparaginase. To immobilize the enzyme within biocompatible capsules would appear to be a potential solution to the problems. Unfortunately, immobilized enzymes are not generally good catalysts at low concentration of reactants. Since the tumor cells can survive at low levels of asparagine, effective enzyme therapy requires reduction of the asparagine concentration to very low levels. The consequence is that leukemia therapy based on the use of immobilized asparaginase remains in the experimental stage.

G. P. Royer

Production

A unique property of enzymes, both proteinaceous and RNA-derived, is the ability to catalyze chemical reactions in a mild aqueous environment, that is, with relatively low temperatures (58–212°F or

15–100°C), atmospheric pressure, and moderate pH. Thus enzymes are indispensable where harsh chemical treatment or extreme conditions are unacceptable, such as in the baking, brewing, food processing, detergent, leather, and textile industries. Novel enzymes have been developed which catalyze chemical reactions in nonaqueous organic solvents, at elevated temperatures, under conditions of extreme pH. Moreover, the chemical specificity with which enzymes catalyze chemical reactions in mild or relatively harsh conditions has led to the use of enzymes for industrial-scale biochemical processes, including synthesis of peptides, chemical feedstock compounds, and specialty compounds. *See* FOOD MANUFACTURING; LEATHER AND FUR PROCESSING.

Sources. Enzymes for traditional industrial use (for example, food processing, laundry detergents, and leather processing) must be produced in large quantities and at modest cost. Consequently, they frequently are obtained as by-products from biological materials that are available in large quantities. Three common sources are plants (malt diastase, papain, bromelain, and ficin), animals (pancreatic enzymes, and digestive enzymes such as lipases, pepsin, and rennin), and microbes (amylases, proteases, cellulases, lipases, pectinases, invertase, and glucose oxidase). The number of plant- and animal-derived enzymes tends to be restricted by the availability of raw material and cost of recovery. For example, animal-derived enzymes are a low-cost by-product of the meat packing industry. However, the supply of animal-derived enzymes is restricted by the number of animals processed at meat packing plants.

To ensure a consistent supply of low-cost enzymes, industrial enzyme manufacturers have relied on microbes as a source. In fact, even specialty industrial enzymes used in the biomedical and biotechnology industries are obtained from microbes.

The use of bacteria, fungi, and yeast can be attributed directly to the fact that they can be grown rapidly in large quantities on inexpensive culture media. Moreover, microbes can be manipulated by traditional mutagenesis and modern genetic engineering techniques to increase the yield of specific enzymes. Thus, enzymes derived from microbes are available in virtually unlimited amounts at high yield and lower cost than those derived from plants and animals.

Fermentation. The process by which microbes are cultivated is called fermentation. When used by bacterial physiologists, this term refers to the anaerobic metabolism of growth substrates. However, in industry the term is applied to any process of cultivating microbial cells in a liquid medium. Microbial fermentation takes two forms, submerged fermentation and semisolid culture. In recent times, the former has predominated, and the latter has been used to only a very limited extent. Production can be divided into three stages: ferme recovery, by which a concentrated solution of enzyme is produced, enzyme recovery, and purification. *See* BACTERIAL PHYSIOLOGY AND METABOLISM.

A pure culture of the organism is inoculated into a sterilized liquid broth (250–1000 milliliters) contained in Erlenmeyer flasks. The growth medium

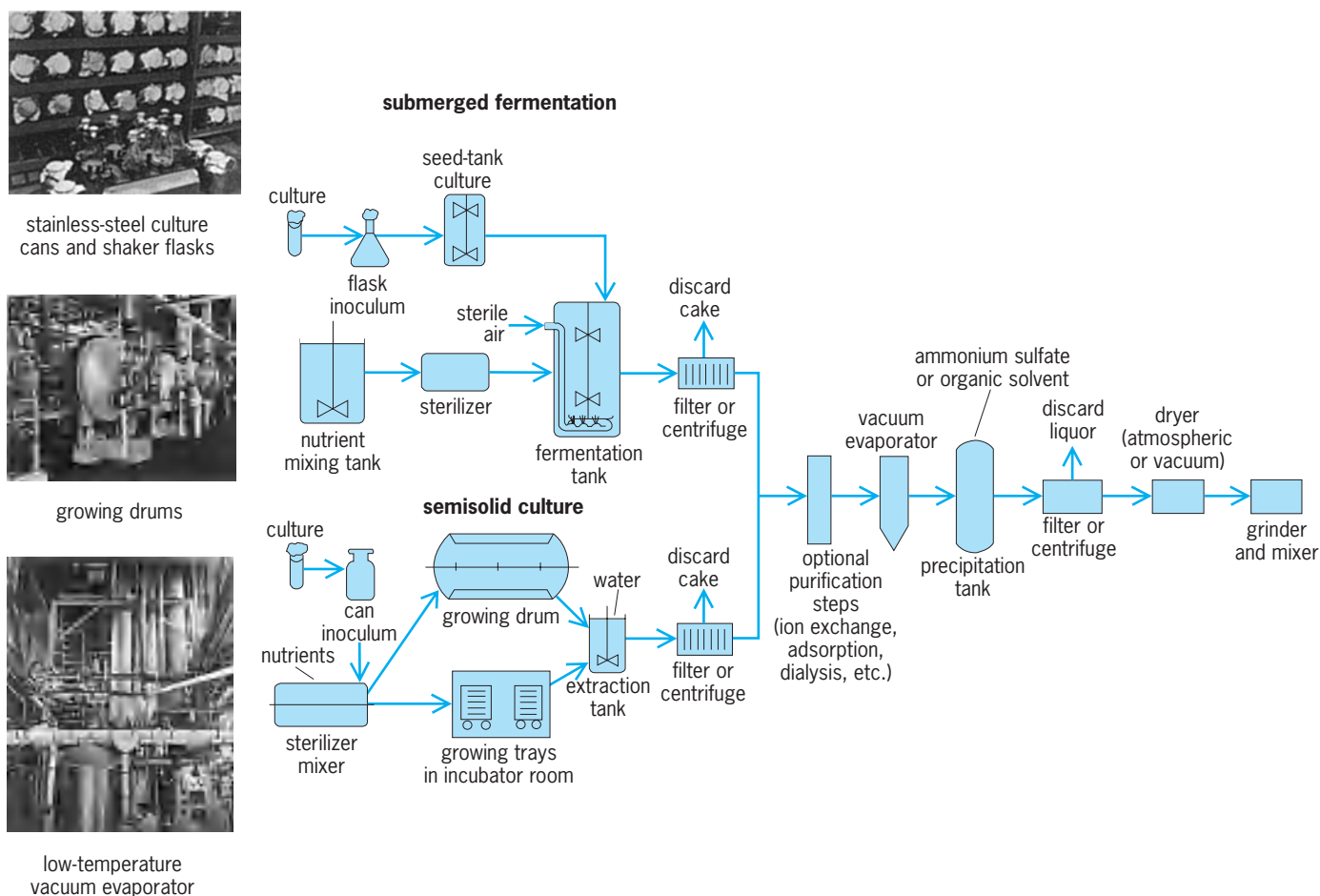


Fig. 3. Procedures for enzyme production by industrial fermentation. (Wallerstein Co., Inc.)

typically employed in the initial and final stages of fermentation is composed of an inexpensive carbon and energy source (potato starch, corn syrup), a nitrogen and sulfur source (hydrolyzed casein or beef extract, ammonia and sulfur salts), minerals, and growth stimulants. The growth in flasks is used to inoculate seed tanks containing 10 to several hundred liters of sterilized culture broth. In the final stage of submerged fermentation, the contents of the seed tanks are transferred to a culture medium in production tanks (4000–100,000 liters). These are stainless steel tanks equipped for the passage of air sterilized by filtration, propeller agitation, temperature control by water jacket or thermal coils, pH control, and nutrient addition (Fig. 3).

Depending on the microorganism, enzyme production is maximal after 1 day to a week. The fermentation tank is then emptied for enzyme recovery. See CULTURE; FERMENTATION.

Enzyme recovery. The first step in recovery is to concentrate the enzyme preparation. The procedure depends on whether the enzyme is contained within the cell cytoplasm (glucose oxidase, invertase) or is secreted into the medium (protease, amylase, pectinase). Enzymes contained within cells can be concentrated by centrifugation. This is followed by disintegration of cells by autolysis or a mechanical method such as high-frequency sonic vibration. In the former procedure, enzymes produced by cells degrade the

microbial cell wall, resulting in cell lysis and release of the contents. Enzymes are recovered by vacuum evaporation after intact cells and other insoluble constituents have been removed by centrifugation or filtration. See CENTRIFUGATION; FILTRATION.

Purification. Enzyme concentrates prepared as described above are frequently marketed without further purification after addition of suitable buffering and stabilizing agents. Depending on the application, further steps may be required before marketing. Additional purification may employ standard biochemical methods such as fractional precipitation by solvents (methanol, acetone, isopropanol) or salts (ammonium sulfate, calcium acetate). Fractional precipitation is rapid, inexpensive, and applicable to very crude enzyme suspensions. It is often the first step in purification of crude broths, primarily to remove low-molecular-weight compounds with objectionable odors or tastes.

Techniques such as adsorption, ion-exchange chromatography, and gel filtration chromatography, although uncommon for enzymes used for traditional industrial applications, are typical for the purification of enzymes used in the biomedical and biotechnology industries. The added purification steps are often included to eliminate cross-contamination by enzymes or other products that would render the enzymes unsuitable for use. This is particularly true for enzymes used in the biomedical field (such

as streptokinase, cereadase, and DNase) that require intravenous administration. Enzyme preparations contaminated with endotoxin, a natural by-product of gram-negative bacteria, result in septic shock upon intravenous injection. This potentially life-threatening reaction can be eliminated only by rigorous purification of enzyme preparations prior to marketing. See BACTERIA; ENDOTOXIN.

Genetic engineering. The availability of industrial enzymes has been limited by identification of suitable, naturally occurring sources of the enzyme. Genetic engineering, including recombinant DNA cloning, site-directed mutagenesis, and the polymerase chain reaction, has made it possible to extend the range of naturally occurring sources and to create completely novel enzymes that are tailored to meet the specific requirements of industrial processes.

The recombinant DNA cloning and expression of foreign genes from animal, plant, and microbial cells in high-yielding bacterial cells has generated a virtually unlimited supply of enzymes. The cloning of enzyme-encoding genes on autonomously replicating extrachromosomal DNA elements (plasmids) frequently results in the high-level expression of the cloned gene product. High-level expression results because certain plasmids exist in hundreds of copies per bacterial cell. The presence of hundreds of copies of a gene that encodes a single enzyme enables the cell to synthesize an enormous quantity of the cloned gene product compared to the source from which the gene was obtained originally. Consequently, recombinant DNA cloning of an enzyme-encoding gene can provide, in certain instances, high yields of the enzyme irrespective of the gene's source.

Recombinant DNA manipulations have been used to further enhance enzyme production from cloned genes by genetically manipulating promoters. These are DNA sequences, located upstream from the region of the DNA encoding the gene product, that are required for transcription and translation of the gene product. Enzyme manufacturers are, therefore, no longer restricted to the use of enzymes from high-yield natural sources. It is possible to identify an enzymatic activity that suits the needs of an existing process, clone the gene encoding the desired enzyme, and manipulate the gene's expression in a bacterial or other suitable host to provide a continuous and relatively low cost supply of the enzyme.

Recombinant DNA methods also provide the means to create novel enzymes. The structure of individual genes can be manipulated by fusion or by the rearrangement of distinct genes from different organisms, thereby changing the characteristics of the product encoded by the gene. In combination, these methods are used to create and modify the activities of enzymes used in industrial processes.

Another extremely powerful technology is the use of random mutagenesis of cloned genes by the polymerase chain reaction technique combined with screening of the mutagenized gene product under se-

lective conditions. Although the approach is referred to as random mutagenesis, the mutagenesis is restricted to a specific region of a cloned gene bounded by polymerase chain reaction primers. By restricting the location of mutations, it is possible to target mutations to a general region of cloned DNA, for instance the region corresponding to the active site of an enzyme. After generating DNA that has been randomly mutated within a defined region, the DNA is cloned into appropriate sites within the gene of interest. The mutated genes are introduced into a host organism that can express the altered enzymes derived from the mutated genes. Cells containing the mutated genes are then screened under selective conditions to reveal those cells that are expressing altered gene products with novel activities. This method has been used to obtain enzymes with altered activities compatible with industrial uses. The development of enzymes with catalytic activity in nonaqueous solvents has profound implications for industrial use, particularly since many industrial processes are performed in organic rather than aqueous solvents. Moreover, combining the use of random polymerase chain reaction mutagenesis with screening under selective conditions circumvents the need for time-consuming and expensive x-ray crystallographic analysis of enzyme three-dimensional structure. See GENETIC ENGINEERING; MOLECULAR BIOLOGY.

Edmund J. Stellwag

Bibliography. P. B. Chock et al. (eds.), *Enzyme Dynamics and Regulation*, 1987; P. A. Frey (ed.), *Mechanisms of Enzymatic Reactions: Stereochemistry*, 1986; P. Gacessa and J. Hubble, *Enzyme Technology*, 1987; C. P. Hollenberg and H. Sahm (eds.), *Microbial Genetical Engineering and Enzyme Technology*, 1987; T. Palmer, *Understanding Enzymes*, 1985.

Enzyme inhibition

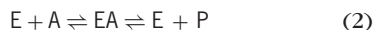
The prevention of an enzymatic reaction due to interaction of an enzyme with a substance that decreases the rate of catalysis by the enzyme. The compound causing the effect is termed an inhibitor. Enzyme inhibitors are important regulators in the normal control of enzymatic reactions in living organisms, are useful in the study of enzyme mechanisms and cellular reactions, and are important chemotherapeutic agents. Inhibitors may have chemical structures that are similar to one of the substrates of the enzyme but may also be quite different.

Classification of inhibitors. Inhibitors bind reversibly to different forms of the enzyme during the catalytic reaction. They include dead-end inhibitors, which may resemble structurally one of the substrates but do not undergo the catalytic reaction; and reaction products, which can act to make the reaction proceed in the reverse direction. In reaction (1), substrate A is converted to product



P. The simplest enzymatic mechanism is shown in

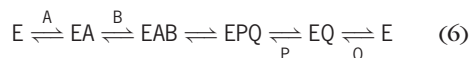
reaction (2), where the enzyme E forms an enzyme-



substrate complex EA, which then converts A to P and releases the product. A dead-end inhibitor (I_1 , where subscripts 1, 2, 3 refer to different chemical compounds) may bind to the free enzyme, reaction (3), or to the EA complex, reaction (4). In principle,



an inhibitor also could bind to both forms of enzyme. Most enzymatic reactions have more than one substrate and product, and the reaction of two substrates to form two products is shown in reaction (5). One possible bi-substrate enzyme mechanism, reaction (6) shows that the enzyme binds the substrates



and releases the products in an ordered sequence. For such a mechanism, an intermediate form of enzyme may bind inhibitors, as in reaction (7).



Three kinds of inhibitors are classified and distinguished experimentally by a study of the effects on the rate of the enzymatic reaction at different concentrations of the inhibitor and varied concentrations of one of the substrates.

Competitive inhibition. Competitive inhibition results when the inhibitor binds reversibly to the same form of enzyme to which the substrate binds. For instance, the dead-end inhibitor I_1 in reaction (3) is a competitive inhibitor against varied concentrations of substrate A. The product P in reaction (2) or the product Q in reaction (6) are also competitive inhibitors against substrate A, as these compounds bind to enzyme form E.

Uncompetitive inhibition. Uncompetitive inhibition results when a dead-end inhibitor binds to a form of the enzyme after the substrate binds, for instance to EA, as shown in reaction (4). Uncompetitive inhibition can also result when the inhibitor binds to a subsequent intermediate form of enzyme, such as in reaction (7) after substrates A and B bind.

Noncompetitive inhibition. Noncompetitive inhibition can result from a variety of reactions, for instance, if a dead-end inhibitor binds to both free enzyme and to the enzyme-substrate complex, as in reactions (3) and (4), where I_1 and I_2 are the same compound. Inhibitor I_1 and product P are also noncompetitive inhibitors against substrate B in reaction (6). See CATALYSIS.

Kinetics of inhibition. The inhibition effects are described quantitatively by a general equation (8),

$$v = VS/[K_m(1 + I/K_{is}) + S(1 + I/K_{ii})] \quad (8)$$

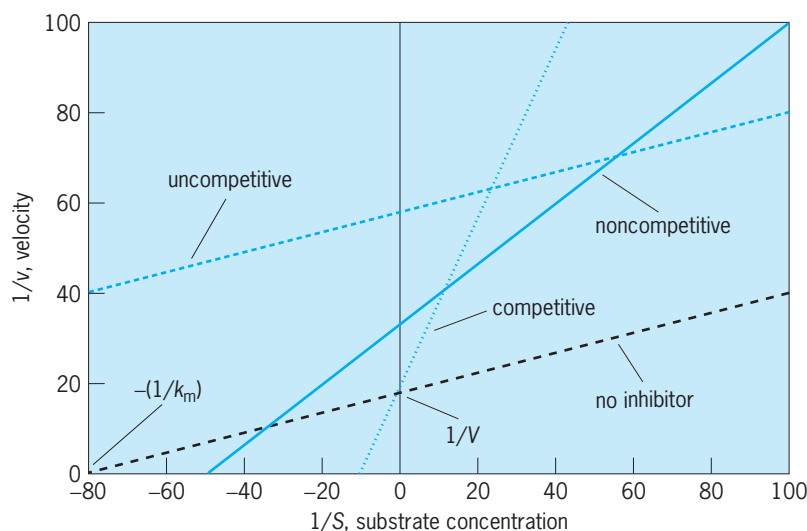


Fig. 1. Different kinds of inhibition.

which is a modification of the Michaelis-Menten equation and which can be written in its double reciprocal form (9). Here v is the initial velocity of

$$1/v = (K_m/V)(1 + I/K_{is})(1/S) + (1/V)(1 + I/K_{ii}) \quad (9)$$

the reaction at various concentrations of a substrate (S), V is the maximum velocity of the reaction when the substrate concentration is much higher than the Michaelis constant (K_m), and K_{is} and K_{ii} are dissociation constants named for the effects of the inhibitors on slopes and intercepts, respectively, of a plot of $1/v$ against $1/S$. Thus, the different kinds of inhibition can be distinguished experimentally by studying the effects of one or more concentrations of inhibitor on the velocity as represented by a double reciprocal (Fig. 1).

A competitive inhibitor affects the slope but not the intercept ($1/V$) at infinite concentration of S of a plot of Eq. (9), because K_{is} has a finite value, whereas the term $(1 + K_{ii})$ is not present in the equation. Since the competitive inhibitor and the substrate bind to the same form of enzyme [such as reactions (2) and (3)], when the substrate concentration is sufficiently high to saturate the enzyme with substrate, the inhibitor has no effect, and the same maximum velocity is obtained as in the absence of inhibitor. In contrast, an uncompetitive inhibitor affects only the intercept of the plot, because K_{ii} has a finite value and the term $(1 + K_{is})$ is not present in the equation. Since the uncompetitive inhibitor and the substrate bind to different forms of enzyme [see reactions (2) and (4)], high substrate concentrations will not prevent the action of the inhibitor; the apparent maximum velocity will be lower, and the apparent K_m value will be higher than in the absence of the inhibitor. A noncompetitive inhibitor affects both the slopes and intercepts of the plot, because K_{is} and K_{ii} have finite, and generally different, values. In general, the apparent K_m value for the substrate, determined from the reciprocal of the intercept on the abscissa ($1/K_m$), will be altered by an inhibitor. If an enzyme is

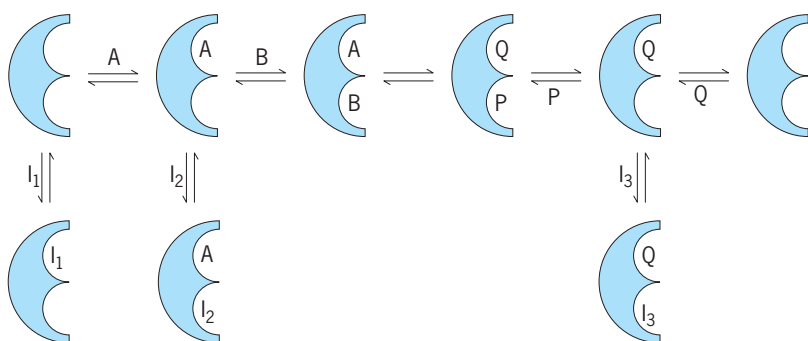


Fig. 2. A typical enzyme may have two binding sites and be subject to be different kinds of inhibition.

inactivated by an irreversible chemical reaction, the kinetics resemble noncompetitive inhibition with equal values for K_{is} and K_{ii} , simply because there is less enzyme present and the apparent K_m value is not changed by the inhibitor. See ENZYME.

Binding modes. Identifying where an inhibitor binds to an enzyme usually requires that the structure of a complex of the enzyme and the inhibitor be determined by a physical method such as x-ray crystallography or nuclear magnetic resonance (NMR). Defining the binding site by the type of inhibition is ambiguous. A competitive inhibitor may bind either to the same site as the substrate or to a different site. A noncompetitive inhibitor may also bind to the same or different sites.

With reference to reaction (6), a typical enzyme may have two binding sites (Fig. 2). Substrates A and B bind to adjacent sites and react to form products P and Q. As a product inhibitor, Q is competitive against varied concentrations of substrate A, since both compounds bind to free enzyme. In contrast, product P is a noncompetitive inhibitor against substrate B, since B binds to EA and P binds to EQ, and the presence of P reverses the reaction sequence. Thus, the inhibition is noncompetitive even though B and P bind to the same site. For the dead-end inhibitors, I_1 is competitive against A and noncompetitive against B, I_2 is competitive against B and uncompetitive against A, and I_3 is competitive against P and uncompetitive against either A or B.

Binding of an inhibitor or substrate may also change the conformation of the enzyme and prevent binding of the other, as shown in Fig. 2, where the inhibitor binds at an allosteric site, which is different from the site to which substrate binds. In this case, the inhibition is competitive, because the substrate and inhibitor bind to the same form of enzyme in a mutually exclusive manner. In general, an inhibitor binding at an allosteric site can be noncompetitive, because it can bind to both the free enzyme and the enzyme-substrate complex. See ALLOSTERIC ENZYME.

Inhibitors and metabolism. Enzymes catalyze the reactions in the metabolism of foodstuffs that yield energy and the materials of living cells. The reactions are controlled by the activities of the enzymes, which are inhibited by products of the reactions and inhibited or activated by various allosteric effectors. The controls are required to maintain the optimal steady-

state concentrations of metabolites. For instance, the three dehydrogenases in the citric acid cycle produce the reduced form of nicotinamide adenine dinucleotide and are inhibited by this product when its concentration becomes high.

Another important control is feedback inhibition of an enzyme that catalyzes a rate-limiting step in a metabolic pathway. The end product of the pathway inhibits the enzyme and prevents production of compounds that are used only for the synthesis of the final product. A textbook example is a bacterial aspartate transcarbamoylase, which produces carbamoylaspartate from carbamoyl phosphate and aspartic acid, and is inhibited by cytidine triphosphate, the end product of the ten-step pathway. Interestingly, another nucleotide, adenosine triphosphate, activates the enzyme, and both nucleotides affect the Michaelis constant for aspartic acid and not the maximum velocity. Thus the nucleotides are "competitive," but they bind at an allosteric site, which is structurally distant from the active site where the substrate aspartic acid binds.

Antimetabolites and antibiotics. A major effort in the battle against disease is the synthesis and evaluation of new enzyme inhibitors. Many enzymes are targets of chemotherapy.

Inhibitors that are structurally similar to the substrate can bind to the same site on the enzyme as the substrate. For instance, malonic acid resembles succinic acid and competitively inhibits succinate dehydrogenase, an enzyme in the citric acid cycle (Fig. 3). As in Fig. 1, products bind to the same site as the corresponding substrate. However, substrate binding sites in enzymes are flexible and can accommodate structurally related molecules. This has led to the development of many kinds of antimetabolites that have been used to inhibit particular steps in metabolic pathways and to exploit the inhibition for therapeutic purposes. The antibiotic sulfanilamide competitively inhibits the incorporation of *p*-aminobenzoic acid into folic acid, an essential vitamin. The antibiotic is efficacious because bacteria cannot survive when they cannot synthesize folic acid, but humans do not make folic acid and normally consume sufficient amounts in the diet. See ANTIBIOTIC; CITRIC ACID CYCLE; SULFONAMIDE.

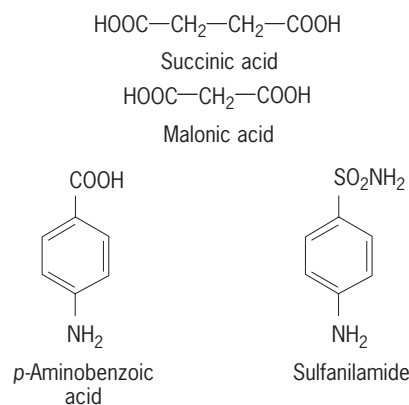


Fig. 3. Structural formulas for two essential metabolites (succinic acid and *p*-aminobenzoic acid) and their competitive inhibitors (malonic acid and sulfanilamide).

The anticancer drug methotrexate inhibits dihydrofolate reductase, which is in the pathway for the biosynthesis of thymidylate and is essential for synthesis of DNA. Methotrexate is a competitive inhibitor against dihydrofolate and binds in the same site of the enzyme where the folate binds. In contrast, nonnucleoside inhibitors of the reverse transcriptase encoded by human immunodeficiency virus bind at an allosteric site some distance from the active site. See CHEMOTHERAPY.

Uncompetitive and noncompetitive inhibitors offer a potential kinetic advantage for inhibiting a metabolic step, because high concentrations of the substrate cannot abolish the effect of the inhibitor. Examples are formamides and sulfoxides, which resemble chemically the product aldehyde that forms when alcohol dehydrogenase oxidizes an alcohol. The mechanism is given in reactions (6) and (7), where the alcohol is substrate B and the inhibitor I_3 binds to EQ after the aldehyde product P has dissociated. These inhibitors effectively inhibit alcohol metabolism.

Inhibitors that chemically react at a substrate binding site irreversibly inactivate the enzyme and are also very effective inhibitors. These may be called active-site-directed or suicide inactivators. Thymidylate synthetase is inactivated by a derivative of 5-fluorouracil, and synthesis of nucleotides required for DNA replication is inhibited, which is beneficial in treatment of certain cancers. Aspirin (acetylsalicylic acid) acetylates a serine residue in the active site of prostaglandin H synthetase, and inhibits the cyclooxygenase activity that initiates the conversion of arachidonic acid to physiologically active eicosanoids. Penicillin reacts with an enzyme that is required to make cross-links in bacterial cell walls and thus causes cell lysis in actively growing bacteria. These inactivators are also used to locate the binding site on the enzyme. See EICOSANOIDS. Bryce V. Plapp

Bibliography. P. F. Cook and W. W. Cleland, *Enzyme Kinetics and Mechanism*, Garland Science, 2007; S. J. Teague, Implications of protein flexibility for drug discovery, *Nat. Rev.*, 2:527–541, 2003; T. H. Venkataramaiah and B. V. Plapp, Formamides mimic aldehydes and inhibit liver alcohol dehydrogenases and ethanol metabolism, *J. Biol. Chem.*, 278:36699–36706, 2003; D. Voet and J. G. Voet, *Biochemistry*, Wiley, 2004.

Eoacanthocephala

An order (or a class according to some taxonomists) of the Acanthocephala characterized by the presence of a small number of giant subcuticular nuclei which are similar to the embryonic nuclei. Body spines may or may not be present and the chief lacunar vessels are dorsal and lateral. The cement gland of the male is a single syncytial organ with a specialized cement reservoir. Ligament sacs persist in the female. Proboscis hooks are few in number and arranged in circles. Eggs are ellipsoidal and thin shelled. These worms are parasitic in cold-blooded vertebrates (turtles, fish). The cystacanth occurs in crustaceans.

Neoechinorhynchus emydis is a parasitic species in North American turtles. The body is white and typically curved in extended specimens. The females are 0.4–0.9 in. (10–22 mm) long, whereas the males are 0.3–0.6 in. (6–14 mm). The body is 0.2–0.4 in. (0.5–1 mm) wide. There are typically six giant subcuticular nuclei of which five are middorsal and one midventral. The male reproductive organs are in the posterior half of the body, the testes are elongate, contiguous in tandem, and the large syncytial cement gland contains eight nuclei. The proboscis is globular with three circles of six hooks, each in quincuncial arrangement. The lemnisci are unequal in length; the longer contains two nuclei, the shorter only one nucleus. Eggs are oval with three membranes; the middle membrane does not completely encircle the acanthor. The first intermediate host is an ostracod (*Cypria maculata*); the second intermediate host is a snail (*Campeloma rufum*), with the cystacanth encysting in the foot of the snail. The definitive host is the map turtle (*Graptemys geographica*). In this life cycle the snail appears to be an essential second intermediate host and not simply a transport host. See ACANTHOCEPHALA. Donald V. Moore

Eocene

The second oldest of the five major worldwide divisions (epochs) of the Tertiary Period, that is, the epoch extending from the end of the Paleocene Epoch to the beginning of the Oligocene Epoch; also referred to as the middle epoch of the older Tertiary (Paleogene of some authors, Nummulitic of earlier French authors). See CENOZOIC; OLIGOCENE; PALEOCENE; TERTIARY.

CENOZOIC	QUATERNARY	Holocene
		Pleistocene
	TERTIARY	Pliocene
		Miocene
		Oligocene
		Eocene
Paleocene		
MESOZOIC	CRETACEOUS	

The Eocene/Oligocene boundary has been codified at the 63-ft (19-m) level in the Massignano section of the northern Apennines (Italy) and corresponds to the extinction level of certain planktonic foraminiferal groups. This boundary has an estimated age of 33.7 million years (Ma), based on radioisotopic ages in underlying and overlying ash beds. The Paleocene/Eocene boundary has been formally defined at the 5.2-ft (1.6-m) level in the Dababiya Quarry near the village of Dababya, approximately 22 mi (35 km) south of Luxor in the Upper Nile Valley, Egypt. This level coincides with a global carbon isotope excursion associated with significant climatic warming and biotic changes and is about 1 million years older than the base of the classic Ypresian Stage, normally considered the oldest stage of the

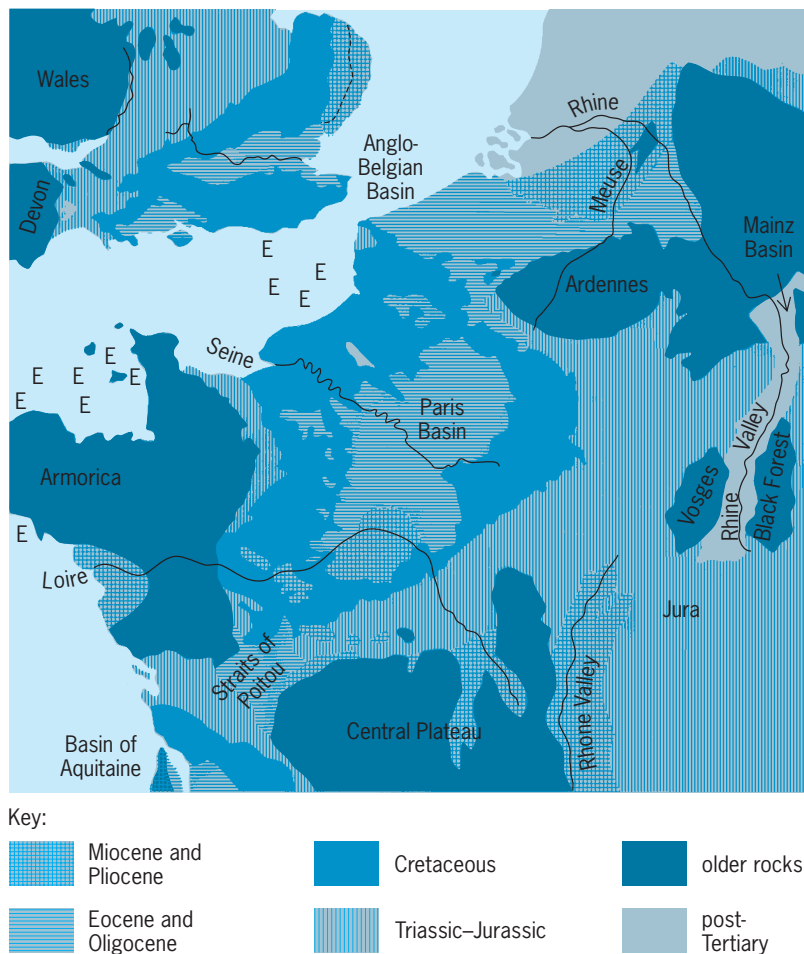


Fig. 1. The Paris and London (Anglo-Belgian) basins and their framework. E = Eocene outcrop. (After A. M. Davies, *Tertiary Faunas*, vol. 2, George Allen and Unwin, 1934)

Eocene. Currently there are two opinions regarding what to do with the associated stage boundaries. One proposes lowering the base of the Ypresian Stage to coincide with the newly relocated and defined Eocene base at the carbon isotope excursion with an estimated age of about 55 Ma. The other proposes retaining the Ypresian Stage in its present position, with an estimated age for its base of about 54 Ma based on the radiometric ages of volcanic ash beds in northwestern Europe, and inserting the Sparnacian Stage as the lowest stage of the newly redefined Eocene. Lowering the Ypresian Stage by nearly 1 million years is counter to normal chronostratigraphic procedure. See PALEOMAGNETISM; ROCK AGE DETERMINATION.

Eocene strata are widespread throughout the world and on the deep ocean floor. They include the common sedimentary types and vary from terrestrial, to marginal (estuarine), to normal marine pelagic origin. Igneous activity, while not as extensive as in the later part of the Cenozoic, was notable in some areas such as East Greenland, Oregon, Washington, and British Columbia.

In 1833, Charles Lyell provided the first basic subdivision of the Tertiary Period into three basic units: Eocene, Miocene, and Pliocene. The basis was the percentage of living mollusk species present in each

subdivision. The Eocene was said to have less than 5% of living forms in it. See MIOCENE; PLIOCENE.

The content (and extent) of the Eocene was subsequently modified with the creation of the Oligocene from the upper part of Lyell's original Eocene and the creation of the Paleocene from the lowermost part of Lyell's Eocene. Rocks of Eocene age are typically well developed in the Anglo-Belgian-Paris basins (Fig. 1). The Eocene is developed as a series of relatively shallow-water marls, limestones, and sands, with gypsum (typically developed in the Montmartre district of Paris) characterizing the uppermost part of the Eocene in the Paris Basin.

Subdivisions. The Anglo-Belgian-Paris basins have served as the type area for the establishment of a sequence of standard stage subdivisions of the Eocene. The lower Eocene corresponds to the Ypresian (Cuisian is synonymous), the Lutetian and the English Bartonian to the middle Eocene, and the Italian Priabonian to the upper Eocene. The Bartonian, previously thought to be the boreal, upper Eocene equivalent of the Mediterranean Priabonian, has since been shown actually to lie stratigraphically between the Lutetian and Priabonian and to be a part of the middle Eocene; while the French terms Auversian, Marinesian, and Ludian are seldom used anymore. In the American Gulf and Atlantic Coastal Plain, the terms Wilcoxian, Claibornian, and Jacksonian are used, although they do not correspond precisely to their European counterparts. A local stage classification is used in California based on benthic foraminiferans: Penutian, Bulitian, Ulatisian, and Narizian.

Terrestrial subdivisions include the Wasatchian, Bridgerian, Uintan, and Duchesnean land mammal ages. New Zealand marine stages include the Waipawan (upper part), Mangaorapan, Heretaungan (lower Eocene), Porangan, Bortonian, Kaiatan (middle Eocene), and Runangan (upper Eocene).

Climate and stable isotopes. Early Paleogene temperatures, including those of high latitudes, were the warmest of the Cenozoic; peak warming occurred in the early Eocene. The Earth was in a greenhouse state, with partial pressure of carbon dioxide ($p\text{CO}_2$) levels in the early Eocene estimated to have been six times higher than present-day values. During the late Paleocene to the early Eocene, deep-sea temperatures at high southern latitudes warmed by some 7–9°F (4–5°C), from about 50–52°F (10–11°C) to about 57–61°F (14–16°C), while surface temperatures increased by some 9–11°F (5–6°C), with maximum temperatures in excess of 68°F (20°C). At low latitudes, surface-water temperatures remained relatively constant and comparable to values of the present-day ocean. Superimposed on this long-term trend was a relatively abrupt (<10,000 years) 2.5–3% drop in $\delta^{13}\text{C}$ (the difference in isotopic ratios $^{12}\text{C}/^{13}\text{C}$ between a sample and a standard) and concomitant marine productivity that has been associated, in turn, with a major turnover (extinction of almost 50%) of the deep-sea benthic (bottom-dwelling) foraminiferal fauna. This drop in $\delta^{13}\text{C}$ has been identified both in marine organisms and in mammalian

bone enamel and paleosol carbonates in terrestrial sections in the Big Horn Basin of the western interior of North America and in the type Sparnacian (that is, earliest Eocene) in the Paris Basin. *See* EXTINCTION (BIOLOGY); GEOLOGIC THERMOMETRY; PALEOSOL.

Theoretical calculations and empirical evidence (decrease in eolian grain size) suggest a decrease in effective atmospheric transport and concomitant upwelling and transport of nutrient flux to the oceanic photic zone during the early Eocene. With the advent of the middle Eocene (Lutetian; about 49 Ma), a long-term, stepwise 13–14°F (7–8°C) decline of low-latitude bottom-water temperatures set in; this change bears a strong resemblance to the fossil record and to stable isotope records of surface-water temperatures at high latitudes over the same interval, suggesting a strong linkage in the thermal evolution of these two areas. As the southern oceans cooled, a more vigorous thermohaline circulation developed and culminated in an accelerated cooling of deep waters in the late Eocene and relatively abrupt drop of 7–9°F (4–5°C) in earliest Oligocene time. This latter event was associated with the growth and expansion of global ice volume (most of which was probably present in the form of glaciers on Antarctica) estimated to have been in excess of 70% of present-day ice volume. The Earth had clearly passed into the icehouse state, in which it remains today.

Sea-level changes. Although the evidence for changes in sea level through time can be clearly read from the alternating history of transgressions and regressions in the stratigraphic record, the causal processes considered responsible for such changes are still unclear. An explanation may relate to global, or eustatic, changes of sea level; more local or regional tectonic events centered in the Earth's continental crust; or the major ocean ridges on the sea floor. A key method for distinguishing between these explanations is by precise stratigraphic correlation; synchrony between unconformities bounding sedimentary sequences in geographically and positionally separated basins would strongly suggest that global, or eustatic, control of depositional sequences is responsible, whereas nonsynchronous or overlapping unconformities would suggest local or regional tectonic control. During the Neogene, glacioeustasy (globally synchronous and essentially quantitatively equivalent changes in sea level) is generally accepted as the driving force controlling sea-level oscillations; during ice-free periods of Earth's history, and during the early Paleogene (Paleocene-Eocene), this explanation does not account for observations in the rock stratigraphic record. Recent evidence suggests that glacioeustasy may have controlled sea-level fluctuations of 66–82 ft (20–25 m) during the Late Cretaceous. *See* CRETACEOUS; MID-OCEANIC RIDGE.

The history of sea level during the Cenozoic Era is one of gradual and inexorable withdrawal of the sea from the continents, punctuated by significant lowering of the sea level in the late Paleocene and the earliest Oligocene (Fig. 2). The latter event is generally considered to be linked with the first major buildup of ice on the Antarctic continent. Major sea-level

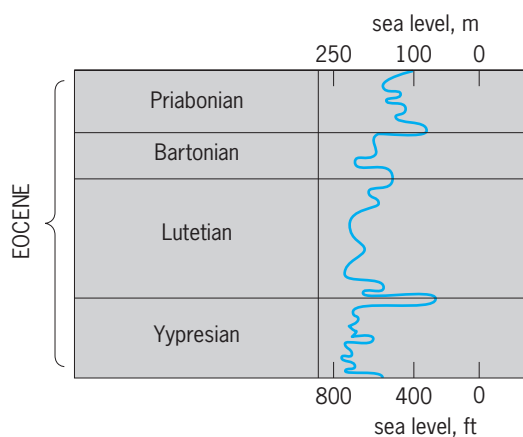


Fig. 2. Eustatic sea-level curve for the Eocene. (After B. U. Haq, J. Hardenbol, and P. R. Vail, *The chronology of fluctuating sea level since the Triassic*, *Science*, 205:1156–1167, 1987)

risks occurred during the early Eocene (Ypresian Age; around 55–54 Ma) and the early middle Eocene (Lutetian Age; around 49–47 Ma), whereas major falls in sea level took place at the early/middle Eocene (Ypresian/Lutetian) boundary (around 49 Ma) and at the middle/late Eocene (Bartonian/Priabonian) boundary (around 40 Ma).

A series of apparently correlative unconformities separating the synchronous alternation of shallow marine and brackish-water facies in the upper Paleocene–lower Eocene of the Paris-London-Hampshire and Belgian basins suggests a predominantly eustatic control on deposition in this passive margin setting. Comparable stratigraphic sequences have been delineated on the eastern and Gulf Coastal Plain of North America. *See* CONTINENTAL MARGIN; SEISMOLOGY; UNCONFORMITY.

Paleoslope modeling of the east coast of North America suggests that early and middle Eocene relative sea level stood 180–400 ft (55–120 m) higher than the present-day level and that the coastline would have shifted about 42 mi (70 km) during this change. Paleoslope modeling refers to the use of benthic foraminiferans to reconstruct ancient environments on the sea floor along the slope of the continental shelf and slope. Different assemblages of benthic foraminiferans occupy different parts of the slope (with a certain degree of overlap), and each depth-related association is termed a biofacies. The faunal compositions of ancient biofacies are generally similar to those of the present day, so they can be used to approximate conditions in the past. On the assumption of a similar gradient of the slope in the past compared with today, the successional changes in benthic foraminiferal biofacies are used to reconstruct the ancient sea-floor environment.

The Eocene was a period of repeated transgressions and regressions of the sea; in particular during the middle Eocene, one of the most extensive global marine transgressions took place, leaving extensive shallow-water limestones with rich and diversified molluscan faunas in the Gulf Coast and the Paris Basin; the latter faunas became a focal point of

extensive studies in the nineteenth century and played a major role in the formulation of Lyell's concept of the Eocene.

Life. Early Paleogene temperatures, including those of high latitudes, were the warmest of the Cenozoic; peak warming occurred in the early Eocene. Equatorial surface waters may have been several degrees cooler than at present, and latitudinal thermal gradients were lower; this resulted in a poleward expansion of tropical life-forms. However, dramatic changes in climatic conditions appear to have taken place at the end of the Paleocene. A significant increase in the temperature of bottom waters of the ocean (about 5°F or 3°C) may have been coupled with a decrease in wind intensity, resulting in a decrease in surface productivity and lowered oxygen levels in bottom waters which may have triggered a major turnover in the deep-water benthic foraminiferans.

The diversification of life seen in the Paleocene continued in the Eocene, a reflection of the poleward expansion of the tropics, particularly during early Eocene time. In the oceanic realm, microplanktonic animals (foraminiferans) and plants (calcareous nannoplankton) flourished and diversified, as did true bony fishes and siphonate gastropods. In shallow, tropical waters, larger foraminiferans—*Nummulites* and discocyclinids—extended their geographic range to latitude 50° north, but the latter group disappeared at the end of the Eocene owing to cooling temperatures. Indeed, microplanktonic animals and plants experienced a gradual but inexorable decline in diversity starting in the late middle Eocene. Succeeding Oligocene faunas and floras were much reduced in diversity and much more uniformly distributed. See FORAMINIFERIDA; NUMMULITES.

On land, subtropical floras extended as far north as southern England and the North American Pacific coast of Puget Sound and southern Alaska. Indeed, the floras of southern England resembled those of modern-day China, Malaysia, and Australia. In the humid interior, thick and extensive mud deposits (the Green River Shale) in Colorado contain a beautifully preserved fresh-water fish fauna eagerly sought after by fossil collectors, as well as vast quantities of untapped oil comparable to that in the oil-rich tar sands of Alberta Province, Canada.

Europe was separated from the Eurasian landmass east of the Urals by a north-south seaway extending from the Arctic to the Tethys Sea—the Turgai Straits. Following the elimination of the elevated corridor that allowed transatlantic poleward migration between Europe and North America in late early Eocene time, middle and late Eocene time witnessed the development of extensive endemic animal evolution. Bats, flying lemurs, creodont carnivores, artiodactyls (cloven-hoof mammals, such as cattle, deer, and camels) and perissodactyls (odd-toed, hooved mammals, such as rhinoceroses and horses), notoungulates (predominantly South American), and edentates reflect the diversification of primitive placental forms. The massive, rhinoceroslike herbivores called titanotheres and uintatheres appeared alongside the

small early progenitors of the modern horse, *Hyracotherium* (known popularly as *Eobippus*). See MAMMALIA; PERISSODACTYLA.

In the Eocene, some mammals turned toward life in the sea; sea cows appeared in the middle Eocene, while the earliest whales (zeuglodonts) appeared in the North Atlantic–Gulf of Mexico region and the aquatic ancestors of the proboscideans appeared in the late Eocene.

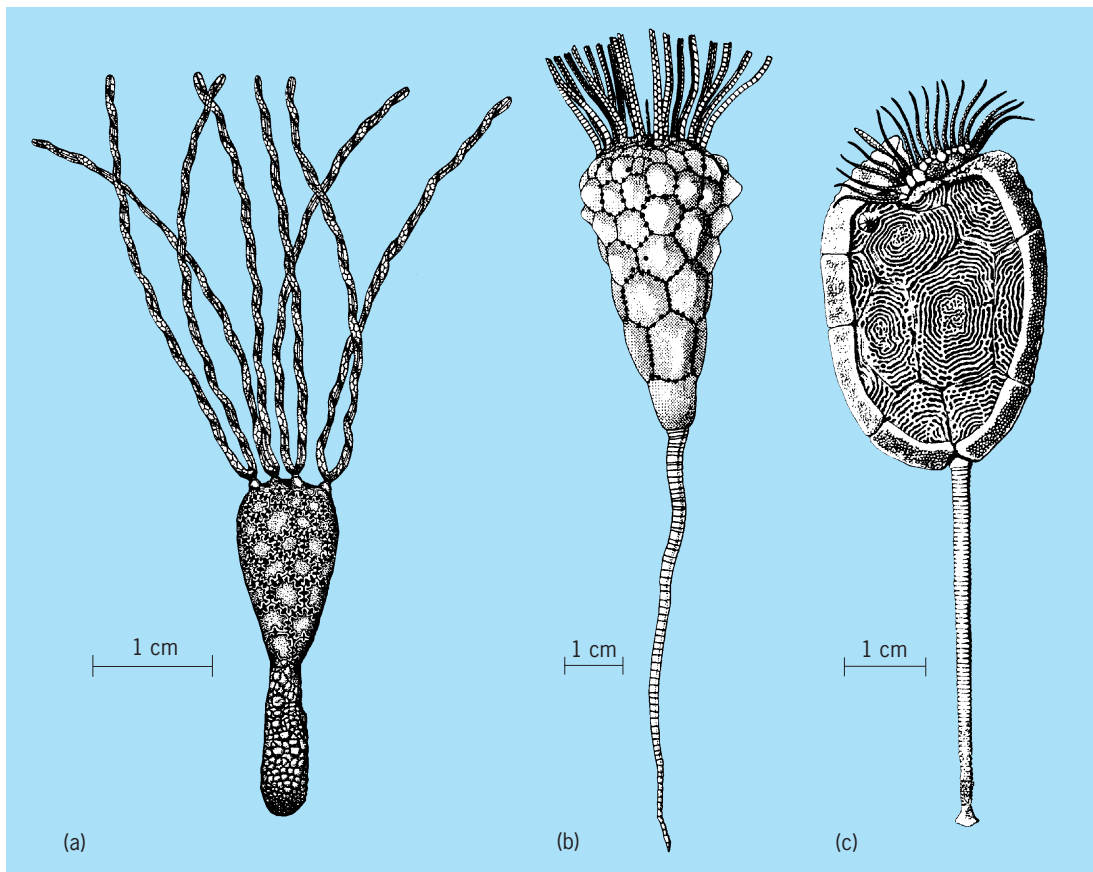
The presence of warm southern ocean currents between Australia and Antarctica, which were in the process of separating, led to increased evaporation and precipitation on the great landmass of Antarctica and to gradual cooling, which, in turn, led inexorably to the initiation of continental ice on that continent in the succeeding epoch, the Oligocene. The Panamanian Isthmus was below sea level, and North and South America were separated from each other, resulting in the separate but parallel evolution of many mammals on the two continents. Marsupials, ancestral anteaters, armadillos, ground sloths, large rodents, and other now-extinct placentals evolved separately until the two landmasses were reunited by the elevation of the Panamanian Isthmus in the Pliocene some 3 Ma. See PALEOBOTANY; PALEONTOLOGY.

W. A. Berggren

Bibliography. M.-P. Aubry et al., Chronostratigraphic terminology at the Paleocene/Eocene boundary, in S. L. Wing (ed.) et al., *Causes and Consequences of Globally Warm Climates in the Early Paleogene* (Special Papers), Geological Society of America, vol. 369, pp. 551–566, 2003; M.-P. Aubry, S. G. Lucas, and W. A. Berggren (eds.), *Late Paleocene–Early Eocene Biotic and Climatic Events in the Marine and Terrestrial Records*, Columbia University Press, 1998; K. Ouda and M.-P. Aubry (eds.), *The Upper Paleocene–Lower Eocene of the Upper Nile Valley: Stratigraphy*, Micropaleontology Special Issue, vol. 49(1), 2003; D. R. Prothero, *The Eocene–Oligocene Transition*, Columbia University Press, 1993; D. R. Prothero, L. C. Ivany, and E. Nesbitt (eds.), *From Greenhouse to Icehouse*, Columbia University Press, 2002.

Eocrinoidea

A medium-sized class of primitive, brachiole-bearing, blastozoan echinoderms of the class Crinozoa that ranged from the Early Cambrian to the Middle Silurian, although few eocrinoids survived past the Middle Ordovician. About 32 eocrinoid genera have been described from North America, Europe, North Africa, and Australia; other occurrences of distinctive plates that may belong to eocrinoids have also been noted. Eocrinoids are the most diverse class of echinoderms known from the Cambrian with about 15 genera, and different members appear to have been ancestral to nearly all of the more advanced brachiole-bearing echinoderm classes that appeared in the Early or Middle Ordovician, such as rhombiferans, parablastoids, and coronoids. Eocrinoids have a globular, conical, or flattened theca or body, with



Reconstructions of Cambrian and Ordovician eocrinoids. (a) *Gogia spiralis*, from the Middle Cambrian of Utah, showing the primitive holdfast; sutural pores between the numerous thecal plates; and long, spiraled brachioles attached to the summit (after R. A. Robison, *Middle Cambrian eocrinoids from Western North America*, *J. Paleont.*, 39:360–363, 1965). (b) *Rhopalocystis destombesi*, from the Early Ordovician of Morocco; note the more advanced columnal-bearing stem, a theca that still retains sutural pores, and incomplete brachioles on summit (after G. Ubaghs, *Eocrinoidea*, in R. C. Moore, ed., *Treatise on Invertebrate Paleontology*, pt. S, 1968). (c) *Mandalacystis dockeri*, from the Middle Ordovician of Oklahoma, showing the stem with its distal attachment, flattened theca with more massive plates around the margins, and short brachioles attached to the summit (after R. D. Lewis et al., *Mandalacystis, a new rhipidocystid eocrinoid from the Whiterockian Stage (Ordovician) in Oklahoma and Nevada*, *J. Paleontol.*, 61:1222–1235, 1987).

many irregularly arranged to partly organized, imbricate or adjacent plates (see **illus.**). Most Cambrian genera have sutural pores on the plate margins, apparently for respiration. Many early eocrinoids have a multiplated, cylindrical to slightly inflated holdfast for attaching the theca to objects lying on the sea floor. Holdfasts apparently evolved into a true columnal-bearing stem in late Middle Cambrian eocrinoids, and most later genera have that advanced type of attachment structure. Eocrinoids typically have two to five short ambulacral grooves radiating from the mouth on the summit to many long, erect, biserial brachioles that were used for feeding. Most eocrinoids were attached, low- to medium-level suspension feeders that used the brachioles to collect small food particles drifting by the theca. Most researchers have argued that eocrinoids are a valid class containing genera that did not develop the fold-like respiratory structures found in more advanced blastozoan classes. However, others have recently proposed that eocrinoids are a paraphyletic stem group that should be discarded with the included genera reassigned to other classes. See ECHINODERMATA.

James Sprinkle

Bibliography. R. C. Moore (ed.), *Treatise on Invertebrate Paleontology*, pt. S, 1968; C. R. C. Paul and A. B. Smith, The early radiation and phylogeny of echinoderms, *Biol. Rev.*, 59:443–481, 1984; J. Sprinkle, *Morphology and Evolution of Blastozoan Echinoderms*, Mus. Compar. Zool., Harvard Univ., Spec. Publ., 1973.

Eolian landforms

Topographic features formed mainly by wind in regions with much loose sand, silt, or clay. An estimated 500 million metric tons (556 million tons) of dust per year are transported by the wind. Airflow in wind can be either laminar or turbulent, but is mostly turbulent. Wind blowing across the land surface contains many swirls and eddies, some of which produce upward components of movement. Generally, the higher the mean wind velocity, the greater the amount of upward turbulence. The ratio of the velocity of upward turbulent gusts near the ground to mean wind velocity varies, but it averages about 1.5. Thus, particles whose settling velocities are lower

than one-fifth of the mean wind velocity are likely to be conveyed upward by wind gusts and transported downwind. Larger particles remain mostly close to the ground until unusually high gusts lift them up. Small eddies occur under most atmospheric conditions, unless the air is unusually thermally stable. As wind velocity rises, corkscrew eddies develop with axes parallel to the mean wind direction; and larger, horizontal eddies, measured in hundreds of meters, are superimposed on the smaller, corkscrew eddies. *See* CLAY; DUST STORM; LAMINAR FLOW; SAND; VELOCITY; WIND.

Movement of particles by wind. Three modes of particle movement by the wind are suspension, saltation, and surface creep. Turbulent eddies are capable of transporting fine silt and clay particles upward and may suspend particles in the air. However, if the terminal velocity of sand grains is greater than the vertical component of the turbulent velocity, the particle trajectory path is a smooth curve with a trajectory height measured in centimeters, rather than in kilometers as is typical for fine particles in suspension. Sand-size particles traveling in low, smooth trajectories and bouncing along on the surface move by saltation. Particles that are too large to be lifted from the surface by the wind may be rolled or pushed along the surface by impact creep, produced by the impact of saltating particles. Of these, saltation is the most significant mode of transportation because fine sand grains are most easily moved by the wind. Smaller grains are more difficult to move because of particle cohesion and aerodynamic effects, and larger particles are more difficult to move because of their greater mass.

When a saltating grain returns to the surface, it may crash into a group of sand grains, driving them into the air. Saltation impact is a significant process in putting dust and other small grains into the air by winds that would otherwise not be strong enough to initiate grain movement. Saltating grains that strike surfaces composed of larger grains, too large to be moved by wind alone, push them downwind a short distance.

Sand grains saltating along a hard, rocky surface bound high into the air upon impact and are thus more easily kept in motion than particles hitting loose sand, where momentum is dissipated among the loose grains, often trapping the impacting grain. As golfers are aware, golf balls hit into a sand trap do not bounce nearly as much as balls landing on a green. Accumulation of sand grains in one place thus leads to the trapping of still more sand, rather than spreading out evenly over the ground surface. *See* MOMENTUM.

Coarse sand grains on the ground surface that are too large to be directly set in motion by the wind may be nudged forward by the impact of bombarding sand grains. The impact of a high-velocity, saltating grain can move a grain six times its diameter and more than 200 times its own weight. This type of downwind sand movement is known as impact creep or surface creep. Creep may account for about one-fourth of the total sand movement.

Wind is considerably limited in the size of particles it can transport. Not only does settling velocity determine the particle size of suspended grains, but also it is important in limiting the height to which saltating grains rise above the ground. In addition, the nature of the surface from which saltating grains rebound is important. The rebound height above a sand-mantled surface may be only 0.5 m (1.6 ft) or less, whereas above a rock surface it can be several meters. Abrasion by windblown sand rarely exceeds about 1 m (3.2 ft) above the ground. Wooden and steel telegraph posts in the Sahara severed about 1 m above the ground show the effective height of saltating grains there.

Wind is an unusually selective sorting agent. Eolian sediments are typically among the best-sorted deposits in nature. The most common grain sizes found by sieve analyses of dune sands are 0.30–0.15 mm in diameter. Diameters less than 0.08 mm are rare because most grains of this size are carried off in suspension.

Eolian abrasion and erosion. Wind erosion takes place by abrasion, the wearing away of material by frictional impact, and by deflation, the blowing away of loose material. Wind abrasion is caused by the collision of sand grains with surfaces. Spectacular examples occur when automobiles caught in sandstorms are stripped of paint and their windshields frosted in a matter of minutes. In nature, saltating sand grains, bouncing along the ground near the surface, collide with rocks or other sand grains that they abrade. Windblown sand blasts exposed rock faces and creates small-scale abrasion features whose details depend on the lithology, inclination, and orientation of the rock; the velocity and direction of the wind; and the size, hardness, and shape of the abrasive material. The results of wind abrasion include the polishing, pitting, grooving, fluting, and faceting of rocks and the rounding of sand grains in transport. *See* EROSION.

Among the most common features of wind abrasion are ventifacts, etched and pitted faces of exposed rock. If a rock face is normal to the prevailing wind direction, a faceted face is abraded at right angles to the wind, with sharp edges between the facet and the original surface. Rocks that periodically rotate so as to expose different sides to the direction of prevailing wind develop multiple-faceted sides. The abraded facets are typically pitted, fluted, grooved, and polished.

Larger-scale erosional effects of wind include yardangs, elongate ridges aligned parallel to the direction of prevailing wind with rounded windward faces that taper in the downwind direction. Yardangs vary in length from several meters to 1 km (0.6 mi), and from a few meters to 200 m (656 ft) high. Bedding and other structures in the material are truncated by the troughs, displaying their erosional origin. That they are wind-eroded is shown by their alignment parallel to the direction of prevailing wind and by their composition of easily eroded material. The elongate ridges are separated by intervening troughs, whose origin may also include erosion

by running water and weathering. *See* WEATHERING PROCESSES.

Where strong winds blow across surfaces lacking coarse material, removal of small debris by the wind forms deflation hollows. The size and shape of deflation hollows vary considerably, depending on the nature of the underlying material, strength of the wind, and local weathering conditions.

Sand dunes. A patch of sand makes an effective trap, which continues to grow in size until it becomes a mound of sand. Protrusion of the sand mound high enough into the air to disrupt laminar airflow over its surface significantly affects further accumulation and movement of sand, and a dune is formed. Once a sand grain becomes part of a dune, it advances downwind only at the rate of advance of the entire dune, rather than at its previous saltation rate. *See* DUNE.

When a critical height is reached, a slip face forms on the lee side of the dune. Sand moves up the windward side of the dune by saltation and is blown over the brink of the crest, increasing the slope until it exceeds the angle of repose for loose sand. An avalanche forms down the lee side, creating a slip face standing at the angle of repose of loose sand, about $30\text{--}34^\circ$. Successive slip faces become progressively buried until they are once again unearthed by erosion on the windward side of the dune as the dune gradually advances downwind. The most important factors are dune height and rate of sand flow. Measured rates of dune migration are 15–20 m/y (50–65 ft/y).

Barchans are crescent-shaped dunes with a steep slip face concave in the downwind direction. The average width of barchan dunes, measured from horn to horn, is about 37 m (121 ft), and the heights of the dunes are consistently about one-tenth the width. The horns of the crescent taper in the downwind direction because, even with constant lateral supply of sand, saltating grains will move more quickly over the hard ground adjacent to the dune patch and most slowly where the sand is thickest. Because the edges of the dune are lower and the rate of advance of the slip face is inversely proportional to the dune height, the edges move downwind at a faster rate than the thicker, central part of the dune, giving the dune its crescent shape. Barchans are typically best developed on barren desert floors where the supply of sand is meager, the prevailing wind is relatively uniform, and vegetation is scarce. They may occur as single dunes or in groups where individual barchans coalesce to form complex shapes.

Transverse dunes are elongate dunes that form perpendicular to the direction of prevailing winds. They originate in several ways. Some are clearly transitional forms with barchan dunes and consist of several barchans connected along single lines.

Parabolic dunes are crescent-shaped dunes in which the slip face of the dune is convex downwind and the horns of the crescent point upwind. Some are clearly transitional to barchan dunes and form as a result of anchoring of the horns by vegetation while the central portion continues to migrate downwind.

Vegetation is an important factor in eolian processes because it is effective in breaking up airflow near the ground and serves to anchor thin sand deposits. The result is that the horns of a parabolic dune, anchored by vegetation, trail the thicker central part of the dune. In some cases, the horns of the parabolic dune may be left far behind, giving the dune a hairpinlike shape.

Longitudinal dunes are elongate ridges of sand parallel to the direction of prevailing winds. One of the places that they may form is where wind funnels sand through a gap in a ridge. As the wind blows through the gap, its velocity increases due to the Bernoulli effect; but when the wind passes through the gap, its velocity decreases to that of the general flow, and it deposits some of the sand carried by the enhanced velocity through the gap. *See* BERNOULLI'S THEOREM.

Seif dunes are longitudinal dunes that occur in long, straight, parallel chains over exceptionally long distances, sometimes 60–190 km (37–118 mi). Perhaps the most intriguing aspect of these unusual dunes is not only their length but also the regularity of their spacing. Most longitudinal dunes are separated by barren rock floors or desert pavements of rather uniform width. No definitive explanation has been demonstrated for these extraordinary dunes.

Star dunes typically consist of three or more sharp-edged ridges extending radially from a high, pointed, central peak. They seem to originate in areas where winds blow from more than one prevailing direction.

Loess. Loess consists of usually well sorted, wind-blown silt and clay deposited as a relatively homogeneous, unstratified blanket over the Earth's surface. Typical loess contains about 40–50% silt, up to 30% clay, and up to 10% fine sand. *See* LOESS.

An estimated 30% of the United States is mantled with loess, and as much as 10% of the Earth's land area is covered by loess. The sources of most loess deposits seem to be related to areas affected by Pleistocene glaciations and climatic changes, especially downwind from glacial outwash plains. Thick loess deposits occur in the Mississippi Valley region of the Central Plains; in Alaska; in the Columbia Plateau of eastern Washington; in portions of Europe south of the Scandinavian Ice Sheet; and in China, Mongolia, and the former Soviet Union. The largest known loess-covered region is 635,000 km² (245,110 mi²) in north-central China. *See* CLIMATE HISTORY; GLACIAL EPOCH; PLEISTOCENE.

The loess deposits of China and the Central Plains of the United States are especially thick, reaching 300 m (984 ft). Elsewhere, loess thicknesses generally are several tens of meters. During the dust-bowl storms of the Central Plains in the 1930s, immense volumes of silt became airborne. In 1935 a dust storm near Wichita, Kansas, was estimated to hold about 5 million metric tons (5.6 million tons) of dust in suspension over a 78-km² (30-mi²) area, and 300 metric tons/km² (333 tons/mi²) was deposited in a single day near Lincoln, Nebraska. Deposition of loess in China presently occurs at a rate of several millimeters per year.

Don J. Easterbrook

Bibliography. R. Cooke, A. Warre, and A. S. Goudie, *Desert Geomorphology*, 1993; A. S. Goodie, S. Stokes, and I. Livingstone, *Aeolian Environments, Sediments and Landforms*, 2000; V. P. Tchakerian (ed.), *Desert Aeolian Processes*; 1995.

Ephedra

A genus of low, leafless, green-stemmed shrubs belonging to the plant class Ginkgoopsida of the division Pinophyta (see **illus.**). They grow in dry,



A species of *Ephedra*, growing in Utah. (Courtesy of Tony Gauba, from National Audubon Society)

alkaline soils around the world. In the southwestern United States these plants are called Mormon tea and jointfir. The drug ephedrine is extracted from the Asiatic species *Ephedra sinica* and *E. equisetina*. It is used medicinally in the treatment of colds, hay fever, and asthma. See EPHEDRALES; GINKGOOPSIDA; PINOPHYTA; PLANT KINGDOM.

Perry D. Strausbaugh; Earl L. Core

Ephedrales

An order of the class Ginkgoopsida having about 35 species in the genus *Ephedra*. These plants are mostly of arid regions. They include freely branched, low shrubs with reduced, scalelike leaves and green photosynthetic twigs. The species are dioecious (seldom gynodioecious). The genus is known from Asia to Europe, in northern Africa, in the United States and Mexico, and from Bolivia to Patagonia. It is not known in the fossil record. See GINKGOOPSIDA; PINOPHYTA; PLANT KINGDOM.

Thomas A. Zanoni

Ephemeris

A source of data pertaining to features of various objects, usually astronomical, that depend on time. The ephemeris can be a printed or computer-readable table such as on a compact disk, or it can be a set of numerical data accessible by computer programs. Most ephemerides pertain to positions of celestial objects as seen from any place in the solar system. The term may also apply to other parameters, such as those describing the state of rotation or orientation of a body in space.

A basic ephemeris is usually derived from the known laws of physics that, in turn, are described by equations of motion. Before computers, these equations were solved and evaluated using elaborate analytical formulas requiring extensive manual effort. With computers, the process is done numerically. Initial ephemerides are compared with positional measurements (observations), the initial parameters are adjusted to produce better agreement (between ephemerides and observations), and the whole process is then iterated. The resultant ephemerides, stored in the computer, are accessible to produce tabular listings or to give other types of positional information at various requested times. See NUMERICAL ANALYSIS.

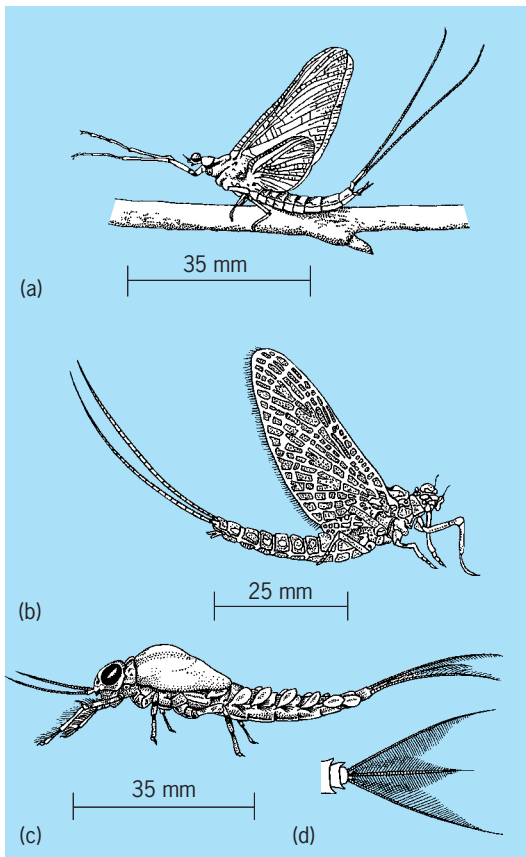
Traditional printed ephemerides are available in many national astronomical almanacs. On the Internet, there are interactive Web sites that enable customizing one's own ephemeris. Also available are computer subroutines for approximate accuracy that users can employ in their self-written programs. Finally, there are precision ephemerides for full accuracy, available in the form of computer data files and reading routines, to be used as a part of a user-written main computer program. Computer links to all of these sources are available on the Web site of the International Astronomical Union's Commission 4 (Ephemerides). See ALMANAC; CELESTIAL MECHANICS; INTERNET.

E. Myles Standish

Bibliography. G. Beutler, *Methods of Celestial Mechanics*, Springer, 2005; J. M. A. Danby, *Fundamentals of Celestial Mechanics*, Willmann-Bell, 1988; *Explanatory Supplement to the Astronomical Ephemeris and the American Ephemeris and Nautical Almanac*, H. M. Stationery Office, London, 1961; P. K. Seidelmann (ed.), *Explanatory Supplement to the Astronomical Almanac*, University Science Books, Mill Valley, CA, 1992.

Ephemeroptera

An order of insects commonly known as mayflies. They are aquatic and live in clean, fresh waters during their immature, or nymphal, lives. Nymphs are adapted to aquatic environments ranging from ponds to mountain streams. Most mayfly nymphs are vegetarians, and most species are present throughout the year as nymphs. The adult stage is very brief, lasting from a few hours to several days (see **illus.**). As



Ephemeroptera. (a) Mayfly, adult male, *Hexagenia* sp., in usual resting position (after A. H. Morgan, *Field Book of Ponds and Streams*, G. Putnam's Sons, 1930). (b) Mayfly, subadult male, *Callibaetis* sp., newly emerged from water, standing on surface film (after A. H. Morgan, *Kinships of Animals and Man*, McGraw-Hill, 1955). (c) Mayfly nymph, *Isonychia* sp., from rapid stream; (d) segmented tailpieces (after A. H. Morgan, *Field Book of Ponds and Streams*, G. Putnam's Sons, 1930).

immature insects, they constitute a year-round basic food supply for carnivores of their communities, especially for fishes, and trout in particular. Nymphs and adults are models for artificial bait flies.

The short-lived adult female mayflies drop their eggs in packets into the water or scatter them over its surface. Certain small species fix the eggs to surfaces of rocks partly submerged in rapid currents. The egg stage may last from a few days to several months and, after hatching, the nymph may mature rapidly (in a few weeks) or may grow slowly (requiring as much as 2 years to complete its development), according to the species. At the end of their growth period, most nymphs rise to the water surface, shed their skins, and become subadults, or subimagoes. Some nymphs may crawl from the water onto vegetation, rocks, or sticks prior to emergence. In a short time, varying from minutes to a day or two, they shed the subadult skin and become adults, or imagoes.

Mayfly nymphs are distinguished from all other aquatic insects by the paired tracheal gills on the back of each of the first seven abdominal segments. The body ends in two or three finely segmented tailpieces. The mouthparts are highly interesting

examples of adaptations for scraping, cutting, and crushing the various plant cells, which are the chief food.

As subadults, mayflies have the general adult form and useless mouthparts. They are clothed by a thin, furry, grayish skin which accounts for the common name duns. The location of the subadult stage, however, between the nymphal and adult stages, suggests the sequence of stages in the metamorphosis of hemimetabolous insects such as larva, pupa, and adult of the butterfly. The subimago stage ends with the final molt of the mayfly's life. No other insects molt in their winged form.

The adult emerges with transparent wings, shining body, and useless mouthparts; the males generally have extremely large eyes, which may be divided, and long front legs. In a few mayflies, the legs of the female, and those of the male, except for the front ones, are nonfunctional and the brief adult life is spent entirely in flight. The quivering of the tailpieces during the mating flight gave them the name spinners. Thousands join in the twilight mating swarms, rhythmic dancing flights unequalled by other insects. Sometimes, the adults, which may be attracted by lights, accumulate in such masses that they become a hazard to vehicular traffic on highways and bridges over large rivers. See INSECTA.

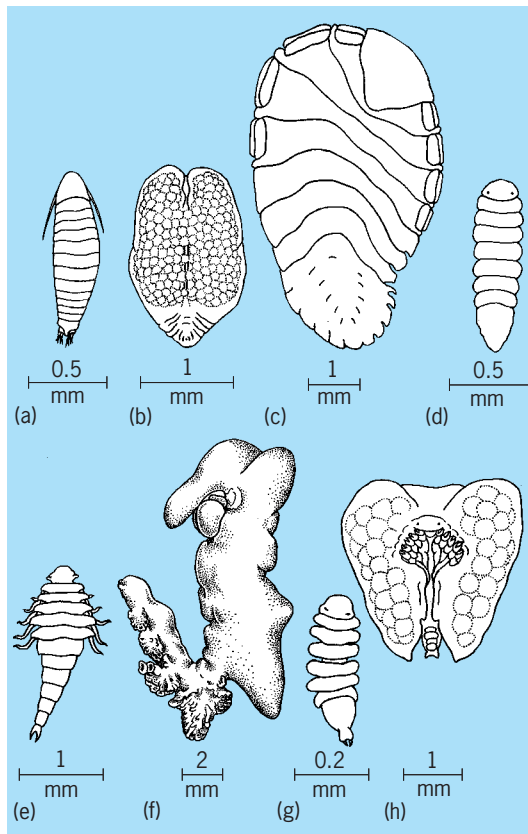
Anne H. Morgan; Lewis Berner

Bibliography. L. Berner, *The Mayflies of Florida*, 1950; B. D. Burks, *The Mayflies, or Ephemeroptera, of Illinois*, 1953; J. W. Fannie and A. Leonard, *Mayflies of Michigan Trout Streams*, 1962; A. H. Morgan, *Field Book of Ponds and Streams*, 1930; J. G. Needham, J. R. Traver, and Y. C. Hsu, *Biology of Mayflies*, 1935; S. P. Parker (ed.), *Synopsis and Classification of Living Organisms*, 2 vols., 1982.

Epicaridea

A suborder of the Isopoda which are parasitic on various crustaceans, mainly marine forms. The females are sometimes modified so strongly as to leave almost no indication of their isopod nature. The female pierces the host's skin with the aid of styliform mandibles to suck blood. It seizes the host's skin by means of its prehensile pereopoda. The pleopoda are respiratory in function. The dwarf male attaches to the female, from which it takes nourishment, and retains its isopod structure. The first larva, the epicardium, is free-living, while the second, the microniscium, is temporarily ectoparasitic on Copepoda. The third larval stage, the cryptoniscium, is free-swimming and seeks a final host, on which it undergoes final development to the adult. The suborder is divided into two tribes, the Cryptoniscina and Bopyrina.

Cryptoniscina. These crustaceans live on Entomostraca and are protandrous hermaphrodites. The cryptoniscium larva attaches to a suitable female, develops testes without any change in its external appearance, and acts as a male. After the death of the host, it metamorphoses into a female, ultimately



Epicaridea. (a) *Cypronicus ovalis* male and (b) female. (c) *Bopyrus squillarum* female and (d) male. (e) *Portunium flavidus* male and (f) female. (g) *Prodajus bilobatus* male and (h) female.

undergoing anatomical degradation that leaves it a mere bulky sac distended with developing embryos. The brood chamber is internal and arises from paired invaginations of the sternum. *Hemioniscus* lives in the mantle cavity of Cirripedia, *Cypronicus* (illus. a and b) lives in the shell of Ostracoda, and *Liriopsis* infests Rhizocephala, which is in turn a parasite of crabs. *Asconiscus*, *Podascon*, and *Cabirops* live in the brood pouch of Mysidacea, Amphipoda, and Isopoda, respectively.

Bopyrina. These isopods are dioecious, but it has been shown experimentally in some forms that the cryptoniscium is ambipotent, and a presumptive young male removed from the female develops into either a female or intersex. The tribe includes three families, rich in species. Bopyridae parasitize Decapoda. *Bopyrus* and its allies (illus. c and d) live in the branchial cavity of prawns, crabs, and hermit crabs; *Pbryxus* lives under the abdomen of prawns; and *Athelges* lives above the abdomen of hermit crabs. The female shows an asymmetrical development of the segments, which are often partially fused, and it shows frequent disappearance of some of the pleopoda or one of their rami. Similar degeneration also occurs in the male, which is nevertheless symmetrical. The brood pouch is external, as in other Isopoda. Branchial parasites cause a swelling of the branchial cavity of the host. In some hosts the parasites produce an atrophy of the gonad

and modification of secondary sexual characters. Entoniscidae, such as *Portunium* (illus. e and f), live in the visceral cavity of crabs and porcellanids. The female is vermiform, lacks metamerism, and is surrounded by a thin membrane derived from the host. Pereiopoda are atrophied, whereas pleopoda are excessively developed. Dajidae, for example, *Prodajus* (illus. g and h), live on Schizopoda, clinging to various parts of the host. The female is symmetrical, scarcely showing any metamerism. The fourth, fifth, and sixth pereiopods are reduced and pleopoda are abortive. See ISOPODA.

Sueo M. Shino

Bibliography. J. G. Baer, *Ecology of Animal Parasites*, 1951; M. Caullery, *Parasitism and Symbiosis*, 1952; S. P. Parker (ed.), *Synopsis and Classification of Living Organisms*, 2 vols., 1982.

Epicycloid

A curve traced by a point on a circle that rolls on the convex side of a fixed circle. The term is also occasionally applied to the curve generated by a point on the prolongation of the radius of a circle as the circle rolls on a straight line. If the rolling circle has radius r and the fixed circle radius R , as in the illustration,

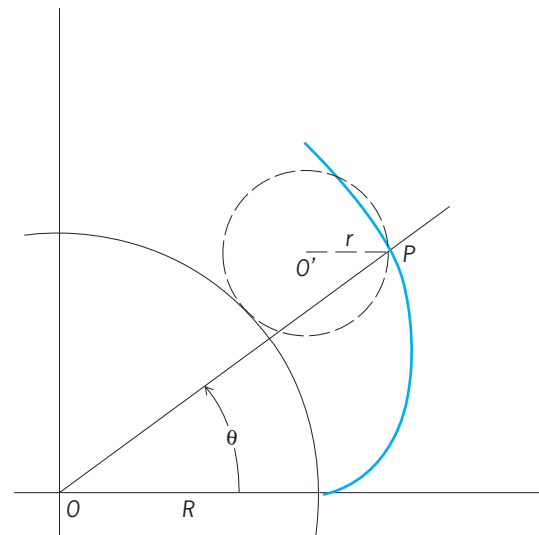


Diagram of an epicycloid.

parametric equations of the epicycloid are as shown below. The curve is closed provided R/r is a ratio-

$$x = (R + r) \cos \theta - r \cos (1 + R/r)\theta$$

$$y = (R + r) \sin \theta - r \sin (1 + R/r)\theta$$

nal number. It is a cardioid when $r = R$. The curve was known to the Greek astronomer Hipparchus about 140 B.C., and was investigated by Gerard Desargues in 1639 and Leonhard Euler in 1781. Desargues made the first applications of epicycloids to the design of gear teeth. See CARDIOID.

Leonard M. Blumenthal

Epidemic

The occurrence of cases of disease in excess of what is usually expected for a given period of time. Epidemics are commonly thought to involve outbreaks of acute infectious disease, such as measles, polio, or streptococcal sore throat. More recently, other types of health-related events such as homicide, drownings, and even hysteria have been considered to occur as "epidemics."

Confusion sometimes arises because of overlap between the terms epidemic, outbreak, and cluster. Although they are closely related, epidemic may be used to suggest problems that are geographically widespread, while outbreak and cluster are reserved for problems that involve smaller numbers of people or are more sharply defined in terms of the area of occurrence. For example, an epidemic of influenza could involve an entire state or region, whereas an outbreak of gastroenteritis might be restricted to a nursing home, school, or day-care center. The term cluster may be used to refer to noncommunicable disease states.

In contrast to epidemics, endemic problems are distinguished by their consistently high levels over a long period of time. Lung cancer in males has been endemic in the United States, whereas the surge of lung cancer cases in women in the United States represents an epidemic problem that has resulted from increase in cigarette smoking among women in general. A pandemic is closely related to an epidemic, but it is a problem that has spread over a considerably larger geographic area; influenza pandemics are often global.

Origins of epidemics. Disease and epidemics occur as a result of the interaction of three factors, agent, host, and environment. Agents cause the disease, hosts are susceptible to it, and environmental conditions permit host exposure to the agent. An understanding of the interaction between agent, host, and environment is crucial for the selection of the best approach to prevent or control the continuing spread of an epidemic.

The importance of the agent-host-environment relationship is exemplified by the problem of epidemic influenza, which can result in high attack rates, or illness levels, and substantial mortality among certain population groups. The agent for epidemic influenza is the highly infectious influenza A or B virus. The host is someone who has never before been infected by the specific strain of influenza virus and lacks protective antibodies to it. If the environment is a closed one, such as a nursing home, the efficiency of virus transmission from a source, most likely another person with a case of influenza, to the susceptible host is heightened. *See* INFLUENZA.

For infectious diseases, epidemics can occur when large numbers of susceptible persons are exposed to infectious agents in settings or under circumstances that permit the spread of the agent. Spread of an infectious disease depends primarily on the chain of transmission of an agent: a source of the agent, a route of exit from the host, a suitable mode of

transmission between the susceptible host and the source, and a route of entry into another susceptible host. Modes of spread may involve direct physical contact, such as touching or sexual intercourse, between the infected host and the new host, or airborne spread, such as coughing or sneezing. Indirect transmission takes place through vehicles such as contaminated water, food, or intravenous fluids; inanimate objects such as bedding, clothes, or surgical instruments; or a biological vector such as a mosquito or flea.

Responses to epidemics. The detection of epidemics in the United States is dependent on public health surveillance systems that monitor the incidence of different types of diseases. Surveillance for infectious disease epidemics or outbreaks relies heavily on reports of disease cases by physicians, hospitals, day-care centers, and other sources for health information. In the United States, state and local public health agencies bear primary responsibility for investigating disease outbreaks, implementing control measures to halt their spread, and conducting long-term follow-up monitoring to detect recurrences. *See* PUBLIC HEALTH.

When potential outbreaks are detected, a public health agency may initiate an epidemic investigation. The investigation is a systematic procedure that combines aspects from the disciplines of epidemiology, clinical medicine, laboratory science, and communications, as well as common sense. The steps involved in an epidemic investigation help to determine who became ill, where the individuals were affected, and when they were affected. By examining those person, place, and time features, the investigator may be able to explain why the epidemic occurred and to identify appropriate control measures. However, because of the urgency associated with many epidemics, public health officials may need to take steps to control the problem while the investigation continues.

The primary purpose of an epidemic investigation is identification of the source or mode of disease transmission in order to recommend control measures. However, investigations of epidemics are conducted for other reasons. They include evaluation of intervention effectiveness, epidemiologic research in situ, and on-site training of public health workers and others engaged in epidemic control. In addition, such studies help fulfill the requirements of specific disease-control programs and provide more-immediate information in response to political or public concerns. *See* EPIDEMIOLOGY; INFECTIOUS DISEASE.

Bibliography. D. W. Fraser et al., Legionnaires' disease: Description of an epidemic of pneumonia, *N. Engl. J. Med.*, 297:1189-1197, 1977; W. W. Holland, R. Detels, and G. Knox (eds.), *Oxford Textbook of Public Health*, 3d ed., 1997; J. L. Kelsey et al., *Methods in Observational Epidemiology*, 2d ed., 1996; J. J. Sacks et al., A nurse-associated epidemic of cardiac arrests in an intensive care unit, *J. Amer. Med. Ass.*, 259:689-695, 1988; P. A. Schulte, R. L. Ehrenberg, and M. Singal, Investigation of

occupational cancer clusters: Theory and practice, *Amer. J. Public Health*, 77:52-56, 1987.

Epidemic viral gastroenteritis

A clinical syndrome characterized by acute infectious gastroenteritis with watery diarrhea, vomiting, malaise, and abdominal cramps with a relatively short incubation period (12-36 h) and duration (24-48 h). A viral etiology is suspected when bacterial and parasitic agents are not found. In the United States, no etiologic agent can be found in 70% of the outbreaks of gastroenteritis. Most of these may be due to viral agents, such as the Norwalk, Snow Mountain, and Hawaii agents, astroviruses, caliciviruses, adenoviruses, nongroup A rotaviruses, and paroviruses. Epidemics are common worldwide, and have occurred following the consumption of fecally contaminated raw shellfish, food, or water, although the virus may be spread by airborne droplets as well. Epidemics are most frequent in residential homes, camps, institutions, and cruise ships. Many individual cases of mild diarrhea may in fact occur in epidemics for which the source of the infection cannot be found. Epidemic viral gastroenteritis is distinct from rotavirus diarrhea, a seasonal disease in winter that is the most common cause of diarrhea in young children, and affects virtually all children in the first 4 years of life. See DIARRHEA; INFANT DIARRHEA.

Although many different agents may cause epidemic viral gastroenteritis, most of these viruses are difficult to detect even with the most sophisticated tests. Since the diarrhea is often mild and of short duration, attention should be given to rehydration therapy and prevention by identification of the source. Some epidemics have been noteworthy, such as those of group B rotavirus that occurred in China in 1982-1983, affecting more than 1 million people. Fatalities have been associated with severe dehydration and loss of fluids and electrolytes in the stool. Research may improve the ability to diagnose these agents, understand the epidemiology of the diseases they cause, and develop protective strategies. Some of these agents, such as astroviruses and caliciviruses, appear to affect young children, whereas other viruses such as Norwalk, Snow Mountain, Hawaii, and Taunton agents have been identified in older children and adults. The degree to which these agents induce an immune response is not known. See ANIMAL VIRUS.

Roger Glass

Epidemiology

The study of the distribution of diseases in populations and of factors that influence the occurrence of disease. Epidemiology examines epidemic (excess) and endemic (always present) diseases; it is based on the observation that most disease does not occur randomly, but is related to environmental and personal characteristics that vary by place, time, and subgroup of the population. The epidemiologist at-

tempts to determine who is prone to a particular disease; where risk of the disease is highest; when the disease is most likely to occur and its trends over time; what exposure its victims have in common; how much the risk is increased through exposure; and how many cases of the disease could be avoided by eliminating the exposure.

In the course of history, the epidemiologic approach has helped to explain the transmission of communicable diseases, such as cholera and measles, by discovering what exposures or host factors were shared by individuals who became sick. Modern epidemiologists have contributed to an understanding of factors that influence the risk of chronic diseases, particularly cardiovascular diseases and cancer, which account for most deaths in developed countries today. Epidemiology has established the causal association of cigarette smoking with heart disease; shown that acquired immune deficiency syndrome (AIDS) is associated with certain sexual practices; linked menopausal estrogen use to increased risk of endometrial cancer but to decreased risk of osteoporosis; and demonstrated the value of mammography in reducing breast cancer mortality. By identifying personal characteristics and environmental exposures that increase the risk of disease, epidemiologists provide crucial input to risk assessments and contribute to the formulation of public health policy.

Epidemiologic studies, based mainly on human subjects, have the advantage of producing results relevant to people, but the disadvantage of not always allowing perfect control of study conditions. For ethical and practical reasons, many questions cannot be addressed by experimental studies in humans and for which observational studies (or experimental studies using laboratory animals or biomedical models) must suffice. Still, there are circumstances in which experimental studies on human subjects are appropriate, for example, when a new drug or surgical procedure appears promising and the potential benefits outweigh known or suspected risks.

Descriptive studies. Descriptive epidemiologic studies provide information about the occurrence of disease in a population or its subgroups and trends in the frequency of disease over time. Data sources include death certificates, special disease registries, surveys, and population censuses; the most common measures of disease occurrence are (1) mortality (number of deaths yearly per 1000 of population at risk); (2) incidence (number of new cases yearly per 100,000 of population at risk); and (3) prevalence (number of existing cases at a given time per 100 of population at risk).

These measures can be calculated for all causes combined or for specific causes of disease and can pertain to an entire population or to specific subgroups (age or racial). When a rate for a specific subgroup is calculated, the number of people in the subgroup with the disease is the numerator and the population is the denominator. For example, to calculate the incidence rate for endometrial cancer in black women aged 50-54 in 1980, the numerator

should include the new cases of endometrial cancer diagnosed in black women in that age group in 1980 and the denominator should give the total number of black women in that age group in the population that year.

The risk of most diseases is related to age. For example, heart disease mortality rates increase markedly with age. Thus, a country with a young population may have low overall rates of heart disease compared to another country whose inhabitants are, on average, much older; still, the risk of heart disease in each specific age group could be similar in the two countries. To compare the risk of disease across populations, differences in the age structure of those populations must be taken into account; similarly, to compare risk in the same population over time, the changing age composition of that population must be taken into account.

The effects of age can be controlled in several ways. For example, with age adjustment by the direct method, the age-specific rates of the populations of interest are applied to the age distribution of a standard population, and the resulting numbers of expected deaths for each population are compared. By the indirect method, the age-specific rates of a standard population are applied to the number of people in each age group in the study population, and the resulting numbers of deaths are compared to the number actually observed. Similar procedures adjust rates for other factors that influence risk of disease, for example, sex or race.

Descriptive measures are useful for identifying populations and subgroups at high and low risk of disease and for monitoring time trends for specific diseases. They provide the leads for analytic studies designed to investigate factors responsible for such disease profiles.

Analytic studies. Analytic epidemiologic studies seek to identify specific factors that increase or decrease the risk of disease and to quantify the associated risk.

Observational studies. In observational studies, the researcher does not alter the behavior or exposure of the study subjects, but observes them to learn whether those exposed to different factors differ in disease rates. Alternatively, the researcher attempts to learn what factors distinguish people who have developed a particular disease from those who have not.

The observational analytic studies most often performed are the cohort study and the case-control study. In the former (also called prospective study), a disease-free group of people, often characterized by a common feature such as occupation, is identified; data on exposures of interest are collected; and the people are observed over time to see if they develop the disease. Occurrence can be compared between the entire cohort (factory workers exposed to certain chemicals) and an external group (the general population) or between subgroups in the study cohort who differ by exposure (workers in the same factory, some exposed and others not exposed to certain chemicals). The relative risk of disease asso-

ciated with exposure is determined by dividing rates of disease in the exposed group by rates in the non-exposed comparison group. Cohort studies are ideal in the respect that data gathered on exposure predate development of disease. However, such studies require large numbers of people and long follow-up periods to ensure reliable estimates of disease occurrence. Cohort studies are typically expensive and administratively complex.

In case-control studies, people who have already developed the disease of interest (cases) and people who are free of the disease (controls) are studied. Data are collected from individuals in both groups on personal characteristics or previous exposure to factors that are suspected to be determinants of the disease, and the two groups are compared. Although the relative risk of disease associated with exposure cannot be calculated directly in case-control studies, it can be closely estimated by a quantity known as the odds ratio, that is, the odds of exposure in the cases divided by the odds of exposure in the controls. Because it is not necessary to observe study subjects over time for the occurrence of disease, case-control studies take less time than cohort studies and require fewer subjects to provide statistically meaningful results; they are therefore generally less expensive to conduct. However, it is not always possible to be certain that the factor suspected of causing the disease preceded the development of the disease; the time sequence is important for establishing causality. Furthermore, the accuracy of exposure information often depends on the recall of past events (dietary practices in early life, dosages of past medications, or occupational exposures to particular chemicals); and this recall may be poor.

Experimental studies. In experimental studies, the investigator alters the behavior, exposure, or treatment of people to determine the impact of the intervention on the disease. Usually two groups are studied, one that experiences the intervention (the experimental group) and one that does not (the control group). Ideally, subjects are randomly assigned to the groups and, where possible, those responsible for assessing outcomes are unaware of (blind to) the assignments. Intervention studies include risk-reduction trials of attempts to alter risk factors (such as cessation of smoking by middle-aged men at risk of heart disease); screening trials of the effectiveness of procedures that detect early subclinical cases (such as the Papanicolaou test to screen for cervical cancer); or clinical trials of the effectiveness of specific treatments (such as studies comparing survival after coronary bypass surgery to survival on medication alone).

Outcome measures include incidence, mortality, and survival rates in both the intervention and control groups. Some studies extend only to rates of compliance (the percentage of individuals assigned to the smoking cessation program who attend the program compared to those not assigned who may seek help on their own) or to risk-factor changes (the percentage of smokers assigned to the smoking cessation program who quit smoking compared to

those not assigned) rather than to measuring disease occurrence per se (heart disease rates in those assigned to the smoking cessation program compared to those not assigned).

Association versus causation. An association of a factor with a disease does not necessarily mean that the two are causally related. For example, people who carry matches may have higher rates of cancer, but matches do not cause the disease. Criteria for establishing causality between a factor and a disease include the strength of the observed relationship, supported by any biological evidence, consistency with results of other epidemiologic studies, and the correct time sequence between exposure to the factor and development of the disease.

Observed associations may be chance findings. Biostatistical methods, which are essential in modern epidemiology, are used to test the possibility that observed results are due to chance and consider the effects of other factors on the associations that are found. The association of a factor with a disease may arise because both are related to another factor (known as a confounding factor). For example, smoking is independently related to both alcohol consumption and oral cancer; thus, individuals who drink alcohol are at increased risk of oral cancer in part because of their smoking. The relative risk of oral cancer attributed to alcohol consumption would be artificially high if smoking (a confounding factor) were not taken into account. Design and analysis techniques deal with the effects of potential confounding factors. See DISEASE; EPIDEMIC.

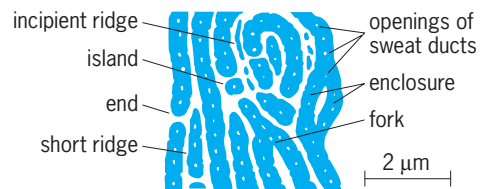
Virginia L. Ernster

Bibliography. G. D. Friedman, *Primer of Epidemiology*, 4th ed., 1994; C. H. Hennekens and J. E. Buring, *Epidemiology in Medicine*, 1987; S. B. Hulley, *Designing Clinical Research*, 1988; J. L. Kelsey, W. D. Thompson, and A. S. Evans, *Methods in Observational Epidemiology*, 2d ed., 1996.

Epidermal ridges

Minute corrugations of the skin. They compose a sculpturing, termed dermatoglyphics, which characterizes the palmar and plantar surfaces of humans and the nonhuman primates. These areas lack hair and sebaceous (oil) glands, but sweat glands are numerous. In certain kinds of monkeys, a portion of the undersurface of the tail bears similarly specialized skin.

Epidermal ridges do not course in a uniform direction, like the ribs of corduroy, but are arranged in configurations of variable types, as exemplified in the familiar arches, loops, and whorls of the fingers. Ridges are dissimilar to the ribs of corduroy also in that they present frequent interruptions, branchings, and other irregularities. All these details are unchanged in lifetime except as they may be destroyed or damaged by accident or disease. Even a small area of skin has a unique individuality of these detailed features, which provide a basis for positive personal identification.



Characteristics of individual ridges and their configurations, seen in a print. (After H. Cummins and C. Midlo, *Fingerprints: Palms and Soles*, McGraw-Hill, 1943)

The ridges over the human hand as a whole, in young adult males, average 0.019 in. (0.48 mm) in breadth. They are slightly narrower in females, 0.017 in. (0.43 mm).

Ridges serve two functions: (1) They increase security of contact with objects, in the manner of the milling of a tool handle. Ducts of sweat glands open on the summits of ridges, and moistening of the skin augments the security of contact. (2) They enhance the sense of touch. In passing the fingers or palm over an object for judging its texture, the slight displacement of ridges heightens stimulation of the underlying nerve endings.

The characteristics of individual ridges and their collective configurations (see *illus.*) are not ordinarily studied by direct inspection of the skin. Instead, they are examined in prints, usually impressions of the inked surface, made for the purpose of record, or natural imprints made by chance contact. See FINGERPRINT; INTEGUMENTARY PATTERNS. Harold Cummins

Epidermis (plant)

The outermost layer (occasionally several layers) of cells on the primary plant body. Its origin and structure are variable. In addition to a general description of the epidermis of plants, this article singles out five structural components of the tissue: (1) cuticle; (2) stomatal apparatus (including guard cells and subsidiary cells); (3) bulliform (motor) cells; (4) trichomes; and (5) root hairs.

In most roots the epidermis has a common origin with other tissues and therefore becomes established a short distance back of the apex. In shoots that have stratified meristems, the outermost layer (first layer of the tunica) usually gives rise to the epidermis. In many grasses, however, this layer divides and contributes to internal tissues of the leaf primordia before the epidermis is organized. In plants that have unstratified shoot meristems, the embryonic epidermis (protoderm) may originate on the flanks of the apical cone or on organ primordia. See APICAL MERISTEM.

Leaves, herbaceous stems, and floral organs usually retain the epidermis through life. Most woody stems retain it for one to many years, after which it is replaced. In roots it is usually short-lived. See PERIDERM.

Ordinarily the epidermis is compact and devoid of intercellular spaces, aside from the stomata. In surface view the individual cells commonly appear

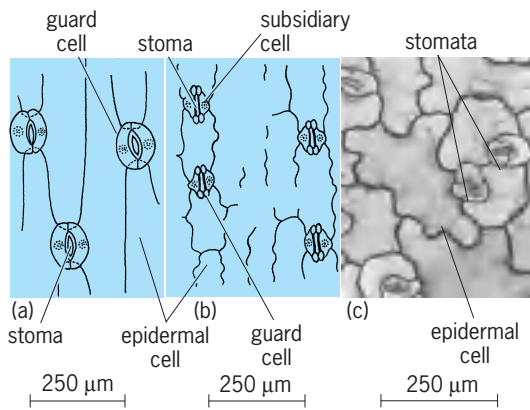


Fig. 1. Plant epidermis in surface view. (a) *Iris* leaf showing kidney-shaped guard cells. (b) *Zea* leaf showing dumbbell-shaped guard cells; two subsidiary cells are associated with each stoma. (c) *Sedum* leaf; each stoma is associated with three or four subsidiary cells.

polygonal or have wavy outlines (**Fig. 1**). Those on stems and linear leaves may be conspicuously elongated parallel to the axis of the organ. Specialized epidermal cells include the silica and cork cells in the leaves of some grasses, the lithocysts in *Ficus* (**Fig. 2a**) and *Cannabis*, and fiberlike cells of *Raphia* and *Stylidium*. In the dry bulb scales of garlic the epidermis consists partly or entirely of heavily thickened sclereids. In the seed coats of many legumes it consists of a compact layer of macrosclereids arranged in columnar fashion. The epidermis of certain Polypodiaceae is photosynthetic. There are also plants in which the entire leaf, excluding the veins, consists of a double (*Elodea*) or even a single (*Hymenophyllum*) layer of epidermis. See SCLERENCHYMA.

Differential thickening of the walls characterizes most epidermal cells. In the more delicate types of

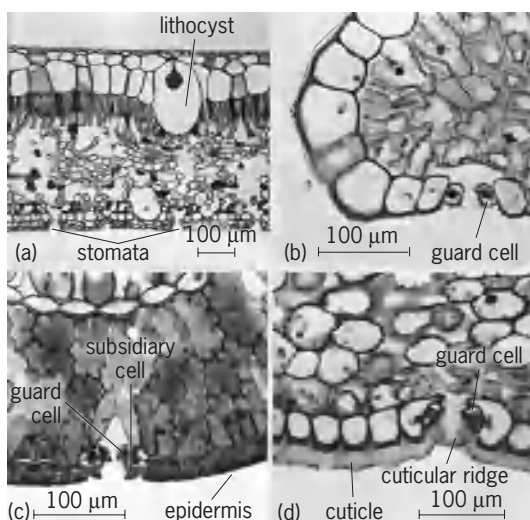


Fig. 2. Specialized epidermal cells in transections of leaves. (a) Leaf of *Ficus elastica*, each side having approximately three layers of cells in the multiple epidermis; lithocyst containing cystolith whose calcium deposits were removed during preservation. (b) Margin of leaf of *Lilium* with relatively thin-walled epidermis. (c) Leaf of *Pinus* with thick-walled epidermis and sunken guard cells. (d) Leaf of *Clivia* with heavy cuticle.

epidermis, the outer tangential wall is commonly the thickest (**Fig. 2b**). In the leaves of many conifers, however, the walls of epidermal cells are so massively thickened that the cell cavities are almost obliterated (**Fig. 2c**).

Multiple epidermis. This tissue is several layers of cells in depth. In its formation the protoderm undergoes periclinal (parallel with the circumference) divisions. A multiple epidermis occurs in many species of *Ficus*, *Begonia*, and *Peperomia*. The outermost layer of a multiple epidermis resembles a uniseriate (single series) epidermis, whereas the other layers constitute a water-storage tissue. In some peperomias, the water-storage tissue is thicker than the rest of the leaf. The velamen (parchmentlike covering) of orchid roots is a multiple epidermis.

Cutin, cuticle, and waxes. Cutin is a mixture of fatty substances characteristically found in epidermal cells. It impregnates the outer cell walls and occurs as a continuous layer (cuticle) on the outer surface (**Fig. 2d**). The cuticle covers the surfaces of young stems, leaves, floral organs, and even apical meristems. The cuticle can be stripped from some organs in relatively large sheets, a feature indicating its continuity. Cutin has been detected in the walls of cells adjacent to the intercellular spaces within leaves. There are some reports that it occurs in the epidermis of young roots.

Waxes appear as a deposit on the outside of the cuticle in many plants; the bloom on purple grapes and plums is an example. Most often the waxes are present in small quantity, but the leaves of some plants may be almost white with wax (*Echeveria subrigida*). The waxes of a few species are of great commercial value in the manufacture of polishes for floors, furniture, automobiles, and shoes. Carnauba wax is obtained from the leaves of a palm (*Copernicia cerifera*); candelilla wax, from the pencil-like stems of *Euphorbia antisyphilitica*. Other substances, such as gums, resins, and salts, usually in crystalline form, may be deposited on the outside of the cuticle.

Stomata. The apertures in the epidermis which are surrounded by two specialized cells, the guard cells, are known as stomata. The singular form, stoma, is derived from the Greek word for mouth. However, some authorities prefer to include both aperture and guard cells within the concept of stoma.

Stomata occur most commonly on the leaves and young stems of vascular plants, with the exception of certain submerged water plants (*Elodea*). They are sometimes found on sporangia (*Botrychium*) and on the sporophytes (spore-producing generation) of certain nonvascular plants (mosses, *Anthoceros*). The apertures of stomata are contiguous with the intercellular space system of underlying tissues and thus permit gas exchange between internal cells and the external environment. See PLANT PHYSIOLOGY.

The number of stomata per square centimeter of epidermis varies from none to more than 100,000. In many plants stomata are restricted largely to the lower surface of the leaf. Occasionally they occur

only on the upper surface (*Castalia*, *Ammophila*). They may also occur in equal numbers on both sides of the leaf (*Abronia*, *Atriplex*). In oleander, stomata are restricted to special crypts on the lower leaf surface. In many dicotyledons stomata are arranged at random; in many monocotyledons and in the leaves of gymnosperms they are arranged with their long axes parallel to that of the leaf. The guard cells are frequently associated with subsidiary cells, which differ from other epidermal cells in size and form.

The guard cells of many plants are at about the same level as the other epidermal cells; in others they project above them. Xerophytes, plants which can subsist with a small amount of moisture, often have a thick-walled epidermis and the guard cells may be deeply sunken. In the majority of plants guard cells are kidney-shaped in surface view; in grasses, however, they resemble a pair of dumbbells. The guard cells undergo changes in form during the changes in turgor (water pressure), a phenomenon associ-

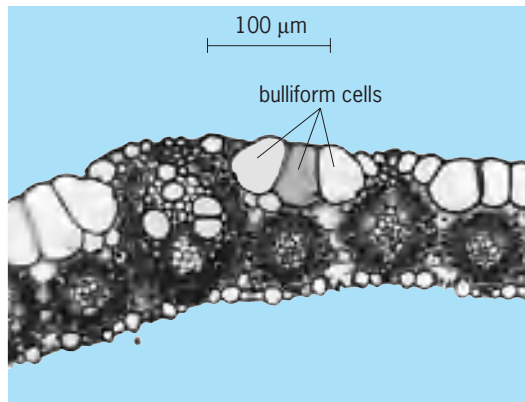


Fig. 3. Transection of leaf of *Andropogon* showing large bulliform cells in the upper epidermis.

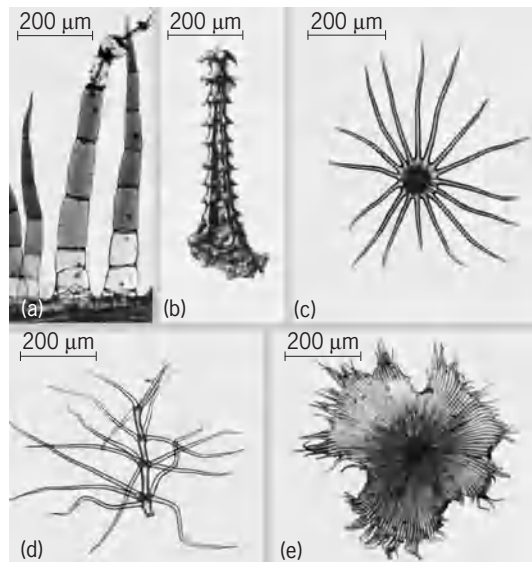


Fig. 4. Trichomes. (a) Uniseriate hairs of *Gynura*. (b) Unicellular anchor hair of *Mentzelia*. (c) Stellate hair of *Solanum*. (d) Branched (candelabra) hair of *Platanus*. (e) Peltate scale of *Shepherdia*.

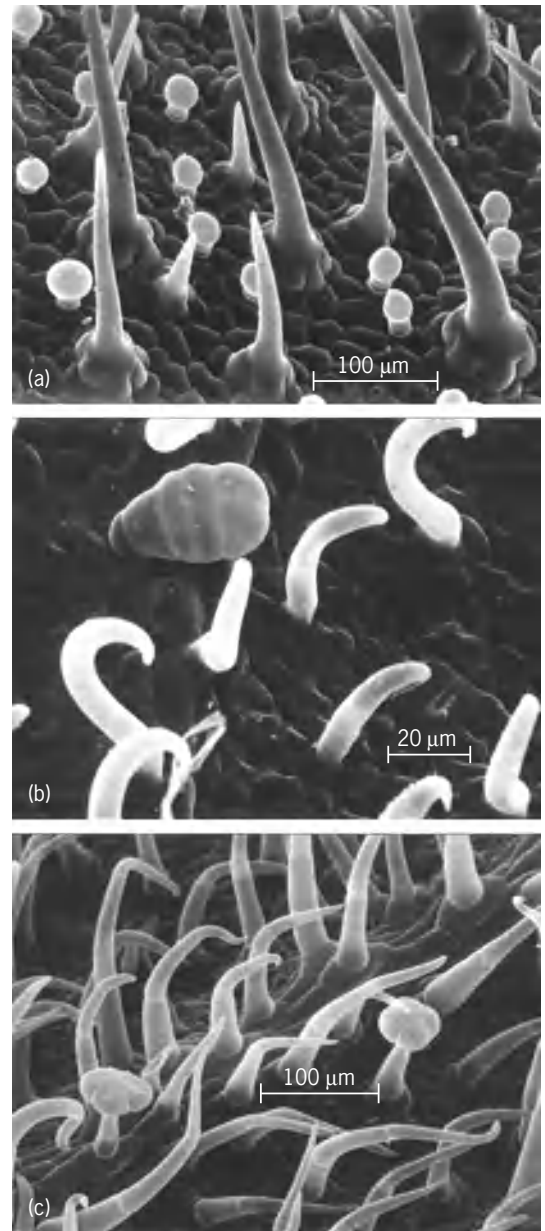


Fig. 5. Trichomes. (a) Unicellular and glandular (colleters) hairs of the geranium (*Pelargonium*). (b) Unicellular-hooked and uniseriate, club-shaped hairs of the bean (*Phaseolus*). (c) Uniseriate and glandular (colleters) hairs of the tomato (*Lycopersicon*).

ated with the characteristic uneven thickening of the guard cell walls.

The opening and closing of the stomatal aperture is caused by relative changes in turgor between the guard cells and surrounding epidermal cells. It is accompanied by hydrolysis, or condensation, of starch within the guard cells. See PLANT-WATER RELATIONS.

Bulliform (motor) cells. These large, highly vacuolated cells occur on the leaves of many monocotyledons but are probably best known in grasses (Fig. 3). They are thought to play a role in the unfolding of developing leaves and in the rolling and unrolling of mature leaves in response to alternating wet and dry periods. Not all investigators, however, subscribe to

this idea. The water vesicles in the epidermis of the ice plant (*Mesembryanthemum crystallinum*) are similar in appearance to bulliform cells.

Trichomes. Appendages derived from the protoderm are known as trichomes; the simplest are protrusions from single epidermal cells. Included in the concept, however, are such diverse structures as uniseriate hairs (Fig. 4a), multiseriate hairs (*Begonia*, *Saxifraga*), anchor hairs (Fig. 4b), stellate hairs (Fig. 4c), branched (candelabra) hairs (Fig. 4d), peltate scales (Fig. 4e), stinging hairs, and glandular hairs (Fig. 5). Many of these structures arise on developing plant organs and are most conspicuous at this stage. They may be persistent or short-lived. Persistent trichomes are responsible for the silver-gray color of sagebrush (*Artemisia*), Spanish moss (*Tillandsia*), and Russian olive (*Elaeagnus*). Mature trichomes may or may not retain their protoplasts. The assignment of functions to trichomes is often uncertain, although they are useful in plant identification. Cotton and kapok fibers are unicellular epidermal hairs. Cotton fibers grow from the seed coat, those of kapok from the pod epidermis.

Root hairs. These are thin-walled extensions of certain root epidermal cells. They develop only on growing root tips and may arise from any epidermal cell, or from specialized cells known as trichoblasts. Root hairs increase the absorbing area of the root tip manifold. The life of a given root hair is usually numbered in days. Persistent root hairs, reported for some plants, probably play no role in absorption. See PHLOEM; PLANT ORGANS; PLANT TISSUE SYSTEMS; ROOT (BOTANY); SECRETORY STRUCTURES (PLANT); XYLEM.

Norman H. Boke

Epidote

The group name for a family of minerals of general composition $\text{Ca}_2(\text{Fe}^{3+}, \text{Al}, \text{Mn}^{3+})\text{Al}_2\text{O}[\text{SiO}_4][\text{Si}_2\text{O}]$ -(OH) which occur widely in metamorphic and igneous rocks (Fig. 1). Epidote [octahedral ferric iron (Fe^{3+}) dominant] and clinozoisite [aluminum (Al) dominant] represent the most common compositions among the epidote group; a third composition, piemontite [manganese (Mn^{3+}) dominant], is less abundant. Allanite refers to compositions displaying significant rare-earth (such as lanthanum or cerium) substitution for calcium (Ca^{2+}), with corresponding replacement of Fe^{3+} by ferrous iron (Fe^{2+}). A fifth member, zoisite, is equivalent to clinozoisite, but it has a different crystalline system. Rare epidote-clinozoisites abundant in chromium (Cr), vanadium (V), and lead (Pb) and allanites rich in fluorine (F), beryllium (Be), and phosphorus (P) also exist. See CRYSTAL STRUCTURE; SOLID SOLUTION.

Structure. The epidote structure (Fig. 2) consists of chains of octahedra (M sites) elongated parallel to the y crystallographic axis and cross-linked by the silicate groups SiO_4 and Si_2O_7 . There are three distinct kinds of octahedra (M1, M2, and M3) arrayed along two types of chains. One chain type consists entirely of edge-sharing octahedra (M2), whereas the

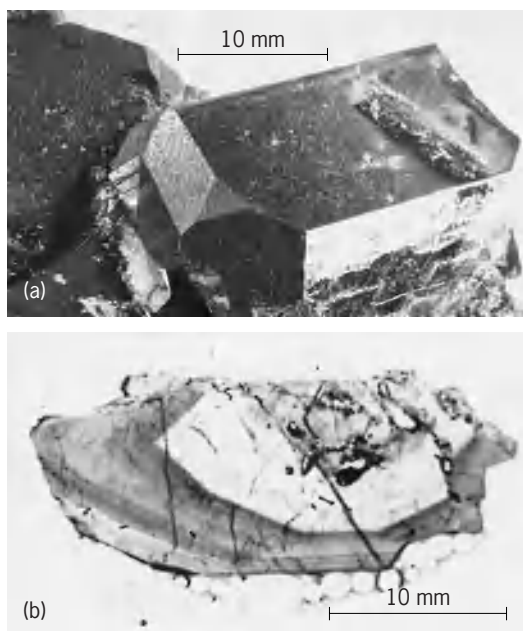


Fig. 1. Epidote. (a) Crystals from Prince of Wales Island, Alaska (Yale University). (b) Oscillatory zoning of these crystals.

other contains M1 octahedra alternating with M3 octahedra along its length. The resulting framework structure contains large cavities (A sites) where the Ca^{2+} cations are housed in 9-10-fold coordination. Like the chains of octahedra, epidote crystals themselves are also elongated parallel to the y axis. Such a relationship between internal crystal structure and external crystal form is a common theme among minerals in general. Epidotes exhibit cleavage on the $yz(100)$ and $xy(001)$ planes, and their hardnesses on the Mohs scale vary between 6 and 7.

In typical epidote-clinozoisites, Fe^{3+} and Mn^{3+} substitute for Al principally in M3 sites; M2 sites house only Al. M1 sites are usually dominated by Al but

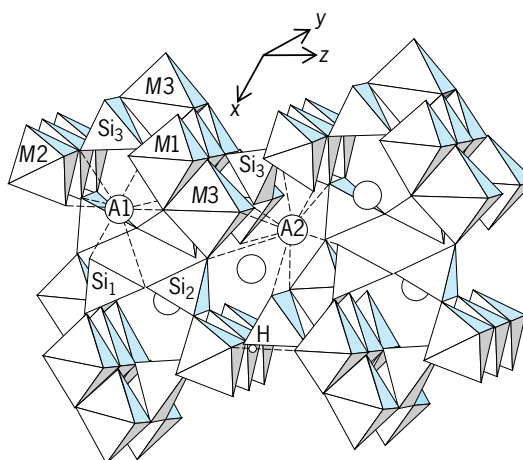


Fig. 2. Crystal structure of the epidote-group mineral clinozoisite. Si_3 represents the silicate group SiO_4 ; Si_1 and Si_2 represent the silicate group Si_2O_7 . M1, M2, and M3 are distinct kinds of octahedra. A1 and A2 are cavities in the framework structure. (After W. A. Deer, R. A. Howie, and J. Zussman, *An Introduction to the Rock-Forming Minerals*, Halsted Press, 1992)

may contain Mn^{3+} . The explanation for the preference of these cations for the M3 site is twofold. First, both cations are larger than Al^{3+} , and the M3 octahedron is the largest in the structure. Second, the M3 sites are relatively distorted, and this characteristic may permit Mn^{3+} to gain additional stabilization energy here, over the other M sites.

Occurrence. Epidote group minerals, particularly epidote and clinozoisite, are common and widespread in regional- and contact-metamorphic rocks, both as primary and secondary (that is, alteration) minerals. They occur together as individual grains, as intergrowths, or as zoned crystals (Fig. 1a). Epidote and (clino)zoisite are found in aluminous limestones with grossularite, anorthite, microcline, quartz, and calcite; in mafic schists and gneisses with hornblende, albite, and chloritoid; in actinolite greenschists with chlorite, sphene, albite, quartz, calcite, and magnetite; in hornfels with diopside, actinolite, grossularite, and albite; in glaucophane schists; in quartzites; and in slates. Approximate depth-temperature conditions of their formation range from 5–25 km (3–15 mi) and 300–500°C (570–1020°F; low-grade) to 5–25 km (3–15 mi) and 450–650°C (840–1200°F; medium-grade). However, their stabilities are also sensitive to the pressure of oxygen in the rock during metamorphism. A zoned crystal will grow in a local environment of fluctuating f_{O_2} during metamorphism. Clinozoisite forms during periods of low f_{O_2} , whereas epidote grows during oscillations to high oxygen pressure. The cause of the oscillations in oxygen pressure is probably diffusion-controlled crystal growth, in which the growth rate exceeds that of oxygen diffusion to the crystallization site. See METAMORPHIC ROCKS.

In igneous rocks, epidote-(clino)zoisite occurs primarily in granitoids, as hydrothermal veins, as an alteration of plagioclase, or as a primary phase; in basalts, in cavities associated with augite, actinolite, grossularite, and sphene; and in mafic igneous rocks, as an alteration of plagioclase (that is, saussurite). See IGNEOUS ROCKS.

Piemontite and allanite are relatively restricted in occurrence. Piemontite is found chiefly in manganese ore deposits, altered volcanic rocks, and some schists. Allanite is an accessory mineral in granitoids, pegmatites, their volcanic equivalent rocks, and mafic schists and gneisses.

Epidote-group minerals occur as excellent crystals in several localities, notably near Salzburg, Austria; Isère, France; Piemonte, Italy; Riverside, California; and Prince of Wales Island, Alaska. See SILICATE MINERALS.

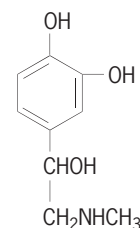
Peter S. Dahl

Bibliography. P. S. Dahl and L. M. Friberg, The occurrence and chemistry of epidote-clinozoisites in mafic gneisses from the Ruby Range, southwestern Montana, *Contrib. Geol.*, 18:77–82, 1980; W. A. Deer, R. A. Howie, and J. Zussman, *An Introduction to the Rock-Forming Minerals*, 2d ed., 1992; W. A. Dollase, Refinement of the crystal structures of epidote, allanite, and hancockite: *Amer. Mineral.*, 56:447–464, 1971; E. W. Heinrich, *Microscopic Identification of Minerals*, 1965; K. Langer

and M. Raith, Infrared spectra of Al-Fe(III)-epidotes and zoisites, $Ca_2(Al_{1-p}Fe^{3+p})Al_2O(OH)[Si_2O_7][SiO_4]$, *Amer. Mineral.*, 59:1249–1258, 1974; C. Klein and C. S. Hurlbut, Jr., *Manual of Mineralogy*, 21st ed., 1993; G. H. Myer, New data on zoisite and epidote, *Amer. J. Sci.*, 264:364–385, 1966; E-an Zen and J. M. Hammarstrom, Magmatic epidote and its petrologic significance, *Geology*, 12:515–518, 1984.

Epinephrine

A hormone that is an important secretion from the adrenal medulla; also known as adrenaline. In 1901, J. Takamine isolated epinephrine in pure form from extracts of adrenal glands. Its structure (shown below) was finally proved by synthe-



sis by F. Stolz in 1904 to be the catecholamine 1-(3,4-dihydroxyphenyl)-2-methyl-aminoethanol. See ADRENAL MEDULLA; ADRENAL MEDULLA HORMONE; HORMONE.

In mammals, the necessary components for the manufacture of the hormone by the body are the amino acids tyrosine or phenylalanine. Tyrosine is converted to dihydroxyphenylalanine (dopa), which is subsequently converted to dopamine, which in turn is modified to form norepinephrine. Norepinephrine is the immediate precursor for epinephrine. This final conversion requires the enzyme phenylethanolamine-*N*-methyltransferase (PNMT). PNMT is present in significant amounts only in the adrenal medulla and in some areas of the brain. Its expression is highly dependent on a supply of glucocorticoids, such as cortisol. Blood flowing from the adrenal cortex drains into the medulla, providing these cells with high glucocorticoid concentrations, sufficient to maintain PNMT production. See AMINO ACIDS; DOPAMINE; NORADRENERGIC SYSTEM.

Epinephrine and norepinephrine are stored in adrenal medullary cells (known as chromaffin cells) in granules bound to proteins called chromogranins. The granules also contain adenosine triphosphate (ATP). In response to release of acetylcholine from nerves innervating the medulla, the granules fuse with the chromaffin cell membrane, and their contents, including epinephrine, are released by exocytosis and quickly enter the blood.

Epinephrine is a sympathomimetic substance; that is, it acts on tissue supplied by sympathetic nerves, and generally the effects of its action are the same as those of other nerve stimuli. Epinephrine effects are mediated through alpha and beta adrenergic receptors, located in many tissues. There exist several

subtypes of alpha and beta receptors. Thus, epinephrine plays an important role in preparing the organism to meet conditions of physiologic emergency. Although medullary cells contain both epinephrine and norepinephrine, it is epinephrine that is released in sufficient quantities to have a greater physiologic impact on tissues.

When injected intravenously, epinephrine causes an immediate elevation in blood pressure and heart rate. Its effects enhance blood flow into muscle and additionally relax bronchial smooth muscle, allowing for greater oxygenation of blood. The chief metabolic changes following the injection of epinephrine are a rise in the basal metabolic rate and an increase of blood sugar. Epinephrine effects the latter action by causing an increase in the rate of glycogenolysis, or breaking down of stored sugar in the liver. These effects of epinephrine are transitory. See CARBOHYDRATE METABOLISM; GLYCOGEN.

Robert Kempainen; Choh Hao Li

Bibliography. W. F. Boron and E. L. Boulpaep, *Medical Physiology*, Elsevier Science/Saunders, Philadelphia, 2005; W. F. Ganong, *Review of Medical Physiology*, 22d ed., McGraw-Hill, New York, 2005; P. R. Larsen et al., *Williams Textbook of Endocrinology*, 10th ed., Elsevier Science/Saunders, Philadelphia, 2003.

Epitaxial structures

Epitaxial interfaces in solids are a special class of crystalline interfaces where the molecular arrangement of one crystal on top of another is defined by the crystallographic and chemical features of the underlying crystal. Royer (1928) systematically showed how the geometry of the atomic arrangement on a substrate could affect the crystallographic orientation of the overgrowth. The term "epitaxy" was introduced to describe the importance of having parallelism between two lattice planes with similar networks of closely similar spacing. Epitaxial phenomena are important to study and understand, as they occur widely in nature (such as oxidation) and are the foundation by which modern semiconductor devices are grown and fabricated. See CRYSTAL GROWTH.

Epitaxial interfaces are a subset of a class of interfaces where lattice planes achieve a correspondence across an interface. If the matching is not perfect, such a correspondence can be achieved by a number of ways, including dilation and contraction of lattice planes; rotation of overgrowth (epilayer relative to the orientation of the substrate) until a set of closely matched lattice spacing can be found; and tilting of the epilayer with respect to the substrate (Fig. 1).

A combination of these effects can occur in some cases. When the epitaxial relationship is only approximate, a perfect matching of the two structures on an atomic scale would in general create prohibitively large elastic stresses in each crystal. Thus, small regions of perfect matching may describe the interface

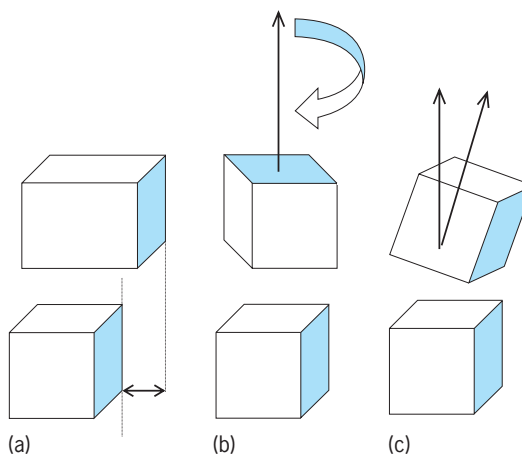


Fig. 1. Matching of lattice planes between the substrate and the epitaxial layer by (a) expansion or contraction of lattice plane, (b) in-plane rotation, and (c) tilting (out-of-plane rotation).

structure, separated by zones of mismatch. Depending upon the accommodation mechanisms between the crystals (that is, dilation, rotation, or tilting), the structural matching at the atomic or molecular level is often accompanied by the formation of interfacial defects which take the form of misfit dislocations. The crystallographic features of these misfit dislocations can be of varying complexity, ranging from one-dimensional features such as dislocations to two-dimensional defects such as atomic ledges and step (Fig. 2). See CRYSTAL DEFECTS; CRYSTAL STRUCTURE.

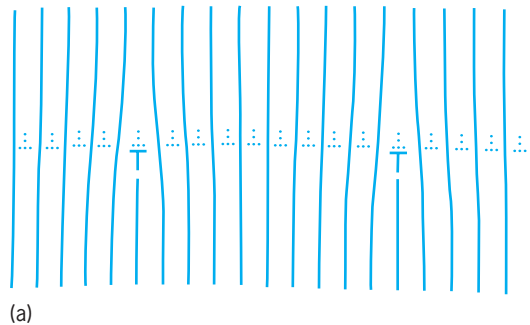
The extent to which these stresses can be accommodated before the breakdown of interfacial matching (or coherency) is very dependent upon the thickness of the film. Hence there exists a "critical thickness" before coherency at the atomic level breaks down. Numerous theoretical and experimental studies exist to understand the details of how these limits are reached. From a technological perspective, this is an important issue as one tries to grow epitaxial films which are coherent and relatively defect-free, and there is a need to know the mechanisms and means by which these limits can be extended. See SEMICONDUCTOR HETEROSTRUCTURES.

A simple but technologically important demonstration of defect accommodation at epitaxial interfaces is the epitaxial growth of cube-cube matching of unit cells where the spacing in the cube direction (or lattice parameter) is slightly different but the mismatch can still be accommodated by dilation processes (Fig. 3). In this example, the epilayer material has a slightly larger lattice parameter compared to the substrate. To achieve coherent matching of the planes that lie perpendicular to the interface, these planes in the epilayer must contract slightly to accommodate the equivalent lattice spacing of the substrate. As the substrate is of far greater thickness or volume than the epilayer, all the accommodation process occurs entirely in the thin film. This process of contraction leads to a distortion of the lattice symmetry of the epilayer, leading to a tetragonal structure

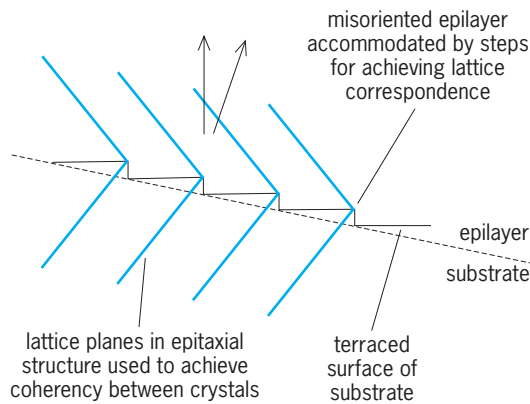
instead of maintaining the original cubic symmetry. The resulting epitaxial heterostructure now consists of a substrate and a strained layer. Multiple layers may be grown on top of the original epitaxial film, leading to a strained-layer superlattice. Such artificially layered architectures, especially in semiconductors, form an important class of materials which possess new and unique electronic properties different from the bulk forms of the same material. Introducing different levels of interfacial strain in these composite structures can alter the band gap. This unique blend of physics, mechanics, and crystal growth is the basis of developing new types of opto-electronic devices through band-gap engineering. See ARTIFICIALLY LAYERED STRUCTURES.

There are numerous modes of lattice accommodation apart from the example given above, including relaxation effects normal to the epitaxial interface; dislocation propagation and regeneration; formation and growth of planar defects; changes in surface morphology of the epitaxial film; and phase transitions such as ordering, and other forms of compositional modulations.

While this discussion has focused on issues governed by crystallography and structural issues, it should be noted that chemical effects play an equally important role. The concept of “matching” extends not only to geometrical matching but also to chemical compatibility. The extent to which the epitaxial films are mechanically stable due to coherency



(a)



(b)

Fig. 2. Misfit accommodation occurs (a) between the epitaxial layer and the substrate through the presence of interfacial misfit dislocations and (b) by atomic steps at an interface where the epitaxial film is tilted with respect to the substrate.

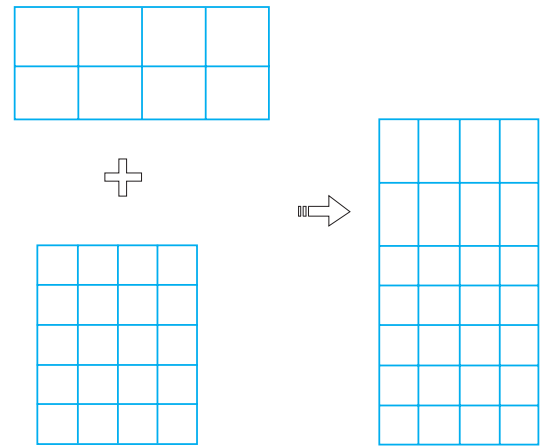


Fig. 3. Schematic illustration of the tetragonal distortion effect to accommodate purely dilational strained-layer epitaxial structure.

stresses is governed not only by the extent of lattice misfit but also by the strength of the chemical bond between the epilayer and the substrate. This property of adhesion is manifested by the extent to which the overlayer wets the substrate. There are two extreme cases by which one may describe this adhesion or wetting behavior at the molecular level. In one, proposed by F. C. Frank and J. H. van der Merwe, the adhesion between the epilayer and the substrate exceeds the binding energy between the crystal atoms, resulting in complete wetting (Fig. 4a). The nucleation of an epitaxial layer occurs by the deposition of a monolayer which is easily deformed elastically and can adhere to the substrate crystal perfectly without the formation of dislocations up to a certain level of misfit (which in some cases can be as high as 10–15%). When the thickness of the layer increases, and hence also its rigidity, this state becomes unstable and epitaxial dislocations are introduced. This mode of growth is often sought in a variety of technological applications, which require the development of large-area, defect-free epitaxial layers. It should be noted that the earlier example for growing strained layer structures is based on the Frank-van der Merwe mechanism of wetting behavior. Extremely thin layers that are only a few atoms thick can be produced in this manner. Such thin layers form the microstructural foundation for the fabrication of quantum wells, which are extremely important in semiconductor device applications. See QUANTIZED ELECTRONIC STRUCTURE (QUEST).

The other extreme is known as the Volmer-Weber mode, which is typical of weak adhesion, resulting in poor wetting characteristics. The resulting film is not continuous but consists of small islands. This type of nanoscale architecture, however, is of value when trying to fabricate quantum dot structures, which involve having nanoscale epitaxial particles rather than continuous layers. Some systems exhibit an intermediate behavior known as the Stranski-Krastanov mode, where there is good adhesion for the first monolayer while the adhesion of the subsequent layers is weaker (Fig. 4b and c). These

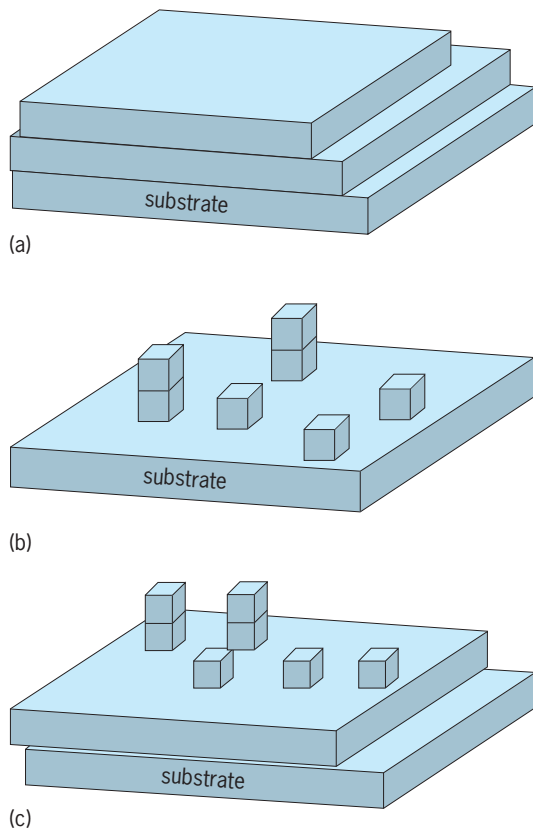


Fig. 4. Modes of epitaxial growth: (a) Frank-van der Merwe mode producing a complete wetting of epitaxial layers. (b) Volmer-Weber mode with totally incomplete wetting resulting in an island growth structure. (c) Stranski-Krastanov mode with initial layer growth accompanied by island growth.

latter modes of growth are important governing mechanisms in developing selective epitaxial growth on geometrically patterned substrates. See NANO-STRUCTURE. Krishna Rajan

Bibliography. A. A. Chernov, *Modern Crystallography*, vol. III: *Crystal Growth*, Springer-Verlag, Berlin, 1984; P. Kordos and J. Novak (eds.), *Heterostructure Epitaxy and Devices*, NATO ASI Ser., vol. 48, Kluwer Academic, 1998; M. G. Lagally (ed.), *Kinetics of Ordering and Growth at Surfaces*, NATO ASI Ser., vol. 239, Plenum, 1990; F. R. N. Nabarro (ed.), *Dislocations in Solids*, Elsevier, Amsterdam, 1983; K. Rajan, J. Narayan, and D. Ast (eds.), *Dislocations and Interfaces in Semiconductors*, TMS, Warrendale, PA, 1988; H. G. Schneider, V. Ruth, and T. Kormany (eds.), *Advances in Epitaxy and Endotaxy*, Elsevier, Amsterdam, 1990.

Epithelium

One of the four primary tissues of the body, which constitutes the epidermis and the lining of respiratory, digestive, and genitourinary passages. The major characteristic of epithelium is that the cells are close together, separated by a very small amount of intercellular substance. The epithelium (endothelium) lining the inner cavities of the heart, blood vessels, and lymphatics differs from that of all the other

groups because under abnormal conditions it may be transformed into a different type of cell characteristic of connective tissue, the fibroblast. Epithelium may be derived from any of the three primary germ layers of the very early embryo—ectoderm, endoderm, or mesoderm.

With very few exceptions (specifically, the stria vascularis of the cochlea, and the hypertrophied thyroid gland), epithelium is free of blood vessels. Nutrients reach the epithelium, and waste products leave it after passing through the ground substance of adjacent connective tissue.

Functions. The functions of epithelium are varied and include (1) protective function, by completely covering the external surface (including the gastrointestinal surface—and the surface of the whole pulmonary tree including the alveoli); (2) secretory function, by secreting fluids and chemical substances necessary for digestion, lubrication, protection, excretion of waste products, reproduction, and the regulation of metabolic processes of the body; (3) absorptive function, by absorbing nutritive substances and preserving water and salts of the body; (4) sensory function, by constituting important parts of sense organs, especially of smell and taste; and (5) lubricating function, by lining all the internal cavities of the body, including the peritoneum, pleura, pericardium, and the tunica vaginalis of the testis.

Structural integrity. The forces which hold the epithelial cells together are not satisfactorily understood. The intercellular substance between the cells, also called cement substance, is undoubtedly important. When the cement substance is weakened by removing calcium (with versene) or protein (with the enzymes trypsin or papain), the epithelial cells may be readily separated. The cement substance also contains a carbohydrate moiety. Interdigitation of adjacent cell surfaces and occurrence of intercellular bridges in certain cells may be important in holding the cells together. Finally, in certain cells local modifications of contiguous surfaces and the intervening intercellular substances, which together form the terminal bars, may be effective in the same way.

Terminal bars are visible with the light microscope after suitable staining and appear to close the spaces between the epithelial cells of the intestine at their free surface. The terminal bars are resolvable with the electron microscope into three components: (1) a very short outer tight junction, where the adjacent cells are so close that they obliterate the intercellular space; (2) an adherent zone immediately below this, where the cell membranes of adjacent cells are separated by a space of 15–20 nanometers; and (3) a desmosomal region. In the desmosomal region there is a short stretch of the cell membranes of adjacent cells where the cell membrane is slightly thickened and in direct contact with fine filaments which run into the adjacent cytoplasm. In the intercellular space between these thickened portions of the cell membranes is a thin layer of densely staining material. The intercellular bridges referred to above are essentially short processes between contiguous cells connected to each other by desmosomes. In

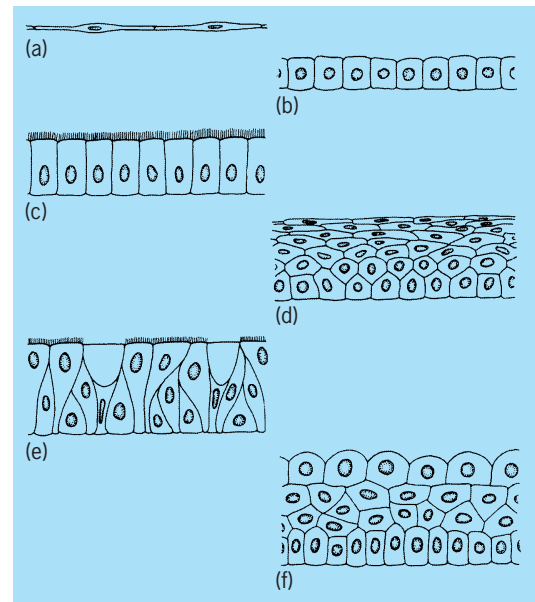
addition in certain sites (as in the skin), hemidesmosomes (thickened cell membrane plus cytoplasmic filaments) may be found in close contact with the subjacent basement membrane.

Structural polarity. A common feature of epithelial cells is their structural, as well as functional, polarity. In the intestinal epithelium, for example, the nucleus is somewhat toward the base of the cell. Between the nucleus and the base of the cell is a set of filamentous mitochondria, and between the nucleus and the surface is a set of granular mitochondria. Between the granular mitochondria and the nucleus is the Golgi apparatus. The striated border occurs only on the free surface of the epithelial cell. This border is visible with the light microscope and is resolvable into an enormous number of cell processes or cylinders 80–90 nm in diameter and up to about 1 micrometer long. Each process comprises a protoplasmic core enclosed in a cell membrane about 12.5 nm thick. After certain methods of preparation, the cell membrane appears as three sheets—two stainable layers separated by a nonstaining space. This trilaminar structure is now considered to be modified by the insertion of various sorts of protein molecules in the lipid coat. These proteins may be involved in antigenicity and cell recognition, permeability, cellular adhesion, and other cellular activities. Outside of this triple-layered cell membrane is another one which may contain protein or polysaccharide of some sort. It should be noted that, although the cell and the cell processes are polarized, the trilaminar structure of the cell membrane appears to be the same around the whole cell except for the specialized zones mentioned above (desmosomes, close junctions, and terminal bars). This cell polarity is equally marked in most epithelial cells and is even more marked in certain gland cells. For submicroscopic structure of cell organelles see CELL (BIOLOGY).

Cell and tissue affinities. It has long been known that many connective tissue cells which may be separate in their original site, except for small areas of contiguity, may in tissue cultures grow out as sheets of cells or epithelia. The factors that control this transformation are unknown but are believed to depend on some property of the cell surface not yet recognized. See CELL ADHESION; TISSUE CULTURE.

It is not altogether clear what forces prevent the sheet of epithelial cells from falling away from the underlying connective tissue. Undoubtedly the basement membrane of the connective tissue, when it occurs, has a cohesive property. When there is no basement membrane, the epithelium rests directly on the ground substance of the connective tissue, which has similar properties. The area of contact between epithelium and this ground substance is frequently increased by irregularities and processes of the basal surface of the epithelial cells.

Surface specialization. The outstanding property of the arrangement of most of the epithelium of the body is the economy of space achieved in the face of a broad exposure of the cell surfaces. The efficiency is achieved by the presence of numerous folds, which may be gross or microscopic and temporary or permanent. A part must also be attributed to



Cellular arrangements in epithelial tissues. (a) Squamous. (b) Cuboidal. (c) Columnar. (d) Stratified squamous. (e) Pseudostratified. (f) Transitional.

the surface specialization of the epithelial cells themselves, such as their minute, fingerlike processes. Another specialization of the surface or epithelial cell is the occurrence of motile cilia, each of which has a complicated structure. Especially prominent in the cilia are the nine pairs of filaments peripherally disposed in each cilium, with two central filaments. The peripheral, and frequently the central, filaments join at the base of the cilium to form a basal body, which is in turn connected with a rootlet that extends into the cytoplasm and may have a very complicated fine structure. See CILIA AND FLAGELLA.

Classification. Classification of epithelia is based on morphology, that is, on the shape of the cells and their arrangement (see *illus.*):

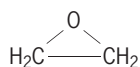
- I. Single-layered.
 - A. Squamous (mesothelium, descending loop of Henle in the kidney) flat.
 - B. Cuboidal (duct, thyroid, choroid plexus)—cubelike.
 - C. Columnar (intestine), sometimes ciliated (Fallopian tube, or oviduct)—tall.
- II. Multiple-layered or stratified.
 - A. Squamous (skin, esophagus, vagina)—superficial cells thin and flat, deeper cells cuboidal and columnar.
 - B. Columnar (pharynx, large ducts of salivary glands), sometimes ciliated (larynx)—two or more layers of tall cells.
- III. Pseudostratified (male urethra), sometimes ciliated (respiratory passages)—all cells reach to basement membrane but some extend toward the surface only part of the way, while others reach the surface.
- IV. Transitional (urinary bladder)—like stratified squamous in the fully distended bladder, in the empty bladder, superficial cells rounded, almost spherical.

Special properties. An important property of epithelium is the ability of its cells to glide over surfaces. This allows replacement of dead cells to take place in the normal state, while presenting a closed surface to the external environment; replacement is especially important in wound repair, when it is necessary that a gap in the surface be filled quickly. Gliding ability is also manifested normally in the movement of cells which slide over each other in transitional epithelium, for example, when the urinary bladder is being distended or contracted. *See* GLAND. Isidore Gersh

Bibliography. P. S. Amenta, *Histology and Human Anatomy*, 6th ed., 1995; T. S. Leeson, C. R. Leeson, and A. A. Paparo, *Text-Atlas of Histology*, 1988.

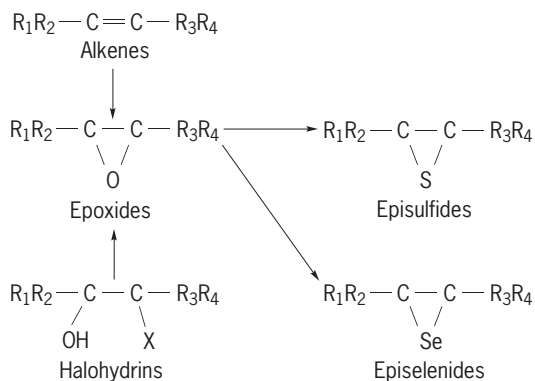
Epoxyde

A member of a class of three-membered ring cyclic ethers that are also known as oxiranes or alkylene oxides. The basic structure of an epoxyde is analogous to that of the first member, ethylene oxide:



See ETHYLENE OXIDE.

Epoxydes are made primarily by select oxidation reactions of alkenes; however, another classical preparation results from the ring closure of halohydrins by way of intramolecular nucleophilic displacement; that is, the reaction occurs within the same molecule, the alkoxide ion displacing the halogen to form a ring. The interest in epoxydes results from their ease in preparation, and their usefulness as a reactive functional group that can give a variety of products after treatment with either electrophilic or nucleophilic reagents, or on occasion after treatment with oxidizing (for example, periodic acid) and reducing (for example, titanocene dichloride) agents. The ease of opening of the strained three-membered ring epoxydes with attack of reagents in a stereospecific manner gives one or two stereochemical products (when applicable), usually in good yield. Epoxydes are also used to prepare monomers, prepolymers, polymers, and copolymers, and to promote polymerization reactions, as shown in the scheme below, where the R terms are functional



groups and X represents a halogen. *See* ALKENE;

ELECTROPHILIC AND NUCLEOPHILIC REAGENTS; HETEROCYCLIC COMPOUNDS.

Preparation. Epoxydes can be readily prepared by treatment of alkenes with peroxides (for example, hydrogen peroxide) or peroxy-acids (for example, *m*-chloroperoxybenzoic acid or magnesium monoperoxyphthalate). Several commercially available epoxydes include ethylene oxide (oxirane, $\text{R}_1 = \text{R}_2 = \text{R}_3 = \text{R}_4 = \text{H}$), propylene oxide (methyloxirane, $\text{R}_1 = \text{CH}_3, \text{R}_2 = \text{R}_3 = \text{R}_4 = \text{H}$), styrene oxide (phenyloxirane, $\text{R}_1 = \text{C}_6\text{H}_5, \text{R}_2 = \text{R}_3 = \text{R}_4 = \text{H}$), cyclohexene oxide ($\text{R}_1/\text{R}_3 = -\text{CH}_2\text{CH}_2\text{CH}_2\text{CH}_2-$, $\text{R}_2 = \text{R}_4 = \text{H}$), and epi-chlorohydrin [1-(chloromethyl) oxirane, $\text{R}_1 = \text{ClCH}_2-$, $\text{R}_2 = \text{R}_3 = \text{R}_4 = \text{H}$].

Specialty epoxyde preparations also include the treatment of compounds such as allyl alcohol with *tert*-butyl hydroperoxide (*t*-BuOOH), titanium tetraisopropoxide, and optically active diethyl tartrate; the treatment of two bromine atoms bonded to the same carbon atom (geminal dibromides; $\text{R}_1\text{R}_2\text{CBr}_2$) with lithium or butyllithium, followed by the condensation with carbonyl compounds (aldehydes or ketones; to variously substituted epoxydes); the condensation of α -haloesters with carbonyl compounds (to epoxyesters); the addition of sulfur ylides (for example, dimethylsulfonium methylide) or diazomethane (CH_2N_2) to aldehydes or ketones (usually to geminally disubstituted epoxydes); or the treatment of aldehydes with hexamethylphosphorous triamide (to symmetrically 2,3-substituted epoxydes). *See* ALDEHYDE; ESTER; KETONE.

Reactions. Nucleophilic reagents used for the ring-opening of epoxydes include hydroxide and water, alcohols or alkoxides, amines, bisulfite, hydrosulfide or mercaptide ion, acetylides, cyanide, dithiane anions, enamines, Grignard reagents, oxazine anions, and other reagents. Treatment of epoxydes with select nucleophilic reagents can lead to other products: reactions of epoxydes with triphenylphosphine form an olefin (epi-oxy-elimination); propylene oxide and certain other epoxydes react with lithium diethylamide to form allyl alcohols; and reaction of sodium azide with an epoxyde, such as styrene oxide, forms a β -azide alcohol, which can then be cyclized to an arizidine (three-membered heterocyclic nitrogen compound) by subsequent treatment with triphenylphosphine.

Common electrophilic reagents used for the ring-opening of epoxydes include hydrochloric and hydrobromic acids (hydroiodic acid with difficulty), which results in the intermediate formation of oxonium intermediates from protonation of the epoxyde oxygen atom. Straightforward acid-catalyzed hydrolysis of ethylene oxide can produce ethylene glycol (also, ethylene glycol, a 1-2 diol, can be cyclodehydrated, that is, lose a molecule of water to form the heterocyclic ethylene oxide). Similarly, treatment of epoxydes with carboxylic acids yields β -hydroxyalkylcarboxylates ($\text{RCOO}-\text{CH}_2\text{CH}_2\text{OH}$). Also, select electrophilic reagents (for example, boron trifluoride etherate) rearrange epoxydes to carbonyl compounds. Other reagents cause rearrangements of epoxy compounds. *See* REACTIVE INTERMEDIATES.

Uses. Poly(ethylene oxide) $[-(\text{CH}_2\text{CH}_2\text{O})_n]$ was one of the first polymers to be prepared in the laboratory. The use of numerous epoxides (as monomers, or for the synthesis of other epoxide monomers or prepolymers) in polymer synthesis has been quite extensive, and these polyethers have been prepared by both cationic and anionic polymerization techniques. The molecular weights of resulting polyethers generally range from 500 to 10,000, usually because of interfering chain-transfer reactions. An especially interesting synthesis is the opening of propylene oxide to give optically active poly(propylene oxides). Epoxides are used in the preparation of copolymers. *See* COPOLYMER; FURAN; POLYETHER RESINS; POLYMERIZATION.

Epoxides are promoters; that is, they are easy to polymerize, and their polymerization causes more difficult cyclic ethers such as tetrahydrofuran to polymerize. Epichlorohydrin has also been used for the preparation of prepolymers by condensations with the sodium salt of bisphenols. The epichlorohydrin/bisphenol A and other epoxide prepolymers have been used for the preparation of epoxy resins (thermosets), which are cross-linked (3D) polymers, linear polymers that are connected by cross-linked molecular units. They are insoluble, infusible, and intractable, and sometimes are called epoxy resins. For example, a low-molecular-weight epoxy prepolymer can be treated and cross-linked with a polyamine or smaller molecules such as diethylenetriamine, or cross-linked by the addition of carboxylic acid anhydrides, such as maleic anhydride. Other polymers can also be incorporated or participate in the cross-linking process, and include phenolics, ureas, and melamines. The products resulting from these prepolymer techniques are used for surface coating materials, molding, pipes, laminating, repair of damaged automobile bodies, manufacture of articles reinforced with glass fibers, durable and tough epoxy resin adhesives (glues), and many other applications.

Naturally occurring polymers such as rubber can be treated with hydrogen peroxide or peracetic acid to form an epoxidized polymer. Also, another naturally occurring polymer, cellulose, can be transformed into hydroxyethylcellulose by treating alkali cellulose with ethylene oxide. *See* CELLULOSE; RUBBER.

Related compounds. The three-membered sulfur analogs of epoxides are episulfides (also known as alkylene sulfides or thiiranes). Episulfides can be prepared from the treatment of diazoalkanes with sulfur, condensation of thioketones with sulfur ylides, and the reaction of epoxides with phosphene sulfide or thiourea/titanium tetraisopropoxide. These cyclic thioethers have been used as monomers, especially in copolymer preparations. There is some evidence for the preparation of several episelenides (replacement of the oxygen of epoxide with selenium; also known as ethylene selenides or seleniranes); but heating solutions of these less stable compounds results in the elimination of selenium and yields the alkene.

Charles F. Beam

Bibliography. T. L. Gilchrist, *Heterocyclic Chemistry*, 3d ed., 1997; R. B. Seymour and C. E. Carraber, Jr., *Seymour and Carraber's Polymer Chemistry*, 5th ed., 2000; M. Smith and J. March, *March's Advanced Organic Chemistry*, 5th ed., 2000.

Epsomite

A mineral with the chemical composition $\text{MgSO}_4 \cdot 7\text{H}_2\text{O}$. Epsomite, or epsom salt, occurs in clear, needlelike, orthorhombic crystals. More commonly it is massive or fibrous, although crystals from salt lakes on Kruger Mountain near Orville, Washington, are reported to be several feet long. Fracture is conchoidal. Luster varies from vitreous to silky. Hardness is 2–2.5 on Mohs scale and specific gravity is 1.68. The mineral has a salty bitter taste and is soluble in water.

Epsomite is found as a capillary coating in limestone caves and in coal or metal mine galleries. It loses readily some of its water of crystallization in dry air. It is also found associated with gypsum and in thin layers in salt deposits of oceanic origin or from salt lakes. *See* MAGNESIUM. Edward C. T. Chao

Epstein-Barr virus

An antigenically distinct member of the herpesvirus group of viruses, whose genome is deoxyribonucleic acid (DNA). Under electron-microscopic observation, the mature virus is indistinguishable in size and structure from the other human herpesviruses; that is, it has a nucleocapsid about 110 nanometers in diameter, and when an envelope is present, the particle diameter is about 180 nm. The 78-nm central core, which is often observed in other herpesviruses, is not clearly visible in Epstein-Barr (EB) particles within infected cells, and few of the intracellular particles possess an envelope. Only about 10% of the extracellular virus particles are enveloped.

EB virus is the cause of one benign disease (infectious mononucleosis), and is associated with certain types of cancer; however, the great majority of EB virus infections are clinically inapparent.

The virus was detected initially by electron microscopy in a small proportion of cells in continuous lymphoblastoid cell lines derived from Burkitt's lymphoma (but particles have not been seen in cells of the tumor itself). This lymphoma is a tumor indigenous to children in central Africa. The virus also has been detected in cell lines derived from nasopharyngeal carcinomas, a type of cancer found with high frequency in persons from southern China. The virus is found in peripheral blood leukocytes from normal individuals and from patients with infectious mononucleosis. *See* INFECTIOUS MONONUCLEOSIS; LYMPHOMA.

With some strains of EB virus, attempts to infect human adult and cord blood lymphocytes and marmoset lymphocytes in cultures have resulted in the

establishment of continuous cell lines, suggesting that these cells have been transformed (that is, “immortalized”) by the virus.

The high prevalence and high titers of EB antibody (detectable by immunofluorescence, complement-fixation, and gel-diffusion tests) in patients with Burkitt’s lymphoma and those with postnasal carcinoma led to the assumption that the association might be etiologic. However, subsequent seroepidemiological surveys showed that infection with EB virus in normal populations is widespread not only in Africa and in Asia but also all over the world. In some areas, including urban areas of the United States, about 50% of children 1 year old, 80–90% of children over age 4, and 90% of adults have antibody to EB virus. Although the mechanism of virus transmission is unknown, it could take place through oropharyngeal excretion, which now has been found to be very common, especially in individuals with overt infectious mononucleosis but also in persons without symptoms. Furthermore, EB virus has an extreme predilection for cells of lymphoid origin, and hence could quite possibly be merely a passenger virus in the lymphoma cells. The question of whether EB virus is etiologically related to lymphoma or to postnasal carcinoma, or to both, remains to be answered. Indeed, other viruses (for example, herpes simplex virus and reovirus) have also been isolated from Burkitt’s lymphomas. On the other hand, the possibility remains that EB virus is the agent responsible for the induction of these human cancers, or is one of the cofactors causing the malignancies.

The techniques for searching for viruses in human cancers have become quite sophisticated, but technical problems still make it difficult to determine whether a virus is merely a passenger inhabiting the tumor cells. Even more difficult is the identification of traces of a virus which may have caused the cancer but which is no longer present in its complete or infective form in the cancer it initiated. Such traces may include only selected portions of the viral genome residing in the cancer cell and molecular probes have been developed for this search. Since EB virus now has been shown capable of inducing lymphomas in the cottontop marmoset (a South American monkey), studies of this virus–primate host system should help to clarify the relationship between EB virus and human disease.

If EB virus is indeed confirmed as having a role in the development of human malignancies, then one major question to be resolved is how a virus so ubiquitous can be involved in so wide a variety of responses—ranging from asymptomatic infection through infectious mononucleosis to production of cancers. However, it should be recalled that many virus infections (for example, polio virus, hepatitis viruses, certain of the arboviral encephalitides) have a wide spectrum of outcomes, ranging from inapparent infection to severe syndromes. See ANIMAL VIRUS; ONCOLOGY.

Joseph L. Melnick

Bibliography. D. V. Ablashi et al. (eds.), *Epstein-Barr Virus and Human Disease*, 1990; *Peak Immunity*:

How to Fight Epstein-Barr Virus, Candida, Herpes Simplex, and Other Immuno-Depressive Disorders and Win, Medical Sciences, General Medical Series, 1989.

Equalizer

An electronic filter that modifies the frequency response (amplitude and phase versus frequency) of a system for a specific purpose. While filters in some applications perform the conceptually simple operations of rejecting specific bands of frequencies and passing other bands of frequencies, equalizers typically realize a more complicated frequency response in which the amplitude response varies continuously with frequency, amplifying some frequencies and attenuating others. An equalizer may have a response fixed in time or may be automatically and continuously adjusted. However, its frequency response is usually matched to some external physical medium, such as an acoustic path or communication channel, and thus inherently needs to be adjustable.

Equalizers can be used in many applications. In music and sound reproduction, equalizers can compensate for artifacts of the electrical-to-sound conversion or for unwanted characteristics of the acoustic environment such as sound reflections or absorption. Sound-recording and sonar systems can use equalizers to reduce unwanted interference. Most analog recording and playback devices, such as audio and video tape recorders, incorporate equalizers to compensate for the undesirable aspects of the recording medium, such as high-frequency roll-off, as well as to reduce noise and maximize dynamic range. See ELECTRICAL INTERFERENCE; ELECTRICAL NOISE; SONAR; SOUND RECORDING.

An important application of equalization is to enhance the performance of systems that communicate or record digital signals (streams of bits). The transmission of digital (as opposed to analog) information is rapidly becoming dominant, and equalization is a critical component of almost all such systems. All communications and recording systems utilize a physical medium, such as wires; coaxial cables; radio, acoustic, or optical-fiber waveguides; or magnetic and optical recording media. These media cause distortion; that is, the output signal is different from the input signal. For example, on radio or acoustic channels there are often multiple paths from transmitter to receiver, each having a slightly different delay and superimposed at the receiver. An equalizer is an electrical device that compensates for this distortion, reversing the effect of the channel and returning a waveform approximating the input signal. The channel output signal in response to a particular input signal (. . . , 0, 0, 1, 0, 0, . . .) may differ from the input, but the equalizer output reproduces the channel input, at least to close approximation (**Fig. 1**). See DISTORTION (ELECTRONIC CIRCUITS).

If the characteristics of the channel are well known, the equalizer can be fixed or nonadaptive. More commonly, the detailed characteristics of a

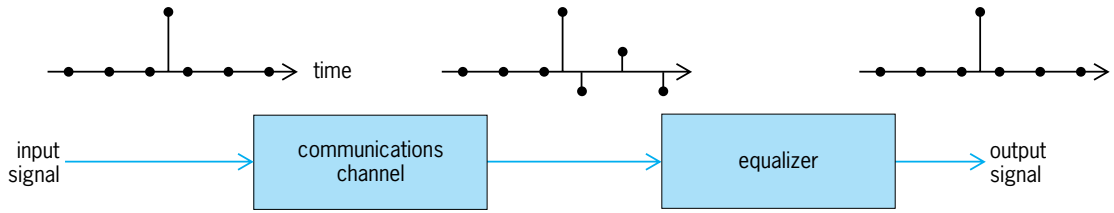


Fig. 1. Communications channel with an equalizer placed at the output. The equalizer recovers an accurate replica of the channel input signal.

channel are not known in advance. For example, an equalizer may be required to compensate for any length of wire, from very short up to a maximum (such as the connection from a telephone central office to the subscriber's telephone instrument). In other cases, the channel may be varying with time, as is characteristic of the radio channel from a fixed transmitter to a moving vehicle. (On a radio channel, the signal arriving at the vehicle may represent multiple superimposed paths, each path corresponding to a reflection from a different building or mountain and having a different delay. The net effect depends on the vehicle location, and hence changes as the vehicle moves.) An adaptive equalizer does not need to know the channel in advance, but is nevertheless able to adjust itself to compensate. Adaptive equalizers are important for achieving high bit rates in telephone computer modems, and also for digital communications over radio channels. See DATA COMMUNICATIONS; MODEM.

Prior to about 1970, equalizers were usually realized by continuous-time filters constructed from interconnected resistors, capacitors, and inductors. Using modern electronic technologies, equalizers are almost always realized in discrete time. That is, the continuous-time channel output, say $x(t)$ where t is time, is first sampled at uniformly spaced time intervals, yielding a sampled signal given by Eq. (1),

$$x_k = x(kT) \tag{1}$$

where T is the interval between samples. As long as T is sufficiently small, no essential information is lost. The discrete-time signal x_k , consisting of samples of the input signal, is input to the equalizer. Time is represented by the discrete variable k rather than the continuous variable t .

Structure. A common internal structure for a discrete-time equalizer is the transversal filter (Fig. 2). The output at any given time k is a linear combination of the current and a set of past input samples, given mathematically by Eq. (2) where the

$$y_k = \sum_{m=0}^{N-1} a_m x_{k-m} \tag{2}$$

current input and output samples are x_k and y_k respectively, and the set of N weighting factors, a_m , $1 \leq m \leq N - 1$ are known as the filter coefficients. The coefficients are chosen or adjusted to control the equalization. A transversal filter can be thought of as introducing its own distortion controlled by the filter coefficients, and this distortion is chosen to re-

verse the distortion on the channel. This is usually possible as long as N is sufficiently large.

Modes of adaptation. In order to make a transversal filter adaptive and, in particular, to determine the channel distortion by looking at the output of the channel alone, the equalizer is normally operated in two modes: the training mode and the tracking mode. In the training mode, the channel characteristics are learned for the first time (for example, with a telephone computer modem at the beginning of a phone call), with very little prior knowledge of the channel necessary. In the tracking mode, the characteristics of the channel are followed, assuming those characteristics do not change too quickly. If the channel characteristics are fixed, the tracking mode is unnecessary.

In the training mode (Fig. 3a), for purposes of learning the channel, a specially chosen training signal u_k is applied to the channel input, and is presumed to be known to the equalizer by prior agreement. (Normal communications over the channel is precluded during the training period.) The difference between the equalizer output and the training signal, known as the error signal e_k , is used to adapt the equalizer. The filter coefficients are adjusted so as to force this error signal, given by Eq. (3), to be

$$e_k = u_k - \sum_{m=0}^{N-1} a_m x_{k-m} \tag{3}$$

small. The e_k being consistently nearly zero is an indication that the equalizer is correctly compensating for the channel distortion. The channel characteristics are not learned directly; rather, the filter coefficients that cause the equalizer to reverse the channel distortion are determined.

The tracking mode (Fig. 3b) proceeds while actual communication is occurring, and hence it uses whatever channel input signal is to be communicated, that signal being unknown to the equalizer. A tracking mode is possible only for digital communications, where the transmitted signal u_k assumes only a

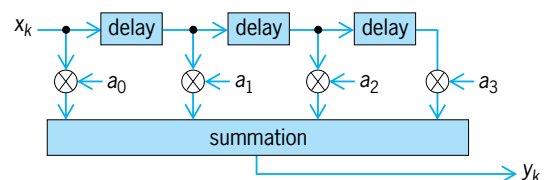


Fig. 2. Transversal filter, which forms a weighted summation of the current and past input signal samples (shown for $N = 4$ coefficients).

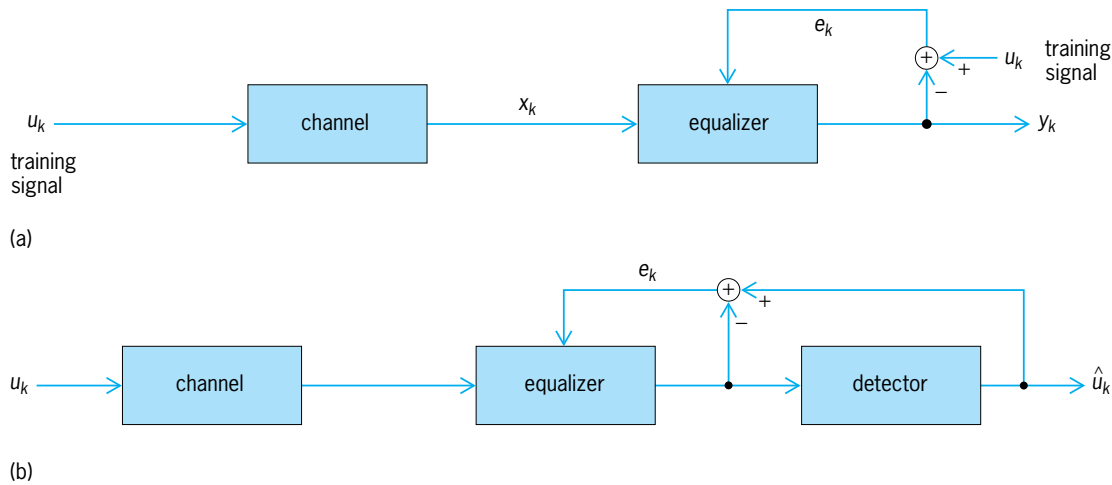


Fig. 3. Equalizer modes of operation. (a) Training mode. The channel input signal \hat{u}_k is known to the equalizer, the adaptation of which has available the error signal e_k . (b) Tracking mode. The error signal is formed by using a detected version \hat{u}_k of the channel input signal \hat{u}_k .

finite and discrete set of values at any given instant of time. (For example, input signal values may carry only one bit of information. That one bit can assume one of two values, 0 or 1, which can be represented by one of two signal values, say $u_k = +1$ or $u_k = -1$.) Tracking assumes that the equalization is always accurate enough that any residual distortion caused by incorrect coefficients does not cause errors in the detection of the input signal values. Thus, a detector placed at the output of the equalizer forms an accurate estimate \hat{u}_k of the input signal u_k , giving a reference for equalizer adaptation. As in training, an error signal is formed, replacing the actual input signal u_k by the detected signal \hat{u}_k . Correct operation depends on $\hat{u}_k = u_k$, at least most of the time. It is quite possible for this scheme to fail if the channel changes too quickly, since the uncorrected distortion may become so large that the detector does not work properly. However, adaptation works properly in the presence of occasional detection errors (such as those caused by random noise) since the adaptation averages out these errors.

Adaptation algorithms. Adaptation chooses the filter coefficients a_m , $1 \leq m \leq N - 1$, to minimize, in some sense, the squared error e_k^2 . The goal is not to minimize this squared error at a particular time k , but to use a criterion that ensures that the squared error is consistently small. For tracking, the adaptation should have a finite memory; that is, it should not depend on the far past (when the channel was different) but the recent past. An example is given by the quantity in Eq. (4) for

$$\xi_k = \sum_{m=0}^{\infty} \beta^m e_{k-m}^2 \quad (4)$$

$0 < \beta < 1$. Since β^m is larger for small m and falls off rapidly for larger m , if ξ_k is minimized at time k , the recent channel characteristics are weighted much more heavily. It is a straightforward exercise in calculus to find the coefficients a_m , $1 \leq m \leq N - 1$, that minimize ξ_k . A slowly changing channel can be

tracked by re-solving for the coefficients at each time instant k .

There are a number of available algorithms for adapting the transversal filter coefficients, differing in the criterion that is minimized and the algorithmic formulation of the minimization. If greater computational complexity is acceptable, a greater rate of change in the channel characteristics can be tracked. Adaptive equalizers are usually implemented in a software program running on a special-purpose microprocessor known as a programmable digital signal processor. In an application such as digital communication with a moving vehicle by radio, where the channel characteristics can change rapidly, the adaptation algorithms are quite sophisticated and require a very high rate of computation. See ALGORITHM; DIGITAL FILTER; ELECTRIC FILTER; MICROPROCESSOR; OPTIMIZATION.

David G. Messerschmitt

Bibliography. C. Cowan and P. Grant, *Adaptive Filters*, 1985; K. Feher, (ed.) *Advanced Digital Communications Systems and Signal Processing Techniques*, 1987; E. Lee and D. Messerschmitt, *Digital Communication*, 1993; B. Mulgrew and C. Cowan, *Adaptive Filters and Equalizers*, 1988; B. Widrow and S. Stearns, *Adaptive Signal Processing*, 1985.

Equation of continuity

An equation of continuity appears in many branches of physics. In dynamic field theory, it is essentially a statement that charge is conserved or that the rate of increase of charge in any region equals the current \mathbf{i} flowing into that region. If v is the volume enclosed by the surface S , this statement may be expressed in integral form, as given in Eq. (1), where ρ is the

$$\int_S \mathbf{i} \cdot \mathbf{n} \, dS = \int_v \nabla \cdot \mathbf{i} \, dv = -\frac{\partial}{\partial t} \int_v \rho \, dv \quad (1)$$

charge density and \mathbf{n} is a unit vector normal to S . The second integral comes from the first by Gauss' divergence theorem. When currents and charges are

confined to a surface S bounded by the curve s , this becomes Eq. (2), where θ is the angle between ds and

$$\oint i \sin \theta ds = -\frac{\partial}{\partial t} \int_S \rho dS \quad (2)$$

\mathbf{i} which is directed into the area S . This states that the current crossing the boundary of an area equals the rate of increase of charge in the area. The volume in Eq. (1) is arbitrary, so the integrands in the last two integrals are equal. This leads to the differential form in Eq. (3), where i_x , i_y , and i_z are the rectangular

$$\nabla \cdot \mathbf{i} = \frac{\partial i_x}{\partial x} + \frac{\partial i_y}{\partial y} + \frac{\partial i_z}{\partial z} = -\frac{\partial \rho}{\partial t} \quad (3)$$

lar components of \mathbf{i} . Maxwell's equations satisfy Eq. (3). For application to the motion of charged particles, Eq. (3) would be written as Eq. (4). When \mathbf{i} and

$$\nabla \cdot \rho \mathbf{v} = \rho \left(\frac{\partial v_x}{\partial x} + \frac{\partial v_y}{\partial y} + \frac{\partial v_z}{\partial z} \right) = -\frac{\partial \rho}{\partial t} \quad (4)$$

ρ vary sinusoidally with time, they may be written as phasors, which are complex numbers such that when multiplied by $e^{j\omega t}$, the real part of the product gives the amplitude, phase, and time dependence. The only change in the preceding equations when \mathbf{i} and ρ are phasor quantities is the replacement of $\partial/\partial t$ by $j\omega$. See FLUID-FLOW PRINCIPLES; GAUSS' THEOREM; MAXWELL'S EQUATIONS. William R. Smythe

Equation of time

The annual, cyclic variation between mean solar time shown on uniformly running clocks and apparent solar time displayed on sundials.

In the course of the Sun's daily east-to-west transit of the sky, the Sun crosses the meridian, an imaginary line running from north to south that passes overhead and divides the sky into equal halves. An observer in the middle of a time zone generally thinks of noon as being the moment that the Sun reaches the meridian. This event, however, corresponds to noon recorded by mechanical or electronic clocks on only four dates each year (approximately April 16, June 14, September 1, and December 25). On all other dates the Sun reaches the meridian either

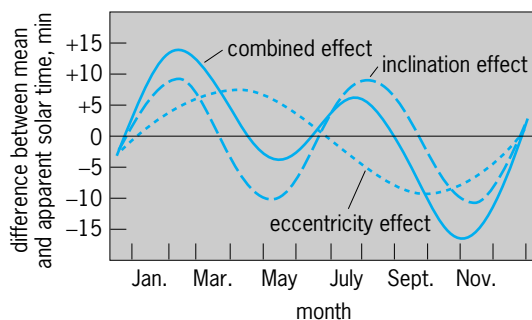


Fig. 1. Graph of the equation of time, showing how the equation results from combining the effects of the inclination of the Earth's axis and the eccentricity of its orbit. (After B. M. Oliver, *The shape of the analemma*, *Sky Telesc.*, 44:20-22, July 1972)

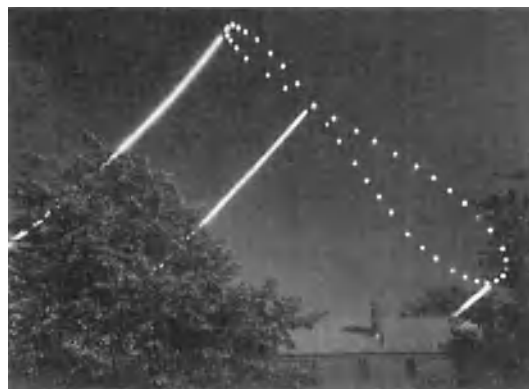


Fig. 2. Analemma recorded in a multiple-exposure, year-long photograph obtained with an east-facing, permanently mounted camera that made an exposure approximately once a week at exactly 8:30 a.m. Eastern Standard Time. On three occasions the shutter was left open from dawn until shortly before 8:30 to show the Sun's path across the sky. All images of the Sun were made through a special dense filter. The foreground was recorded without the filter one afternoon when the Sun was in the western sky and roughly behind the camera (as indicated by the shadow of the chimney). (Photograph by Dennis di Cicco)

early or late, with the extremes being 16.3 min early around November 3 and 14.3 min late around February 12. This difference is the equation of time, and results from the combined effects of Earth's axis of rotation being tipped 23° relative to Earth's orbital plane and the elliptical rather than circular shape of the orbit. See MERIDIAN.

The Earth's orbital motion causes the Sun to move eastward against the background stars along a path called the ecliptic. This motion is about 1° each day. Although the Earth turns once on its axis every $23^h 56^m 4^s$, it must turn slightly more in order for the eastward-moving Sun to return to the meridian. Thus, the length of time from one meridian passage of the Sun to the next averages 24 h over the course of a year.

If the Earth had a perfectly circular orbit and its axis of rotation was perpendicular to the orbital plane, the Sun's motion would be along the celestial equator and uniform throughout the year. It is this idealized case of a mean Sun that shows the mean solar time kept by mechanical and electronic clocks and watches. However, the elliptical orbit brings the Earth closest to the Sun in January. This proximity, as J. Kepler discovered in the sixteenth century, causes the Earth to move more rapidly in its orbit than in July, when the Earth is farthest from the Sun. The changing orbital speed varies the Sun's apparent rate of motion along the ecliptic and is in part responsible for the equation of time. See CELESTIAL MECHANICS; KEPLER'S LAWS.

Slightly more influential is Earth's tipped axis, which varies the Sun's position north and south of the celestial equator according to the season. Around the time of the spring and autumn equinoxes the Sun moves at a steep angle relative to the celestial equator. Its daily motion projected onto the equator is less than at the solstices, when the Sun travels parallel to

the equator. This situation also creates a departure between the Sun's actual position and that of a mean Sun moving uniformly along the celestial equator. The effects of the inclination of the Earth's axis and the eccentricity of its orbit are combined (Fig. 1).

If the Sun's daily position is considered relative to the meridian at the moment a clock reads noon, from the equation of time it is seen that the Sun will sometimes be west of the meridian (early) and sometimes east of it (late). The coupling of this variation with the north-south movement of the Sun along the ecliptic causes the Sun to mark out a large figure-eight known as the analemma (Fig. 2). This pattern is sometimes shown on the tropical zone of Earth globes (usually in the Pacific Ocean). The analemma shows the latitude at which the Sun passes directly overhead on any given date and the equation of time. See EARTH ROTATION AND ORBITAL MOTION; TIME. Dennis Di Cicco

Bibliography. B. M. Oliver, The shape of the analemma, *Sky Telesc.*, 44:20-22, July 1972.

Equations, theory of

The branch of mathematics concerned with finding facts concerning the roots of algebraic equations and finding methods for obtaining them. The most important type of algebraic equation is the polynomial equation in one unknown which is an expression of the form $f(x) = a_n x^n + a_{n-1} x^{n-1} + \dots + a_1 x + a_0 = 0$, where x is called the unknown, or variable; n is a positive whole number; and the a_i , with $i = 0, 1, \dots, n$, are constants, or fixed numbers, called coefficients of the equation. The left member of the equation is called a polynomial in one variable of degree n . A root of such an equation is a number which, when substituted for the variable x , makes the left member zero. For example, 3 is a root of the equation $x^3 + 2x^2 - 13x - 6 = 0$. In addition, systems of equations in one or more variables are considered, and here the problem is to find values for the variables which simultaneously satisfy each equation of the system. See LINEAR SYSTEMS OF EQUATIONS; POLYNOMIAL SYSTEMS OF EQUATIONS.

The topics covered in a systematic study of the theory of equations can be placed in the following principal subdivisions: properties of a polynomial which do not depend on the particular number system containing the coefficients of the polynomial; factorization of polynomials; equations with coefficients which are rational, real, or complex numbers; determination of bounds for real roots, and systematic methods for approximating real roots of equations; the solution of quadratic, cubic, and quartic equations by radicals.

This article is limited to polynomials and to equations which have rational, real, or complex numbers as coefficients. Each of these number systems constitutes a number field.

Let $f(x) = a_n x^n + a_{n-1} x^{n-1} + \dots + a_1 x + a_0$ and $g(x) = b_m x^m + b_{m-1} x^{m-1} + \dots + b_1 x + b_0 \neq 0$ be polynomials with coefficients in a number field F . Then the division algorithm, which is

a formal statement of the ordinary division process for polynomials, states that there exist unique polynomials $q(x)$ and $r(x)$, with coefficients in F , such that $f(x) = q(x)g(x) + r(x)$, where either $r(x) = 0$ or $r(x)$ has degree less than $g(x)$. If $r(x) = 0$, then $g(x)$ divides $f(x)$. Immediate corollaries of the division algorithm are the remainder and factor theorems. If $g(x)$ is the linear polynomial $(x - a)$, then $f(x) = q(x)(x - a) + r$, where r is a constant. Then, substituting $x = a$, the result is $f(a) = r$. This is the remainder theorem. If $f(a) = r = 0$, then a is a root of the equation $f(x) = 0$, giving the factor theorem, which states that a is a root of $f(x) = 0$ if, and only if, $(x - a)$ is a factor of $f(x)$. It follows from the factor theorem that the equation $f(x) = 0$ has at most n roots in F , where $f(x)$ is of degree n .

The greatest common divisor (gcd) of two polynomials $f(x)$ and $g(x)$ is a polynomial $d(x)$ which divides both $f(x)$ and $g(x)$, and is divisible by every other polynomial which also divides $f(x)$ and $g(x)$. A process called the euclidean algorithm, based on successive application of the division algorithm, is used to find the gcd $d(x)$ having coefficients in the smallest number field that contains the coefficients of $f(x)$ and $g(x)$. The polynomials $f(x)$ and $g(x)$ are called relatively prime if their gcd is a constant.

The factorization of a polynomial $f(x)$ depends on the particular number field F under consideration. For example, the polynomial $x^5 - 1/2x^4 - x^3 + 1/2x^2 - 2x + 1$ having coefficients which are rational numbers, factors as $(x^2 + 1)(x^2 - 2)(x - 1/2)$ over the rational numbers; as $(x^2 + 1)(x - \sqrt{2})(x + \sqrt{2})(x - 1/2)$ over the real number; and as $(x + i)(x - i)(x - \sqrt{2})(x + \sqrt{2})(x - 1/2)$ over the complex numbers. A polynomial $f(x)$ with coefficients in F is irreducible over F if it cannot be expressed as a product of polynomials of lower degree. Every polynomial can be expressed in essentially one way as a product of irreducible factors, although there is no general algorithm which enables one to obtain this expression. For polynomials with rational coefficients, there is such an algorithm, devised by L. Kronecker. There are methods for finding the repeated factors of a polynomial which are often useful as a first step in factoring a polynomial.

If a polynomial equation $f(x) = 0$ with rational coefficients is multiplied by a suitable whole number, an equation with whole number coefficients is obtained. If r/s , a rational number in lowest terms, is a root of such an equation, then r divides the constant term and s divides the leading coefficient. Hence, the rational roots of $f(x) = 0$ can be found by a finite number of trials.

The fundamental theorem of algebra states that a polynomial equation with complex coefficients has a complex root. From this it follows immediately that a polynomial of degree n with complex coefficients factors into n linear factors over the complex numbers.

If $f(x) = 0$ has real coefficients, then it can be shown that, if the complex number $a + bi$ is a root of $f(x)$, the conjugate complex number $a - bi$ is also

a root. Thus, the real, irreducible factors of a polynomial with real coefficients are linear or quadratic. A further consequence is that, if $f(x)$ with real coefficients has odd degree, then $f(x) = 0$ has at least one real root.

The property of a polynomial $f(x)$ with real coefficients, which is basic for the study of the equation $f(x) = 0$, is that $f(x)$ defines a continuous real function. The location principle which follows from this states that if there exist real numbers $a < b$ such that $f(a)$ and $f(b)$ have opposite signs, then $f(x) = 0$ has a real root r such that $a < r < b$. The location principle is used to isolate the real roots and is fundamental in systematic schemes such as Horner's method and Newton's method for approximating the real roots. A numerical method, known as Graeffe's method, can be used to approximate both the real and the complex roots.

Important results used in finding the real roots of $f(x) = 0$ with real coefficients include the following: Rolles' theorem from differential calculus, which states that, between two real roots of $f(x) = 0$, there is at least one real root of the derivative $f'(x) = 0$; Sturm's theorem, which gives an exact count of the number of real roots in an interval between two real numbers $a < b$; and Descartes' rule of signs, which states that the number of positive roots of $f(x) = 0$ equals the number of variations in sign of the coefficients of $f(x)$ minus a nonnegative even number. The negative roots of $f(x) = 0$ are positive roots of $f(-x) = 0$.

The results on the bounds for the real roots of $f(x) = 0$ are based on the fact that, for sufficiently large values of x , the sign of $f(x)$ is the same as the sign of the leading term $a_n x^n$ of $f(x)$. In particular, $f(x) = 0$ has no real roots for $|x| \geq |a_k/a_0| + 1$, where a_k is the coefficient of $f(x)$ having the greatest numerical value.

Polynomial equations of degree 2, 3, and 4 are solvable by radicals. This means that there are formulas which give the roots in terms of the coefficients of the equation and that these formulas involve only the rational operations and the operation of extraction of roots. The principal methods for the solution of the cubic and quartic equations were devised by J. Cardan and L. Ferrari, respectively. By use of the Galois theory of equations, it can be proved that for $n > 4$ there cannot exist a formula involving only rational operations and root extractions for expressing the roots of every polynomial equation of degree n in terms of the coefficients. See DETERMINANT; MATRIX THEORY.

Ross A. Beaumont

Bibliography. R. A. Barnett and T. J. Kearns, *Elementary Algebra: Structure and Use*, 6th ed., 1994; D. Dobbs and R. A. Hanks, *Modern Course on the Theory of Equations*, 2d ed., 1992.

Equator

The great circle around the Earth, equally distant from the North and South poles, which divides the Earth into Northern and Southern hemispheres. It is

the greatest circumference of the Earth because of centrifugal force from rotation, and resultant flattening of the polar areas.

The Earth's rotational axis is vertical to the plane of the Equator, and because the inclination of the axis is 66.5° (66.55°) from the plane of the ecliptic, the plane of the Equator is always inclined 23.5° from the ecliptic.

At noon on the days of the vernal and autumnal equinoxes (March 21 and September 23), the Sun on the Equator is 90° above the horizon. The lowest angle of the Sun is 66.5° . This occurs at noon, at the winter and summer solstices (June 21 and December 22), when the Sun is vertical at the Tropic of Cancer and Tropic of Capricorn, respectively. Days and nights are always equal because the plane of the ecliptic intersects the equatorial plane at the Earth's center. Consequently, the Sun at the Equator rises and sets at approximately 6 A.M. and 6 P.M. throughout the year.

The celestial equator in astronomy is equally distant from the celestial poles and is the great circle in which the plane of the terrestrial Equator intersects the celestial sphere. See ASTRONOMICAL COORDINATE SYSTEMS; MATHEMATICAL GEOGRAPHY.

Van H. English

Bibliography. A. H. Strahler and A. Strahler, *Introducing Physical Geography*, 4th ed., 2005.

Equatorial currents

Ocean currents near the Equator. The westward trade winds that prevail over the tropical Atlantic and Pacific oceans drive complex oceanic circulations characterized by alternating bands of eastward and westward currents (**Fig. 1**). The intense currents are confined to the surface layers of the ocean; below a depth of approximately 100 m (330 ft) the temperature is much lower, and the speed of ocean currents is much slower. The westward surface currents tend to be divergent—they are associated with a parting of the surface waters—and therefore entrain cold water from below. The water temperature rises as the currents flow westward, so that temperatures are low in the east and high in the west, except between 3° and 10°N where eastward surface currents create a band of warm water across the Pacific and Atlantic oceans. The distinctive sea surface temperature pattern in which surface waters are warm in the west and cold in the east, except for the warm band just north of the Equator, reflects the oceanic circulation (**Fig. 2**). A dramatic change in this pattern every few years during El Niño episodes, when the temperature of the eastern tropical Pacific Ocean rises, is associated with an intensification of the eastward currents and a weakening (sometimes reversal) of the westward currents. See EL NIÑO.

South Equatorial Current. This current flows westward in the upper ocean, has its northern boundary at approximately 3°N , and attains speeds in excess of 1 m/s (3.3 ft/s) near the Equator. It is directly driven by the westward trade winds and has its

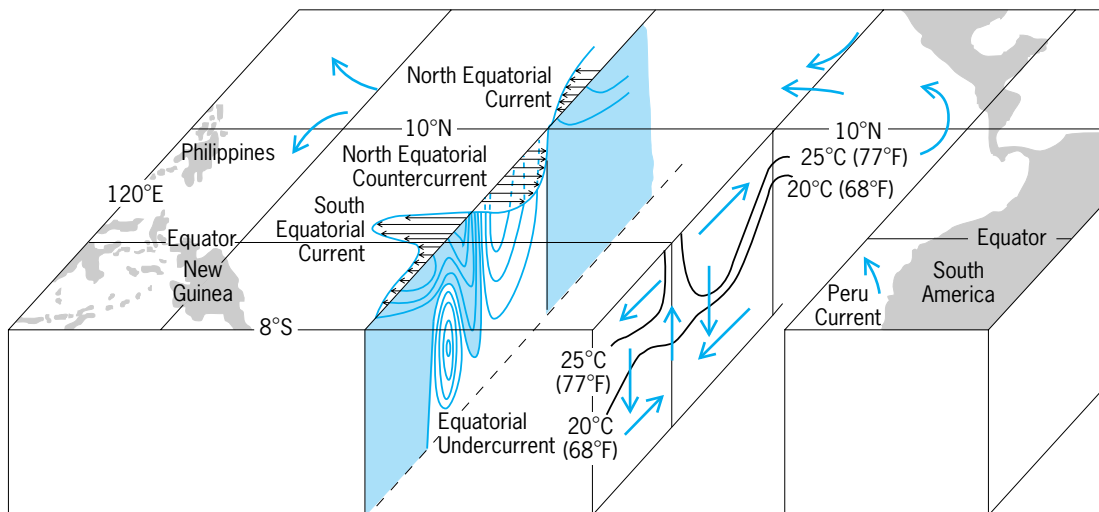


Fig. 1. Circulation of the tropical Pacific Ocean. The shaded areas indicate regions of westward flow. The principal eastward currents are the Equatorial Undercurrent centered on the Equator, and the North Equatorial Countercurrent to the north of the Equator. The heavy arrows in the vertical plane show intense upwelling at the Equator, poleward flow in the surface layers, and sinking motion off the Equator, sustained by equatorward flow at depth of 100 m (330 ft) approximately. This motion affects the thermal field as is evident from the two isotherms.

origins in the cold, northwestward-flowing Peruvian coastal current. Because the Coriolis force deflects water parcels to their right in the Northern Hemisphere and to their left in the Southern Hemisphere, this current is divergent at the Equator. As a consequence, cold water from below wells up along the Equator. See CORIOLIS ACCELERATION; UPWELLING.

North Equatorial Countercurrent. This current flows eastward immediately to the north of the South Equatorial Current. The boundary between these two currents is a sharp thermal front that is clearly evident in satellite photographs. The front can literally be a green line, hundreds of yards wide, because of the abundance of phytoplankton. (It frequently has westward-propagating undulations with a wavelength on the order of 1000 km or 600 mi. These are instability waves associated with the enormous shear of the two currents.) This current, which is counter to the wind, is driven by the torque (curl) that the wind exerts on the ocean. To its north is a colder westward current known as the North Equatorial Current. See OCEAN WAVES; PHYTOPLANKTON.

Equatorial Undercurrent. This current, which in the Pacific Ocean was originally known as the Cromwell Current, is an intense, narrow, eastward, subsurface jet that flows precisely along the Equator across

the width of the Pacific. Its core, where speeds can be in excess of 1.5 m/s (5 ft/s), is at an approximate depth of 100 m (330 ft); its width is approximately 200 km (120 mi). A similar current exists in the Atlantic Ocean. In the Indian Ocean it is often present along the Equator, in the western part of the basin during March and April when westward winds prevail over that region. Such winds (including the trade winds over the Pacific and Atlantic oceans) pile up warm surface waters in the west while exposing cold waters to the surface in the east. Isotherms along the Equator therefore slope downward to the west (Fig. 1). Thus a column of water has a higher average temperature in the west than in the east. Furthermore, if horizontal pressure gradients near the ocean floor are negligible—the absence of strong currents near the ocean floor suggests that such is indeed the case—then the sea surface must have a higher elevation in the west than in the east. (The difference in height is on the order of 1 m or 3 ft.) Hence, westward winds, by driving the surface waters westward, also maintain this slope of the sea surface and thus create an eastward pressure. Off the Equator, the subsurface eastward flow driven by the pressure force is deflected equatorward, in both hemispheres, by the Coriolis force. At the Equator, where the

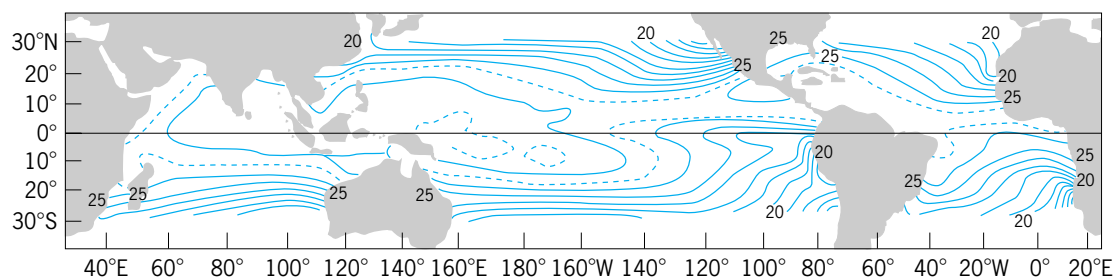


Fig. 2. Isotherms at the ocean surface. The contour interval is 1°C (1.8°F). Broken-line contours are 27°C (81°F) and 29°C (84°F).

Coriolis force vanishes, the convergent water flows in the direction of the pressure force and becomes the Equatorial Undercurrent. *See* ATLANTIC OCEAN; INDIAN OCEAN; PACIFIC OCEAN.

Numerical models that simulate the oceans are capable of reproducing the observed circulation when it is forced with the observed winds. From their results it is then possible to infer aspects of the circulation that are difficult to measure, for example, the extremely small vertical component of the currents. The vertical and meridional flows link the various eastward and westward equatorial currents (Fig. 1). A water parcel that flows eastward in the Equatorial Undercurrent gradually rises to the surface until it joins the westward South Equatorial Current, a divergent surface current in which the parcel flows poleward. If it is in the Northern Hemisphere, the parcel sinks in the neighborhood of 3°N, near the thermal front mentioned earlier. At depth, it flows back toward the Equatorial Undercurrent. *See* SIMULATION; WIND STRESS.

Time dependence of the currents. The winds in the tropics are subject to considerable seasonal and interannual fluctuations that force variations in the oceanic circulations and, in particular, in patterns of sea surface temperature. Changes in those patterns have a strong influence on the atmospheric circulation, and hence they can cause climate fluctuations such as El Niño. This circular argument—changes in sea surface temperature patterns both cause and are caused by changes in atmospheric conditions, especially the winds—implies that interactions between the ocean and atmosphere are the basis of climate variability. There is an asymmetry in the roles of the ocean and atmosphere because, whereas the atmosphere adjusts rapidly—that is, within a week or two—to a change in sea surface temperatures, the ocean has far more inertia and takes many months to adjust to a change in the winds. It is this memorylike character of the ocean that sets the time scale for the recurrence of El Niño phenomena.

The sudden onset of westerly winds over the equatorial Indian Ocean, occurring twice a year at the time of the equinoxes, provides valuable information about the memory of the oceans, and how equatorial oceans adjust to a sudden change in the winds. *See* EQUINOX.

The winds not only generate oceanic currents but also excite waves. Those waves propagate back and forth across the ocean basin and are instrumental in bringing the ocean into a state of equilibrium. The time they take to propagate across an ocean basin is a measure of the memory of the ocean—of the time it takes the ocean to be in equilibrium with the winds. These subsurface waves are most prominent along the layer that separates the warm surface waters from the cold deep water, and their speed increases with decreasing latitude. The memory of the ocean is therefore shortest near the Equator, and is shorter for a small ocean such as the Indian or Atlantic Ocean than for a larger one such as the Pacific. Thus, the seasonally varying winds force an equilibrium response in the equatorial Atlantic

Ocean—the winds and currents fluctuate practically in phase—but not in the Pacific Ocean. The Pacific Ocean is so large that, even on interannual time scales, oceanic conditions at a given time are not simply determined by the forcing at that time but in part depend on the forcing at earlier times. Hence, at the height of El Niño, waves excited much earlier are still propagating around the basin. Those waves signal the termination of El Niño and set the stage for complementary conditions that prevail until waves excited during El Niño signal its return at a later date. *See* MARITIME METEOROLOGY; OCEAN CIRCULATION; TROPICAL METEOROLOGY.

S. George Philander
Bibliography. Open University, *Ocean Circulation*, 2d ed., 2001; J. Pedlosky, *Ocean Circulation Theory*, 2004; S. G. Philander, *El Niño, La Niña and the Southern Oscillation*, 1990.

Equilibrium of forces

In a mechanical system the condition under which no acceleration takes place. Newtonian mechanics today is based upon two definitions which modify, but are essentially equivalent to, Newton's three fundamental laws. These definitions postulate the action of forces on particles. A particle is defined as a conceptual volume element that has mass and is sufficiently small to have point location. A body is defined as a system of particles. To develop the mechanics of a body, these definitions are applied to each of its particles and their influences summed. *See* ACCELERATION.

Newtonian particle laws. The law of motion is that, in a newtonian frame of reference, a particle of mass m acted on by resultant force \mathbf{F} has acceleration \mathbf{a} in accordance with the equation $\mathbf{F} = kma$. Therein, k is a positive constant whose value depends upon the units in which \mathbf{F} , m , and \mathbf{a} are measured.

The action-reaction law states that when one particle exerts force on another, the other particle exerts on the one a collinear force equal in magnitude but oppositely directed.

Equilibrium of a particle. A particle at rest or in uniform (unaccelerated) motion in a newtonian reference frame is in equilibrium. With few exceptions, a frame of reference fixed with respect to Earth is considered to be newtonian and the equation $\mathbf{F} = kma$ applied. Accordingly, when a particle is in equilibrium on Earth, $\mathbf{a} = 0$ and thus $\mathbf{F} = 0$. Also, when $\mathbf{F} = 0$, $\mathbf{a} = 0$ because k and m are not zero. These considerations lead to a theorem and its corollary: The resultant of forces exerted on a particle in equilibrium is zero; and if the resultant of forces exerted on a particle is zero, the particle is in equilibrium.

Body equilibrium. A body acted upon by force is in equilibrium when its constituent particles are in equilibrium. The forces exerted on its particles (and therefore on the body) are either internal or external to the body. An internal force is one exerted by one particle on another in the same body.

An external force is one exerted on a particle or the body by a particle not of the body.

Being particle actions, internal forces obey the action-reaction law. They occur in pairs whose individual and combined resultants are zero. Hence a further theorem is that the resultant of all forces internal to a body is zero.

The resultant of forces on each particle and therefore on all particles of a body in equilibrium is zero. Because the forces are either internal or external and the resultant of internal forces is zero, it follows as a theorem that the resultant of all external forces acting upon a body in equilibrium is zero.

Equations. Forces whose resultant is zero have zero vectors of total force and total moment. Consequently, both the algebraic sum of their components in any direction and their sum of moments about any line are also zero. Specifically, relative to an *OXYZ* reference frame, the components and moments of forces externally applied to a body in equilibrium are related through the equations given below.

$$\begin{aligned}\Sigma X &= 0 & \Sigma M_x &= 0 \\ \Sigma Y &= 0 & \Sigma M_y &= 0 \\ \Sigma Z &= 0 & \Sigma M_z &= 0\end{aligned}$$

See RESULTANT OF FORCES.

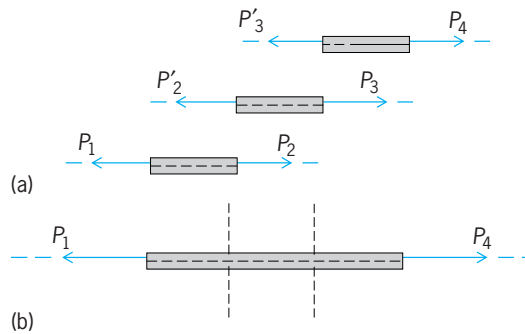
Forces can be collinear, coplanar, parallel, concurrent, and so forth. According to the class of a particular external force system and its orientation in *OXYZ*, certain of these equilibrium equations may be either trivial (containing no terms on their left sides) or dependent (derivable by arithmetic operations on one or more previously written equations of the set). The number of equations that are independent (contain different information) is the number necessary to define the resultant of the particular class force system which acts.

Force interaction of bodies. Under certain conditions, particles of one body exert pertinent forces on particles of another. Such forces, being particle actions, obey the action-reaction law. Thus, the entire force system exerted by one body on another coexists with an equal but opposite system exerted by the other on the one. The resultants of the co-existent systems are also equal and opposite. These conditions lead to the theorem that the mutual force actions of bodies are equal and opposite.

Contact and body forces. External forces are either contact or body forces according to their manner of application.

Contact forces are those exerted between the contacting particles of bodies which touch physically. The resultant of such forces is a force, couple, or force and couple, depending partly upon the nature of the contacting surfaces. For example, rigid bodies whose contacting surfaces do not penetrate, interlock, or adhere exert resultant compressive (push) forces on one another which act at the point or within the area of contact. Additionally, if either surface is considered to be smooth (incapable of resisting or applying force tangent to the surface), the force is normal to the surfaces at the contact.

Body forces are those exerted upon certain or all particles of a body by matter which need not contact



Examples of equilibrium forces. (a) Action forces and reaction forces acting on particles of a body. (b) Balanced external forces on a body.

the body. Gravitational, magnetic, and electrostatic forces are examples. A body's weight is the force resultant of essentially parallel downward forces resulting from the gravitational pull of the Earth on each particle of the body. The weight force acts downward through the geometric center of a homogeneous body. See CENTER OF GRAVITY.

Problem of equilibrium. To determine the forces which act upon matter at rest on Earth is the usual problem of equilibrium. The problem is statically determinate if solution is possible by employing the equations of equilibrium only. These equations relate all external forces; hence accurate solution requires that external forces be recognized and described according to their agencies of application. To summarize this information the free-body diagram is conveniently employed (see *illus.*). A free-body diagram is the sketch of a body of interest with all external forces acting thereon symbolized by arrows drawn on the sketch. See DYNAMICS; STATICS. Nelson S. Fisk

Equine infectious anemia

A lentivirus-induced disease of the horse family with an almost worldwide distribution. It is characterized by recurring fever, platelet reduction, weight loss, edema, and anemia. Although death can occur, there is usually an eventual cessation of clinical signs. However, host defenses are unable to completely eliminate the virus, and the animal remains a persistently infected inapparent carrier.

Virus. The equine infectious anemia virus is in the *Lentivirus* genus of the family *Retroviridae*. Although it is closely related to the human immunodeficiency viruses (HIV-1, HIV-2), its genetic organization is the least complex of this group of viruses. The equine infectious anemia virus particle (approximately 115 nanometers in diameter) is composed of a lipid-bilayer membrane surrounding a cone-shaped core, which in turn encapsulates a diploid ribonucleic acid (RNA) genome. Transmembrane proteins protrude through the lipid bilayer to support surface unit proteins on the exterior of the particle. These surface unit proteins are believed to be responsible for attachment to host cells, which in the case of the equine infectious anemia virus are predominantly of

the monocyte-macrophage lineage (although some strains may also infect endothelial cell types). This function of the surface unit proteins makes them an important target for neutralizing antibodies. Following entry into the host cell, the single-stranded RNA genome is converted into a double-stranded, deoxyribonucleic acid (DNA) copy (termed the provirus) by the viral enzyme reverse transcriptase. Integrase, another viral enzyme, then inserts the proviral DNA into the DNA of the host, where it directs the synthesis of progeny viral particles.

Mode of transmission. Equine infectious anemia virus infects only members of the horse family. The mechanical transfer of blood between animals by large blood-feeding insects (mainly horse flies and deer flies) is probably the most important mechanism of natural transmission. However, the virus can be efficiently transmitted by humans if sterile techniques are not observed during veterinary procedures.

Symptoms and causes. Clinical responses to equine infectious anemia virus infection can range from an extremely severe disease resulting in death to an absence of obvious signs. This variation appears to be determined by both the virus strain and the genetic composition of the host. However, when disease occurs it is usually observed shortly after exposure to the virus (incubation periods of 10–45 days are common) and consists of fever, platelet reduction (thrombocytopenia), lethargy, loss of appetite, and petechial hemorrhages. Most horses survive this acute episode but in many cases progress to the chronic form of the disease, characterized by recurring fever, thrombocytopenia, anemia, weight loss, edema, and hemorrhage. Each fever episode is associated with massive amounts of viral replication that occurs predominantly in the mature tissue macrophages of the liver and spleen and results in the release of millions of viral particles into the bloodstream. This extensive replication triggers the release of powerful molecules called cytokines, such as tumor necrosis factor alpha (TNF α), that are normally associated with inflammatory reactions. As TNF α is both an inducer of fever and an inhibitor of platelet production, it appears that such molecules may account for many of the symptoms associated with equine infectious anemia. However, specific immune responses also contribute significantly to clinical signs such as anemia and thrombocytopenia. This is due to immune complexes that bind to receptors on erythrocytes and platelets, targeting these cells for destruction by complement or phagocytic cells.

Pathology. There are no pathological lesions unique to equine infectious anemia, although several abnormalities are typically associated with the chronic form of the disease. These abnormalities include enlargement of the spleen, liver, and lymph nodes along with hemorrhages on mucosal and serosal surfaces. At the microscopic level, immune complexes are frequently deposited in the glomeruli of the kidneys, Kupffer cells in the liver may contain stainable iron granules from the breakdown of ery-

throcytes, and the number of hematopoietic cells in the bone marrow may be reduced.

Immune responses and viral persistence. In most cases, immune responses eventually control the extensive viral replication associated with each fever episode and bring about a cessation of clinical signs. However, these responses require a lengthy period to fully develop and even then are incapable of eliminating the virus completely. Equine infectious anemia virus employs numerous mechanisms to avoid removal by host defenses. These mechanisms range from infecting and possibly compromising the efficiency of a key cell type in the immune system to possessing a surface unit protein with a fundamental structure that confers partial resistance to neutralizing antibodies. Another property of this protein is its ability to withstand considerable alterations in amino acid content without loss of function. This, in combination with a high mutation rate resulting from the possession of an error-prone reverse transcriptase, facilitates the emergence of variants that are completely resistant to the strain-specific neutralizing antibodies produced within the first months of infection. Although additional strain-specific neutralizing antibodies may be generated against these variants, the mutational capabilities of this virus permit still more resistant types to arise. This process explains why each clinical episode of equine infectious anemia is caused by a different antigenic variant and why termination of this cycle appears to coincide with the eventual production of neutralizing antibodies that are broadly cross-reactive instead of strain-specific. In addition to evading neutralizing antibodies, the ability to produce a constant stream of viral variants because of rapid mutation rates may also be valuable in combating cytotoxic lymphocyte responses. However, this has yet to be fully evaluated. *See* IMMUNITY.

Control. Several inactivated virus or genetically engineered viral protein vaccines have been tested, but these have been completely successful only against the strain of virus from which they were derived. In China, an attenuated live virus vaccine has been used which has reportedly induced a broader form of immunity, but only after multiple doses and a postvaccination incubation period of several months. The delay in the development of full immunity, along with the fact that this approach would confound the current protocols for the diagnosis and control of equine infectious anemia, has prevented the adoption of attenuated virus vaccines outside China.

A safe and effective vaccine that rapidly induces protection against all strains of the equine infectious anemia virus is not yet available; this virus is controlled in most areas of the world by serological testing. The most commonly used method is the agar immunodiffusion, or Coggins, test that detects the presence of antibodies against the major core protein of the virus. If a horse tests positive, further transmission of the equine infectious anemia virus can usually be prevented by segregation or quarantine. The application of these techniques resulted in more than a tenfold decrease in the reported rate of positive tests

in the United States between 1975 and 1998. See RETROVIRUS; VIRUS INFECTION, LATENT, PERSISTENT, SLOW.

R. Frank Cook; Charles J. Issel

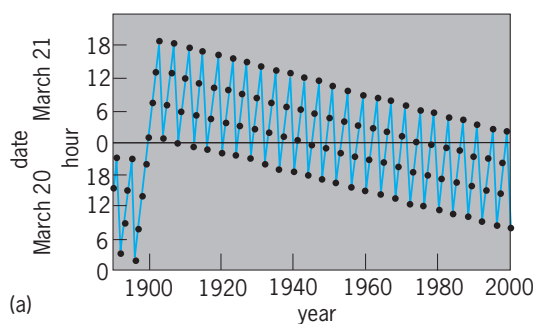
Bibliography. R. F. Cook, C. J. Issel, and R. C. Montelaro, Equine infectious anaemia, in M. J. Studdert (ed.), *Virus Infections of Equines*, pp. 297–323, Elsevier Science, Amsterdam, 1996; W. Maury, Regulation of equine infectious anemia virus expression, *J. Biomed. Sci.*, 5:11–23, 1998; R. C. Montelaro, J. M. Ball, and K. E. Rushlow, Equine retroviruses, in J. A. Levy (ed.), *The Retroviridae*, vol. 2, pp. 257–360, Plenum Press, New York, 1993; D. C. Sellon, F. J. Fuller, and T. C. McGuire, The immunopathogenesis of equine infectious anaemia virus, *Virus Res.*, 32:111–138, 1994.

Equinox

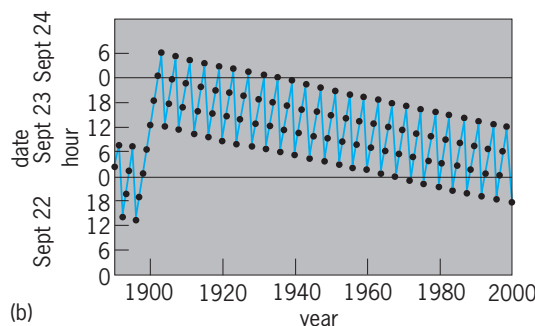
One of the two places in the sky where the ecliptic crosses the celestial equator; or one of the two times of the year when the Sun crosses these points. The ecliptic is the great circle across the sky that marks the mean path of the Sun; the celestial equator is the great circle that is an extension into the sky of the Earth's mean Equator. These two great circles meet at two points, one of which is the vernal equinox and the other the autumnal equinox. The Sun passes the vernal equinox each year about March 20, and the autumnal equinox about September 22. The vernal equinox can occur as early as March 19 and as late as March 21; for most of the twenty-first century it will be on March 20. The autumnal equinox can occur as early as September 21 and as late as September 24; for most of the twenty-first century it will be on September 22. The dates and times drift with the difference between the actual solar years and 365 days, and are corrected by leap years. This results in a 4-year variation superimposed on a negative 11-min-per-year slope (illus. a and b). Since 2000 was an ordinary leap year in the Gregorian calendar, unlike 1900, the dates will continue to decline through 2100 (illus. c). See ASTRONOMICAL COORDINATE SYSTEMS; CALENDAR; TIME.

At the vernal equinox, the Sun crosses from southern to northern declinations, marking the beginning of Northern Hemisphere spring and Southern Hemisphere autumn. At the autumnal equinox, the Sun crosses from southern to northern declinations, marking the beginning of Northern Hemisphere autumn and Southern Hemisphere spring. At the equinoxes, the Sun is directly above the Earth's Equator.

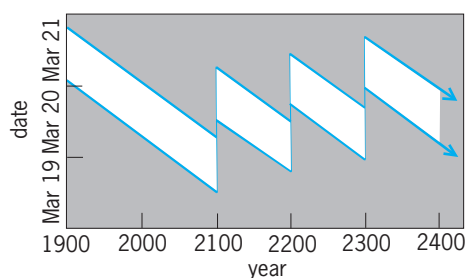
The term equinox is derived from the Latin for equal nights, indicating that the day and night are of equal duration. However, the actual duration of daylight is several minutes longer on the days of the equinoxes, and the actual dates of equal days and nights follow the autumnal equinox and precede the vernal equinox by a few days. The equinoctial dates are geometrical constructions in which the Sun is treated as a point; in actuality the top of its disk rises a few minutes ahead of its center. Further-



(a)



(b)



(c)

Date and time of the equinoxes (Universal Time or Greenwich mean time). (a) Vernal equinox, 1890–2000. (b) Autumnal equinox, 1890–2000. (c) Approximate 4-year range of times of the vernal equinox, 1900–2400. (After R. L. Reese and G. V. Chang, *The date and time of the vernal equinox: A graphical analysis of the Gregorian calendar*, *Amer. J. Phys.*, 55:848–849, 1987)

more, refraction in the Earth's atmosphere makes the Sun appear higher in the sky than it actually is, an effect that also lengthens daylight by several minutes. At sunrise, the top of the Sun is actually below the horizon, over which it is elevated by refraction. See METEOROLOGICAL OPTICS; REFRACTION OF WAVES.

Since precession of the equinoxes changes the equatorial plane over time, astrometrists distinguish between the dynamical equinox, defined by the intersection of the mean equatorial plane and the ecliptic, and the catalog equinox, which sets the zero value for the right ascensions in a star catalog for the particular epoch of that catalog. The precession, caused by the gravitational pull of the Sun and Moon on the Earth's equatorial bulge, moves the equinoxes along the ecliptic about 50 seconds of arc per year. Thus, the equinoxes occur 20 min earlier each year. See EARTH ROTATION AND ORBITAL MOTION; PRECESSION OF EQUINOXES.

Since the vernal equinox was in the constellation

Aries when Hipparchus studied it 2000 years ago, it is known as the first point in Aries. It is, however, in Pisces. Jay M. Pasachoff

Bibliography. R. L. Reese and G. Y. Chang, The date and time of the vernal equinox: A graphical analysis of the Gregorian calendar, *Amer. J. Phys.*, 55:848-849, 1987; P. K. Seidelmann (ed.), *Explanatory Supplement to the Astronomical Almanac*, rev. ed., 1992.

Equisetales

An order of the division Sphenophyta, subkingdom Embryobionta. The Equisetales, commonly known as horsetails, is represented by a single living genus, *Equisetum*, with about 25 species found both in moist and dry habitats. These plants grow throughout the world, except in Australia and New Zealand. The plants range from herbaceous to shrubby and rarely exceed 3 or 4 ft (0.9 or 1.2 ft) in height, although some tropical species grow much taller.

The plant body is commonly composed of perennial underground stems (rhizomes) with various types of aerial stems (Fig. 1) in different species. Some of these are perennial, others annual; some are unbranched and reproductive, others much branched and vegetative. The bushy structure of the latter is suggestive of the common name, horsetail.

The leaves are reduced to scales which appear fused in whorls at the nodes. In the young leaf, the central region contains stomates and chloroplasts which function in photosynthesis for a short time. The burden of photosynthesis is carried on by cells in the cortex of the stem. There are no true leaf gaps in the stele. See LEAF.

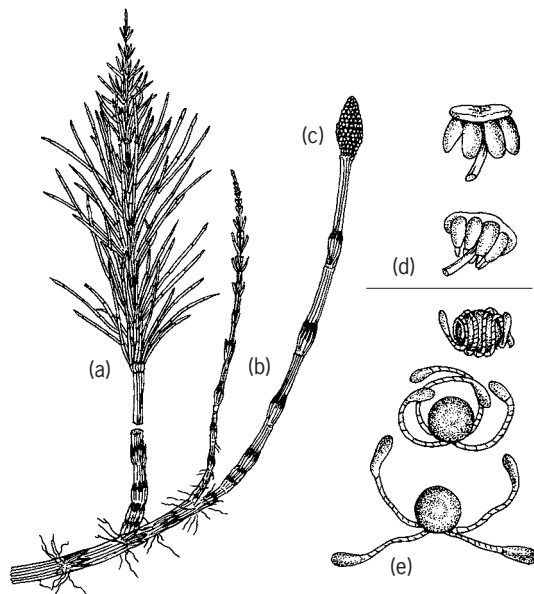


Fig. 1. *Equisetum arvense*. (a) Sterile shoot. (b) Fertile shoot growing from an underground rootstock. (c) Cone. (d) Two views of shield-shaped sporangiophores. (e) Spores, greatly enlarged. As spore dries, elaters expand. (After E. W. Sinnott and K. S. Wilson, *Botany: Principles and Problems*, 6th ed., McGraw-Hill, 1963)

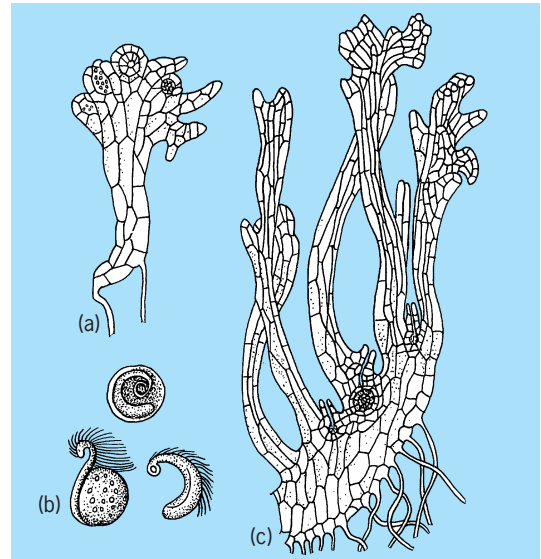


Fig. 2. *Equisetum* gametophyte. (a) Small gametophyte, showing antheridia only. (b) Male gametes. (c) Larger gametophyte; branching lobes have archegonia at bases. (After E. W. Sinnott and K. S. Wilson, *Botany: Principles and Problems*, 6th ed., McGraw-Hill, 1963)

The stem is composed of conspicuous nodes with long, strongly ribbed internodes. The cells of the stem originate in a single, pyramidal apical cell, giving rise eventually to a ring of vascular bundles with numerous openings or canals present except at the node. There is no cambium, but the basal part of each internode functions for some time as an intercalary meristem. The walls of sclerenchyma fibers in the cortex and of the epidermal cells are more or less impregnated with silica and aid in giving mechanical support to the stem. Early settlers, noting this, made use of the vegetative plants for scouring pots, pans, and floors, and as a result gave them the common name of "scouring rushes." See APICAL MERISTEM; CORTEX (PLANT); EPIDERMIS (PLANT); SCLERENCHYMA; STEM.

Except for the first-formed primary root, all roots are adventitious. Each adventitious root is derived from a single apical cell and develops an exarch protosteles. There is no true cambium, but an intercalary meristem (region of cell division) has been observed.

The strobili (cones) are solitary and terminal, either on the main stem or on lateral branches. They are composed of a central axis with a whorled series of compact lateral branches known as sporangiophores. Sporophylls (spore-bearing leaves) are entirely suppressed. Because of the pressure resulting from their development, each sporangiophore appears compressed and hexagonal in surface view. At maturity the several sporangia ripen on the lower or abaxial portion of the sporangiophore. The spores produced are alike (homosporous). Each possesses a complex wall structure having four long, slender appendages known as elaters. These elaters are hygroscopic (coil and uncoil in response to changes in humidity) and thus aid in dissemination of spores.

Spores which reach a suitable substrate germinate immediately, producing a flat plate or mound

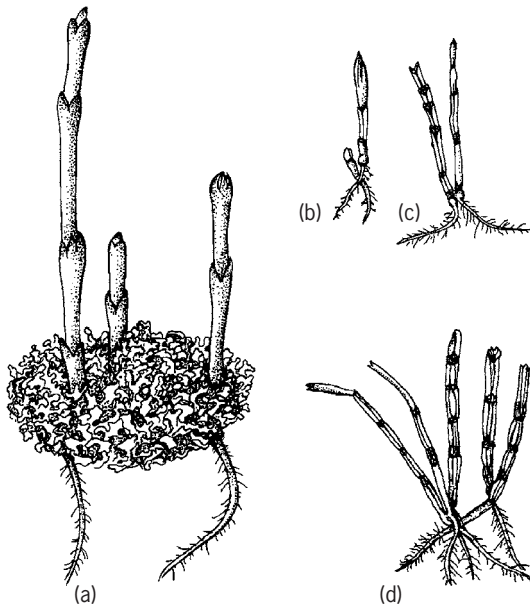


Fig. 3. Sporophytes. (a) *Equisetum* gametophyte with sporophytes. (b, c, d) The young sporophytes of *Equisetum hiemale*. (After G. M. Smith, *Cryptogamic Botany*, vol. 2, 2d ed., McGraw-Hill, 1955)

of green tissue with numerous filamentous branches called the gametophyte (Fig. 2) which usually gives rise to the sex organs, antheridia (male) and archegonia (female). The motile sperms produced in the antheridia are relatively large, spirally coiled cells with numerous apical flagella. The archegonia are formed in the cushionlike portion of the gametophyte, with the chambers containing the eggs well buried in the thallus. Several eggs on the same gametophyte may be fertilized, following which each may develop an embryonic sporophyte (Fig. 3) without a suspensor (a chain of cells which serves to put the embryo in a favorable position in relation to its food supply). With the formation of the primary root and stem, the young sporophyte produces chlorophyll and becomes independent (self-sustaining). See EMBRYOIONTA; REPRODUCTION (PLANT); ROOT (BOTANY).

Paul A. Vestal

Bibliography. L. Benson, *Plant Classification*, 2d ed., 1979; H. C. Bold, *Morphology of Plants and Fungi*, 5th ed., 1987; H. C. Bold and J. La Claire, *The Plant Kingdom*, 5th ed., 1987; G. M. Smith, *Cryptogamic Botany*, vol. 2, 2d ed., 1955.

Equivalent circuit

A representation of an actual electric circuit or electronic device by a simple circuit whose behavior is very near to that of the actual circuit over a specified range of conditions. When these conditions are satisfied, the equivalent circuit may be said to constitute a macromodel of the actual circuit. The use of equivalent circuits is important in many of the analysis and design operations associated with electronic circuits. For example, a frequently used electronic device is the operational amplifier. One of the better-

known operational amplifiers, the μA 741, is realized as an integrated circuit containing over 20 transistors. Many electric circuits, for example, active- RC (resistance-capacitance) filters, require the use of large numbers of operational amplifiers. In such applications, the use of equivalent circuits greatly simplifies analysis and design operations. See ELECTRIC FILTER; INTEGRATED CIRCUITS; OPERATIONAL AMPLIFIER; TRANSISTOR.

Two of the best-known equivalent circuits are the Thévenin equivalent circuit and the Norton equivalent circuit. The Thévenin equivalent circuit consists of the series connection of a voltage source and a two-terminal circuit (Fig. 1a); the Norton equivalent circuit consists of the parallel connection of a current source and a two-terminal circuit (Fig. 1b). In both circuits, the output of the voltage or current source is usually dependent on the electric signals applied to the input terminals. See NETWORK THEORY; THÉVENIN'S THEOREM (ELECTRIC NETWORKS).

Two common applications of equivalent circuits are the modeling of large circuits (such as the operational amplifier mentioned above), and the modeling of individual electronic solid-state devices such as transistors. An example of an equivalent circuit for an operational amplifier is shown in Fig. 2a. The Thévenin equivalent circuit is used. The gain constant K is usually very large (of the order of 10^5). The output resistor R_{out} is usually very small and frequently is omitted. Solid-state devices such as transistors are inherently nonlinear in their behavior. For such nonlinear components, the use of equivalent circuits is usually restricted to modeling the actual circuit behavior under the conditions that the circuit variables have such small variations that the nonlinear relations may be approximated by idealized ones. In such a scenario, the equivalent circuits are usually referred to as small-signal equivalent circuits. An example of such a device is a first-order

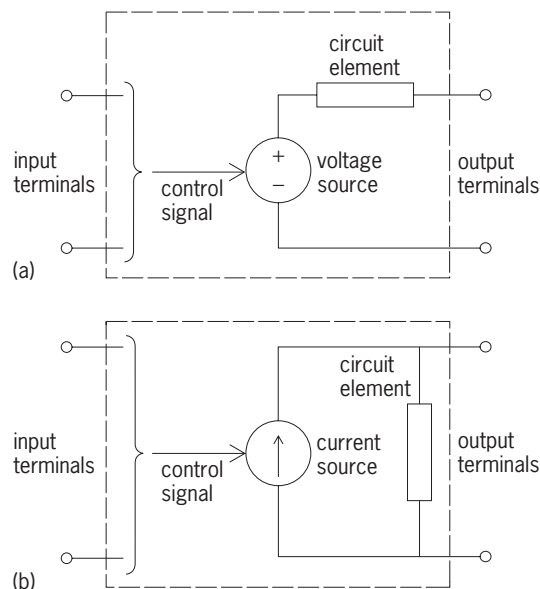


Fig. 1. Basic equivalent circuits. (a) Thévenin equivalent circuit. (b) Norton equivalent circuit.

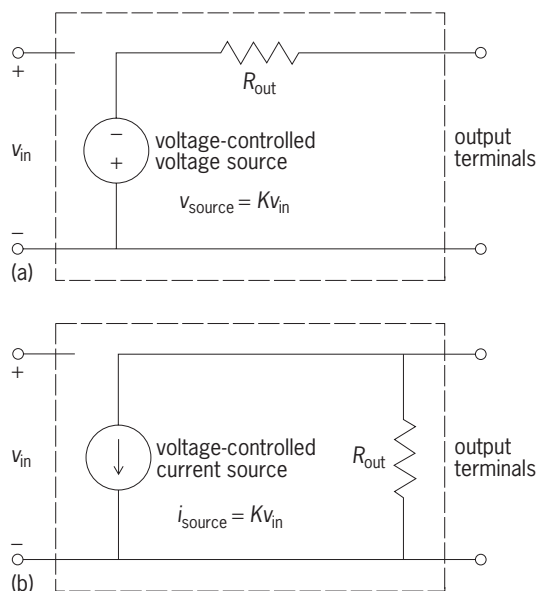


Fig. 2. Equivalent circuit applications. (a) Operational amplifier equivalent circuit applications. (b) MOSFET equivalent circuit.

small-signal model for a metal-oxide-semiconductor field-effect transistor (MOSFET; Fig. 2b). The Norton equivalent circuit is used. The output resistance R_{out} is usually very large, and frequently is omitted from the equivalent circuit. See CIRCUIT (ELECTRICITY); CIRCUIT (ELECTRONICS). Lawrence P. Huelsman Bibliography. L. P. Huelsman, *Basic Circuit Theory*, 3d ed., 1991; D. E. Johnson, J. R. Johnson, and J. L. Hilburn, *Electric Circuit Analysis*, 1989; J. W. Nilsson and S. A. Riedel, *Electric Circuits*, 7th ed., Prentice Hall, 2004.

Equivalent weight

The number of parts by weight of an element or compound which will combine with or replace, directly or indirectly, 1.008 parts by weight of hydrogen, 8.00 parts of oxygen, or the equivalent weight of any other element or compound. The term equivalent weight comes from the law of equivalent proportions, which states that the weights of two elements A and B which combine separately with identical weights of another element C are either the weights in which A and B combine together, or are related to them in the ratio of small whole numbers. A standard weight of 8.000 parts is chosen for oxygen. For all elements, the atomic weight is equal to the equivalent weight times a small whole number, called the valence of the element. See ATOMIC MASS; COMPOUND (CHEMISTRY); ELEMENT (CHEMISTRY); VALENCE.

An element can have more than one valence and therefore more than one equivalent weight. The use of the terms is explained below.

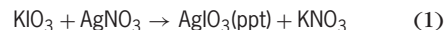
1. Ammonia, NH_3 , contains 1 atom of nitrogen combined with 3 atoms of hydrogen. Since the equivalent weight of hydrogen is equal to its atomic weight, the equivalent weight of nitrogen is one-third its atomic weight and its valence is 3.

2. Magnesium oxide, MgO , contains 1 atom of magnesium combined with 1 atom of oxygen. Since the equivalent weight of oxygen is one-half its atomic weight, the equivalent weight of magnesium is also one-half its atomic weight and its valence is 2.

3. Phosphorus forms phosphorus trichloride, PCl_3 , and phosphorus pentachloride, PCl_5 . Since the equivalent weight of chlorine is equal to its atomic weight, in the trichloride the equivalent weight of phosphorus is one-third its atomic weight and its valence is 3, and in the pentachloride the equivalent weight is one-fifth its atomic weight and its valence is 5.

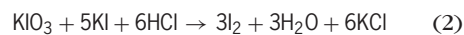
The equivalent weight of a compound depends on the reaction in which it takes part. Thus:

1. In the reaction between potassium iodate, KIO_3 , and silver nitrate, $AgNO_3$, one molecule of silver iodate, $AgIO_3$, is precipitated for every molecule of silver nitrate. This is represented by reaction (1).



Since the equivalent weight of silver is its atomic weight, the equivalent weight of potassium iodate, in this reaction, is its molecular weight.

2. When potassium iodate, KIO_3 , is reduced to iodine, I_2 , by potassium iodide, KI , three molecules of iodine are produced per molecule of potassium iodate. This is represented by reaction (2). Since the



equivalent weight of iodine is one-half its molecular weight, the equivalent weight of potassium iodate, in this reaction, is one-sixth its molecular weight.

This concept, together with that of gram-equivalent weight, tends to have been abandoned, and relations are expressed in terms of balanced stoichiometric chemical equations and relative numbers of moles reacting. See ELECTROCHEMICAL EQUIVALENT; MOLE (CHEMISTRY); OXIDATION-REDUCTION; STOICHIOMETRY. Thomas C. Waddington

Erbium

A chemical element, Er, atomic number 68, atomic weight 167.26, belonging to the rare-earth group. The naturally occurring element is made up of the six

1																	18																		
1	2											13	14	15	16	17	2																		
3	4											5	6	7	8	9	10																		
Li	Be											B	C	N	O	F	Ne																		
11	12	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18																		
Na	Mg	Al	Si	P	S	Cl	Ar	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36										
K	Ca	Sc	Ti	V	Cr	Mn	Fe	Co	Ni	Cu	Zn	Ga	Ge	As	Se	Br	Kr	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52	53	54
Rb	Sr	Y	Zr	Nb	Mo	Tc	Ru	Rh	Pd	Ag	Cd	In	Sn	Sb	Te	I	Xe	55	56	57	58	59	60	61	62	63	64	65	66	67	68	69	70		
Cs	Ba	La	Hf	Ta	W	Re	Os	Ir	Pt	Au	Hg	Tl	Pb	Bi	Po	At	Rn	87	88	89	90	91	92	93	94	95	96	97	98	99	100	101	102		

lanthanide series 57 58 59 60 61 62 63 64 65 66 67 68 69 70
La Ce Pr Nd Pm Sm Eu Gd Tb Dy Ho Er Tm Yb

actinide series 89 90 91 92 93 94 95 96 97 98 99 100 101 102
Ac Th Pa U Np Pu Am Cm Bk Cf Es Fm Md No

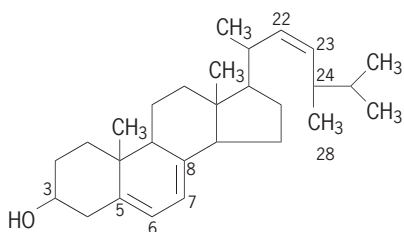
stable isotopes. The rose-pink oxide, Er_2O_3 , dissolves in mineral acids to give rose-colored solutions. The salts are paramagnetic and the ions are trivalent. At low temperatures the metal is antiferromagnetic and at still lower temperatures becomes strongly ferromagnetic. For properties of the metal see RARE-EARTH ELEMENTS

Frank H. Spedding

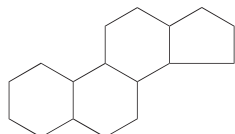
Bibliography. F. A. Cotton et al., *Advanced Inorganic Chemistry*, 6th ed., Wiley-Interscience, 1999; K. A. Gschneidner Jr., J.-C. Bünzli, and V. K. Pecharsky (eds.), *Handbook on the Physics and Chemistry of Rare Earths*, 2005.

Ergosterol

A steroid, with the structure shown below, that be-



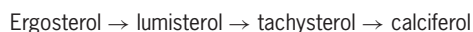
longs to the class of unsaponifiable lipids. It was first isolated from ergot bodies (black or dark-purple club-shaped bodies of hardened mycelium from the ergot fungus, which replace the seeds of various grasses, such as rye and wheat). Ergosterol is a white crystalline compound, insoluble in water and soluble in organic solvents. It is derived from the cyclopentaperhydrophenanthrene skeleton:



See ERGOT AND ERGOTISM.

The hydroxyl group at carbon atom 3 is of the β configuration. Ergosterol differs from cholesterol by having three positions of unsaturation, at carbon atoms 5-6, 7-8, and 22-23, and containing a methyl group (carbon atom 28) substituted for a hydrogen atom at carbon atom 24. See CHOLESTEROL.

Ergosterol makes up to 90-100% of the steroids in certain yeasts and can be extracted from these yeasts in commercial quantities. Ultraviolet irradiation of ergosterol leads to the formation of calciferol (vitamin D_2) by way of a series of irreversible reactions:



Under controlled conditions, a more than 50% yield of vitamin D_2 can be obtained. The process is used commercially. See VITAMIN D.

The biosynthetic pathway of ergosterol in yeast is probably similar to that of cholesterol in animal tissues, since both acetate and squalene are converted to ergosterol by yeast:



The amino acid methionine serves as the donor of the methyl group (carbon atom 28) of ergosterol. This reaction probably takes place after the steroid nucleus has been formed. See STEROID.

Willem J. Van Wagtenonck

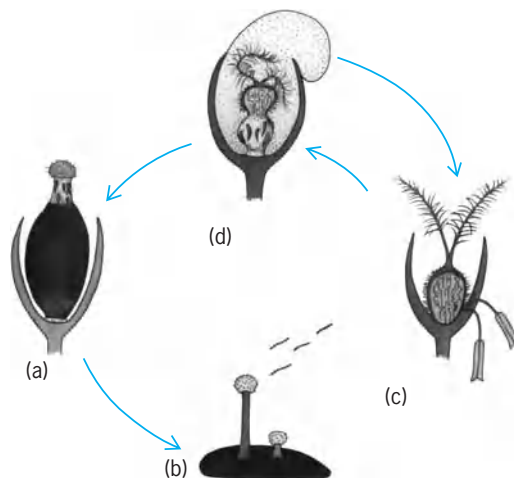
Bibliography. H. Eriksson and J. A. Gustafsson, *Steroid Hormone Receptors: Structure and Function*, 1984; B. Green and R. E. Leake (eds.), *Steroid Hormones: A Practical Approach*, 1987; H. L. Makin, *Biochemistry of Steroid Hormones*, 2d ed., 1984.

Ergot and ergotism

Ergot is the seedlike body of fungi (molds) of the genus *Claviceps*; ergotism is a complex disease of humans and certain domestic animals caused by ingestion of grains and cereals infested with ergot. Ingestion of these long, hard, purplish-black structures called sclerotia (illus. a) may lead to convulsions, abortion, hallucinations, or death. During the Middle Ages, hundreds of thousands of people are believed to have died from this disease, often referred to as holy fire, St. Anthony's fire, or St. Vitus's dance. Epidemics in humans, although less prevalent in modern times, last occurred in 1951, and the potential danger is always present, as shown by annual livestock losses due to ergot poisoning. See POISON.

Fungus. Ergot of rye derives its name from *l'ergot*, the French word for cockspur, owing to the similarity in size and shape. Later, the word ergot was applied to all sclerotia of *Claviceps*, even though some were of globose or irregular shape. There are 32 recognized species, most of which infect members of the grass family. Only three species are parasitic on the rushes and sedges. See PLANT PATHOLOGY.

The other stages of the ergot fungus are less conspicuous and often go unnoticed. After the sclerotia overwinter or persist through drought, as is the case in some desert and tropical species, they germinate under favorable, moist conditions and give rise to small mushroomlike fruiting bodies (illus. b). The spiny spherical heads of the fruiting bodies contain



Ergot life cycle. Parts a-d explained in text.

long filamentous ascospores (sexual spores) which are ejected and blown about by the wind. Ascospores landing on a healthy stigma of a grass flower stick and enter, as does pollen, the ovary in the form of a fine mycelial network (illus. *c*). Soon after infection, a special convoluted structure called a sphacelium develops between the ovary and the base of the floral cavity, at which point further invasion of host tissue is restricted (illus. *d*). The sphacelium produces hundreds of thousands of conidia (asexual spores) embedded in a sugary matrix called honeydew. This sticky, mildly fetid substance attracts insects which feed on the sugars. When the insects move to other flowers, the conidia adhering to the legs and mouthparts infect healthy ovaries of other plants (illus. *c*). After several weeks the sclerotium develops between the sphacelium and the flower base, completing the cycle (illus. *a*). See ASCOMYCOTA.

Both honeydew and sphacelium are nontoxic and do not contain the poisonous alkaloids responsible for ergotism. The honeydew itself, composed almost entirely of sugars, is quite nutritious. The major sugars are glucose, fructose, sucrose, and fructose-containing oligosaccharides; in some species sugar alcohols are the major honeydew constituents. Honeydew sugars are known to be products of the parasite rather than products of the host, as was previously believed. By converting host sucrose into some of the above-mentioned sugars, the sphacelium creates a sucrose deficit in the floral cavity. The resulting flow of host sucrose is continuous, since fungal enzyme activity will not allow a sucrose equilibrium to exist.

Sclerotia have an unusual chemical makeup. They carry only 10% water by weight, and 50% of the dry weight is composed of fatty acids, sugars, and sugar alcohols, which make the ergot a storehouse of energy. Unfortunately, they also contain the poisonous alkaloids, ranging from 0 to 1.2% of the dry weight.

Disease. There are three types of ergotism (gangrenous, convulsive, and hallucinogenic). Their symptoms often overlap; the hallucinogenic form is usually observed in combination with one of the other two. The unusual combination of gangrenous and convulsive symptoms is sometimes observed in the Balkans and areas near the Rhine River.

The onset of gangrenous ergotism is generally characterized by lassitude, nausea, and pains in the limbs followed by alternating sensations of intense heat and cold. As bodily extremities become numb, livid watery vesicles may appear on the affected parts (usually arms and legs). Finally, the diseased area turns black, dry, and becomes mummified.

In convulsive ergotism, various parts of the body become grossly deformed as a result of clonic or tonic convulsions, or both. Generalized neurological stimulation causes epileptiform seizures, whereas specific stimulation might involve ravenous hunger, violent retching, tongue biting, or abnormal breathing patterns. This form involves a longer recovery period and often results in permanent nerve damage and subsequent sensitization.

The hallucinogenic form often includes symptoms of one of the other types. In its more pure form, it is referred to as choreomania, St. Vitus's dance, or St. John's dance. Vivid hallucinations are accompanied by psychic intoxication reminiscent of the effects of many of the modern psychedelic drugs (nervousness, physical and mental excitement, insomnia, and disorientation). Early reports state that the disease usually manifested itself in the form of strange public dances that might last for days or weeks on end. Dancers made stiff jerky movements accompanied by wild hopping, leaping, screaming, and shouting. They were often heard conversing with devils or gods, and danced compulsively, as if possessed, until exhaustion caused them to fall unconscious or to lie twitching on the ground. High mortality rates were associated with severe epidemics involving any of the three forms of ergotism.

History. The earliest undisputed occurrence of ergotism was documented in Germany in A.D. 857, about 700 years before ergot was described in any of the medieval herbals. However, the most serious epidemics were recorded between the late 900s and the 1800s. Although both gangrenous ergotism (mainly in western Europe) and hallucinogenic ergotism were reported throughout this period, the convulsive form was not described until the late 1400s. The epidemics came in times of famine, usually during the years of heavy rain following severe winters. Under these conditions rye became heavily infected with ergot.

Convulsive ergotism is little understood even today. It occurred mainly in northern Europe east of the Rhine, with the greatest number of fatalities among children and weaned babies. However, it did not affect nursing babies, even those with severely afflicted mothers. The disease was much less prevalent on farms and dairies or among individuals consuming bacon, eggs, or milk (or other vitamin A-rich food). These and clinical observations have led to the hypothesis that ergot intake combined with avitaminosis A was the cause of convulsive ergotism. This does not, however, explain the absence of the gangrenous form in eastern Europe.

Ergot was known to cause ergotism in France as early as the late 1500s. After this period, with the cleaning of grain and the increased use of wheat (wheat is a poor host for ergot), the incidence of ergotism fell rapidly in the West. In the East, ergot was not accepted as the cause of ergotism until the terrible 1771 epidemic; continued resistance to such ideas by many peasants caused widespread epidemics even into the twentieth century. Once ergot was known to cause the disease, many governments set safety limits on ergot ranging in grain from 0.1 to 0.3%. A level of 8 to 10% ergot in grain is considered to be fatal. In modern times ergotism has occurred seldom, and then usually from overdoses of the drug in medicine or in attempts to procure abortion.

The success with which the disease is controlled in humans has been brought about by (1) agricultural inspection, (2) use of wheat, potatoes, and maize instead of rye, (3) limited control of ergot, (4) reserves

of sound grain, and (5) forecasting severe ergot years. The most recent and best-recorded epidemic was in southern France in 1951 when an unscrupulous miller used moldy grain to make flour.

Ergot in medicine. Although the first description of the use of ergot in childbirth appeared in 1582, earlier reports indicate its use by midwives far back into European and Chinese history to stimulate uterine contractions and to control uterine bleeding. Its formal introduction into medicine did not come until the early 1800s. After commercial production began, powdered ergot preparations were found to be variable and unstable. Later in the century, ergot was successfully used to treat migraine, but this effect, too, was variable.

In the early twentieth century two ergot alkaloids (ergotamine and ergotamine) were isolated. Unfortunately, they caused significant side effects and were not as specific or active as some of the crude aqueous preparations. Shortly after its discovery, ergotamine was found to be effective in the treatment of migraines. Both ergotamine and ergotamine cause vasoconstriction that can lead to gangrene with chronic use. *See* HEADACHE.

In 1935 a new water-soluble ergot alkaloid, ergonovine, was synthesized. It is a small lysergic acid derivative, whereas ergotamine is a large, tetracyclic peptide linked to lysergic acid. Ergonovine is used to facilitate childbirth by stimulating uterine contractions. Ergotamine was found to be a mixture of similar large alkaloids (ergocorine, ergokryptine, and ergocristine). Many other important lysergic acid derivatives have been produced by means of semisynthesis. Of these derivatives, LSD-25 (*d*-lysergic acid diethylamide) is the most famous. LSD has been used experimentally, mainly in psychiatry and neurophysiology. *See* PSYCHOTOMIMETIC DRUGS.

Through extensive research, many other uses for ergot alkaloids has been found. Ergotamine and dihydroergotamine are used to treat migraines, and methysergine is used in migraine prophylaxis. Dihydroergotamine is prescribed for hypertension, cerebral diffuse sclerosis, and peripheral vascular disorders. Ergocorine and the less toxic agroclavine have been reported as unusual experimental birth-control agents. The drugs appear to inhibit implantation of the ovum. Several semisynthetic alkaloids are also active implantation inhibitors.

Alkaloid production. For more than a century, scientists have tried to produce sclerotia in pure culture, tissue culture, and sophisticated extractions of host tissues, but only the mycelial stage has been grown. Neither the sphaecelia nor the sclerotia can be cultured away from the host. These failures have been due to a lack of understanding about the physiological requirements that are met during parasitism.

In 1951 the first ergot alkaloids were found in pure liquid cultures of mycelial mutants. The products were clavine alkaloids, which were interesting but of no medical value (except as possible birth-control agents). These compounds are ergoline derivatives (lysergic acid minus carbonyl group) and are not eas-

ily converted to lysergic acid derivatives. High yields of the more valuable alkaloids have been produced in liquid culture by new mycelial strains. Such production is not without its problems. Most important, the strains are unstable and may revert to nonalkaloid production with subsequent culture. However, bromocryptine and other ergot alkaloids have been used to alleviate some of the symptoms of Parkinson's disease. *See* ALKALOID; PARKINSON'S DISEASE.

R. L. Mower

Bibliography. E. L. Backman, *Religious Dances in the Christian Church and in Popular Medicine*, 1952; G. Barger, *Ergot and Ergotism*, 1931; F. J. Bové, *The Story of Ergot*, 1970; M. Goldstein et al. (eds.), *Ergot Compounds and Brain Function: Neuroendocrine and Neuropsychiatric Aspects*, 1980.

Ericales

An order of flowering plants, division Magnoliophyta (angiosperms), in the large asterid assemblage (often Asteridae in previous systems of classification). Nearly all of the 24 families assigned to the order have previously been considered members of several orders in the subclass Dilleniidae. Ericales are a diverse group that have general asterid characters: tenuinucellate ovules; flowers with fused sepals, and often, fused petals with the anthers fused at least basally to the petals, even when the petals are apparently free; cellular endosperm formation; and tegumentary tapeta.

The largest families are Ericaceae (1350 species), Primulaceae (1000 species), Myrsinaceae (1000 species), Sapotaceae (800 species), and Balsaminaceae (600 species). The two biggest families are largely temperate herbs (Primulaceae), and the next three are largely tropical trees or herbs (Balsaminaceae). A relationship between the first four families has been known for many years, but the last has nearly always been considered related to the rosoid family Geraniaceae.

Familiar plants belonging to Ericales include rhododendrons (*Rhododendron*, Ericaceae), camellias (*Camellia*, also the genus of tea; Theaceae), primroses (*Primula*, Primulaceae), phlox (*Phlox*, Polemoniaceae), and impatiens (*Impatiens*, Balsaminaceae). *See* BLUEBERRY; DILLENIIDAE; MAGNOLIOPSIDA; PLANT KINGDOM.

Mark W. Chase

Eriocaulales

An order of flowering plants, division Magnoliophyta (Angiospermae), subclass Commelinidae of the class Liliopsida (monocotyledons). The order consists of the single family Eriocaulaceae, with about 1200 species. The Eriocaulales are Commelinidae with a reduced (or no) perianth and with unisexual flowers aggregated into a dense, involucrate head that is elevated above the clustered, basal leaves on a long peduncle. Although the individual

flowers are small and inconspicuous, the heads are more or less showy and pollination is usually by insects, in spite of the absence of nectar and nectaries. The perianth typically consists of two series of similar, white-hyaline tepals, three members in each series. Each locule of the ovary has a single, pendulous, orthotropous, tenuinucellate ovule. The order is of negligible economic importance. *See* COMMELINIDAE; LILIOPSIDA; MAGNOLIOPHYTA; PLANT KINGDOM. Arthur Cronquist

Erosion

Erosion is the result of processes that entrain and transport earth materials along coastlines, in streams, and on hillslopes. Wind and water are common agents through which forces are applied to resistant rocks, soils, or other unconsolidated materials. Erosion types often are designated on the basis of the agent: wind erosion, fluvial erosion, and glacial erosion. Fluvial erosion usually has been regarded as the most effective type in shaping the land surface during recent geologic time. Under certain environmental conditions, however, wind erosion moves considerable quantities of earth materials, as demonstrated during the "dust bowl" years in the United States. Glacial erosion shaped much of the land surface during the Quaternary Period of geologic time. Each type of erosion produces distinctive landforms, contributing to the diversity of terrestrial landscapes. *See* DESERT EROSION FEATURES; EOLIAN LANDFORMS; GEOMORPHOLOGY; GLACIOLOGY; MASS WASTING; QUATERNARY; STREAM TRANSPORT AND DEPOSITION.

Many terms are used in association with erosion processes, although not with consistency. Erosion involves the entrainment and transportation of earth materials from source areas to deposition sites. Soil loss refers to material actually transported from a particular hillslope or field, and may be less in volume than erosion due to on-site deposition in microtopographic depressions. Sediment is the product of erosion, and sediment yield is the amount of sediment transported from a watershed. Other terms are found in the literature that differ somewhat in connotation within particular contexts and for particular scientific disciplines. *See* SEDIMENTOLOGY.

Erosion includes various processes, depending upon perspective, ranging from in-situ rock weathering, mass-movement, sheetwash on hillslopes, channel bed and bank scour, work of fauna, and deposition. Most scientists regard weathering, mass-movement, and deposition as separate processes and faunal erosion to be of minor importance except in special circumstances. *See* WEATHERING PROCESSES.

Forces exerted by erosion processes must exceed resistances of earth materials for entrainment and transportation to occur. Environmental conditions determine the magnitude of the forces, the resistances, and the relations among them. Erosion rates are highly variable in time and space due to changing relations between forces and resistances. The major factors governing wind-erosion rates are wind veloc-

ity, topography, surface roughness, soil properties and soil moisture, vegetation cover, and land use. The major factors governing fluvial-erosion rates on hillslopes are rainfall energy, topography, soil properties, vegetation cover, and land use. The major factors governing fluvial-erosion rates in stream channels are depth and velocity of water flow, together with the size and cohesiveness of the bed and bank materials. The major factors governing glacial-erosion rates are the depth and velocity of ice flow, together with the hardness of the bed and side-wall materials.

During the twentieth century, public and research interest focused on fluvial erosion. The erosion rate under natural, undisturbed conditions is known as geologic, normal, or natural erosion. The erosion rate following land disturbance, usually due to human activities such as agriculture, mining, or construction, is known as accelerated or anthropic erosion. Human activities are capable of changing erosion rates by two or three orders of magnitude. Conversely, soil-conservation practices can reduce erosion rates to approximate the natural rate if the properly designed, applied, and maintained. Effective erosion control is based upon either decreases in forces impinging on surfaces or increases in the resistances of surface materials. *See* SOIL CONSERVATION.

Accelerated erosion by fluvial processes may be the most important environmental problem worldwide because of its spatial and temporal ubiquity. Erosion rates commonly exceed soil-formation rates, causing depletion of soil resources. The effects of erosion are insidious due to the removal of the fertile topsoil horizon, compromising food production. Sediment frequently is transported well beyond the source area to degrade water quality in streams and lakes, harm aquatic life, reduce the water-storage capacity of reservoirs, and increase channel-maintenance costs. *See* SOIL.

T. J. Toy

Bibliography. R. P. C. Morgan, *Soil Erosion and Conservation*, 3d ed., 2005; T. J. Toy, Accelerated erosion, process, problems, and prognosis, *Geology*, 10:524-529, 1982; T. J. Toy et al., *Soil Erosion: Processes, Prediction, Measurement, and Control*, 2002; F. R. Troeh, J. A. Hobbs, and R. L. Donahue, *Soil and Water Conservation*, 4th ed., 2003.

Errantia

A group of 34 families of Polychaeta in which the anterior, or cephalic, region is more or less fully exposed and the body is often long, linear to short, and depressed. The segments of these worms are similar or change gradually. The pharynx is often heavily muscularized and eversible, and its inner walls are fortified with calcified or chitinized plates or jaws. Some families are benthic throughout their lives, others are entirely pelagic, and some have pelagic larval and reproductive stages and benthic trophic development. Errantia occur in all seas, at all depths, and in inland seas or lakes. They range in length from a

few millimeters to 6 ft (2 m). See POLYCHAETA.

Scale bearers and allies. There are 6 scale-bearing and 3 allied families included in this large artificial category.

Aphroditidae. With 9 genera and about 85 species, they include the sea mouse (*Aphrodita*) (Fig. 1), the giant cold-water *Laetmonice*, and several other deep-water genera.

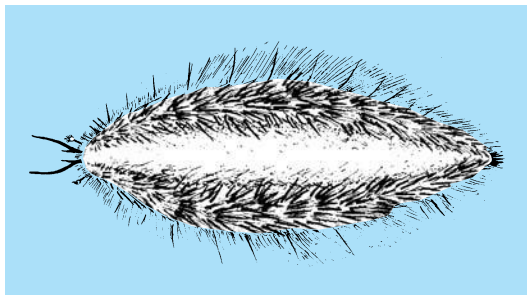


Fig. 1. *Aphrodita*, of the Aphroditidae, dorsal view.

Polynoidae. This is the largest of polychaete families, with about 75 genera and 600 species. The body may be short or long and large or minute; it is covered with elytra (Fig. 2). Representatives of this family are in all seas and at all depths. Most are benthic, some are pelagic (*Drieschia*, *Nectochaeta*), and others are commensal (*Arctonoe*, *Hesperonoe*, *Polynoe*); some are known only from abyssal depths (*Macellicephalo*). Common genera are *Harmothoe*, *Lepidonotus*, and *Halosydna*.

Polyodontidae. The 8 or 9 genera and about 36 species are tubicolous and often large-bodied, with characteristic cephalic and parapodial structures. Some have spinning parapodial glands, which secrete the fibers to construct tubes. They are worldwide, occurring chiefly in tropical latitudes. Typical genera are *Euphantalis*, *Pantbali*, and *Polyodontes*.

Sigalionidae. Known through 10 genera and about 130 species, they are sometimes large and long and seldom short and small (*Pholoe*). The best-known genera are *Leanira*, *Psammolyce*, *Sigalion*, *Stbenelais*, and *Thalenessa*. They are worldwide and occur in all latitudes, from littoral to abyssal depths.

Other families. The other 5 families are small, each with distinctive characteristics: The Peisidicidae have 1 genus and 3 species; the Pareulepididae, 1 genus and 8 species; the Pisionidae, 5 genera and 12 species; the Chrysopetalidae, 4 genera and 18 species; and the Palmyridae, 1 genus and 2 species.

Amphinomorpha families. There are 3 families included in this group. The Amphinomidae, with 18 genera and about 120 species, comprise the stinging or fire worms; a burning sensation results from handling them. Common genera are the large, brilliantly colored *Amphinome* and *Hermodice*—the “scorpions of the sea”—which may be 10–16 in. (25–40 cm) long and as thick as a finger. Most have a characteristic caruncle. They abound in circumtropical intertidal zones among corals and corallines. The chalky white, very numerous bristles can penetrate

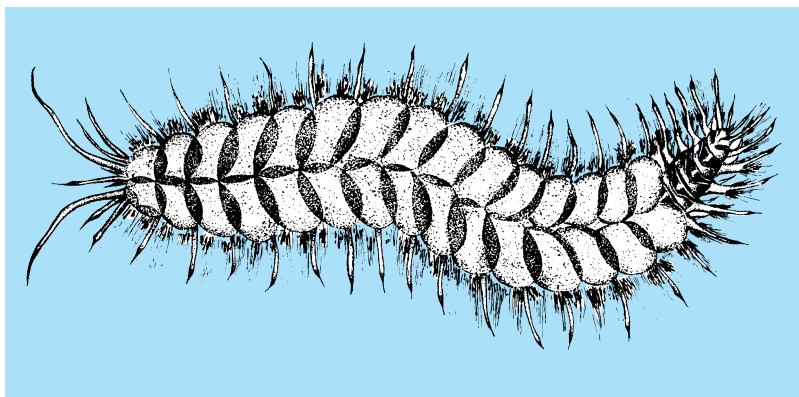


Fig. 2. *Harmothoe*, of the Polynoidae, dorsal view.

human skin and cause discomfort which may last for hours. The Euphrosinidae (Fig. 3) and the Spintheridae (Fig. 4) are smaller; the first family has 2 genera and 45 species and the second 1 genus and 8 species. They are short and depressed and are associated with sponges and other colonial animals.

Leaf-bearing and pelagic families. The leaf-bearing and pelagic groups are composed of 8 families. The Phyllodocidae, with about 30 genera and more than 240 species, are often brilliantly iridescent and are highly motile, with arresting colors and ornamentation. Most are long and linear; others are short

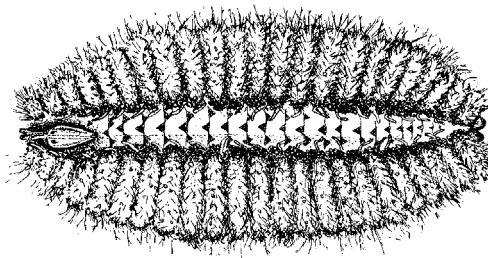


Fig. 3. *Euphrosine*, of the Euphrosinidae, dorsal view.

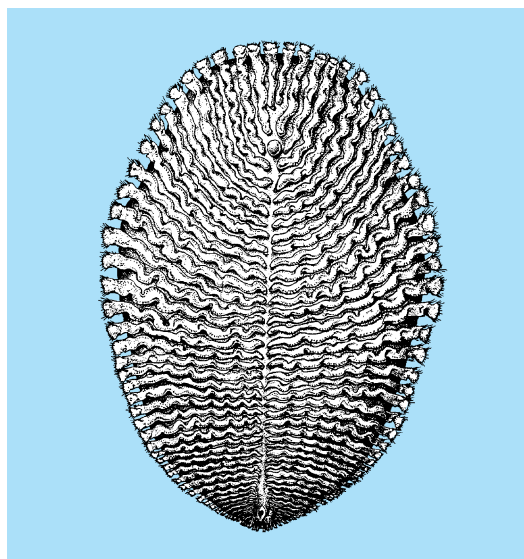


Fig. 4. *Spinther*, of the Spintheridae, dorsal view.

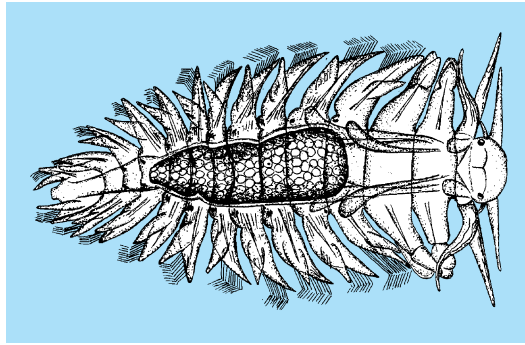


Fig. 5. *Lopadorrhynchus*, of the Lopadorrhynchidae, dorsal view.

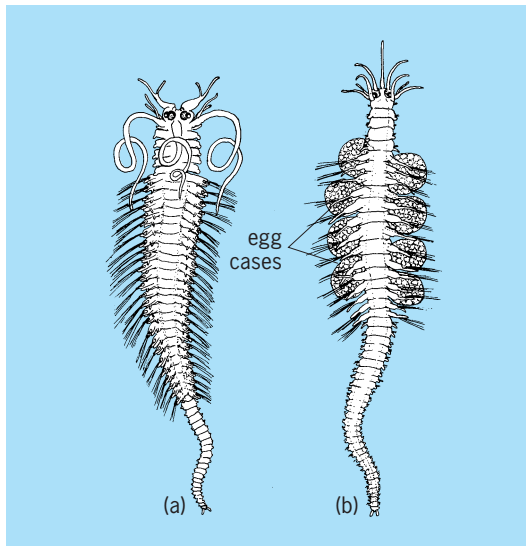


Fig. 6. *Procerastea* of the Syllidae (Autolytinae). (a) Epitokous male, dorsal view. (b) Epitokous female with attached egg cases, dorsal view.

and tumid. They occur in all seas and at all depths. Common genera are *Anaitides*, *Eulalia*, *Eteone*, and *Phyllodoce*.

The pelagic families allied to the Phyllodoceidae are the Alciopidae, with 10 genera and 41 species;

the Lopadorrhynchidae (Fig. 5), with 2 genera and 5 species; the Iospilidae, with 4 genera and 7 species; the Pontodoridae, with a single genus and species; and the Typhloscolecidae, with 3 genera and 14 species. The benthic Lacydonidae have 3 genera and 6 species, and the aberrant pelagic Tomopteridae, or glass worms, 2 genera and about 40 species. The development of *Tomopteris* demonstrates that its affinities with other polychaetes are doubtful, the resemblances being a result of evolutionary convergence.

Hesionidae and Pilargidae. The Hesionidae have 22 genera and more than 82 species, and the Pilargidae, 6 genera and about 25 species. They differ from other families chiefly in cephalic parts. Most are small, short, and depressed; a few are long and linear. Some are free-living; others commensal. Because of their small size, they escape collection unless special efforts are exerted. The best-known genera are *Hesionie*, *Leocrates*, *Oxydromus*, and *Podarke* in the Hesionidae and *Ancistrosyllis* and *Pilargis* in the Pilargidae.

Syllidae. These worms comprise the second largest of polychaete families; they have 55–60 genera and more than 400 species. Most occur along intertidal rocky shores among epiphytic plants or animals or in encrusting sponges. They are chiefly foraging predators or scavengers, identified by their long, linear, translucent bodies with articulated cirri. Some are very minute and are often found bearing eggs or juveniles on body segments. The largest syllid, *Trypanosyllis*, found in Antarctic seas, may be 4 in. (100 mm) long. Most syllids are a few millimeters long.

The Syllidae are easily divisible into four subfamilies: the Autolytinae, with sinuous pharynx and without ventral cirri (*Autolytus* and *Myriamida*); the Eusyllinae, with a thick body and nonarticulated cirri (*Eusyllis*, *Pionosyllis*, and the luminescent *Odontosyllis*); the Syllinae, with many genera and species (*Typosyllis*, *Haplosyllis*, *Syllis*, and *Trypanosyllis*); and the *Exogoninae*, with a short, small body of few segments.

Reproduction is extremely diversified, both sexual and asexual, and sexes are sometimes dimorphic (Fig. 6). Schizogamy or reproduction by fission results in individuals with heads but no pharyngeal structures, as in *Chaetosyllis*, *Ioda*, *Tetraglene*, *Polybostrichus*, and *Sacconereis*.

Nereidae. A family known for about 27 genera and more than 360 species, they include some of the largest and best-known Errantia and are frequently used as laboratory animals. Most are marine, occurring in worldwide seas and at all depths; a few are euryhaline, entering tidal streams or inland seas, or occur in thermal springs and high mountain lakes. They have been studied for their neurosecretory cells in reproductive processes. Reproduction is usually sexual, with individuals similar or heteromorphid as epitokous adults. Development proceeds through pelagic, planktonic larvae; through lecithotropic (yolky) ova without pelagic

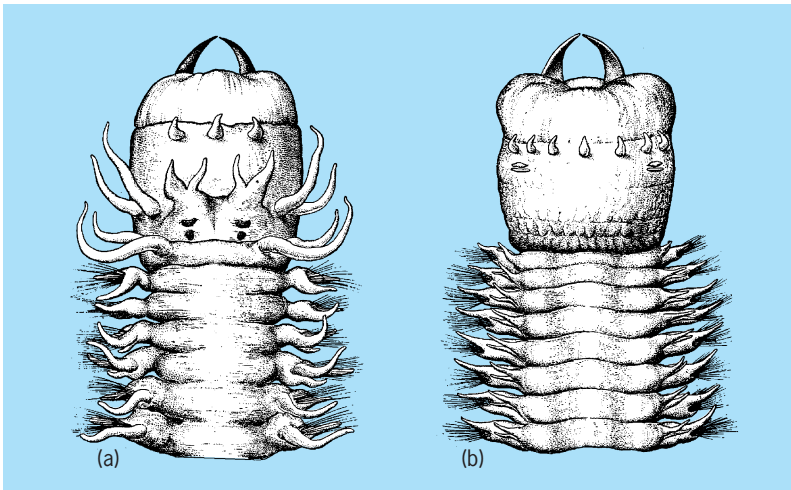


Fig. 7. *Ceratocephale*, of the Nereidae, with everted pharynx. (a) Dorsal. (b) Ventral.

larva; or by viviparity, in which juveniles are released from the female body. Common genera are *Nereis*, *Neantbes*, *Ceratocephale* (Fig. 7), *Platynereis*, and *Perinereis*. The subfamily Namanereinae occurs largely in fresh water. *Tylorrhynchus heterochaetus* is a pest in rice paddies in southeastern Asia.

Nephtyidae. With 4 genera and about 85 species, they are sometimes called sandworms for their occurrence in intertidal sands, where their highly opalescent colors and rapid movements render them almost invisible. The eversible pharynx is distinctive (Fig. 8). They occur in worldwide seas at all depths, and a few are euryhaline. Common genera are *Nephtys* and *Aglaophamus*.

Sphaerodoridae. This family has 2 genera and about 18 species; all are small (up to 1.4 in. or 35 mm long), short- to long-bodies, and usually papillated; usually their occurrence is rare.

Glyceridae. Known from 3 genera and about 66 species the Glyceridae are characterized by the enormous eversible proboscis used in food capture and forward progression. The largest genus, *Glycera*, has about 58 species and occurs in the shallowest to greatest ocean depths throughout the world. *Glycerella* is a polar genus, and *Hemipodus*, known for 16 species from shallow depths, occurs chiefly in the Western Hemisphere.

Goniadidae. Closely allied to the Glyceridae, the family Goniadidae has 7 genera and more than 60 species, from all seas and at all depths. *Glycinde*, *Goniada*, *Goniadides*, and *Ophioglycera* are worldwide in occurrence; *Goniadella* and *Progoniada* are limited to the North Atlantic; and *Goniadopsis* occurs in the Red Sea and Indian Ocean.

Eunicea. This is a superfamily which comprises 6 families. The Onuphidae are known for 10 genera and about 100 species; they are tubicolous, grazing, herbivorous, scavenging, are worldwide in occurrence, and exist in intertidal to abyssal depths. Conspicuous representatives are the quill worm (*Hyalinoecia*) and *Diopatra*, which is chiefly intertidal or in shelf depths, sometimes forming vast colonies. *Notbria* and *Paronuphis* are common in the shallow to greatest ocean depths. *Rbamphobrachium* is unique for its greatly prolonged setae in the first segment, which extend far back into the body and

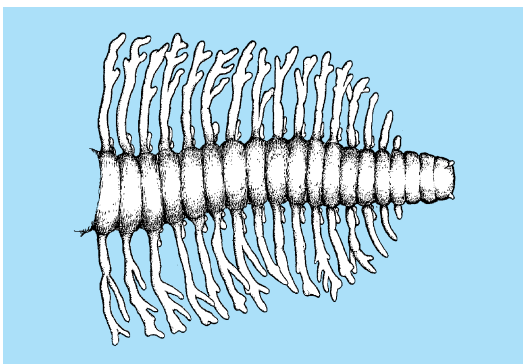


Fig. 8. Anterior end of *Aglaophamus*, of the Nephtyidae, with everted proboscis, dorsal view.

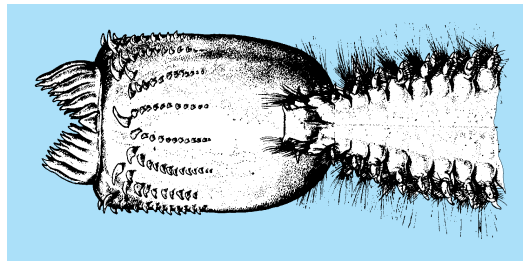


Fig. 9. *Iphitime*, of the Lysaretidae, dorsal view.

are capable of unusual forward extension, like the long shafts of knitting needles. Some onuphids display unusual care of the brood.

The Eunicidae have 7 genera and about 225 species, coming mainly from tropical seas. Included are the coral-eating *Lysidice*; the edible, swarming palolo of the South Pacific islands; and a bait worm, *Marphysa*. Some of its members are errantiate; others tubicolous. All have a highly characteristic pharyngeal armature consisting of maxillae and mandibles.

The Lumbrineridae have 3 genera and about 130 species; *Lumbrineris* is by far the largest genus and is known in all seas at all depths. *Ninoe* and *Cenogenus* are small, chiefly deep-water genera.

The Arabellidae are known for 9 genera and about 66 species. *Arabella*, *Drilonereis*, and some others are chiefly free-living, or partly parasitic, whereas *Haematocleptes*, *Oligognathus*, *Labidognathus*, and *Ophiuricola* are wholly parasitic. They retain the typical pharyngeal armature of the superfamily Eunicia.

The Lysaretidae have 3 genera and 8 species; some are large fish-bait worms (*Lysarete* and *Halla*) in tropical regions where they occur. *Iphitime* (Fig. 9) is an ectoparasite on large crabs.

The Dorvilleidae have 3 genera and about 40 species. All are small to minute and are chiefly intertidal in epiphytic growth. The largest genus is *Dorvillea*, occurring in worldwide seas; *Protodorvillea* and *Ophryotrocha* are much smaller forms.

Two parasitic families allied to the Eunicea are the Histriobdellidae, ectoparasites of crayfishes, with 2 genera and 5 species, and the Ichthyotomidae, parasites of fishes, with a single genus and species. See ANNELIDA.

Olga Hartman

Bibliography. O. Hartman, *Catalogue and Index of the Polychaetous Annelids of the World*, University of Southern California, 1959-1965; S. P. Parker (ed.), *Synopsis and Classification of Living Organisms*, 2 vols., 1982.

Erysipelothrix

A genus of gram-positive bacteria comprising at least two species, *E. rhusiopathiae* and *E. tonsillarum*. *Erysipelothrix rhusiopathiae* is the more pathogenic and causes disease in a variety of animal species, including humans. It is a rod-shaped

organism, $0.3 \times 1\text{--}2$ micrometers, that may form filaments in old cultures and in chronically infected tissues. It occurs in the tonsil of healthy swine and on the surfaces of fish and other aquatic species; it is shed in the urine, oral secretions, and feces of infected animals. Organisms deposited in the soil may survive for years and may be a source of infection for susceptible domestic and wild animals and birds.

Epidemiology and transmission. *Erysipelothrix rhusiopathiae* causes disease in swine (erysipelas), turkeys, lambs, and a variety of other domestic and wild mammals and birds. Infections are acquired via skin abrasions, or occur endogenously when organisms residing in the tonsil invade the bloodstream of swine stressed by excessive heat. The organism spreads via contaminated excretions and saliva. Infections in humans are associated with occupational exposure to the organism (veterinarians, butchers, pathologists, fish-handlers). They are usually acquired via minor cuts and abrasions on the hands or arms, leading to a painful local purplish erythema (reddening of the skin) termed erysipeloid. This infection may occasionally lead to invasion of the bloodstream or joints.

Clinical findings and pathology. In swine, the acute disease develops suddenly with high fever, prostration, and conjunctivitis, and death within a day or two of septicemia. A less severe form is associated with rhomboid, deep-red skin lesions and is called diamond skin disease. Chronic swine erysipelas takes the form of verrucous (wartlike) growths on the heart valves and a progressive immune-mediated arthritis resembling rheumatoid arthritis in humans. Polyarthritis is a common manifestation in lambs. The disease is often devastating in turkey flocks, causing mass mortality due to septicemia. A thin capsule on *E. rhusiopathiae* confers resistance to phagocytosis and intracellular killing by host phagocytes. Hyaluronidase and neuraminidase also contribute to virulence. Multiplication of the organism in, and subsequent killing of, phagocytes are critical in determining the outcome of infection. In the terminal stages, thrombi form in the small blood vessels, which may rupture and hemorrhage.

Diagnosis and treatment. *Erysipelothrix rhusiopathiae* is easily demonstrated microscopically and by culture of blood and tissues of affected animals. Penicillin is the antibiotic of choice for treatment of infected animals or humans.

Control and prevention. Protective immunity is stimulated by live attenuated strains or by inactivated virulent *E. rhusiopathiae*. A protein is involved in stimulation of protective antibodies. Live vaccines stimulate both enhanced killing activity of phagocytes and protective opsonizing antibodies. See MEDICAL BACTERIOLOGY.

John F. Timoney

Bibliography. J. F. Timoney et al., *Hagan and Bruner's Microbiology and Infectious Diseases of Domestic Animals*, 8th ed., 1988; J. F. Timoney and M. H. Groschup, Properties of a protective protein antigen of *Erysipelothrix rhusiopathiae*, *Vet. Microbiol.*, 37:381, 1993; R. L. Wood and A. D. Lemman et al. (eds.), *Diseases of Swine*, 1992.

Escalator

A continuous, moving stairway that transports passengers between two levels. Most escalators are installed to service floors less than 20 ft (6 m) apart, but some escalators have been constructed to rise 100 ft (30 m) or more between floors.

The modern escalator consists of a series of steps, fastened to an endless chain running in a track system that extends up an incline. The steps continue around a turnaround at the upper landing, pass back down the incline under the steps available to the passengers, and progress to another turnaround at the lower landing. This chain, known as a step chain, is powered by a machine that is usually connected by a gear, chain, or mechanical coupling to a shaft located at the upper landing turnaround. An alternating-current induction motor, selected for its simplicity and constant speed characteristics, is used to move the steps. Most modern escalators travel at 100 ft/min (0.5 m/s). See INDUCTION MOTOR.

The track system is designed to permit the steps to approach the upper and lower landings horizontally and to maintain a horizontal tread surface wherever the steps will be holding passengers. The treads of the steps are provided with cleats that match a pattern of combs at the landings, and each riser (the vertical front face of the step) is fitted with cleats that mesh with the adjacent step. The cleats are a safety feature, reducing the possibility of objects becoming caught at the landings or between steps.

The escalator has handrails on each side of the moving steps. They move at the same speed as the steps and extend a distance beyond the point where the steps pass under the combs at the landings. This provides support for passengers as they move onto or off the moving steps.

The escalator mechanism is covered with a balustrade that creates a safe environment for the passengers and enhances the product architecturally. There are two types of balustrades, glass and opaque. The glass type has glass panels below the handrails and gives the escalator an open appearance. The opaque type uses nontransparent panels instead of glass. Glass is usually preferred in installations in open atriums found in building lobbies or shopping malls.

The use of escalators is widespread because of their continuous accessibility and their ability to transport large numbers of people. The capacity, in passengers per hour, is determined by the step width and the speed at which the escalator travels. The most common step widths are 24 and 40 in. (60 and 100 cm). The capacity is based on one passenger per 40-in. (100-cm) step or one passenger on every other 24-in. (60-cm) step. These values are reliable for planning the number of escalators required in a building, even though theoretically a larger capacity can be achieved over a short time. Based on such step loading, an escalator running at 100 ft/min (0.5 m/s) can transport 4500 or 2250 passengers per hour with 40- or 24-in. (100- or 60-cm) step widths, respectively.



Cutaway view of an escalator. (Otis Elevator Co.)

Movement of people by escalators is very economical. The escalator is a balanced machine when operating empty. As passengers board in the down direction, they offset the friction in the system. When the down load exceeds approximately 10% capacity, power is generated by the escalator so that the total power consumption of the building containing the unit is reduced. Although power is required to carry passengers up, the cost to carry a passenger up one floor on a fully loaded escalator is extremely small, typically less than a fraction of a cent. D. L. Steel

Escape velocity

Minimum speed away from a parent body that a particle must acquire to escape permanently from the gravitational attraction of the parent. Escape velocity is also termed parabolic velocity. See ORBITAL MOTION.

Earth retains an atmosphere because the escape velocity is considerably higher than the mean velocity of the gas molecules in its atmosphere. For a space ship to escape from Earth and travel to another planet or orbit about the Sun, it must reach escape velocity. This velocity v can be calculated by equating the kinetic energy of the moving body m to the work necessary to overcome the gravity g_0 at the surface of the parent whose radius is r_0 in rising to a height b above the surface of the parent; thus Eq. (1) holds. From this Eq. (2) is derived, where for Earth $r_0 = 6.38 \times$

$$\frac{mv^2}{2} = m \int_0^b g_0 \left(\frac{r_0}{r_0 + b} \right)^2 db = \frac{mg_0 r_0 b}{r_0 + b} \quad (1)$$

$$v = \sqrt{2g_0 r_0} \sqrt{\frac{b}{r_0 + b}} \quad (2)$$

10^6 m and $g_0 = 9.8 \text{ m/s}^2$ and for escape $b = \infty$ so that the term under the second radical is unity; thus $v_{\text{escape}} = 37.0 \times 10^3 \text{ ft/s}$ ($11.2 \times 10^3 \text{ m/s}$). For comparison, the mean velocity of a gas with molecules of mass m at absolute temperature T is $v = \sqrt{3kT/m}$, where k is Boltzmann's constant. Thus hydrogen has a mean velocity in the vicinity of $6.6 \times 10^3 \text{ ft/s}$ ($2 \times 10^3 \text{ m/s}$) and heavier gases in an atmosphere have lower velocities. Escape velocities for other bodies in the solar system are Mercury, $12.5 \times 10^3 \text{ ft/s}$ ($3.8 \times 10^3 \text{ m/s}$); Venus, $34.3 \times 10^3 \text{ ft/s}$ ($10.4 \times 10^3 \text{ m/s}$); Mars, $16.8 \times 10^3 \text{ ft/s}$ ($5.1 \times 10^3 \text{ m/s}$); Jupiter, $201.3 \times 10^3 \text{ ft/s}$ ($61.0 \times 10^3 \text{ m/s}$); Saturn, $121.1 \times 10^3 \text{ ft/s}$ ($36.7 \times 10^3 \text{ m/s}$); Uranus, $72.6 \times 10^3 \text{ ft/s}$ ($22 \times 10^3 \text{ m/s}$); and Neptune, $79.2 \times 10^3 \text{ ft/s}$ ($24 \times 10^3 \text{ m/s}$). See CELESTIAL MECHANICS; SATELLITE (SPACECRAFT); SPACE FLIGHT. R. L. Duncombe

Escapement

A mechanism in which a toothed wheel engages alternate pallets attached to an oscillating member. The escapement is found principally in timepieces but may be employed wherever oscillating motion is required. Its origins are ancient and obscure.

In the simple form of escapement (Fig. 1), oscillating member cc' is an open bar arranged to slide longitudinally in bearings CC , which are attached through a frame (not shown) to the bearing for toothed wheel a . Wheel a turns continuously in the direction of the arrow, and is provided with three teeth b , b' , and b'' . The oscillating member has two pallets c and c' . In the position shown, tooth b is just ceasing to drive pallet c to the right, and is escaping, while tooth b' is just coming in contact with pallet c' , which it will drive to the left.

Although escapements are generally used to convert circular into reciprocating motion, as in the above example, the wheel being the driver, in many cases the action may be reversed. In Fig. 1, if the open slide bar were driven with reciprocating motion, the wheel would be made to turn in the opposite direction from its rotation as the input member. Also, there would be a short interval at the beginning of each stroke of the bar in which no motion would be given to the wheel. The wheel *a* must have one, three, five, or any other odd number of teeth upon its circumference.

Escapements are also adapted for use in automatic factory equipment to regulate the flow of parts from the magazine or feeder track. Such feeding and spacing devices are of various forms, many being adaptations of the ratchet and pawl. See PAWL; RATCHET.

In a mechanical clock or watch, the escapement intervenes between the energy source (spring or elevated weight) and the regulating device (pendulum or balance wheel). It is acted upon by both. The escape wheel is mounted on the same shaft as the last wheel of the gear train, and impulses are delivered from the escape wheel to operate the regulating device. The regulating device, which has a natural period of oscillation, determines the rate at which it

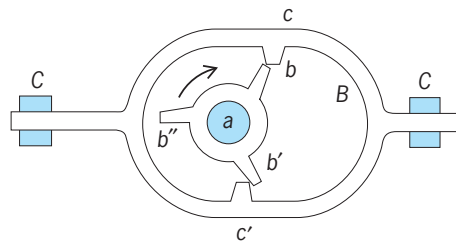


Fig. 1. Simple form of escapement.

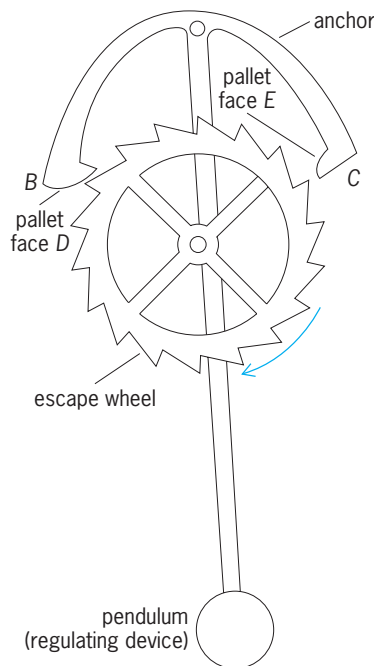


Fig. 2. Anchor recoil escapement.

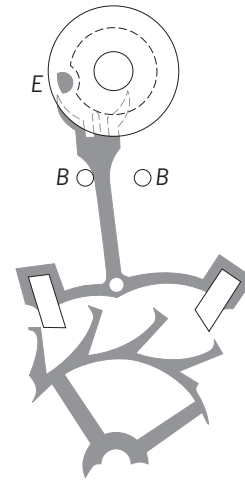


Fig. 3. Detached-lever escapement, a type used in watches.

will receive these impulses, and thus regulates the rate of going of the timepiece.

The anchor recoil escapement (Fig. 2) is used with a pendulum and takes its name from the shape of its oscillating member and its action. This type of escapement appeared late in the seventeenth century and has survived, with modifications, to the present. To simplify the explanation, the pendulum is shown attached directly to the anchor, so that the pendulum and anchor swing as one.

In the position shown, pallet *B* has just received an impulse from the escape wheel, the impulse swinging the pendulum to the left. When pallet *B* has cleared the tooth, allowing the wheel to escape, pallet *C* will be in position to arrest the wheel. Recoil, or momentary reversal of the escape wheel, occurs just after it is arrested because the pendulum has not quite completed its swing. The wheel tooth in contact with pallet *C* will then give the oscillating parts an impulse in the opposite direction.

Deadbeat, an escapement without recoil, has arresting faces of the pallets described by a circular arc whose center is at the pivot point of the anchor. Escape-wheel teeth are contoured to give impulses to these pallet faces, over which they slide without recoil.

Modern watches generally employ a detached-lever escapement (Fig. 3), which has banking pins *B* to limit the oscillation of the anchor and its lever. An escapement is termed detached when the regulating device, in this case the balance wheel, is given an impulse during only a small part of its operating cycle. When the fork reaches the end of its swing, it is lightly locked by a wheel tooth and remains stationary until the returning impulse pin *E* causes sufficient recoil of the escape wheel to release the pallet. The chronometer escapement is a detached escapement that furnishes an impulse in only one direction of the swing of the balance wheel. See CHRONOMETER; CLOCK (MECHANICAL).

Some watches employ a light ratchet powered by a vibrating reed at 360 Hz in place of a balance wheel and escapement. See WATCH. Douglas P. Adams

Bibliography. H. G. Harris, *Handbook of Watch and Clock Repairs*, rev. ed., 1972; H. H. Mabie and F. W. Ocvirk, *Mechanisms and Dynamics of Machinery*, 4th ed., 1987.

Escarpment

A long line of cliffs or steep slopes that break the general continuity of the land by separating it into two level or sloping surfaces. Some very high escarpments, or scarps, may form by vertical movement along faults. Often a whole block of land may be forced upward while the adjacent block is downfaulted. Such scarps are common in the tilted fault-block mountains of eastern California, Nevada and western Utah. In eastern Africa, prominent fault scarps mark the margins of the great rift valleys, whose floors are downfaulted by as much as 8200 ft (2500 m). See FAULT AND FAULT STRUCTURES.

Other types of escarpments form by differential weathering and erosion of contrasted rock types. Less resistant rocks, such as clay or shale, are often eroded from beneath resistant cap rocks, such as sandstone and limestone. With support removed from below, the cap rock fails and the escarpment retreats. Escarpments are often very prominent in arid regions, where hardened weathering products may form extensive cap rocks known as

duricrusts. When inclined strata are eroded, they may produce cuesta escarpments, which have back slopes that approximate the dip of their sedimentary layers and steeper facing slopes (the scarp faces) that truncate the bedding.

Some of the largest known escarpments occur on the planet Mars, where erosion has presumably been much slower than on the Earth in reducing primary structural relief. The great shield volcano Olympus Mons is surrounded by a basal scarp that forms a circle approximately 430 mi (700 km) in diameter. Even more spectacular are the great escarpments that bound the Valles Marineris, a group of steep-walled equatorial canyons more than 2500 mi (4000 km) in length. The escarpments are as much as 4 mi (6 km) high and contain the scars of great landslides that contributed to scarp retreat (see *illus.*).

Victor R. Baker

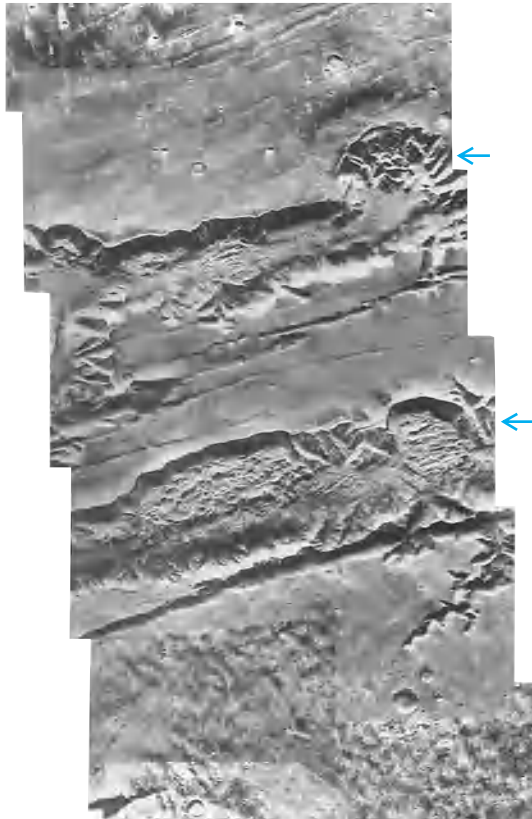
Escherichia

A genus of bacteria named for Theodor Escherich, an Austrian pediatrician and bacteriologist, who first published on these bacteria in 1885. *Escherichia coli* is the most important of the six species which presently make up this genus, and it is among the most extensively scientifically characterized living organisms. *Escherichia coli* are gram-negative rod-shaped bacteria approximately $0.5 \times 1-3$ micrometers in size.

Escherichia coli can be differentiated from other closely related bacteria by a variety of assays for metabolic processes, including the fermentation of specific sugars and the enzymatic modification of amino acids. Modern molecular taxonomic analysis based on the nucleotide sequences of ribosomal ribonucleic acid (RNA) has revealed that *Shigella*, a bacterial genus of medical importance previously thought to be distinct from *E. coli*, is actually the same species.

The natural habitat of *E. coli* is the colon of mammals, reptiles, and birds. In humans, *E. coli* is the predominant bacterial species inhabiting the colon that is capable of growing in the presence of oxygen. The presence of *E. coli* in the environment is taken to be an indication of fecal contamination.

Most strains of *E. coli* are commensals which are harmless to the humans and other animals they colonize, but some strains can cause disease when given access to extraintestinal sites or the intestines of non-commensal hosts. *Escherichia coli* is the most important cause of urinary tract infections. Specialized protein structures on the surface of *E. coli*, known as adhesive fimbriae, permit some strains to colonize the epithelial cells lining the urinary tract, leading to inflammation and tissue damage. Women are more susceptible than men to *E. coli* urinary tract infections due to the proximity of the urethra to the anus. Four out of ten women experience at least one urinary tract infection in their lifetime. *Escherichia coli* urinary tract infections are also associated with catheterization and instrumentation of the urinary



Tithonium Chasma (top) and Ius Chasma (bottom) in the great equatorial canyon system of Mars. The escarpments bounding them are 1.2–3.7 mi (2–6 km) high. The photographed scene is about 155 mi (250 km) across. Arrows indicate scars of the great landslides. (NASA)

tract. Men and women are equally susceptible in this setting. Urinary tract infections may extend into the bloodstream, especially in hospitalized patients whose defenses are compromised by the underlying illness. This may lead to a type of whole-body inflammatory response known as sepsis, which is frequently fatal. Certain *E. coli* strains with a specific polysaccharide capsule on their surface can invade the intestine of the newborn and cause sepsis and meningitis. These strains are acquired at birth from *E. coli* which have colonized the vagina of the mother.

Several different strains of *E. coli* cause intestinal infections. In the developing world, the most important of these are the enterotoxigenic *E. coli*. These strains produce specialized adhesive fimbriae which enable them to colonize the small intestine. Enterotoxigenic *E. coli* produce enterotoxins which act on the epithelial cells lining the small intestine, causing the small intestine to reverse its normal absorptive function and secrete fluid. This leads to a dehydrating diarrhea which can be fatal, especially in poorly nourished infants. Therapy consists of oral or, in serious cases, intravenous rehydration. Enterotoxigenic *E. coli* are transmitted by ingestion of fecally contaminated water and food, and are a common cause of diarrheal disease in travelers in developing countries. Enterotoxigenic *E. coli* are also an important cause of diarrheal disease in livestock animals, especially calves and piglets. These strains do not infect humans.

An important group of pathogenic *E. coli* in developed countries are the enterohemorrhagic strains, especially the serotype known as *E. coli* O157:H7. These strains are normal commensals in cattle but cause bloody diarrhea in humans. A complication of approximately 10% of cases is a potentially fatal disease known as hemolytic uremic syndrome. The virulence of these strains involves the close attachment of bacteria to epithelial cells lining the colon, resulting in alteration of the epithelial cell structure, and the production of Shiga toxin. The toxin enters the bloodstream after being absorbed in the colon and damages the endothelial cells lining the blood vessels of the colon, resulting in bloody diarrhea. In cases of hemolytic uremic syndrome, the toxin circulating in the blood damages blood vessels in the kidney, resulting in kidney failure and anemia. Antibiotics are not recommended at present due to concern that the action of antibiotics on the bacterial cells will increase the release of toxin. Enterohemorrhagic *E. coli* are acquired by the ingestion of undercooked beef, uncooked vegetables, or unpasteurized juices from fruits which have been contaminated with the feces of infected cattle. An infection can also be acquired from contact with a human infected with the organism and from contaminated water. Children and the elderly are at greatest risk of developing hemolytic uremic syndrome.

Other strains which are pathogenic in the human colon include the enteroinvasive *E. coli* (including *Shigella*) and the enteropathogenic *E. coli*. Enteroinvasive *E. coli* induce the epithelial cells lining the

colon to engulf the bacteria. The bacteria then gain access to the interior of the epithelial cells, where they rapidly multiply and destroy the cell. An inflammatory response ensues, and the bacteria are ultimately destroyed by neutrophils. The disease is called bacillary dysentery and is characterized by bloody diarrhea. It can be fatal in poorly nourished children but is not usually associated with serious complications in developed countries, where it remains a common cause of diarrheal disease. Acquisition of enteroinvasive *E. coli* and *Shigella* is by ingestion of food or water contaminated with feces of a person infected with the organism, or by direct contact with infected feces. The disease is usually treated with antibiotics.

Enteropathogenic *E. coli* colonize the small intestine and interact closely with epithelial cells in its lining, altering the structure and function of the cells in a manner similar to enterohemorrhagic *E. coli*. Unlike enterohemorrhagic *E. coli*, enteropathogenic *E. coli* do not release toxins, and damage to the intestine appears to result from structural alterations to the cells. Enteropathogenic *E. coli* have been associated with protracted diarrhea in infants and can occasionally cause severe wasting. There are no routine methods available for the diagnosis of enteropathogenic *E. coli* infections, and the prevalence of disease caused by these organisms is unknown. Antibiotics are of uncertain efficacy. See DIARRHEA; TOXIN. Steve L. Moseley

Bibliography. J. B. Kaper and A. D. O'Brien (eds.), *Escherichia coli O157:H7 and Other Shiga Toxin-Producing E. coli Strains*, ASM Press, 1998; H. L. T. Mobley and J. W. Warren (eds.), *Urinary Tract Infections: Molecular Pathogenesis and Clinical Management*, ASM Press, 1996; A. A. Salyers and D. D. Whitt, *Bacterial Pathogenesis: A Molecular Approach*, ASM Press, 1994.

Esociformes

An order of teleost fishes, also known as Haplomi and Esocae. Fishes of this small order, comprising two families and 10 species, can be identified by a combination of the following characteristics: dorsal and anal fins are posteriorly placed on a rather slender body; maxillae are toothless and in the gape of the mouth; small uniform teeth are on the tongue and basibranchial elements; and absent are adipose fin, breeding tubercles, pyloric caeca, and mesocoracoid bones. Both families occur in freshwater of the Northern Hemisphere. See OSTEICHTHYES; TELEOSTEI.

Esocidae (pikes). Esocids have a duckbill-like snout; a forked caudal fin; large canine teeth on the dentary, vomer, and palatine bones; small embedded cycloid scales; a complete lateral line; eight or more pores in the infraorbital canal; and abdominal pelvic fins (**Fig. 1**). The family consists of one extant genus, *Esox*, and five species: *E. lucius* (northern pike) is circumpolar in distribution; *E. reicherti* is endemic to the Amur River in eastern Asia; and three species—*E. americanus* (with two

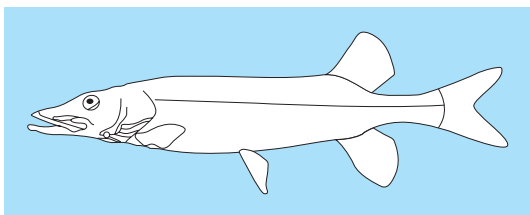


Fig. 1. Example of an esocid (pike). (Courtesy of J. S. Nelson, 2006)

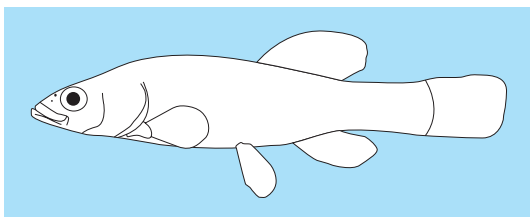


Fig. 2. Example of an umbrid (mudminnow). (Courtesy of J. S. Nelson, 2006)

subspecies, the redbfin pickerel and the grass pickerel), *E. masquinongy* (muskellunge), and *E. niger* (chain pickerel)—are restricted to North America east of the Rocky Mountains.

Esocids are diurnal predators that feed mainly on other fishes, which are attacked from ambush. The muskellunge is the largest species, reaching a length of 1.4 m (4.6 ft).

Umbridae (mudminnows). Mudminnows are distinguished from esocids by the following characteristics: short snout, rounded caudal fin, faint or absent lateral line, and three or fewer pores in the infra-orbital canal (Fig. 2). The family comprises three genera and five species: *Dallia pectoralis* (Alaska blackfish) occurs in Alaska and northeastern-most Siberia; *Novumbra hubbsi* (Olympic mudminnow) is limited to the Olympic peninsula in western Washington; *Umbra limi* (central mudminnow) occurs in east-central North America; *U. pygmaea* occurs in eastern United States; and *U. krameri* inhabits the Danube and Dniester river systems in southeastern Europe.

Umbrids inhabit densely vegetated swamps, bogs, ponds, sluggish creeks, as well as large rivers and lakes that support abundant vegetation. They are usually found on or near a substrate of mud or silt and feed on adult and larval aquatic insects, larval crustaceans, and snails. These hardy fishes are capable of surviving extreme cold and oxygen-depleted water. The maximum total length is 8–33 cm (3–13 in.), depending on the species.

Herbert Boschung

Bibliography. T. Grande, H. Laten, and J. A. Lopez, Phylogenetic relationships of extant esocid species (Teleostei: Salmoniformes) based on morphological and molecular characters, *Copeia*, 2004(4):743–757, 2004; J. A. Lopez, P. Bentzen, and T. W. Pietsch, Phylogenetic relationships of esocoid fishes (Teleostei) based on partial cytochrome b and 16S mitochondrial DNA sequences, *Copeia*, 2000(2):420–431, 2000; J. A. Lopez, W. J. Chen, and G. Ortí, Esociform phylogeny, *Copeia*, 2004(3):449–464, 2004; J. S. Nelson, *Fishes of the World*, 4th ed., Wiley, New York, 2006.

Esophagus

A section of the alimentary canal that is interposed between the pharynx and the stomach. Because of divergent specializations in the various vertebrates, the esophagus cannot be described in general terms and is not always distinguishable.

In humans it is a tube running the full length of the neck and the thorax, held in its position ventral to the vertebral centra by a tunica adventitia of loose connective tissue. It has an inner lining of folded mucous membrane with an exceptionally thick lamina propria, a submucosa of elastic and collagenous connective tissue, and two layers of muscle. The musculature is striated in the anterior third of its length, unstriated in the posterior third, and variably intermixed in the middle. It is supplied with autonomic nerve fibers.

Although normally collapsed, the human esophagus is capable of considerable distension during the rapid passage of swallowed material, under which condition the folds of mucous membrane and lamina propria are temporarily smoothed out. Numerous microscopic esophageal glands open into the lumen, extending their compound tubules out into the submucosa.

In humans the transition from the esophagus to the stomach occurs quite abruptly at the diaphragm. In many vertebrates the distinction between the two is not so clear. The pharynx narrows posteriorly like a funnel and the foregut may thereupon enlarge, but much of what appears to be stomach may have an esophageal character histologically. In a number of fish species there is no stomach enlargement at all. There is a metabranchial foregut anterior to the entrance of ducts from the liver and pancreas, but application to it of the names stomach and esophagus is of questionable value. At the other extreme, in some sharks and bony fishes there is a constricted esophagus with a sparsely glandular or cornified lining. The esophagus of some turtles is studded with harsh horn-tipped papillae. In some of the carnivorous birds the esophagus may be temporarily dilated for the storage of large masses of swallowed food before they can be digested. Many graminivorous (grain- and seed-eating) birds are provided with a permanent ventral esophageal pouch, the crop, for that purpose. Pigeons of both sexes feed their newly hatched young with a thin paste consisting of cells proliferated and sloughed from special areas of the crop under the stimulus of the hormone prolactin. This is pigeon's milk. See DIGESTIVE SYSTEM.

William W. Ballard

Essential oils

Volatile, fragrant oils obtained from plants. Essential oils are distinguished from those known as fixed oils, which are mainly triglycerides of fatty acids. See FAT AND OIL.

Essential oils have been obtained from over 3000 plants and are designated and defined by the plant species and sometimes the geographical location.

The sources of these oils are diverse, including flower petals (for example, rose and jasmine), spices (cinnamon and ginger), pine oil and turpentine, and citrus fruit peels. Compounds present in the juice that may contribute to the distinctive flavor of a fruit or berry are not, strictly speaking, components of the essential oil. Chemically, essential oils are extremely complex mixtures containing compounds of every major functional-group class. The oils are isolated by steam distillation, extraction, or mechanical expression of the plant material; often only certain parts, such as roots, buds, leaves, or flower petals, are used.

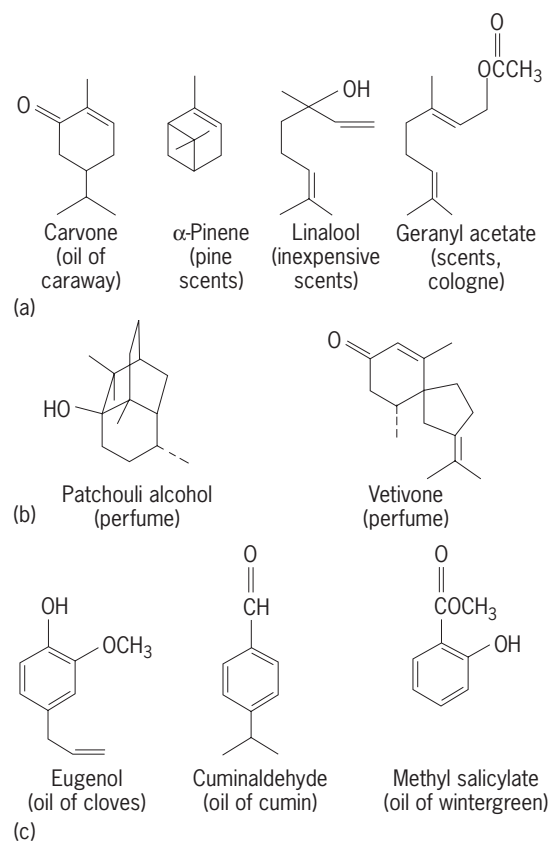
History. Essential oils have been produced and used for flavoring, incense, and medicinal purposes for many centuries. There are indications that crude preparations were known in ancient Egypt and Persia. Distillation of the essence of plants was described about the year 1500.

In the beginnings of modern organic chemistry, essential oils provided a fertile source of compounds for studying the structures and reactions of organic molecules. Early work revealed that several compounds from the oils of aromatic plants had in common a chemically inert structural unit, and these compounds were designated as aromatic. This unit was later determined to be a benzene ring (C_6H_6). From this association the modern term aromatic arose to denote the electronic structure of a large class of compounds of which benzene is the parent. Chemical investigations of essential oils also led directly to the development of terpene chemistry. *See* BENZENE; TERPENE.

Chemistry. Among the hundreds of compounds that have been separated from essential oils are many monoterpenes (10 carbons, giving 2 isoprene units) and a smaller number of sesquiterpenes (15 carbons giving 3 isoprene units). These terpenes may be open-chain, monocyclic, or polycyclic, and usually contain one or more double bonds and hydroxyl groups. Also present, generally in very small amounts, are numerous saturated aliphatic aldehydes and carboxylic esters. The major constituent of several important essential oils is a di- or trisubstituted benzene derivative. A few representative examples are shown in the **illustration**.

Monoterpene hydrocarbons make up the bulk of many essential oils that are used in flavoring foods and beverages. Some of the most valuable oils, used in perfumery, are obtained from tropical plants and contain a high proportion of bi- and tricyclic sesquiterpenes.

At one time essential oils were the principal or only source available for these compounds, but synthetic methods have advanced to the point that any compound that can be isolated from an essential oil can now be prepared by synthesis. Large-scale processes are used to manufacture a number of the major components present in commercially important oils. Linalool and geraniol can be prepared from acetone or isobutene, formaldehyde, and acetylene. Racemic menthol is available from thymol, which is in turn manufactured from toluene. A number of monoterpenes, in optically active form, are derived by simple reactions from the abundant pinenes in turpentine.



Examples of benzene derivatives in some essential oils: (a) monoterpenes, (b) sesquiterpenes, and (c) aromatics.

On the other hand, the most practical source of certain valuable perfume ingredients, such as patchouli alcohol and related complex sesquiterpenes, is the essential oil. *See* MENTHOL.

Characterization and analysis. The importance of essential oils for perfumes or flavors lies in the fact that they are extremely complex mixtures of as many as 200–300 compounds, usually a few major constituents and many minor ones. Distinctive odors and tastes may depend on the blend of all these components. In early work, only a few of the main constituents were isolated by fractional distillation and identified by various chemical methods. The total mixture composing the oil is characterized by such physical properties as color, boiling-point range, specific gravity, refractive index, and optical rotation. These properties and sometimes also information on chemical composition, such as total hydroxyl or aldehyde content, are used to specify essential oils of acceptable medicinal quality in the United States Pharmacopeia. For use as flavors or scents, organoleptic tests are an important part of the evaluation.

The development of gas chromatography in the 1960s and then the use of very long capillary columns, typically 0.001 in. \times 165–330 ft (0.25 mm \times 50–100 m), with a film of absorbent on the wall revolutionized the analysis of essential oils. With these wall-coated open-tube columns, coupled to the inlet of a mass spectrometer, it is possible to separate, identify, and quantitatively estimate over 100 compounds. Others can be separated and characterized

by general type, and still many more can be detected. See GAS CHROMATOGRAPHY.

Other analyses have been based on the infrared spectra of the components separated by capillary gas chromatography columns. An additional physical constant that is used in characterizing the individual terpenes and other components of essential oils is the Kovats index, which is a measure of the retention time of a given compound in a gas chromatogram relative to the retention times of *n*-alkane reference compounds. Carbon-13 nuclear magnetic resonance spectra of an essential oil can also be used for identification of the major components of the mixture. See NUCLEAR MAGNETIC RESONANCE (NMR).

Variability from plant sources. The formation of essential oil in the plant, and consequently the yield and composition of the oil produced, depends on many factors. Genetic differences in plants of the same species that are otherwise indistinguishable (chemotypes) can result in widely different essential oil content. Geographic location and agricultural factors also influence oil production. These include soil, water, nutrients, and climatic variables, such as sunlight, temperature, and day length. The yield of oil from peppermint (*Mentha piperita*) grown in central Washington is threefold greater than that from the same strain grown in heavier soil in the midwestern United States.

Both the quantity and composition of the essential oil can change drastically as the plant matures. In the oil from coriander (*Coriandrum sativum*), the content of aliphatic aldehydes drops and that of the monoterpene alcohols increases tenfold from the stage of full flowering to green fruit.

Processing. The composition and quality of an essential oil is affected significantly by the method of isolation and subsequent processing steps. Enzymatic action (fermentation) in the crude plant material prior to distillation or extraction can bring about hydrolysis of glycosides and release of oil components. The composition of the oil may depend greatly on the isolation method. Steam distillation is the most common process, but sensitive compounds can undergo rearrangement or oxidation on heating. Oils from flower blossoms are extracted by pressing the petals with a purified fat (enfleurage process) or hydrocarbon solvent, and the essential oil is then extracted with alcohol. The alcohol solution is concentrated to give a liquid known as an absolute that is used in the manufacture of perfumes. Extraction of the plant with alcohols can lead to artifacts due to reaction with the solvent. See DISTILLATION; EXTRACTION.

Further processing steps depend on the ultimate use of the essential oil. Fractional distillation (rectification) at reduced pressure can be carried out to remove some of the terpene hydrocarbons and produce a concentrate (folded oil) in which the characteristic flavor or scent is enhanced. Distillation or chemical treatment may also be used to remove an undesired trace component.

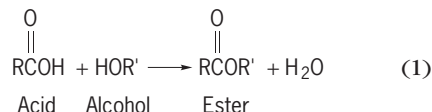
James Moore

Bibliography. E. Guenther, *The Essential Oils*, vols. 1-6, 1947-1952; B. D. Mokerjee and C. J.

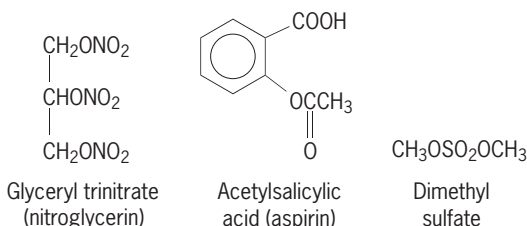
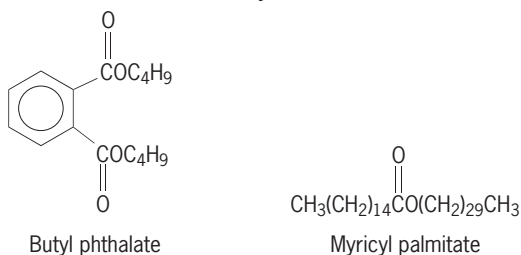
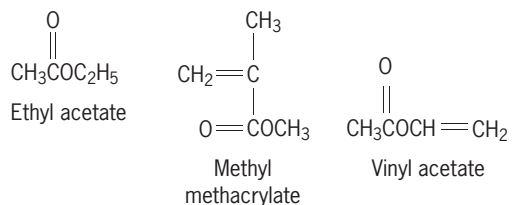
Mussinan (eds.), *Essential Oils*, 1981; *Perfumer and Flavorist*, vols. 1-10, 1976-1986; R. Teranishi, R. G. Buttery, and S. Sugisawa (eds.), *Bioactive Volatile Compounds from Plants*, ACS Symposium Series 525, 1993.

Ester

The product of a condensation reaction (esterification) in which a molecule of an acid unites with a molecule of alcohol with elimination of a molecule of water, reaction (1).



At one time it was thought that esterification was analogous to neutralization, and esters are still named as though they are "alkyl salts" of carboxylic acids, as shown in the following examples:



Properties and uses. Esters are generally insoluble in water and have boiling slightly higher than hydrocarbons of similar molecular weight. An ester may often be characterized by its infrared absorption spectrum. For example, saturated aliphatic esters have a distinctive, strong C=O absorption band at 1750-1735 cm⁻¹, and a second band assigned to C—O stretching at 1300-1000 cm⁻¹.

Ethyl and butyl acetates are volatile industrial solvents, used particularly in the formulation of lacquers. The ethyl acetate produced in the United States is primarily used as a solvent. Higher-boiling esters such as butyl phthalate are used as softening agents (plasticizers) in the compounding of plastics. The natural waxes of biological origin are largely simple esters. For example, a principal component of

beeswax is myricyl palmitate. *See* SOLVENT; WAX, ANIMAL AND VEGETABLE.

Esters of cellulose (cellulose triacetate) are used in photographic film and as a textile fiber (acetate rayon). Cellulose acetatepropionate and cellulose acetatebutyrate have become important as thermoplastic materials. Cellulose nitrate, containing 10.5–11% nitrogen, is called celluloid pyroxylin; with alcohol and camphor (a plasticizer), it forms celluloid. Dynamite cotton is cellulose nitrate of 11.5–12.3% nitrogen content, and gun cotton is cellulose nitrate of 12.5–13.5% nitrogen. Cordite and ballistite are made from gun cotton, which is plasticized with glyceryl trinitrate (nitroglycerin). Dimethyl and diethyl sulfates (esters of sulfuric acid) are excellent agents for alkylating organic molecules that contain labile hydrogen atoms, for example, starch and cellulose. *See* CELLULOSE; EXPLOSIVE.

Esters of unsaturated acids, for example, acrylic or methacrylic acid, are reactive and polymerize rapidly, yielding resins; thus, methyl methacrylate yields a poly(methyl methacrylate). Analogously, esters of unsaturated alcohols are reactive and readily react with themselves; thus, vinyl acetate polymerizes to poly(vinyl acetate). The polyester resins known as glyptals result from the polyesterification of glycerol with phthalic anhydride; the process can be controlled to yield either a fusible or an infusible resin. When the polyesterification is carried out in the presence of a long-chain, unsaturated acid of the drying oil type, the oxidative polymerization of the latter is superimposed upon the polyesterification, resulting in hard, synthetic, weather-resistant enamels, suitable for automobile finishes. Polyesterification of ethylene glycol with terephthalic acid results in a polyester fiber. If the material is formed in sheets, it is a useful photographic film. *See* POLYMER.

Many low-molecular-weight esters have characteristic, fruitlike odors: banana (isoamyl acetate), rum (isobutyl propionate), and pineapple (butyl butyrate). These esters are used to some extent in compounding synthetic flavors and perfumes. *See* ALCOHOL; CARBOXYLIC ACID; DRYING OIL; FAT AND OIL; POLYESTER RESINS; SOLVENT.

Esterification. In the broadest sense, esterification is any reaction in which at least one of the products is an ester. There are many routes to the formation of esters. Some of the more important reactions for preparing esters take place between the following pairs of compounds: (1) an acid and an alcohol, (2) an acid anhydride and an alcohol, (3) an acid chloride and an alcohol, (4) an acid and an unsaturated hydrocarbon such as an olefin or an acetylene, (5) an ester and an alcohol, (6) an ester and an acid, and (7) two different esters. This article treats esterification in only a limited sense—reaction between a carboxylic acid (RCOOH) and an alcohol (R'OH) to give the ester and water. For discussions of reactions of an ester with an alcohol, an acid, or another ester *See* TRANSESTERIFICATION.

Esterification reactions are generally reversible and accompanied by relatively small heat effects of the order of a few kilocalories per mole of ester. Al-

though the reactions generally take place in a single liquid phase in the presence of a catalyst, a limited number of esters have been prepared by passing the reactant vapors over a solid catalyst. In the presence of a catalyst, the reaction is commonly conducted at a temperature of about 100°C (212°F); in the absence of a catalyst, a temperature of about 250°C (480°F) is used to give a reasonable reaction rate. The pressure at which the reaction is conducted is determined only by the volatility of the components of the system. It is usually atmospheric pressure. In order to produce most esters economically, some means must be provided for completing the reaction by removing one or more of the products. *See* ACID ANHYDRIDE; ACID HALIDE.

In a typical industrial procedure for the preparation of ethyl acetate, a mixture of acetic acid, excess ethanol, and sulfuric acid is passed into an esterifying column heated to reflux. A ternary azeotrope containing 70% ethanol, 20% ester, and 10% water separates into layers, one of which contains 85% ethyl acetate. The ester may be purified by fractional distillation, and the recovered starting materials are recycled.

Other commercially important esters are prepared as follows. Dibutyl phthalate is prepared from phthalic anhydride and butanol in a stepwise reaction to form first the monoester and then the diester; cellulose acetate from purified α -cellulose and a mixture of acetic anhydride and acetic acid; alkyd resins from phthalic anhydride, unsaturated fatty acids, and glycerol; nitroglycerine (glycerol trinitrate) from glycerol and the proper mixture of nitric acid and sulfuric acid. Aspirin, the world's most used analgesic, is prepared by the reaction of salicylic acid with acetic anhydride below 90°C (190°F), and is purified by recrystallization. *See* ASPIRIN; AZEOTROPIC MIXTURE.

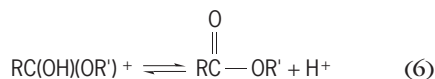
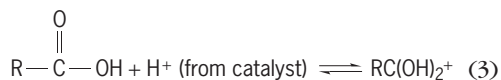
The fact that esterification involves an equilibrium was established in 1862 by M. P. E. Berthelot in his study of the ethyl alcohol-acetic acid system. If 1 mole each of acetic acid and ethyl alcohol react, it is found at equilibrium that $\frac{2}{3}$ mole each of ethyl acetate and water is present at room temperature, along with $\frac{1}{3}$ mole each of alcohol and acid. This can be applied to the equilibrium equation shown as Eq. (2), where K_E is the equilibrium constant and

$$K_E = \frac{[\text{CH}_3\text{COOC}_2\text{H}_5] \times [\text{H}_2\text{O}]}{[\text{CH}_3\text{COOH}] \times [\text{C}_2\text{H}_5\text{OH}]} \quad (2)$$

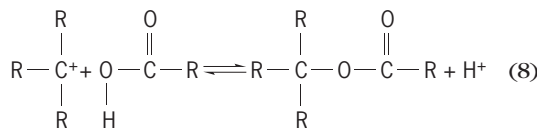
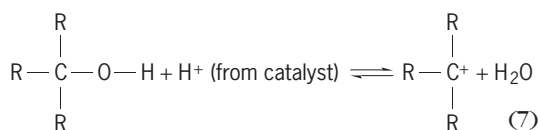
the square brackets signify concentrations in moles per liter of the enclosed reagent. This gives for K_E the value 4. Unless the temperature is deliberately changed, this value is fixed for ethyl alcohol and acetic acid (different alcohol-acid systems have different though characteristic equilibrium constants); indeed, regardless of the starting concentrations of acid and alcohol, the value 4 is maintained.

The mechanism of direct esterification has been much studied. The use of isotopic oxygen (^{18}O) shows that in the reaction of an acid with an alcohol of primary or secondary type, the ester oxygen comes from the alcohol and the acid oxygen

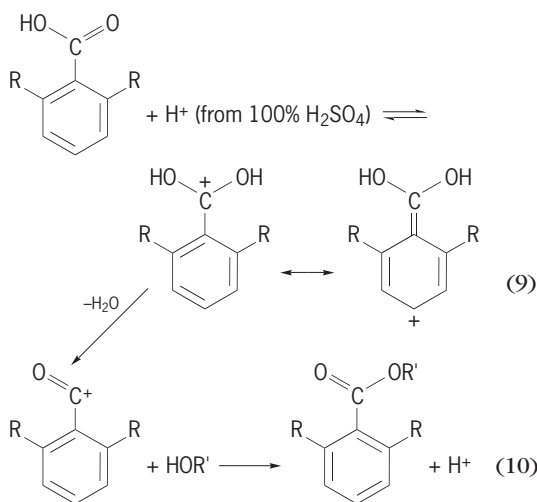
goes to form water. Moreover, in an ordinary acid-catalyzed esterification, the rate of reaction is dependent upon the concentrations of both the carboxylic acid and the alcohol. These observations are accommodated by the mechanistic picture shown as reactions (3)–(6).



In the case of tertiary alcohols, isotopic oxygen studies show that the ester oxygen comes from the carboxylic acid, and the hydroxyl from the alcohol goes to form water, implying the modified picture represented as reactions (7) and (8).



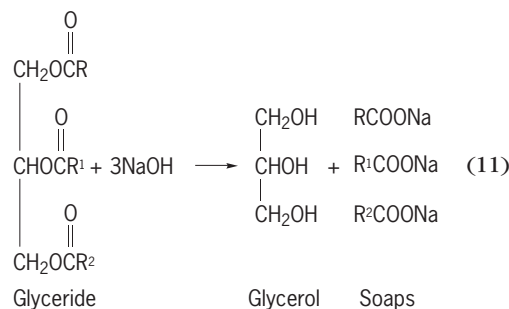
Aromatic acids having substituents in both ortho positions are so hindered in their reaction with alcohols that direct esterification is impracticable. When such acids are dissolved in 100% sulfuric acid and the resulting solution is poured into an alcohol, a good yield of the ester is obtained in a few minutes. The mechanism of this reaction involves an intermediate acyl ion, as shown in reactions (9) and (10).



Hydrolysis. The splitting of esters in such a way as to regenerate the parent acid and alcohol is an example of hydrolysis. It is important, especially in dealing with naturally occurring esters such as those

found in animal and vegetable fats, oils, and waxes. In the presence of dilute mineral acid, hydrolysis of an ester is the reverse of acid-catalyzed esterification; an excess of water is used to ensure complete splitting, and the reaction is carried out at elevated temperatures to speed up the process. Often alcohol is added to solubilize the reactants. Esters formed from glycerol and long-chain carboxylic acids (fats and oils), from long-chain acids and long-chain alcohols (waxes), and simple esters of mono-, di-, or polycarboxylic acids with primary, secondary, or tertiary alcohols, are hydrolyzable under acid conditions, using dilute hydrochloric or sulfuric acids or Twitchell's reagent (prepared from benzene or naphthalene, oleic acid, and concentrated sulfuric acid). However, esters of di-ortho-substituted aromatic carboxylic acids (for example, 2,6-dimethylbenzoic acid) are hindered with respect to hydrolysis, and must be treated according to the Newman technique, which involves first solution in 100% sulfuric acid, and then addition to excess cold water.

The reaction of an ester with a base to form an alcohol and salt of the acid is a type of hydrolysis historically called saponification. Ordinary household soaps are thus made from natural fats and oils of plant or animal origin. They are typically mixtures of the sodium salts of C_{12} and higher fatty acids, and a by-product is glycerol, reaction (11).



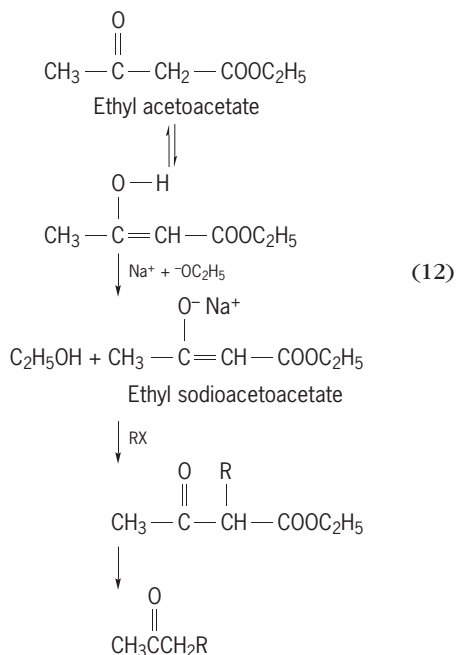
See DETERGENT; SOAP.

Catalytic reduction of esters can be effected at elevated temperature (250°C or 480°F) and pressure (15–20 atm or 1.5–2.0 megapascals), using molecular hydrogen and a copper chromite catalyst; this furnishes a convenient means for the preparation of long-chain mono- or dihydroxy alcohols from ester of the corresponding mono- or dicarboxylic acids. Thus, diethyl succinate ($\text{C}_2\text{H}_5\text{OOC}-\text{CH}_2\text{CH}_2-\text{COOC}_2\text{H}_5$) is reduced to form ethyl alcohol ($\text{C}_2\text{H}_5\text{OH}$) and butylene glycol ($\text{HO}-\text{CH}_2\text{CH}_2\text{CH}_2\text{CH}_2-\text{OH}$). For laboratory reductions, either sodium dissolving in alcohol (Bouveault-Blanc method), or lithium aluminum hydride is preferred.

Esters usually react well with Grignard reagents to yield tertiary alcohols. See GRIGNARD REACTION.

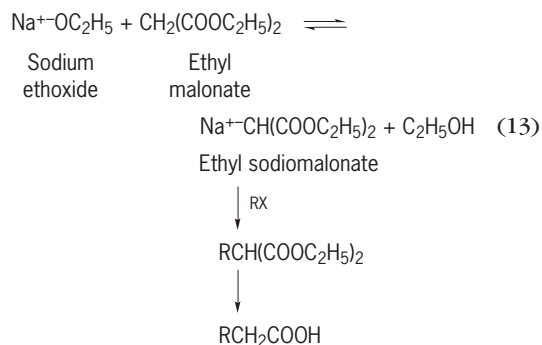
Acetoacetic ester synthesis. An ester of special importance in laboratory synthesis is ethyl acetoacetate, often called acetoacetic ester. Upon treatment with sodium ethoxide, acetoacetic ester forms the salt, ethyl sodioacetoacetate, which reacts with a primary alkyl bromide or iodide (RX) to form the alkyacetoacetate. Hydrolysis and decarboxylation of the latter product provides a general route to the

synthesis of methyl ketones, reaction (12).



See KETONE.

Malonic ester synthesis. Ethyl malonate reacts with a strong organic base such as sodium ethoxide to form the sodiomalonate. Upon treatment with a primary alkyl halide (RX), a substituted malonic ester is



formed, which upon hydrolysis and decarboxylation yields a carboxylic acid. A variety of carboxylic acids have been prepared by reaction sequence (13).

Paul E. Fanta

Bibliography. R. T. Morrison and R. N. Boyd, *Organic Chemistry*, 6th ed., 1992; M. B. Smith and J. March, *March's Advanced Organic Chemistry: Reactions, Mechanisms, and Structure*, 5th ed., 2001.

Estimation theory

A branch of probability and statistics concerned with deriving information about properties of random variables, stochastic processes, and systems based on observed samples. Some of the important applications of estimation theory are found in control and communication systems, where it is used to estimate the unknown states and parameters of the system. For example, the position and velocity of a satellite is estimated from ground radar observations of its range, elevation, and azimuth. These ob-

servations are contaminated with random noise due to atmospheric propagation and radar circuitry. The statistical properties of random noise are assumed known except for some parameters which can be estimated from the data. Generally, the random noise is assumed to have a gaussian distribution, and its mean and covariance may be known or unknown. It is also assumed to be "white," that is, uncorrelated from one time instant to the next. The integral of white noise is a Wiener process or brownian motion process which plays a fundamental role in the theory of stochastic processes. See DISTRIBUTION (PROBABILITY); ELECTRICAL NOISE; STOCHASTIC PROCESS.

The estimation problem for dynamic systems may be divided into two parts: parameter estimation and state estimation. The basic difference between a parameter and the state is that the former either does not change at all or changes slowly in time, whereas the latter continuously evolves in time. For example, the state of a satellite is a six-dimensional vector consisting of three position variables and three velocity variables along the axes of an orthogonal coordinate system. The parameters of the satellite are its mass, inertia, and so on. In many control and communication problems, some of the system parameters are not known with desired accuracy. The problem of estimating these parameters from observed data is called parameter identification, though it is basically a problem of estimation. The more general problem of developing a mathematical model of the system from observed data is called system identification. On the other hand, the problem of state estimation is described by names such as signal processing, filtering, and smoothing. The problem belongs to the theory of stochastic processes and is also commonly known as time series analysis.

Approaches to estimation. Three basic approaches used for estimation are least-squares, maximum-likelihood, and bayesian. An estimator is defined as a function of the observations possessing certain desirable properties such as unbiasedness, consistency, and minimum variance. See STATISTICS.

Least-squares estimation. This was first used by the German mathematician K. F. Gauss in 1821. It is one of the most widely used techniques for estimation. In using the least-squares technique, it is postulated that the observations $\{y(t): t = 1, 2, \dots, T\}$ are related to the unknown parameters, denoted by a vector θ , defined by Eq. (1), where $v(t)$ is a random error which

$$y(t) = f(x_t, \theta, t) + v(t) \quad (1)$$

$$t = 1, 2, \dots, T$$

is white and has zero mean, and x_t is a set of predetermined variables. The least-squares estimate of θ is one that minimizes a scalar function J defined by Eq. (2), where $\sigma_v^2(t)$ is the variance of $v(t)$.

$$J = \sum_{t=1}^T [y(t) - f(x_t, \theta, t)]^2 \frac{1}{\sigma_v^2(t)} \quad (2)$$

See LINEAR ALGEBRA.

The estimate is generally denoted as $\hat{\theta}_{LS}$ and is obtained by numerical minimization of J with respect to θ . The case in which the function $f(x_t, \theta, t)$ is linear

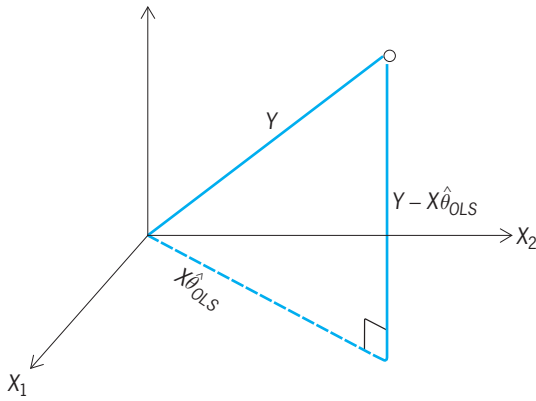


Fig. 1. Least-squares estimation for $T = 3$ and dimension of $\theta = 2$. Error vector is orthogonal to space spanned by columns of X .

in θ and σ_v^2 constant is of special importance and is known as ordinary least-squares. Equation (1) takes the form of Eq. (3).

$$y(t) = x_t^T \theta + v(t) \quad (3)$$

$t = 1, \dots, T$

Let Y denote a $T \times 1$ vector consisting of all observations $\{y(1), \dots, y(T)\}$, where T (the number of observations) is greater than the dimension of θ (the number of unknown parameters). Similarly, defining X and V , Eq. (3) may be written as a single vector equation (4).

$$Y = X\theta + V \quad (4)$$

It is easily shown that for this case J is minimized by the estimator of Eq. (5).

$$\hat{\theta}_{OLS} = (X^T X)^{-1} X^T Y \quad (5)$$

See MATRIX THEORY.

The ordinary least-squares estimator has several interesting properties. It can be shown that in the class of all unbiased linear estimators, $\hat{\theta}_{OLS}$ has the minimum variance. Furthermore, the error vector $(Y - X\hat{\theta}_{OLS})$ is orthogonal to the space spanned by the columns of X (Fig. 1).

The best estimate of Y is $X\hat{\theta}_{OLS}$ which is a projection of Y on the space spanned by the columns of X . The concept of orthogonal projection is more general and carries over to infinite-dimensional spaces that are involved in state estimation. See LEAST-SQUARES METHOD.

Maximum-likelihood estimation. This was first used by R. A. Fisher in 1906 and is based on the concept of a likelihood function. Consider the conditional density function $p(Y|\theta)$, where Y represents the set of observations, and θ the set of unknown parameters. In general, $p(Y|\theta)$ is a function of both Y and θ , but if Y is set equal to its observed value, say Y^* , then $p(Y^*|\theta)$ regarded as a function of θ is known as the likelihood function. A justification for this name comes from the fact that if $p(Y^*|\theta)$ is twice as large for $\theta = \theta_1$ as it is for $\theta = \theta_2$, then the likelihood that θ_1 is the true value is twice as large as θ_2 being the true value. A maximum-likelihood estimate of θ is $\hat{\theta}_{ML}$ if

$p(Y^*|\hat{\theta}_{ML})$ is the maximum of $p(Y^*|\theta)$ with respect to all possible values of θ . See PROBABILITY.

For example, if θ can take only two values, θ_1 or θ_2 , and Y is a scalar, then $p(Y|\theta_1)$ and $p(Y|\theta_2)$ can be plotted as functions of Y (Fig. 2). If the observed value of Y is y , then θ_1 is the maximum-likelihood estimate (MLE) since $p(y|\theta_1) > p(y|\theta_2)$. Similarly, if y' is observed, then θ_2 is the maximum-likelihood estimate. But if y'' is observed, there is an ambiguity. If $p(Y|\theta_1) = p(Y|\theta_2)$ for all values of Y , then θ is said to be unidentifiable from Y .

When $p(Y|\theta)$ is gaussian, then the maximization of $p(Y|\theta)$ is equivalent to the minimization of a quadratic function in Y which is nothing but the least-squares criterion. Thus, for gaussian distributions, the maximum-likelihood estimate and least-squares estimate are equivalent.

Bayesian estimation. In least-squares and in maximum-likelihood estimation, θ is regarded as an unknown but constant parameter. In bayesian estimation, θ is regarded as a random variable and is assigned a prior probability $p(\theta)$. The probability $p(\theta)$ may be objective, that is, based on prior data, or subjective, based on a statistician's judgment. Once Y is observed according to a probability law $p(Y|\theta)$, Bayes' rule is used to obtain the posterior probability $p(\theta|Y)$ according to Eq. (6).

$$p(\theta|Y) = \frac{p(Y|\theta)p(\theta)}{\int p(Y|\theta)p(\theta) d\theta} \quad (6)$$

The observed value of Y is used in Eq. (6) so that $p(Y|\theta)$ is nothing but the likelihood function. The posterior probability $p(\theta|Y)$ summarizes all information regarding θ that can be obtained from Y . It is, therefore, much more informative than a point estimate of θ such as $\hat{\theta}_{LS}$ or $\hat{\theta}_{ML}$. On the other hand, computations of Eq. (6) are more tedious than those of point estimates. But there are cases in which these computations can be simplified.

In the gaussian case, if θ represents the conditional mean of Y and the prior distribution $p(\theta)$ is taken as gaussian, it is easily shown that $p(\theta|Y)$ is also gaussian and that Eq. (6) can be simplified to two equations, one for the mean of $p(\theta|Y)$ and the other for the covariance of $p(\theta|Y)$. If a sequence of observations are made, namely $\{y(1), \dots, y(N)\}$, Eq. (6) can be used repeatedly to obtain $p(\theta|y(1), \dots, y(N))$. The same procedure when used for state estimation in linear dynamic systems leads to a Kalman filter, discussed below. See BAYESIAN STATISTICS.

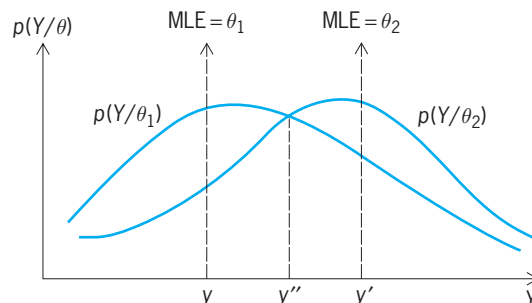


Fig. 2. Maximum-likelihood estimation (MLE).

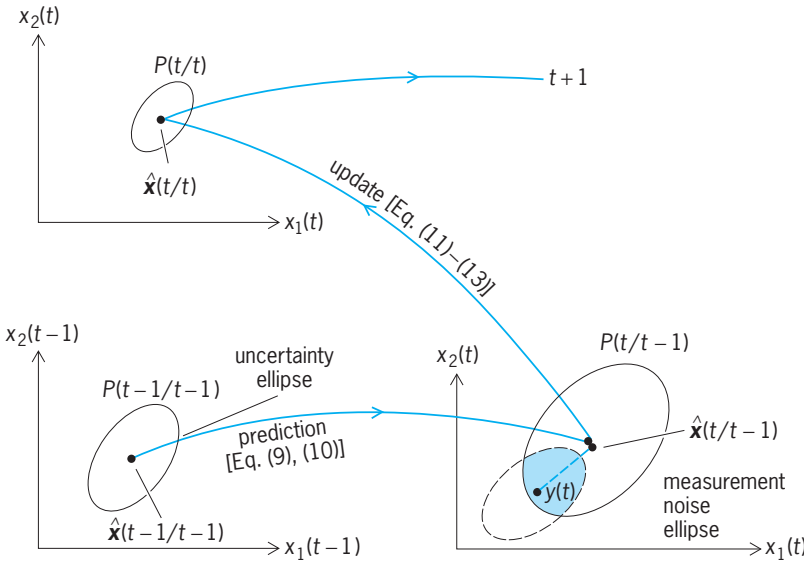


Fig. 3. Discrete-time Kalman filter in two dimensions. The steps involved in obtaining $\hat{x}(t|t)$, $P(t|t)$ from $\hat{x}(t-1|t-1)$, $P(t-1|t-1)$ are shown for a two-dimensional case.

Gauss-Markov processes. A stochastic process $x(t)$ is Markov if given its present state, its future is independent of its past. A fairly general model for Gauss-Markov processes is the Gauss-Markov sequence given by Eqs. (7) and (8), where $x(t)$ is an n -

$$x(t+1) = \Phi(t)x(t) + \Gamma(t)w(t) \quad (7)$$

$$y(t) = H(t)x(t) + v(t) \quad (8)$$

$$t = 0, 1, 2, \dots, N$$

dimensional state vector and $w(t)$ is a q -dimensional zero-mean gaussian white noise sequence with the property that $w(t)$ and $w(\tau)$ are independent if $t \neq \tau$. The sequence $w(t)$ is called process noise, and $v(t)$ which has similar properties is called measurement noise; $y(t)$ is a p -dimensional vector of observations and $\Phi(t)$ is known as the transition matrix of the process. The initial state $x(0)$ also has a gaussian distribution.

The model of Eqs. (7) and (8) is also called a state vector model and arises naturally, in many cases, from a physical model. For example, if $x(t)$ consists of position and velocity components of a moving object subject to random shocks $w(t)$, then for small durations, Eqs. (7) and (8) would describe its motion. The state $x(t)$ takes on different meanings in different contexts, but it always represents the collection of all the information from the past and present behavior of the system that is sufficient to predict the future behavior of the system.

Three state estimation problems arise in connection with Eqs. (7) and (8). Let $\hat{x}(t|\tau)$ represent an estimate of $x(t)$ based on the observation set $\{y(1), \dots, y(\tau)\}$. Based on the values of t and τ , the following names are given to the estimates: filtered estimate, $t = \tau$; predicted estimate, $t > \tau$; smoothed estimate, $t < \tau$.

Kalman filter. A Kalman filter provides estimates that are optimal in the least-squares, maximum-

likelihood, and bayesian senses for the Gauss-Markov model.

Discrete-time case. A discrete-time Kalman filter is a set of recursive equations for the mean and covariance of the gaussian conditional density function $p(x(t)|y(t), \dots, y(1))$ for $t = 1, 2, \dots, N$. Let $\hat{x}(t|t)$ and $P(t|t)$ denote the mean and covariance at time t (Fig. 3). During the prediction phase, $(\hat{x}(t|t-1), P(t|t-1))$ are obtained from Eqs. (9) and (10), where

$$\hat{x}(t|t-1) = \Phi\hat{x}(t-1|t-1) \quad (9)$$

$$P(t|t-1) = \Phi P(t-1|t-1)\Phi^T + \Gamma Q \Gamma^T \quad (10)$$

Q is the covariance of the process noise $w(t)$.

During the update phase, Bayes' rule, Eq. (6), is applied and gives Eqs. (11)–(13), where R is the co-

$$\hat{x}(t|t) = \hat{x}(t|t-1)$$

$$+ K(t)(y(t) - H\hat{x}(t|t-1)) \quad (11)$$

$$K(t) = P(t|t-1)H^T(HP(t|t-1)H^T + R)^{-1} \quad (12)$$

$$P(t|t) = (I - K(t)H)P(t|t-1) \quad (13)$$

variance of the measurement noise $v(t)$. Here $\hat{x}(0|0)$ and $P(0|0)$ are given as the initial distribution of $x(0)$.

Continuous-time case. The continuous-time analog of the model of Eqs. (7) and (8) is a linear differential equation forced by a white noise process, Eqs. (14) and (15).

$$\frac{dx(t)}{dt} = Fx(t) + Gu(t) \quad (14)$$

$$y(t) = Hx(t) + v(t) \quad 0 \leq t \leq T \quad (15)$$

Here $x(0)$ is normally distributed with mean \hat{x}_0 and covariance P_0 . The correlation functions of $u(t)$ and $v(t)$ are Dirac delta functions, $\delta(t - \tau)$, which are zero for $t \neq \tau$ and integrate out to one. These processes are also uncorrelated with $x(0)$ and with each other.

A simple way of deriving the continuous-time Kalman filter is to discretize Eqs. (14) and (15) with a small time-step Δt and express them in the form of Eqs. (7) and (8). Then the discrete-time Kalman filter equations (9)–(13) are used and the limit $\Delta t \rightarrow 0$ is taken. The filtering equations (16)–(19) are obtained.

$$\frac{d\hat{x}(t)}{dt} = F\hat{x}(t) + K(t)(y(t) - H\hat{x}(t)) \quad (16)$$

$$\frac{dP(t)}{dt} = FP(t) + P(t)F^T + GQG^T - P(t)H^T R^{-1}HP(t) \quad (17)$$

$$K(t) = P(t)H^T R^{-1} \quad (18)$$

$$\hat{x}(0) = x_0 \quad P(0) = P_0 \quad (19)$$

A Kalman filter can be represented as a signal-flow block diagram (Fig. 4).

Both the discrete-time and continuous-time models can readily be generalized to the problems of prediction and smoothing.

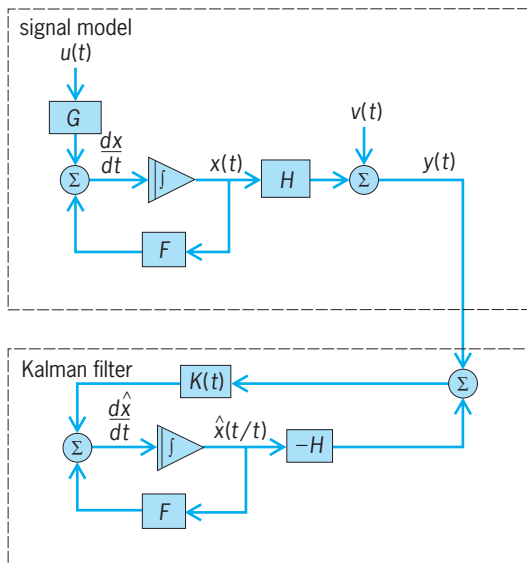


Fig. 4. Block diagram of a continuous-time Kalman filter and the signal model.

Other state estimation methods. The first treatment of the prediction problem was due to A. Kolmogorov (1941) for the discrete-time case. N. Wiener (1942) developed the continuous-time filter independently in the frequency domain. Both of these approaches considered stationary processes and made use only of the correlation or spectral properties of the stochastic processes. R. E. Kalman (1960), on the other hand, used the state vector model which can describe both stationary and nonstationary processes and can incorporate physical information regarding the system. The Kalman filter is equivalent to the Kolmogorov-Wiener filter for the stationary case.

Nonlinear filters. If the system equations (7) and (8) or (14) and (15) are nonlinear, the filtering density function is non-gaussian and it is necessary to evaluate higher-order moments or to propagate the complete density function. In the continuous-time case, this leads to partial differential equations which are very time-consuming to solve. A more practical, though suboptimal, approach is to linearize the nonlinear equations around the latest filtering and prediction estimates. This leads to an extended Kalman filter. More complicated versions are (1) the iterated extended Kalman filter in which the measurement nonlinearity effects are alleviated further by repeated linearization and update, (2) gaussian second-order filters in which certain bias terms due to the nonlinearities in the extended Kalman filter are eliminated, and (3) iterated smoothing filters in which one or more backward-step smoothing estimates are used for linearization.

Applications. Since their introduction in 1960, Kalman filters and their extensions have found numerous applications. Initially, these filters were developed for space applications such as satellite orbit determination, inertial navigation, Apollo lunar landing module guidance, and so on. The applications to power systems and industrial processes were developed shortly thereafter. Kalman filters have been used for forecasting, water quality prediction,

hurricane tracking, and aircraft landing systems. Another area for application of Kalman filters is stochastic control, that is, control of systems with random inputs. It can be shown that an optimal controller for the so-called LQG systems (linear dynamics, quadratic criterion, and gaussian disturbances) consists of a Kalman filter cascaded with a linear-state feedback device. This concept has been used in the design of flight controllers, submarine depth-keeping control, process control, and pilot modeling, and in the study of optimal economic stabilization policies. See FLIGHT CONTROLS; GUIDANCE SYSTEMS; INERTIAL GUIDANCE SYSTEM; OPTIMAL CONTROL (LINEAR SYSTEMS); PROCESS CONTROL; SPACE NAVIGATION AND GUIDANCE; STOCHASTIC CONTROL THEORY.

Raman K. Mehra

System identification. This is the determination of a mathematical model of a system from the observed input and output data. System identification most often refers to dynamical system models, although this term can be extended to static input-output models. A dynamic system model is usually represented by differential equations (for a continuous-time model) or by difference equations (for a discrete-time model). Techniques to determine the model from measured data involve several key steps. The first step is the selection of the model type, including the selection of a continuous-time or discrete-time model; the selection of a linear or nonlinear model; the selection of the order of the model; and deciding whether there is prior knowledge of the system, or such knowledge is lacking, and parameters are to be determined from measurements. The other steps are the design of the experiment, including the choice of the experiment type as batch processing or recursive processing, and the choice of the input signal; the determination of parameter values in the model that provide the best fit to data (the parameter-estimation problem); and model validation. The selection of the model type is quite application dependent, while the remaining three steps involve one of numerous mathematical techniques for parameter estimation. See DIFFERENCE EQUATION; DIFFERENTIAL EQUATION; MODEL THEORY.

Model selection depends largely on the level of prior knowledge of the system. In some situations this knowledge may be sufficient to derive the dynamical model from known principles of physics, with uncertainty only about the values of some parameters in the model. The choice of the model in this case is then dictated by known laws of physics. This is the case of the identification of a partially known system. In other situations the knowledge of the system may be insufficient to derive a system of model equations. In this case the form of the model may be assumed, with all the parameters in it uncertain. The values of all the parameters are then estimated from input and output measurements. In contrast to the previous case, the parameters may have no physical meaning. An example of such a model is the linear discrete-time system given by Eq. (20).

$$y(t) + a_1 y(t-1) + \cdots + a_n y(t-n) = b_1 u(t-1) + \cdots + b_m u(t-m) + v(t) \quad (20)$$

This model has a single input $u(t)$, a single output $y(t)$, and noise $v(t)$, and there is no prior knowledge of the parameters $a_1, \dots, a_m, b_1, \dots, b_m$.

Batch versus recursive identification. If the input and output signals are prerecorded and the parameter-estimation procedure is carried out at some later time, the process is called a batch or off-line identification. The advantage of batch identification is that it is possible to try several models and find the one that produces the best fit to data. However, the identification results are not available in time for active control of the system while the signals are being recorded.

If the identification is carried out during the process operation and the parameter estimates are being modified as new measurements arrive, such a process is called a recursive or on-line identification. Recursive identification can be included in the feedback control algorithm, thus producing an adaptive control law, which can be used to control systems whose models are not completely known initially. See ADAPTIVE CONTROL.

Least-squares identification. As an example of system identification, the parameters of the model given by Eq. (20) can be determined by using the least-squares method in the following manner:

If two vectors are defined by Eqs. (21) and (22),

$$\phi(t) = [-y(t-1), \dots, -y(t-n), u(t-1), \dots, u(t-m)]^T \quad (21)$$

$$\theta = [a_1 \dots a_n b_1 \dots b_m]^T \quad (22)$$

then Eq. (20) can be rewritten as Eq. (23). The vector

$$y(t) = \theta^T \phi(t) + v(t) \quad (23)$$

of parameters θ can then be estimated by the least-squares method, as described above. This would result in a batch identification of the model. A recursive version of the same method is given by Eq. (24),

$$\hat{\theta}(t) = \hat{\theta}(t-1) + L(t)[y(t) - \hat{\theta}^T(t-1)\phi(t)] \quad (24)$$

where $\hat{\theta}(t)$ denotes the estimate of θ at time t .

The term $[y(t) - \hat{\theta}^T(t-1)\phi(t)]$ represents the mismatch between the actual output $y(t)$ and the one predicted by using the previous parameter estimate $\hat{\theta}(t-1)$. The $L(t)$ is a column of $n+m$ time-dependent coefficients which dictate the relative size of the update made in each component of $\hat{\theta}$ from step $t-1$ to step t . The vector $L(t)$ is recursively computed from another formula (not shown here) in such a way that the sum of squared mismatch terms is minimal (the least-squares method).

The above is known as the recursive least-squares (RLS) algorithm. This is one of the most widely known, robust, and easy-to-implement algorithms.

The minimum number of data points required for the algorithm is equal to the dimension of the vector θ , that is, $n+m$. The algorithm requires initialization. The common choice of initial value is $\hat{\theta}(0) = 0$. If the sequences $v(t)$ and $\Phi(t)$ are uncorrelated, the estimate $\hat{\theta}(t)$ will be unbiased.

Choice of input signal. A recursive parameter estimation algorithm may not necessarily converge if the input signal is poorly chosen. A condition known as the persistent excitation condition is sufficient (but not necessary) for the convergence of several algorithms, including the RLS algorithm. In the case of model given by Eq. (20) with $m=n$, the persistent excitation condition requires that the power spectral density of the input signal $u(t)$ be nonzero at $2n$ or more frequencies. This condition is satisfied by a signal that contains at least $2n$ distinct frequencies, and also by white noise.

One particularly useful choice of the input signal is pseudorandom binary noise. This is a waveform switching between two values, for example, $+1$ and -1 , according to a deterministic sequence designed in such a way that the autocorrelation function of the waveform approximates that of a white noise. Pseudorandom binary noise of small amplitude can be added to a constant input, thus providing a signal that does not disrupt normal system operation.

Applications of system identification. System identification is one of the most frequently applied parts of control systems theory. Both batch and recursive estimation algorithms have been used since the 1960s in many industries, in particular in the chemical and petroleum industries, paper-machine control, and ship steering. System identification subsequently became more closely integrated in the overall control loop in the form of adaptive control laws. Adaptive control has been applied to many problems in process control, motion control, and the control of aircraft and spacecraft.

Closely related to the problem of identification is the problem of failure detection in a dynamic system. It is a problem of determining a sudden change in a parameter value caused by a failure of one component in the system. This is usually accomplished by using algorithms which compare residuals in several redundant input-output relations. See CONTROL SYSTEMS; FAULT ANALYSIS; SYSTEMS ENGINEERING.

Andre Z. Manitius

Bibliography. S. M. Bozic, *Digital and Kalman Filtering*, 2d ed., 1995; R. G. Brown and P. Y. Hwang, *Introduction to Random Signal Analysis and Kalman Filtering*, 3d ed., 1996; H. F. Chen and L. Guo, *Identification and Stochastic Adaptive Control*, 1991; M. S. Grewal and A. P. Andrews, *Kalman Filtering: Theory and Practice*, 1993; L. Ljung, *System Identification: Theory for the User*, 2d ed., 1999; J. M. Mendel, *Lessons in Estimation Theory*, 1995; A. Papoulis, *Probability, Random Variables and Stochastic Processes*, 3d ed., 1991; H. Stark and J. W. Woods, *Probability, Random Processes, and Estimation Theory for Engineers*, 2d ed., 1994.

Estrogen

A substance that maintains the secondary sex characters and organs, such as mammary glands, uterus, vagina, and fallopian tubes, of mammalian females. Naturally occurring substances with this activity

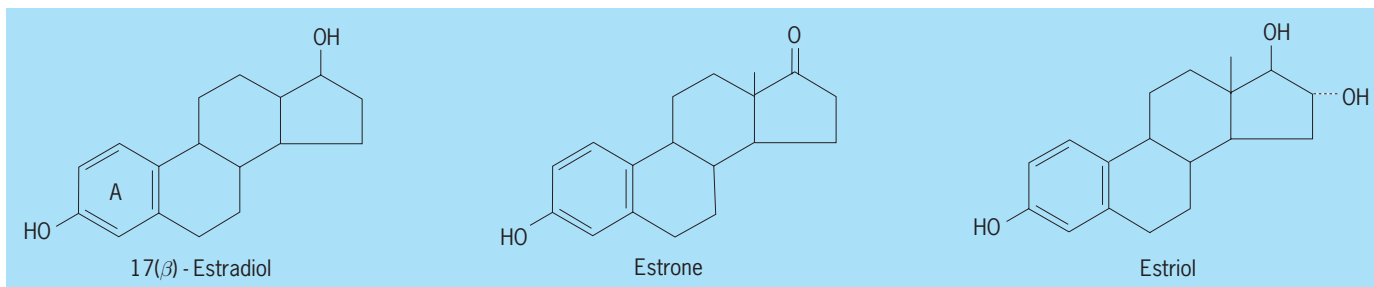


Fig. 1. Structural formulas of estrogenic hormones with ring A phenolic group.

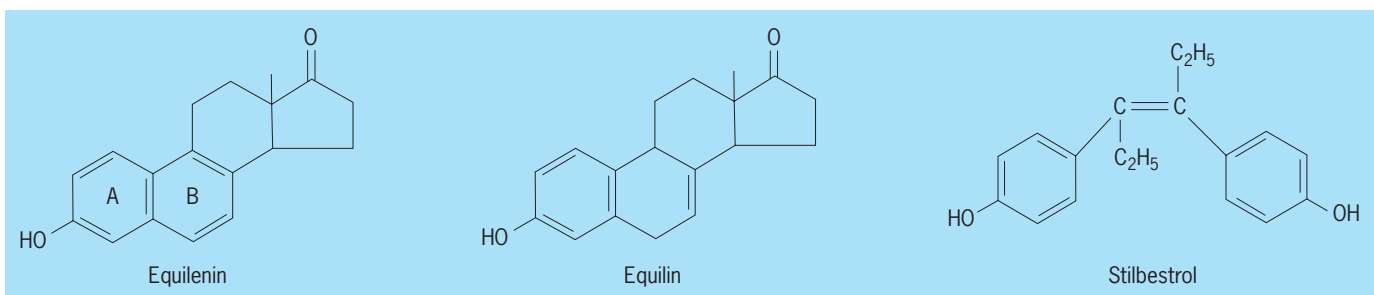


Fig. 2. Structural formulas of estrogenic hormones with ring B unsaturation.

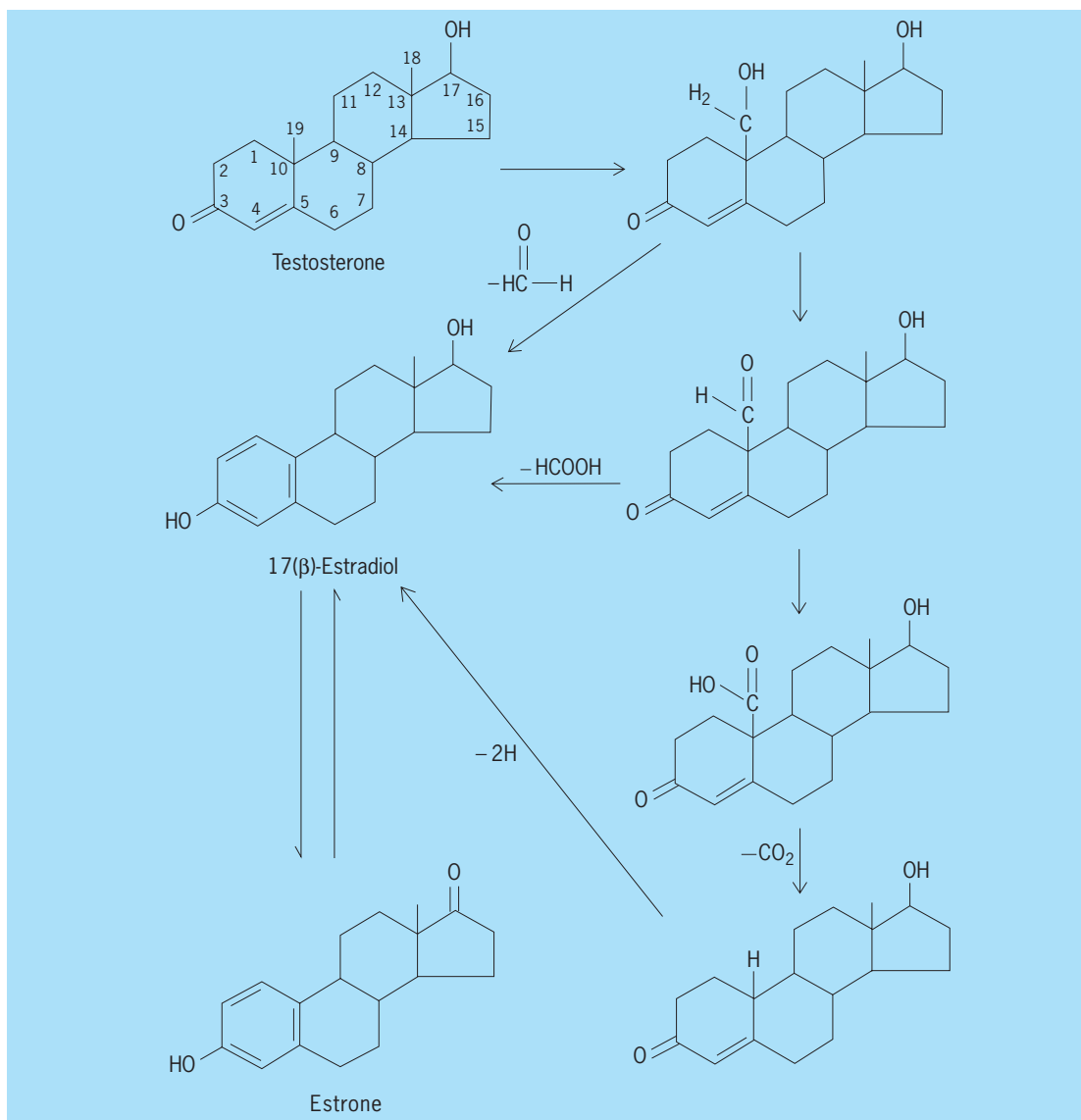


Fig. 3. Pathways for the biosynthesis of estrogens from testosterone.

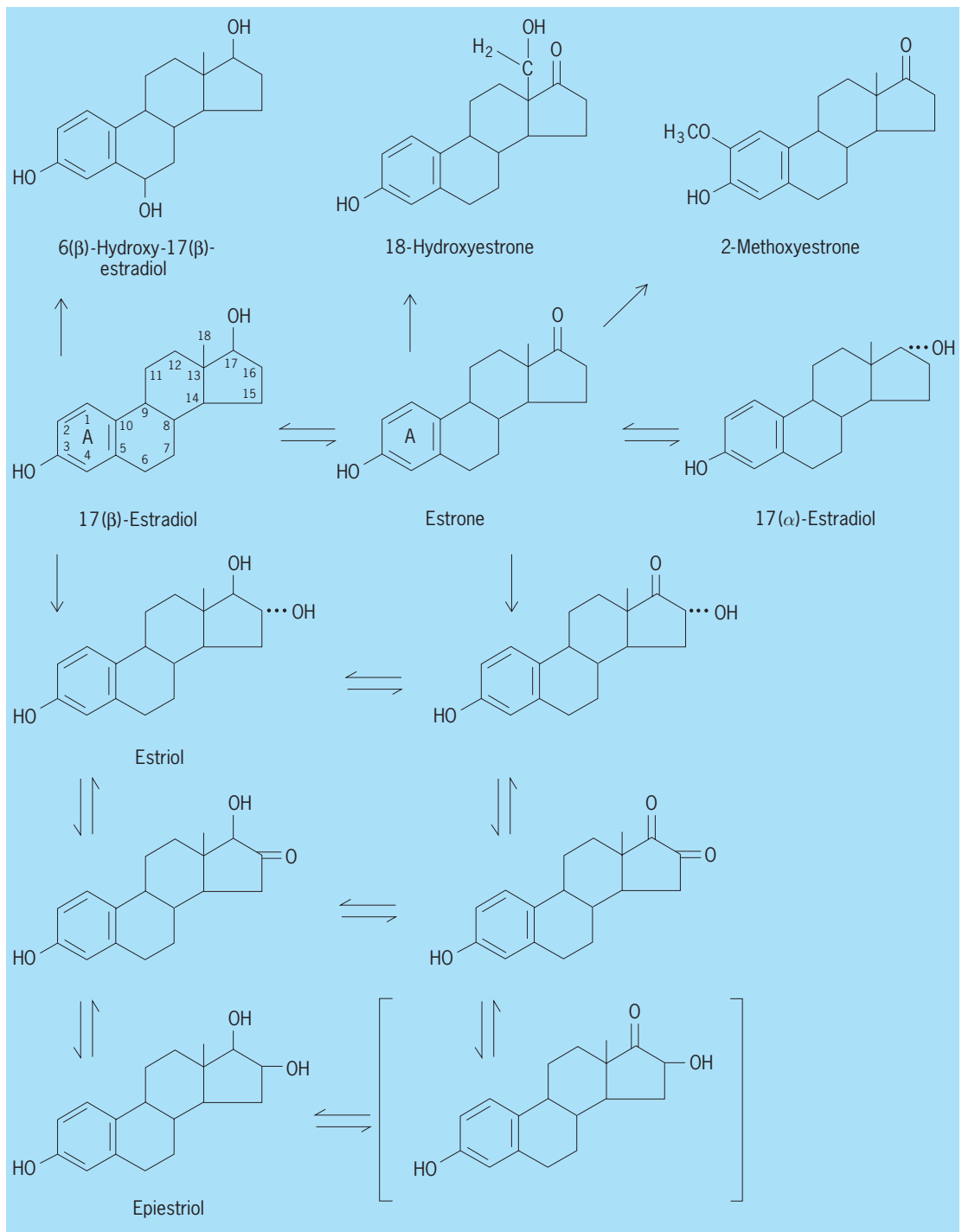


Fig. 4. Reactions involved in the metabolic breakdown of estrogens.

are steroid hormones. The principal estrogenic hormone substances are 17(β)-estradiol, estrone, and estriol. These compounds have a ring A phenolic group (Fig. 1). They are produced and secreted directly into the bloodstream by the ovary, testis, adrenal, and placenta of pregnancy. Two other naturally occurring estrogenic hormones, equilin and equilenin (Fig. 2), have been obtained only from the urine of pregnant mares and are apparently peculiar to that species. Stilbestrol, a synthetic compound with considerable estrogenic activity, has been used exten-

sively in medical practice (Fig. 2).

Ring A phenolic group estrogens. The most active of the estrogens known to be produced by mammalian tissues is 17(β)-estradiol which is interconvertible with estrone. An enzyme present in liver, placenta, and other tissues catalyzes this oxidation-reduction. Estriol is considered to be derived from estrone.

Ring B unsaturated estrogens. Equilenin and equilin are naturally occurring representatives of ring B unsaturated estrogens. Equilenin also has been detected in human adrenal cancer tissue from a patient with a

feminizing syndrome. Stilbestrol is a nonsteroid synthetic estrogen. It has a biological activity equal to that of 17(β)-estradiol by injection and a higher activity by oral administration.

Biosynthesis of estrogens. This involves a series of reactions from androgens such as testosterone and Δ^4 -androstene-3,17-dione. The transformation of androgens to estrogens involves the removal of carbon 19 and two hydrogen atoms, one each at carbons 1 and 2. The loss of the carbon atom most likely occurs by reactions indicated in **Fig. 3**, using testosterone as the example.

Catabolism of estrogens. The metabolic breakdown of estrogens involves various modifications at carbon atoms 2, 6, 16, and 18. These overall reactions are illustrated in the series of equations in **Fig. 4**. See HORMONE; STEROID. Ralph I. Dorfman

Bibliography. B. Green and R. E. Leake, *Steroid Hormones: A Practical Approach*, 1987; V. H. James and J. R. Pasqualini (eds.), *Hormonal Steroids*, 1984; H. L. Makin, *Biochemistry of Steroid Hormones*, 2d ed., 1984; S. Raam (ed.), *Immunology of Steroid Hormone Receptors*, 1988.

Estrus

The period in mammals during which the female ovulates and is receptive to mating. It is commonly referred to as rut or heat. From one estrus period to the next there occurs a series of changes, particularly in the ovary, uterus, and vagina, termed the estrous cycle. With reference to the ovary, the cycle can be divided into a follicular phase, during which the Graffian follicles are ripening, and a luteal phase, during which the corpora lutea develop in the ovulated follicles. During these two phases mainly estrogen and progesterone, respectively, are secreted, and these hormones control the uterine and vaginal changes. The beginning of the follicular phase is termed proestrus and the luteal phase metestrus. Following the latter, there is a period of relatively little change, termed diestrus. In species in which the latter is prolonged, it is termed anestrus.

During proestrus the endometrial lining of both the uterus and vagina thickens. The vagina regresses in metestrus, and vaginal washings exhibit mainly cornified epithelial cells. At the same time uterine glands continue to enlarge and become highly secretory. This is the stage during which the embryo would implant in the uterine lining if the previously ovulated eggs had been fertilized. Toward the end of this period, in the absence of fertilization, the uterine endometrium regresses, gradually in most mammals, but precipitously in primates. In the latter the regression involves considerable loss of tissue and bleeding. See MENSTRUATION.

Some mammals, such as the armadillo, bat, deer, wolf, shrew, and ferret, are monestrous, having a single cycle per year. Most are polyestrous, but considerable variations occur in different species. In dogs there are 2–3 cycles per year. Cats cycle about every 2 weeks during breeding seasons. Cows, horses, and

sheep illustrate seasonal polyestrus with 2½- to 3-week cycles, while laboratory rats and mice maintain a 5- to 6-day schedule throughout the year. In the rabbit there is persistent estrus, and ovulation occurs in response to the stimulus of coitus. This mating-induced ovulation is known to occur also in the cat and ferret during their more restricted period of estrus, and it probably also occurs in a number of other mammals.

The regulation of the estrous cycle, and of breeding seasons, is brought about by feedback-type hormonal interactions between the ovary and pituitary gland, influenced in many cases by environmental factors, such as length of day. Thus, for example, if additional illumination is provided in the fall, ferrets can start breeding during the winter instead of spring; and if the sequence of night and day is changed, mice, which ordinarily mate and ovulate during the night, will adapt to the new conditions. Such environmental stimuli (light, coition) can affect the reproductive cycle through releasing hormones produced by cells in the hypothalamus, which are connected to neural pathways in the brain. Specific releasing hormones control the amounts of gonadotrophic hormones discharged by the anterior lobe of the pituitary into the bloodstream. See ENDOCRINE MECHANISMS; OVUM; PITUITARY GLAND; REPRODUCTION (ANIMAL).

Albert Tyler; Howard L. Hamilton

Estuarine oceanography

The study of the physical, chemical, biological, and geological characteristics of estuaries. An estuary is a semienclosed coastal body of water which has a free connection with the sea and within which the seawater is measurably diluted by freshwater derived from land drainage. Many characteristic features of estuaries extend into the coastal areas beyond their mouths; and because the techniques of measurement and analysis are similar, the field of estuarine oceanography is often considered to include the study of some coastal waters which, by the above definition, are not strictly estuaries. Also, semienclosed bays and lagoons exist in which evaporation is equal to or exceeds freshwater inflow, so that the salt content either is equal to that of the sea or exceeds it. Hypersaline lagoons have been termed negative estuaries, whereas those with precipitation and river inflow equaling evaporation have been called neutral estuaries. Positive estuaries, in which river inflow and precipitation exceed evaporation, form the majority.

Topographic classification. Embayments are the result of fairly recent changes in sea level. During the Pleistocene ice age, much of the seawater was locked up in continental ice sheets, and the sea surface stood about 100 m (330 ft) below its present level. In areas not covered with ice, the rivers incised their valleys to this base level. During the ensuing Flandrian Transgression, when the sea level rose at about 1 m (3.3 ft) per century, these valleys became inundated. Much of the variation in form of the resulting estuaries

depends on the volumes of sediment that the river or the nearby coastal erosion has contributed to fill the valleys.

Where river flow and sediment discharge were high, the valleys have become completely filled and even built out into deltas. Generally, deltas are best developed in areas where the tidal range is small and where the currents cannot easily redistribute the sediment the rivers introduce. They occur mainly in tropical and subtropical areas where river discharge is seasonally very high. The distributaries, or passes, of the delta are generally shallow, and often the shallowest part is a sediment bar at the mouths of the distributaries. The Mississippi and the Niger are examples of this type of delta.

Where sediment discharge was less, the estuaries are unfilled, although possibly they are still being filled. These are drowned river valleys or coastal plain estuaries, and they still retain the topographic features of river valleys, having a branching, dendritic, though meandering, outline and a triangular cross section, and widening regularly toward the mouth, which is often restricted by spits. River discharge tends to be reasonably steady throughout the year, and sediment discharge is generally small. These estuaries occur in areas of high tidal range, where the currents have helped to keep the estuaries clear of sediment. They are typical of temperate regions such as the east coast of North America and northwestern Europe, examples being the Chesapeake Bay system, the Thames, and the Gironde. See COASTAL LANDFORMS.

In areas where glaciation was active, the river valleys were overdeepened by glaciers and fiords were created. A characteristic of these estuaries is the rock bar or sill at the mouth that can be as little as a few tens of meters deep. Inside the mouth, however, the estuaries can be at least 600 m (1800 ft) deep and can extend hundreds of kilometers inland. Fiords are typical of Norway and the Canadian Pacific coast. See FIORD.

Another estuarine type is called the bar-built estuary. These are formed on low coastlines where extensive lagoons have narrow connecting passages or inlets to the sea. Within the shallow lagoons, the tidal currents are small, but the deep inlets have higher currents. Again, a sediment bar is generally present across the entrance. In tropical areas, the lagoons can be hypersaline during the hot season. They are typical of the southern United States and of parts of Australia.

Estuaries are ephemeral features since great alterations can be wrought by small changes in sea level. If the present ice caps were to melt, the sea level would rise an estimated 30 m (99 ft), and the effect on the form and distribution of estuaries would be drastic.

Physical structure and circulation. Within estuaries, the river discharge interacts with the seawater, and river water and seawater are mixed by the action of tidal motion, by wind stress on the surface, and by the river discharge forcing its way toward the sea. The difference in salinity between river water and

seawater—about 35 parts per thousand—creates a difference in density of about 2%. Even though this difference is small, it is sufficient to cause horizontal pressure gradients within the water which affect the way it flows. Density differences caused by temperature variations are comparatively smaller. Salinity is consequently a good indicator of estuarine mixing and the patterns of water circulation. Obviously, there are likely to be differences in the circulation within estuaries of the same topographic type which are caused by differences in river discharge and tidal range. The action of wind on the water surface is an important mixing mechanism in shallow estuaries, particularly in lagoons; but generally its effect on estuarine circulation is only temporary, although it can produce considerable variability and thus make interpretation of field observations difficult. See SEAWATER.

Salt-wedge estuaries. Freshwater, being less dense than seawater, tends to flow outward over the surface of seawater, which penetrates as a salt wedge along the bottom into the estuary (Fig. 1). This creates a vertical salinity stratification, with a narrow zone of sharp salinity change, called a halocline, between the two water masses, which can reach 30 parts per

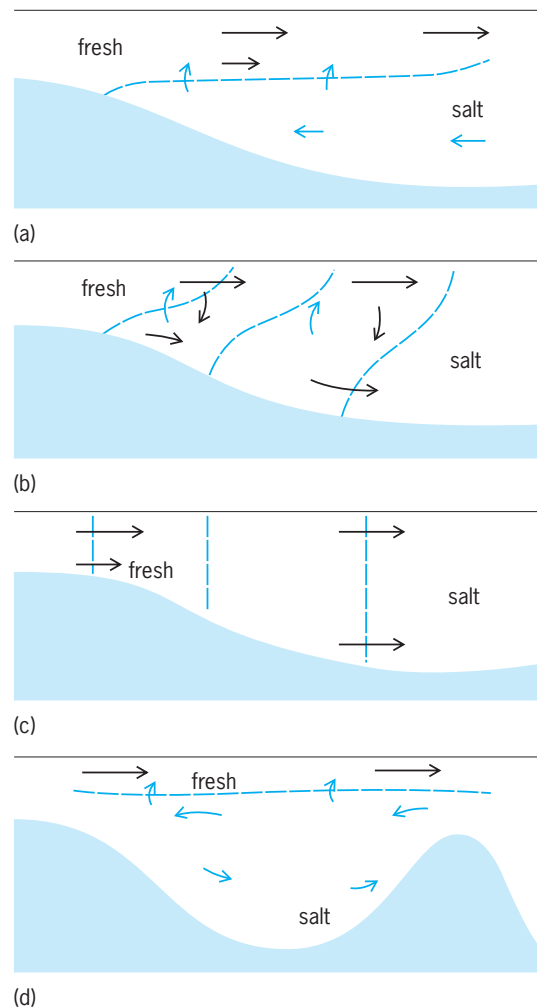


Fig. 1. Diagrams of mixing in estuaries. (a) Salt-wedge type. (b) Partially mixed type. (c) Well-mixed type. (d) Fiord.

thousand in 0.5 m (1.5 ft). If the sea is tideless, the water in the salt wedge is almost motionless. However, if the surface layer flowing toward the sea has a sufficiently high velocity, turbulent mixing can occur through a mechanism known as Kelvin-Helmholtz instability, which is a process where the denser salt water is drawn up into the overflowing freshwater in a coherent "rolled-up" pattern. Ultimately these instabilities break down completely, and the salt water is mixed entirely into the overlying fresh-water mass, increasing the salinity of the upper layer, before eventually being discharged to the ocean. This and similar processes are sometimes referred to as entrainment, which results in a net loss of fluid from the salt wedge. Consequently, for this loss to be replaced, there must be a compensatory flow of salt water toward the head, or landward portion, of the estuary within the salt wedge, but of a magnitude much less than that of the flow in the surface layer. There is a considerable velocity gradient near the halocline as a result of the friction between the two layers. Consequently, the position of the salt wedge will change according to the magnitude of the flow in the surface layer, that is, according to the river discharge. The Mississippi River is an example of a salt-wedge estuary. When the flow in the Mississippi is low, the salt wedge extends more than 160 km (100 mi) inland, but with high discharge the salt wedge extends only a mile or so above the river mouth. Some bar-built estuaries, in areas of restricted tidal range and at times of high river discharge, as well as deltas, are typical salt-wedge types.

Partially mixed estuaries. When tidal movements are appreciable, the whole mass of water in the estuary moves up and down with a tidal periodicity of about 12.5 h. Considerable friction occurs between the bed of the estuary and the tidal currents, and causes turbulence. The turbulence tends to mix the water column more thoroughly than occurs in salt-wedge estuaries, although little is known of the relationship of the exchanges to the salinity and velocity gradients. However, the turbulent mixing not only mixes the salt water into the fresher surface layer but also mixes the fresher water downward. This causes the salinity to decrease toward the head of the estuary in the lower layer and also to increase progressively toward the sea in the surface layer. As a consequence, the vertical salinity gradient is considerably less than that in salt-wedge estuaries. In the surface, seaward-flowing layer, the river discharge moves toward the sea; but because the salinity of the water has been increased by mixing during its passage down the estuary, the discharge at the mouth can be several times the river discharge. To provide this volume of additional water, the compensating inflow must also be much higher than that in the salt-wedge estuary. The velocities involved in these movements are only on the order of a few centimeters per second, but the tidal velocities can be on the order of 100 centimeters per second. Consequently, the only way to evaluate the effect of turbulent mixing on the circulation pattern is to average out the effect of the tidal oscillation, which requires considerable preci-

sion and care. The resulting residual or mean flow will be related to the river discharge, although the tidal response of the estuary can give additional contributions to the mean flow. The tidal excursion of a water particle at a point will be related to the tidal prism, the volume between high- and low-tide levels upstream of that point; and the instantaneous cross-sectional velocity at any time will be related to the rate of change of the tidal prism upstream of the section. In details, the velocities across the section can differ considerably. It has been found that in the Northern Hemisphere the seaward-flowing surface water keeps to the right bank of the estuary, looking downstream, and the landward-flowing salt intrusion is concentrated on the left-hand side (Fig. 2). This is caused by the Coriolis force, which deflects the moving water masses toward the right and is of increased importance for very wide estuaries. Of possibly greater importance, however, is the effect of topography, because the curves in the estuary outline tend to concentrate the flow toward the outside of the bends. Thus, in addition to a vertical circulation, there is a horizontal one, and the halocline slopes across the estuary. Because the estuary has a prismatic cross section, the saline water is concentrated in the deep channel and the fresher water is discharged in the shallower areas. Examples of partially mixed estuaries are the rivers of the Chesapeake Bay system.

Well-mixed estuaries. When the tidal range is very large, there is sufficient energy available in the turbulence to break down the vertical salinity stratification completely, so that the water column becomes

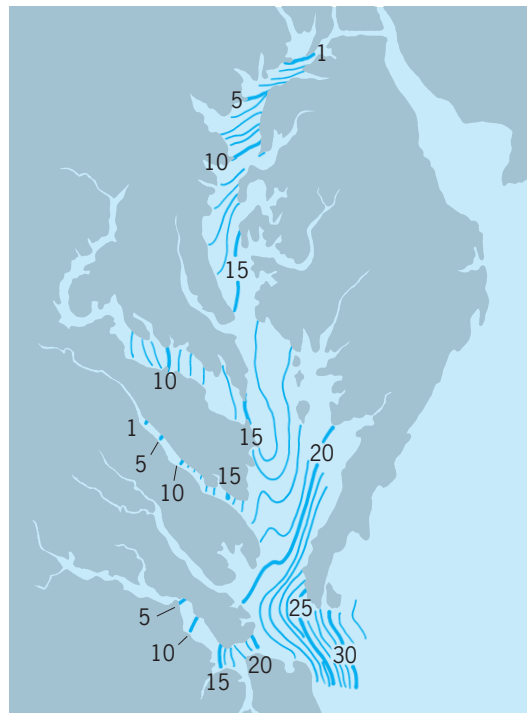


Fig. 2. Typical surface salinity distribution in Chesapeake Bay. Numbers indicate parts per thousand. (After H. E. Landsberg, ed., *Advances in Geophysics*, vol. 1, Academic Press, 1952)

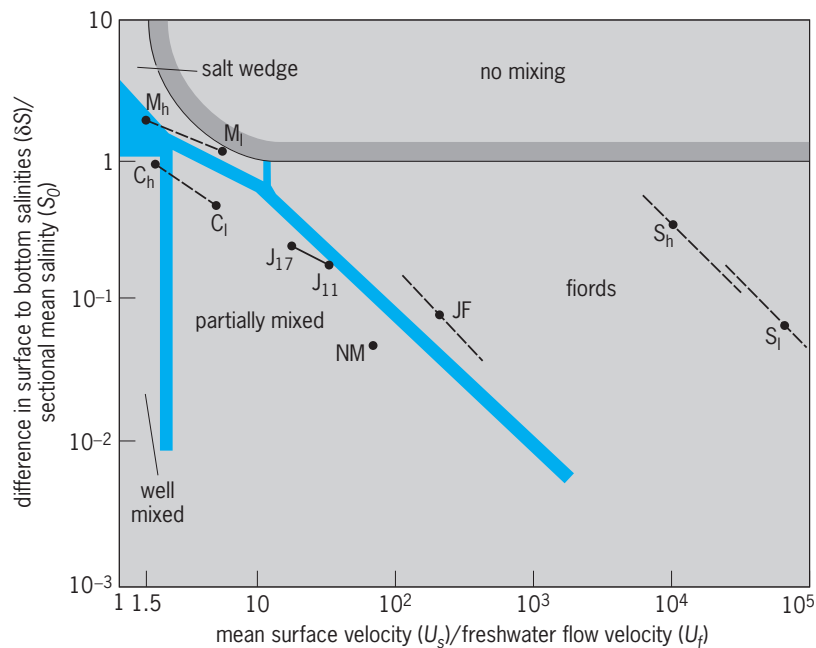


Fig. 3. Classification diagram for estuaries. An estuary appears as a line on the diagram; the upper reaches are less well mixed than the lower sections. Subscript letters refer to high (h) and low (l) river discharge; subscript numbers are distances from the mouth. J = James River; M = Mississippi; C = Columbia River; NM = Narrows of the Mersey; S = Silver Bay; JF = Strait of Juan de Fuca. (After D. V. Hansen and M. Rattray, Jr., *New dimensions in estuary classification*, *Limnol. Oceanogr.*, 11:319–326, 1966)

vertically homogeneous. In this type of estuary there can be lateral variations in salinity and in velocity, with a well-developed horizontal circulation; or if the lateral mixing is also intense, the estuary can become sectionally homogeneous (also called a one-dimensional estuary). Because there is no landward residual flow in the sectionally homogeneous estuary, the upstream movement of salt is produced during the tidal cycle by salty water being trapped in bays and creeks and bleeding back into the main flow during the ebb. This mechanism spreads out the salt water, allowing salt to be maintained within the estuary, but it is probably an effective trapping mechanism for salt for only a small number of tidal excursions landward of the mouth.

Fiords. Because fiords are so deep and restricted at their mouths, tidal oscillation affects only their near-surface layer to any great extent. The amount of turbulence created by oscillation is small, and the mixing process is achieved by entrainment. Thus fiords can be considered as salt-wedge estuaries with an effectively infinitely deep lower layer. The salinity of the bottom layer will not vary significantly from mouth to head, and the surface fresh layer is typically no more than a few tens of meters deep. When the sill is deep enough not to restrict circulation, the inflow of water occurs just below the halocline, with an additional slow outflow near the bottom. When circulation is restricted, the replenishment of the deeper water occurs only occasionally, sometimes on an annual cycle and typically related to severe weather events. Between these replenishment episodes, the bottom layer can become anoxic, with very low dissolved oxygen (DO) levels.

The descriptive classification of estuaries outlined above depends on the relative intensities of the tidal and river flows and the effect that these flows have on stratification. A quantitative comparison between estuaries can be made using the diagram of **Fig. 3**, which is based on a stratification and a circulation parameter.

River plumes. The estuary represents the first stage of the blending of fresh river water into the sea. The water that is discharged from the estuary mouth typically has characteristics (such as salinity and temperature) that lie between those of the fresh river water and seawater, due to mixing and entrainment processes within the estuary. The process of blending with the coastal ocean continues in a region known as a river plume. With regard to the mixing and blending processes, a river plume can be divided into two regions. Certain estuaries may not exhibit a well-defined near-field region; and in some estuaries with particularly low freshwater flow, there may be no evidence of a plume at all.

A near-field plume is often considered an extension of the estuary out into the coastal ocean. It can be broadly defined as the region where the velocity of the outflowing estuarine water (related to river velocity) is sufficient to dominate the physical dynamics of the system. In this region, the less dense estuarine waters flow in a lens near the surface, spreading laterally, as well as mixing vertically with the ocean water underneath. The typical extent of a near-field region is approximately 1–3 km (0.6–2 mi).

The far-field plume begins at the outer boundary of the near-field region, or at the mouth if velocities are not sufficient to generate a near-field region, and often extends several hundreds of kilometers down the coast. In this region, the estuarine water continues to be mixed into the surrounding ocean by wind- and wave-generated processes. Because of effects caused by the Earth's rotation, a far-field plume will typically turn to the right outside of the river mouth and flow adjacent to the coast in the Northern Hemisphere or to the left in the Southern Hemisphere.

The mixing in the near-field region can be much more intense than the wind-driven mixing in the far-field plume, but acts across a very small region, compared to the size of the far field. Thus, several recent studies have suggested that the total amount of mixing occurring in each of the two regions may be roughly comparable.

Flushing and pollution-dispersal prediction. Much research into estuarine characteristics is aimed at predicting the distribution of effluents discharged into estuaries. Near the mouth of a partially mixed estuary, the salinity of the estuarine water is very near that of the adjacent ocean water, implying that the incoming fresh river water has been diluted significantly by mixing with a much larger volume of ocean water. Consequently, estuaries are more effective than rivers in diluting and removing pollutants. It has been observed that increased river flow causes both a downstream movement of the saline intrusion and a more rapid discharge of water to the

sea. The latter effect occurs because increased river discharge increases stratification; increased stratification diminishes vertical mixing and enhances the flow toward the sea in the surface layer. Thus, increased river discharge has the effect of increasing the volume of freshwater accumulated in the estuary, but to a lesser extent than the increase of the discharged volume. Obviously, it takes some time for the freshwater from the river to pass through the estuary. A rough estimate of the flushing time can be determined by dividing the total volume of freshwater accumulated in the estuary by the river flow. For most estuaries the flushing time is 5 to 10 days.

If a conservative, nondecaying pollutant is discharged at a constant rate into an estuary, the effluent concentration in the receiving water will vary with the tidal current velocity and will spread out by means of turbulent mixing. The concentrations will be increased during the next half cycle as the water passes the discharge point again. After several tidal cycles, a steady-state distribution will be achieved, with the highest concentration near the discharge point. Concentrations will decrease downstream but not as quickly as they do upstream. However, the details of the distribution will depend largely on whether the discharge is of dense or light fluid and whether the discharge is into the lower or upper layer. Since its movement will be modified by the estuarine circulation, the effluent will be more concentrated in the lower layer upstream of the discharge point, and it will be more concentrated in the upper layer downstream. To obtain maximum initial dilution, a light effluent would have to be discharged near the estuary bed so that it would mix rapidly as it rose.

For nonconservative pollutants, such as coliform sewage bacteria, prediction becomes more difficult. The population of bacteria dies progressively through the action of sunlight, and concentrations diminish with time as well as by dilution. The faster the mixing, the larger the populations at any distance from the point of introduction, since less decay occurs.

Because of the poor mixing of freshwater into a salt-wedge estuary, an effluent introduced in the surface layer will be flushed from the estuary before it contaminates the lower layer, provided that it is not too dense.

Mathematical modeling. Increasingly, mathematical modeling is being used, with reasonable success in many instances, to predict effluent dispersal with a minimum amount of field data. Although the governing mathematical equations can be stated, they cannot be solved in their full form because there are too many unknowns. To reduce the number of unknowns, various assumptions are made, including some form of spatial averaging to reduce a three-dimensional problem to two dimensions or even one dimension. Mixing parameters, about which little is known, are assumed constant or are considered as a simple variable in space, and are altered so that the model fits the available prototype data.

The first step is usually to model the flow and salin-

ity distribution. Because the density field is important in determining the flow characteristics, density and flow are interlinked problems. Then, for pollutant studies the pollutant is assumed to act in the same manner as fresh or salt water, or the flow parameters are used with appropriate mixing coefficients to predict the distribution. Simple models consider the mean flow to be entirely the result of river discharge, and tidal flow to be given by the tidal prism. Segmentation is based on simple mixing concepts and crude mixing ratios. Salinity and pollutant concentrations can then be calculated for cross-sectionally averaged and vertically homogeneous conditions by using the absolute minimum of field data. These models are known as tidal prism models. One-dimensional models are very similar, but use a finer grid system and need better data for validation. Two-dimensional models assume vertical homogeneity and allow lateral variations, or vice versa. There are difficulties in including the effects of tidally drying areas and junctions; the models become more costly and require extensive prototype data, but they are more realistic. The ideal situation of modeling the flow and salinity distribution accurately simply on the basis of knowledge about the topography, the river discharge, and the tidal range at a number of points is still a long way off.

Ecological environments. Estuarine ecological environments are complex and highly variable when compared with other marine environments. They are richly productive, however. Because of the variability, fewer species can exist as permanent residents in this environment than in some other marine environments, and many of these species are shellfish that can easily tolerate short periods of extreme conditions. Motile species can escape the extremes. A number of commercially important marine forms are indigenous to the estuary, and the environment serves as a spawning or nursery ground for many other species.

River inflow provides a primary source of nutrients such as nitrates and phosphates which are more concentrated than in the sea. These nutrients are utilized by plankton through the photosynthetic action of sunlight. Because of the energetic mixing, production is maintained throughout, in spite of the high levels of suspended sediment which restrict light penetration to a relatively thin surface layer. Plankton concentrations can be extremely high, and then, higher levels of the food web—filter-feeding shellfish and young fish—have an ample food source. The rich concentrations provide large quantities of organic detritus in the sediments which can be utilized by bottom-feeding organisms and which can be stirred up into the main body of the water by tidal action. For a more complete treatment of the ecology of estuarine environments from the biological viewpoint *see* MARINE ECOLOGY.

The stratification present in estuaries tends to produce concentrated regions of detritus and microorganisms, which are attractive to other species as significant food sources. This can be particularly true where a region of strong stratification intersects the

bottom of the estuary, a region known as a front. Higher populations of organisms, from juvenile fish to seals, can often be seen congregating near fronts.

There is a close relationship between the circulation pattern in estuaries and the faunal distributions. Several species of plankton peculiar to estuaries appear to confine their distribution to the estuary by using the water-circulation pattern; pelagic larvae of oysters are transported in a similar manner. The fingerling fish (*Micropogon undulatus*), spawned in the coastal waters off the eastern coast of the United States, are carried into the estuarine nursery areas by the landward residual bottom flow.

Sediments. The patterns of sediment distribution and movement depend on the type of estuary and on the estuarine topography. The type of sediment brought into the estuary by the rivers, by erosion of the banks, and from the sea is also important; and the relative importance of each of these sources may change along the estuary. Fine-grained material will move in suspension and will follow the residual water flow, although there may be deposition and reerosion during times of locally low velocities. The coarser-grained material will travel along the bed and will be affected most by high velocities and consequently, in estuarine areas, will normally tend to move in the direction of the maximum current.

Fine-grained material. Fine-grained clay material, about 2 micrometers in size, brought down the rivers in suspension can undergo alterations in its properties in the sea. Base (cation) exchange with the seawater can alter the chemical composition of some clay minerals; also, because the particles have surface ionic charges, they are attracted to one another and can flocculate. Flocculation depends on the salinity of the water and on the concentration of particles. It is normally complete in salinities in excess of 4 parts per thousand, and with suspended sediment concentrations above about 300 parts per million (1 mg/liter), and has the effect of increasing the settling velocities of the particles. The flocs have diameters larger than 30 μm but effective densities of about 1.1 g/cm³ because of the water closely held within. If the material is carried back into regions of low salinity, the flocculation is reversed and the flocs can be disrupted by turbulence. In sufficiently high concentrations, the suspended sediment can suppress turbulence. The sediment then settles as layers which can reach concentrations as high as 300,000 parts per million and which are visible as a distinctive layer of "fluid mud" on echo-sounder recordings. At low concentrations, aggregation of particles occurs mainly by biological action.

Turbidity maximum. A characteristic feature of partially mixed estuaries is the presence of a turbidity maximum. This is a zone in which the suspended sediment concentrations are higher than those in the river or farther down the estuary. This zone, positioned in the upper estuary around the head of the salt intrusion and associated with mud deposition in the so-called mud reaches, is often related to wide tidal mud flats and saltings. The position of the turbidity maximum changes according to changes in

river discharge, and is explained in terms of estuarine circulation. Suspended sediment is introduced into the estuary by the residual downstream flow in the river. In the upper estuary, mixing causes an exchange of suspended sediment into the upper layer, where there is a seaward residual flow causing downstream transport. In the middle estuary, the sediment settles into the lower layer in areas of less vigorous mixing to join sediment entering from the sea on the landward residual flow. It then travels in the salt intrusion back to the head of the estuary. This recirculation is very effective for sorting the sediment, which is of exceedingly uniform mineralogy and settling velocity. Flocs with low settling velocities tend to be swept out into the coastal regions and onto the continental shelf. The heavier or larger flocs tend to be deposited.

The concentrations change with tidal range and during the tidal cycle, and fluid muds can occur within the area of the turbidity maximum if concentrations become sufficiently high. During the tidal cycle, as the current diminishes, individual flocs can settle and adhere to the bed, or fluid muds can form. The mud consolidates slightly during the slack water period, and as the current increases at the next stage of the tide, erosion may not be intense enough to remove all of the material deposited. A similar cycle of deposition and erosion occurs during the spring-neap tidal cycle. Generally, there is more sediment in suspension in the turbidity maximum than is required to complete a year's sedimentation on the estuary bed.

Mud flats and tidal marshes. The area of the turbidity maximum is generally well protected from waves, and there are often wide areas of mud flats and tidal marshes (Fig. 4). These areas also exchange considerable volumes of fine sediment with sediment in suspension in the estuary. At high water the flats are covered by shallow water, and there is often a long stand of water level which gives the sediment time to settle and reach the bottom, where it adheres or



Fig. 4. Aerial photograph of tidal flats showing the areas of pans, marshes, and vegetation between the channels, Scott Head Island, England. (Photograph by J. K. St. Joseph, Crown copyright reserved)

is trapped by plants or by filter-feeding animals. The ebb flow is concentrated in the winding creeks and channels. At low water there is not enough time for the sediment to settle, and it is distributed over the tidal flats during the incoming tide. Thus, there is a progressive movement of the fine material onto the mud flats by a process that depends largely on the time delay between sediment that is beginning to settle and sediment that is actually reaching the bed. The tidal channels migrate widely, causing continual erosion. Thus, there is a constant exchange of material between one part of the marshes and another by means of the turbidity maximum. As the muds that are eroded are largely anaerobic, owing to their very low permeability, the turbidity maximum is an area with reduced amounts of dissolved oxygen in the water.

Coarse-grained material. Coarser materials such as quartz sand grains that do not flocculate travel along the bed. Those coming down the river will stop at the tip of the salt intrusion, where the oscillating tidal velocities are of equal magnitude at both flood and ebb. Ideally, coarser material entering from the sea on the landward bottom flow will also stop at the tip of the salt intrusion, which becomes an area of shoaling, with consequent decrease of grain size inland. Normally the distribution of the tidal currents is too complex for this pattern to be clear. Especially in the lower part of the estuary, lateral variations in velocity can be large. The flood and ebb currents preferentially take separate channels, forming a circulation pattern that the sediment also tends to follow. The channels shift positions in an apparently consistent way, as do the banks between them. This sorts the sediment and restricts the penetration of bed-load material into the estuary.

Salt-wedge patterns. In salt-wedge estuaries, the river discharge of sediment is much larger, though generally markedly seasonal. Both suspended and bed-load materials are important. The bed-load sediment is deposited at the tip of the salt wedge, but because the position of the salt wedge is so dependent on river discharge, the sediments are spread over a wide area. At times of flood, the whole mass of accumulated sediment can be moved outward and deposited seaward of the mouth. Because of the high sedimentation rates, the offshore slopes are very low, and the sediment has a very low bearing strength. Under normal circumstances, the suspended sediment settles through the salt wedge, and there is a zonation of decreasing grain size with distance down the salt wedge, but changes in river flow seldom allow this process to occur.

Fiord patterns. Sedimentation often occurs only at the heads of fiords, where river flow introduces coarse and badly sorted sediment. The sediment builds out into deltalike fans, and slumping on the slopes of the fan carries the sediment into deep water. Much of the rest of the fiord floor is bare rock or only thinly covered with fine sediment.

Bar-built estuary patterns. Bar-built estuaries are a very varied sedimentary environment. The high tidal currents in the inlets produce coarse lag deposits,

and sandy tidal deltas are produced at either end of the inlets, where the currents rapidly diminish. In tropical areas, the muds that accumulate in the lagoons can be very rich in chemically precipitated calcium carbonate. Daniel G. MacDonald; K. R. Dyer

Bibliography. R. S. K. Barnes and J. Green (eds.), *The Estuarine Environment*, 1972; K. R. Dyer, *Coastal and Estuarine Sediment Dynamics*, 1988; K. R. Dyer (ed.), *Estuarine Hydrography and Sedimentation*, 1980; H. B. Fischer et al., *Mixing in Inland and Coastal Waters*, 1979; B. J. Neilson, J. Brubaker, and A. Kuo (eds.), *Estuarine Circulation*, 1989; B. J. Neilson and E. L. Cronin (eds.), *Estuaries and Nutrients*, 1981; C. B. Officer, *Physical Oceanography of Estuaries and Associated Coastal Waters*, 1976.

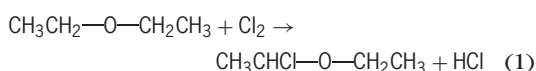
Ether

One of a class of organic compounds characterized by the structural feature of an oxygen atom linking two hydrocarbon groups, R—O—R'. Ethers are used widely as solvents, both in chemical manufacture and in the research laboratory. The most important ether is ethyl ether, C₂H₅OC₂H₅.

The hydrocarbon radicals R and R' may be identical (simple ether) or different (mixed ether). They may be aromatic or aliphatic, and the names of the ethers correspond to the hydrocarbon groups present. Thus, CH₃—O—CH₃ is methyl ether, rarely dimethyl ether, and C₆H₅—O—CH₃ is phenyl methyl ether.

Properties. Ethers are less soluble in water than are the corresponding alcohols, but are miscible with most organic solvents. Low-molecular-weight ethers have a lower boiling point than the corresponding alcohols, but for those ethers containing radicals larger than butyl, the reverse is true. The boiling points approximate those of hydrocarbons of the same molecular weight and geometry, indicating that association of ether molecules in the liquid state is negligible. Inertness at moderate temperatures, an outstanding chemical characteristic of the saturated alkyl ethers, leads to their wide use as reaction media. The organic magnesium compounds known as the Grignard reagents, RMgX, perhaps the most used reagents in organic synthesis, are almost always prepared in ether solutions, and suspensions of alkali metals in ethers are often employed. At higher temperatures, however, ethers are split by the alkali metals and by the halide salts of metalloids.

Ethers may also be split by hydrogen halides. Hydrogen iodide, HI, for example, often reacts at room temperature to form an alcohol and an alkyl iodide. Ethers react with chlorine and bromine considerably more readily than do the corresponding hydrocarbons. The initial reaction involves the formation of hydrogen halide and the substitution of a halogen atom for one of the hydrogens on a carbon adjacent (alpha) to oxygen, as shown by reaction (1).



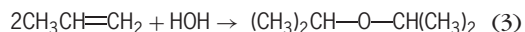
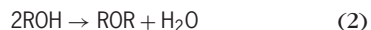
Because such α -halogens are reactive, the halogenated ethers are convenient intermediates for synthesis. Halogenated ethers are known in which the halogen is on a carbon other than that adjacent to the oxygen, but they are relatively inert.

On standing, ethers react with the oxygen of the air to form peroxides. Before distillation, it is essential that any considerable accumulation of peroxides be destroyed, by alkaline hydrolysis or by treatment with a reducing agent, such as ferrous hydroxide. On concentration and heating, ether peroxides detonate with dangerous violence. Some ethers form saltlike addition compounds with Lewis acids, the halogens, or picric acid. These addition compounds are theoretically related to, but usually much less stable than, the corresponding derivatives of amines. This property permits the separation of ethers from inert hydrocarbons by extraction with concentrated sulfuric acid.

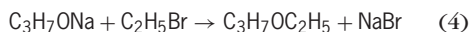
Identification of ethers is difficult. Often the more reactive components of a mixture are removed by chemical reagents, and the residual ethers are identified by a combination of their failure to react and their specific physical properties. The inertness of ethers is utilized in the syntheses of complicated organic molecules, an objectionably reactive alcohol group being protected by converting it to an unreactive ether. Hydrogen iodide may be used to regenerate the alcohol from the ether when the need for protection has passed.

Unsaturated ethers undergo the reactions usually associated with the double bond. Vinyl ethers, in which the double bond is adjacent to the oxygen, are readily polymerized or copolymerized with such monomers as vinyl acetate to yield useful polymers. Vinyl ethers also react in the presence of acid catalysts with compounds that possess active hydrogens. Thus, with alcohols, they form acetals.

Manufacture and preparation. Simple ethers may be considered to be the anhydrides of alcohols and are manufactured from alcohols by catalytic dehydration, as in reaction (2), or from olefins by controlled catalytic hydration, as in reaction (3). Mixed



ethers of definite structure may be prepared by the Williamson synthesis, shown by reaction (4). This



synthesis was of considerable significance historically because a knowledge of the structure of ethers was important in developing the radical theory, a stepping-stone to the present extensive knowledge concerning the arrangements of atoms in the molecules of organic compounds. A closely related reaction is that which takes place between cellulose, alkali, and ethyl chloride to yield an important plastic, the polyethyl ether of cellulose known as ethyl cellulose.

Ethyl ether. The best known of the ethers is ethyl ether, sometimes called diethyl ether, $\text{CH}_3\text{CH}_2\text{OCH}_2\text{CH}_3$. It is used in industry as a solvent and in medicine as an anesthetic.

The older process of manufacture involved heating ethyl (grain) alcohol to moderate temperatures with catalytic quantities of sulfuric acid, H_2SO_4 . Both ethyl ether and ethyl alcohol are manufactured by the controlled catalytic hydration of gasoline. Solid acid catalysts and a flow process are generally used instead of the sulfuric acid methods, and the proportion of alcohol to ether is controlled by variation of temperature and reactant concentrations.

The anesthetic properties were first noticed by Paracelsus (1490-1541) and were rediscovered by Michael Faraday in 1818, but it was not until 1846 that its potential as a surgical anesthetic was demonstrated, by W. T. G. Morton. The ethyl ether intended for anesthetic use differs from the ordinary variety in that possibly injurious impurities are removed. Peroxides are particularly harmful, and storage conditions must inhibit their formation. See ANESTHESIA.

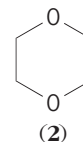
When ether is used as a solvent, its high volatility can cause loss. However this volatility is advantageous in that ether can be readily removed from the concentrated or crystallized product. The toxicity to humans is low, and recovery from overexposure is rapid and complete. It readily forms explosive mixtures with air, and on standing in containers which have been opened, it forms dangerous peroxides. Its freezing point is -117.4°C (-179.3°F); boiling point, 34.6°C (94.3°F); density, 0.7146; and refractive index, 1.35424. The solubility of ether in water is 6.18%, and of water in ether, 1.2%.

Cyclic ethers. Several cyclic ethers are of special importance and interest. The simplest of these is ethylene oxide or oxirane (1), made industrially by



the oxidation of ethylene with air over a silver catalyst. The major portion is used as an intermediate in the hydrolytic manufacture of ethylene glycol. Ethylene oxide is also used in the preparation of non-ionic emulsifying agents, plastics, plasticizers, one type of synthetic rubber, and several important synthetic textiles. Another important use is as a gaseous sterilizing agent.

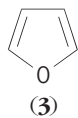
Dioxane or 1,4-dioxane (2) is prepared by the



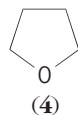
catalytic dimerization of ethylene oxide. It is unusual among substances of low dielectric constant (2.21) in that it is soluble in water in all proportions. Extensively used as a solvent industrially, it readily dissolves fats, waxes, natural and synthetic resins, cellulose

ethers, and lacquers, and it is employed by biologists to prepare paraffin-impregnated tissue sections.

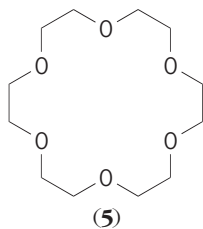
Furan (3), made by the decarbonylation of furfural,



is the most important ether obtained from an agricultural source. Most of it is hydrogenated to form the useful solvent tetrahydrofuran (4).



Crown ethers. Certain large-ring polyethers, the crown ethers, are able to increase the solubility of alkali metal salts in nonpolar organic solvents. Specific metal complexes are formed by these crown ethers with alkali metal cations, the specificity for a given cation depending upon the hole in the middle of the crown ether structure. The 18-member cyclic polyether containing six oxygen atoms is known as 18-crown-6 (5). It readily forms a complex with



the potassium cation. Potassium permanganate dissolves in benzene in the presence of 18-crown-6, forming a purple solution that can oxidize alcohols, alkenes, and alkylbenzenes under neutral conditions. See FURAN; HETEROCYCLIC COMPOUNDS; MACROCYCLIC COMPOUND. Paul E. Fanta

Ethers for oxygenated fuel. In the 1990 the U.S. Clean Air Act was revised to promote the production and use of cleaner-burning fuels. As a result, methyl tertiary butyl ether (MTBE) became a very important petrochemical. Nine United States cities with the most severe ozone pollution have been mandated to have a minimum of 2% oxygen in a reformulated gasoline. Outside the United States, other cities and countries are requiring oxygenated gasoline. Other environmental driving forces such as the phasing out of lead compounds and the reduction of aromatics in gasoline have increased the demand for ethers.

Methyl tertiary butyl ether is the most widely produced ether for oxygenates. It is commonly produced by the dehydrogenation of isobutane and the subsequent reaction of isobutylene with methanol; it is also produced as a coproduct in the propylene oxide technology. Tertiary amyl methyl ether, also used as an oxygenate for fuel, is produced in fluid catalytic crackers and steam naphtha crackers. Because of their high octanes, these and other aliphatic

ethers are good substitutes for aromatics in the gasolines. Also, their relatively low boiling temperature improves the distillation characteristics and drivability performance of the fuel. The oxygen in the ether also reduces the engine exhaust emissions. See GASOLINE; PETROCHEMICAL.

Tetrahydrofuran. Tetrahydrofuran is a solvent and monomer that is widely used. It has many applications in polymers, textiles, and composite materials. The most cost-effective commercial means of production is the oxidation of normal butane to maleic anhydride, which is subsequently reduced to tetrahydrofuran.

Anne Gaffney

Bibliography. R. T. Morrison and R. N. Boyd, *Organic Chemistry*, 6th ed., 1992; M. B. Smith and J. March, *March's Advanced Organic Chemistry: Reactions, Mechanisms, and Structure*, 5th ed., 2001.

Ether hypothesis

James Clerk Maxwell and his contemporaries in the nineteenth century found it inconceivable that a wave motion should propagate in empty space. They therefore postulated a medium, which they called the ether, that filled all space and transmitted electromagnetic vibrations.

During the last half of the nineteenth century, dozens of models were tried, but all broke down at some point. Direct experimental attempts to establish the existence of an absolute ether frame of reference, in which Maxwell's equations hold and light has the velocity c , have failed. The best known of these is the Michelson-Morley experiment, in which an attempt was made to measure the velocity of the Earth relative to the ether. See LIGHT.

Every hypothesis (ether drag, Lorentz contraction, and so on) invented to reconcile some experiment with the ether concept has been disproved by some other experiment. At present, there is no evidence whatever that the ether exists. See MAXWELL'S EQUATIONS. William R. Smythe

Bibliography. J. D. Jackson, *Classical Electrodynamics*, 3d ed., 1998; R. Resnick and D. Halliday, *Basic Concepts in Relativity and Early Quantum Theory*, 2d rev. ed., 1995; L. S. Swenson, Jr., *The Ethereal Aether: History of the Michelson-Morley-Miller Aether-Drift Experiments, 1880-1930*, 1972; E. T. Whittaker, *A History of the Theories of Aether and Electricity*, 2 vols., 1954, reprint 1987.

Ethology

The study of animal behavior. Modern ethology includes many different approaches, but the original emphasis, as expounded by Konrad Lorenz and Niko Tinbergen, was placed on the natural behavior of animals. This contrasted with the focus of comparative psychologists on behavior in artificial laboratory situations such as mazes and puzzle boxes.

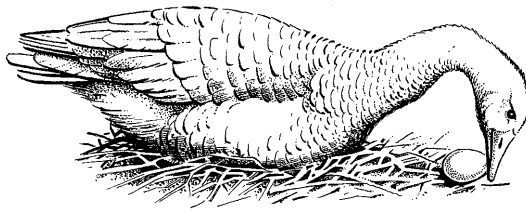


Fig. 1. Egg-rolling response of a greylag goose studied by Lorenz and Tinbergen is evoked by a great variety of shapes and objects near the nest.

Ethologists view the naturalistic approach as crucial because it reveals the environmental and social circumstances in which the behavior originally evolved, and prepares the way for more realistically designed laboratory experiments. The approach goes back to the stress that Charles Darwin placed on hereditary contributions to behavior in all species, including humans. Viewing behavior as a product of evolutionary history has helped to elucidate many otherwise puzzling aspects of its biology and has paved the way for the new science of neuroethology, concerned with how the structure and functioning of the brain controls behavior and makes learning possible.

Innate release mechanisms. A central concept in classical ethology still valuable today is that of the innate release mechanism. This has changed considerably since Lorenz's first characterization of the "innate schemata" of animals, but the major emphasis is similar. If a species has had a long history of experience with certain stimuli, especially those involving survival and reproduction, then to the extent that genes affect the ability to attend closely to such stimuli, natural selection leads to adaptations enhancing responsiveness to them. In the study of innate schemata a common first step was investigation of the development of responsiveness to such stimuli in infancy, focusing on situations that the ethologist as a naturalist knew to be especially relevant to survival. Study of infants, whether animal or human, with innate proclivities that have yet to be submerged by the accumulation of new experiences in growing to adulthood, is a rich source of insights into the biological foundations of behavioral development. Early investigations of the begging of nestling birds, serving to elicit feeding by their parents, indicated the inappropriateness of the term innate schema, which seems to imply that the bird inherits complete mental images. The evidence showed instead that birds are attuned to very specific, isolated "sign stimuli" or "releasers." Only after enrichment and transformation by experience can one confidently begin to speak in terms of mental imagery, which develops out of experiences garnered while responding to sign stimuli.

The later term "innate releasing mechanism," set forth by Tinbergen and Lorenz in an analysis of egg-rolling behavior in the greylag goose (**Fig. 1**), eschews notions of innate mental imagery. The concept has proved fertile in understanding how genes influence behavioral development, and in focusing

attention of neuroethologists on inborn physiological mechanisms that permit learning while encouraging the infant to attend closely to very specific stimuli, the nature of which varies from species to species according to differences in ecology and social organization.

Such specific responsiveness can be imposed at many stages in the transformation of a sensation into a percept—sometimes centrally in the brain and sometimes in the very process of sensing a stimulus, especially in simple organisms. For example, the hearing organ on the antenna of a male mosquito is mechanically tuned to the wing tone of females of the species (**Fig. 2**). In a sense, she is the only thing he hears clearly. Similarly, the visual experience of a toad, although potentially rich and varied, with some capacity for change in the face of individual experiences, appears to be dominated by the "worm-likeness" or "snake-likeness" of shapes and images, the former associated with feeding and the latter with avoidance behavior (**Figs. 3 and 4**). When a human baby responds innately to simple components of a face, its brain is involved as well as its sense organs. It responds to faces as especially captivating and engrossing even to its untutored eyes, but unlike the toad, which learns only with deliberation, the baby quickly proceeds to learn more about detailed configurations of the face, especially its mother's. Soon it is able to distinguish men from women and eventually one stranger from another, displaying skills that can be traced back to its first innate responses to the simple stimuli conveyed even by a crude caricature of the human face (**Fig. 5**).

In birds, such as the herring gull, innate release mechanisms are also better thought of as having evolved to guide processes of perceptual learning, rather than to design animals as though they were automata (**Fig. 6**). Learning is as important in the development of the behavior of many animals as in the human species. Yet as the young organism interacts with social and physical environments with which it has evolved adaptive relationships, the

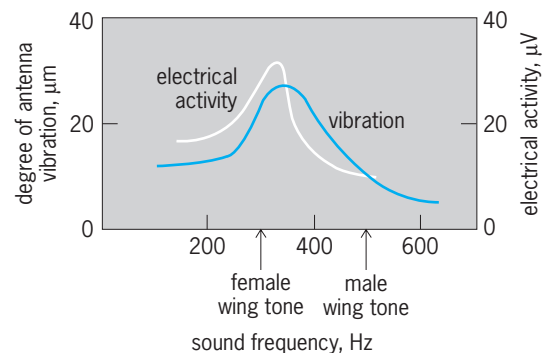


Fig. 2. Electrical activity and vibration of the male mosquito antenna. The hearing organs of the male mosquito, on the antennae, are tuned to resonate to the wing tone of females of the species, at about 300 Hz. This avoids jamming by the male's own wing tone, at about 500 Hz. (After P. Marler and W. H. Hamilton, *Mechanisms of Animal Behavior*, John Wiley and Sons, 1966)

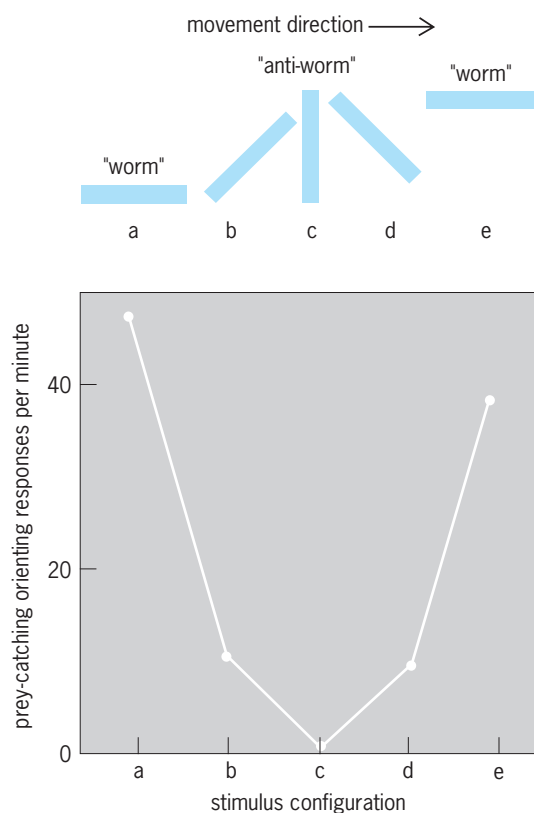


Fig. 3. A toad gives a feeding response to a “worm” stimulus but not to an “antiworm.”

course that learning takes in nature is guided and profoundly influenced by the innate predispositions that the organism brings to bear on dealing with the situation.

It is important, however, to examine the innateness question closely, and to appreciate what the ethologist is *not* claiming by invoking innate responsiveness. Babies are clearly not innately equipped to respond to all characteristics of the human face. The capacity to modify and supplement innate release mechanisms accommodates for individual differences in experience. No two mothers are completely alike in voice, appearance, or in the ways in which they interact with their babies. Similarly in animals innate responsiveness focuses on a few key properties of a situation, with more complete and complex responsiveness learned by experience. Only in the earliest stages of avian and mammalian development is the world of the infant organism completely dominated by simple releasers. In adulthood, recognition may still depend heavily on sign stimuli, but never on them alone. They serve as introductions to the much more complex configurations of sights, sounds, smells, tastes, and textures that come to dominate behavior later in life. Yet the introductory role of innate release mechanisms is crucial in ensuring that some things are attended to closely and others are treated more casually, and in imposing some order on the process of building up impressions, progressing from the simple to the complex.

Fixed-action patterns. Perceptions of the external world provide a basis for both thought and action. It is a fundamental axiom of ethology that each organism’s brain is armed with genetically determined programs of action which, in their own way, are as predictable and controlled as the genetic programs for developing anatomical structures such as a brain or a face. Ethologists have shown that it is possible to reconcile the need to modify patterns of action on the basis of experience with the possession of basic patterns of action that are coordinated by the brain, innately controlled, and often distinct from species to species.

Ethological pioneers such as Oskar Heinroth and Lorenz have provided the insight that what appears to the uninitiated to be a continuously varying series of unrepeatable actions typically proves to have at its core relatively stereotyped “fixed action patterns.” In extreme cases, their structure can be used as a basis for historical reconstructions in much the same way that a paleontologist draws evolutionary inferences from variations in the structure of jaws, crania, and pelvic bones. Analyses of fixed action patterns have

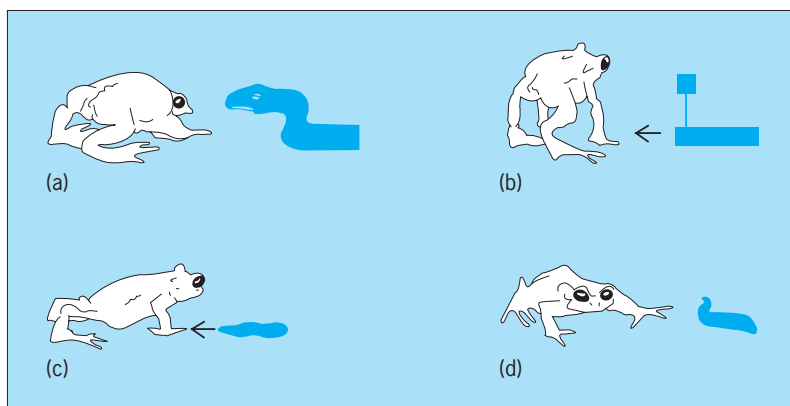


Fig. 4. Enemy images of the common toad. (a) A snake evokes avoidance. (b) A “head-rump” image also evokes avoidance. (c, d) A leech changes from a “worm” to an “enemy” when it rears its head. (After J. P. Ewert, *Neuroethology*, Springer-Verlag, 1980)

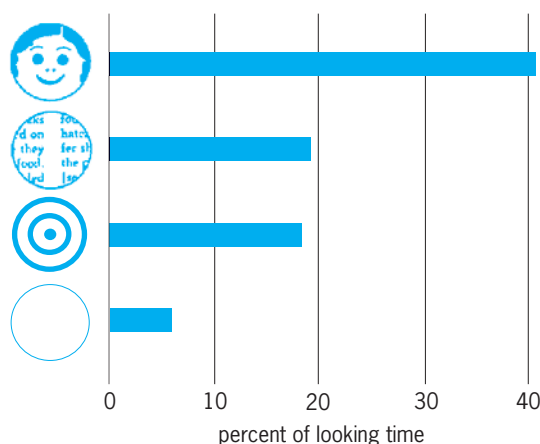


Fig. 5. Babies 2 to 3 months old spend more time looking at a facelike drawing than at newsprint or a bull’s-eye. (After R. L. Fantz, *The origins of form perception*, *Sci. Amer.*, 204:66-72, 1961)

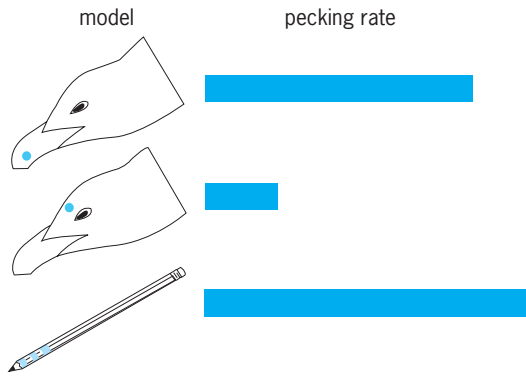


Fig. 6. Herring gull chick pecks at the spot, red in true life, on its mother's bill to get the food that is presented. The spot is critical for getting the newly hatched chick to peck. At this early age it pecks as much or more at a red-banded stick.

been used for reconstructing phylogenetic histories of organisms as diverse as fruit flies, fiddler crabs, fishes, ducks, gulls, and monkeys.

The concept of the fixed action pattern, like that of the innate release mechanism, has gone through many refinements and metamorphoses, but the underlying concept remains fundamental to understanding the development of the ability to act. Innate "motor programs," generated by the brain, are the natural units out of which behavior emerges during development. Each of these programs designates not a single completely stereotyped action, but a range of options which are limited but sufficient that selection among them allows for adjustments through experience. Close study reveals that sometimes the modifiability of actions lies in the potential flexibility of orientation, timing, and sequencing of actions rather than in the basic patterns or coordinations from which complex actions are built up. Thus nervous systems make some behavioral adjustments promptly and easily and others only with much greater difficulty, in harmony with species differences in the requirements for patterns of action as dictated by the species' structure and mode of life. Ethologists have found repeatedly that while social experience is vital in many animals for normal development of actions and responses, animals reared in restricted environments may still develop many units of action that are normal. The animal has to learn, however, how to put them together in an adaptive sequence.

Even when behavior improves as young animals mature, great care must be taken not to confuse effects of experience with consequences of maturation of the nervous system. A classic experiment demonstrated the emergence of normal movement in salamander embryos that were anesthetized through entire phases of development; thus the behavior is not due to practice or learning. In many cases, however, feedback from the performance of actions may be critically important for normal development, as when birds learn to sing.

Imprinting. Modern research on the ethology of learning began when Lorenz discovered imprinting

in geese. He found that if he led a flock of newly hatched goslings himself they became imprinted on him (**Fig. 7**). When mature, they would court people as though confused about their own species identity. Learning occurred very rapidly and tended to be restricted to a short sensitive phase early in life. A newly hatched gosling would normally follow its parents as they led the young away from the nest, and it was only by interposing himself as a parental surrogate that Lorenz discovered that a learning process is involved. The learning is highly focused by genetically determined preferences both to follow a parent-object with particular appearance and emitting species-specific calls, and also to learn most quickly and accurately at a particular stage of development. The interplay between nature (genetic predisposition) and nurture (environmental influence) in learning is displayed especially clearly in imprinting, hence its special interest to biologists and psychiatrists. Indications are that it is not concerned so much with learning about species as with learning to recognize individual parents and kin, both to ensure mating with one's own kind and to avoid incestuous inbreeding.

There are many forms of imprinting. So-called filial imprinting, ensuring that ducklings and goslings follow only their parent, is distinct from sexual imprinting, affecting mate choice in adulthood; the sensitive phases for learning are different in each case. Imprinting-like processes also shape the development of food preferences and abilities to use the sun and stars in navigation.

Song learning has many of the same features as imprinting. There are local dialects in many birdsongs so that, for example, in the San Francisco Bay area one could get a good idea of location just by listening



Fig. 7. Goslings following Lorenz, on whom they are imprinted, adopting him as their mother.

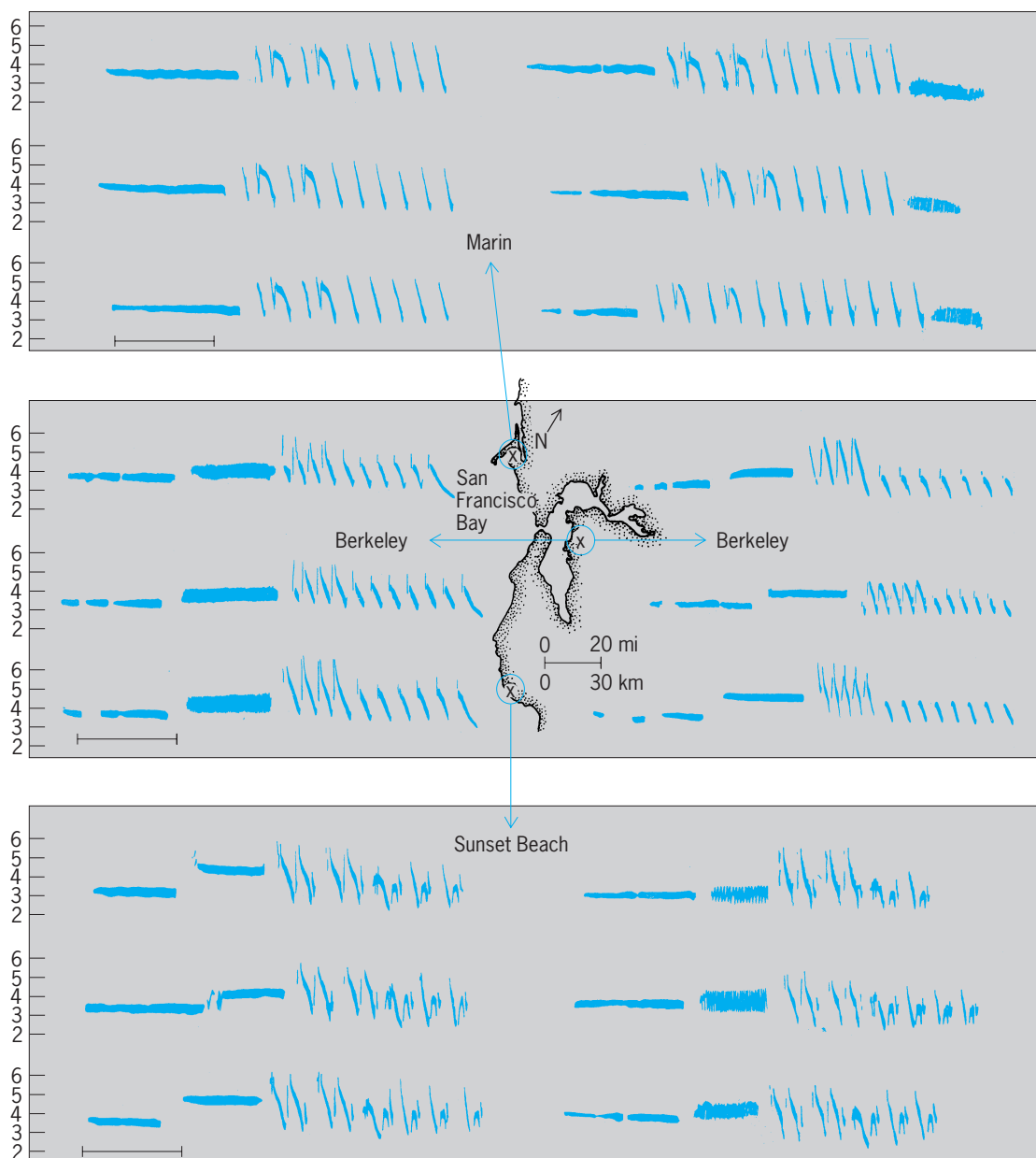


Fig. 8. Sound spectrograms of songs of 18 white-crowned sparrows in the San Francisco Bay area. Marin, Berkeley, and Sunset Beach sparrows each have a different dialect. Sound spectrograms can be read from left to right with sound frequency (in hertz) as a measure of pitch on the vertical axis. (After P. Marler, *A comparative approach to vocal learning: Song development in white-crowned sparrows*, *J. Comp. Physiol. Psychol.*, 71 (no. 2, pt. 2):1–25, 1970)

to the white-crowned sparrow (**Fig. 8**). A sparrow reared away from members of its own kind will sing an abnormal song, lacking any trace of the local dialect. If taped songs are played to a young sparrow, roughly between 2 and 7 weeks of age, it will then memorize the song patterns (**Fig. 9**). Later it will not only sing normally, but will reproduce that particular dialect, even if the dialect is from a different area than the sparrow's birthplace.

This experiment, which has been conducted with many songbirds, shows clearly that songs are nurtured by experience. But nature also intrudes at many points in the learning process. If two species are reared in isolation, both will sing abnormally (**Fig. 10**), but their songs are quite distinct from one

another, preserving enough normal features so that their species identity can easily be distinguished. The deficiencies are clear if the songs are played over a loudspeaker on the territories of wild males. Normal songs are attacked more than those of males reared in isolation, but the latter still get more responses than the songs of another species. The size of individual male song repertoires varies from species to species, ranging from one into the hundreds. Although innate repertoires are always smaller than learned ones, species differences nevertheless persist.

Another way to demonstrate the interaction of nature and nurture is to give the young sparrow a choice of songs of more than one species. Some

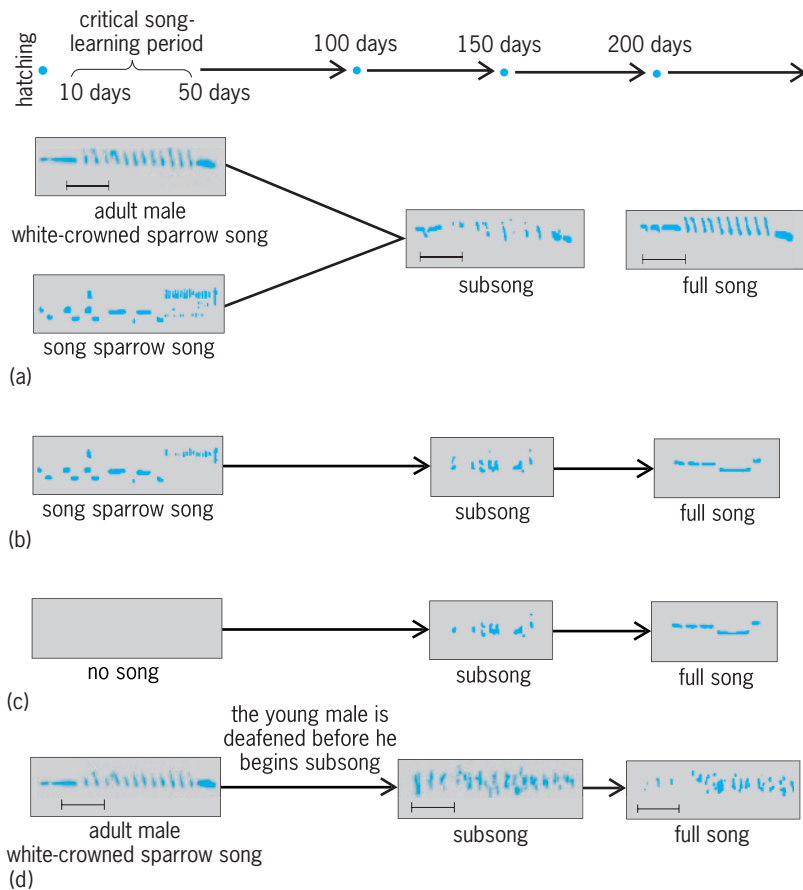


Fig. 9. Diagrams of song development in white-crowned sparrows: (a) when played recordings of normal song; (b) when given songs of a different species, the song sparrow; and (c) when reared in isolation. (d) Song development in a deaf male. Songs of other species may be rejected for learning, as shown with song sparrow vocalization in a and b.

birds, such as starlings and mockingbirds, learn the songs of other species. African wydah birds lay their eggs in other birds' nests, and learn songs from their hosts. But these birds are exceptions. If choices are offered, birds are likely to favor songs of their own kind. Ethologists theorize that songbirds inherit an auditory template in the brain that both helps them in selecting songs of their own kind for learning, and also helps the young male to develop his own song at a later stage. However, he must be able to hear his own voice to achieve this. This is evident from the buzzy, tuneless song of a bird that is deaf, a song so anonymous that it lacks even those characteristics that the bird develops normally in isolation from its own kind. Thus innate auditory templates are responsible for some of the normal features of a birdsong, exerting their influence by way of auditory feedback.

Research on song learning is one of the clearest cases in which neuroethologists have been led to fundamentally new discoveries about how brains support learning. The new insights stem directly from research in the ethological tradition, demonstrating genetic contributions to the song-learning process. Thus unlike psychological studies of animal learning in the laboratory, which have tended to favor the "blank-slate" view of the brain's contribution to learning, ethology emphasizes the need to understand all aspects of the biology of a species under study before one can hope to understand how the animal learns to cope with the many complexities of individual existence and social living. Thus ethology may lead not only to an understanding of how natural behavior evolves, but also to new insights into how

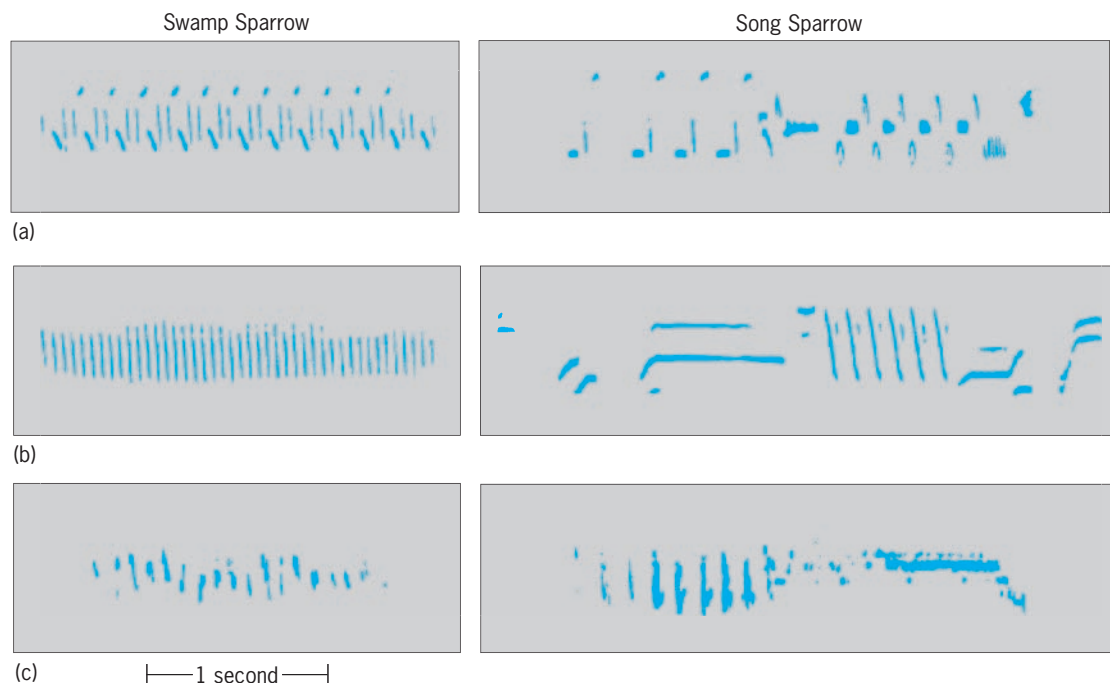


Fig. 10. Comparison of songs of the swamp sparrow and the song sparrow: (a) as heard in nature (normal); (b) as heard from males reared from the egg without any opportunity to hear the normal song of their species (isolated); and (c) as produced by deaf birds. (After P. Marler, *Song learning: innate species differences in the learning process*, in P. Marler and H. Terrace, eds., *The Biology of Learning, Dahlem Conference*, Springer-Verlag, 1984)

brains help organisms learn to cope with social and environmental problems confronting them as individuals. See ANIMAL COMMUNICATION; BEHAVIOR GENETICS; INSTINCTIVE BEHAVIOR. Peter R. Marler

Bibliography. J. L. Gould, *Ethology: The Mechanisms and Evolution of Behavior*, 1982; R. A. Hinde, *Ethology: Its Nature and Relations with Other Sciences*, 1982; J. R. Krebs and N. B. Davies, *An Introduction to Behavioural Ecology*, 3d ed., 1993.

Ethyl alcohol

Probably the best known of the alcohols, ethyl alcohol, formula $\text{CH}_3\text{CH}_2\text{OH}$, is also called alcohol, ethanol, grain alcohol, industrial alcohol, fermentation alcohol, cologne spirits, ethyl hydroxide, and methylcarbinol. Pure ethyl alcohol is a colorless, limpid, volatile liquid which is flammable and toxic and has a pungent taste. It boils at 78.4°C (173°F) and melts at -112.3°C (-170.1°F), has a specific gravity of 0.7851 at 20°C (68°F), and is soluble in water and most organic liquids. It is one of the most important industrial organic chemicals. Ethyl alcohol is produced by chemical synthesis and by fermentation or biosynthetic processes. See BIOSYNTHESIS; ORGANIC SYNTHESIS.

Ethyl alcohol is used as a solvent, extractant, antifreeze, fuel and fuel additive, and intermediate in the synthesis of innumerable organic chemicals. It is also an essential ingredient of alcoholic beverages.

Various grades of ethyl alcohol are produced, depending on their intended use. U.S. Pharmaceutical (USP XV) grade is the water azeotrope of ethyl alcohol and is 95% ethyl alcohol by volume. National Formulary (NFX) grade is 99+% ethyl alcohol by weight; it is also called absolute, or anhydrous, alcohol. This grade is generally prepared by azeotropic dehydration with benzene and therefore usually contains about 0.5% benzene. Denatured alcohol contains a small amount of a malodorous or obnoxious material to prevent the use of this grade of ethyl alcohol for beverage purposes. See AZEOTROPIC MIXTURE.

The major use of ethyl alcohol is as a starting material for various organic syntheses. Bimolecular dehydration of ethyl alcohol gives diethyl ether, which is employed as a solvent, extractant, and anesthetic. Dehydrogenation of ethyl alcohol yields acetaldehyde, which is the precursor of a vast number of organic chemicals, such as acetic acid, acetic anhydride, chloral, butanol, crotonaldehyde, and ethylhexanol. Reaction with carboxylic acids or anhydrides yields esters which are useful in many applications. The hydroxyl group of ethyl alcohol may be replaced by halogen to give the ethyl halides. Treatment with sulfuric acid gives ethyl hydrogen sulfate and diethyl sulfate, a useful ethylating agent. Reaction of ethyl alcohol with aldehydes gives the respective diethyl acetals, and reaction with acetylene produces the acetals, as well as ethyl vinyl

ether. Treatment of ethyl alcohol with ammonia produces acetonitrile, which may be reduced to ethylamine. These and other ethyl alcohol-derived chemicals are used in dyes, drugs, synthetic rubber, solvents, extractants, detergents, plasticizers, lubricants, surface coatings, adhesives, moldings, cosmetics, explosives, pesticides, and synthetic fiber resins.

Herman J. Phaff; Emil M. Mrak
Bibliography. J. R. Critser, Jr., *Biotechnical Engineering: Equipment and Processes*, 1994; M. Roehr (ed.), *The Biotechnology of Ethanol: Classical and Future Applications*, 2001; H. C. Vogel and C. C. Todaro (eds.), *Fermentation and Biotechnical Engineering Handbook: Principles, Process Design, and Equipment*, 2d ed., 1996.

Ethylene

A colorless gas, formula $\text{CH}_2=\text{CH}_2$, with a boiling point of -103.8°C (-155°F) and a melting point of -169.4°C (-273°F). Ethylene is the most important synthetic organic chemical in terms of volume, sales value, and number of derivatives. About half of the ethylene produced is used in the manufacture of polyethylene; ethylene dichloride and vinyl chloride production uses about 20%, synthesis of ethylene oxide and derivatives account for about 12%, and styrene production consumes about 8% of the ethylene. Other important derivatives are ethanol, vinyl acetate, and acetaldehyde.

Production. Thermal cracking of hydrocarbons in the presence of steam is the most widely used process for producing ethylene. Cracking is done at about 1600°C (2912°F) and 30 lb/in.² absolute (21 kilopascals absolute) pressure, followed by rapid quenching to below 1000°C (1760°F). Ethylene is recovered by low-temperature fractionation at 500–550 lb/in.² absolute (340–380 kPa absolute) and purified by very low-temperature (approximately -65°C or -85°F) gas-separation procedures to remove hydrogen, methane, and ethane.

In the United States, hydrocarbon gases—ethane through the butanes—account for about 65% of the feedstocks used to produce ethylene. With these clean raw materials, yields are good (about 50%), and efficiencies and conversion are high (about 80%). Liquid refinery products such as naphthas, kerosines, and gas oils constitute about 25% of process raw materials in the United States. While with these more complex feedstocks the yields and efficiencies are about 30%, the trend has been toward the use of the more available, heavier feedstocks. In some countries where petroleum raw materials are expensive, ethylene is produced by the dehydration of fermentation ethanol.

Derivatives. Polyethylene, the most important derivative of ethylene, is produced by both high- and low-pressure processes to make high- and low-density, high-molecular-weight thermoplastic polymers. Aluminum alkyl catalysts (Ziegler polymerization) are used to polymerize ethylene to relatively low-molecular-weight, straight-chain hydrocarbon

derivatives which are convertible to even-numbered-carbon linear olefins, alcohols, and acids. Commercial processes use palladium-catalyzed oxidation of ethylene to produce acetaldehyde, or if acetic acid is used as the solvent, vinyl acetate. Chlorination and oxychlorination processes are used to make vinyl chloride. Ethylene oxide is produced by silver-catalyzed oxidation of ethylene. Acid-catalyzed hydration of ethylene produces ethanol competitively with fermentation processes. *See* ACETYLENE; ALKENE; ETHYL ALCOHOL; ETHYLENE OXIDE; POLYMER.

Robert K. Barnes

Function as plant hormone. Ethylene exerts a major influence on many aspects of plant growth, development, and senescence. It is produced by plants and acts in trace amounts; the effects on plants are spectacular and commercially important.

The ability of certain gaseous agents existing in smokes, fumes, or illuminating gas to stimulate fruit ripening, to induce pineapple flowering, to hasten the coloring ("degreening") of harvested lemons, or to shed tree leaves prematurely, was recognized long before ethylene was identified as a plant growth regulator. In 1901, D. Neljubow identified ethylene as the agent in illuminating gas which most effectively caused the abnormal growth of dark-grown pea seedlings. Since then, those growth-regulatory effects of illuminating gas and various fumes and smokes have been attributed to the presence of ethylene in the air.

Although fruit shippers were aware that ripe or rotting fruits could hasten the ripening of other fruits stored with them, it was not until 1934 that chemical proof of production of ethylene by ripe apples was provided. Fruit physiologists have since established that ethylene functions as a ripening hormone. Subsequently, ethylene has been identified as a natural product not only of fruits but of all plant parts, including leaves, stems, flowers, and roots. The rate of ethylene production varies with both the type of tissue and its stage of development. In germinating seeds, high rates of ethylene production occur when the radicle (primary root of the seedling) starts to penetrate the seed coat and during the period when the seedling forces its way through soil. In fruits and flowers, ethylene production is induced at the initial stage of fruit ripening and flower fading. The pathway of ethylene biosynthesis in plants has been recently elucidated: ethylene is derived from C-3,4 of methionine, an amino acid, via *S*-adenosylmethionine and 1-aminocyclopropane-1-carboxylic acid as intermediates.

Demonstration of the natural function of ethylene in fruit ripening has stimulated the search for other functions. Among the most diverse physiological effects which have been recognized are: breaking of seed and bud dormancy, stimulation of cortical air space (aerenchyma) development and adventitious root formation, inhibition of root growth, stem strengthening and internode shortening, stimulation of flowering in pineapple, induction of female flower in cucumber, promotion of fruit and leaf degreening, stimulation of fruit ripening, promotion of

fruit and leaf abscission, promotion of fading in various types of flowers, promotion of stem growth in some aquatic and semiaquatic plants, modification of gravitropism and stimulation of latex flow in rubber trees. *See* ABCISSION; DORMANCY.

The use of ethylene gas to regulate plant responses of commercial value is restricted to harvested crops (such as green banana, tomato, or citrus fruits) in special closed chambers. This problem was largely overcome with the introduction and commercialization of ethylene-releasing compounds, such as 2-chloroethylphosphonic acid, which is spontaneously degraded into chloride, ethylene, and orthophosphate upon dilution with water. This "liquid ethylene" can be applied to preharvest or postharvest crops, and is now registered for more than 20 crops. It is used for promoting ripening of fruits such as tomatoes, apples, coffee berries, and grapes; facilitating harvesting of cherries, walnuts, and cotton by accelerating abscission or fruit dehiscence; increasing rubber production by prolonging latex flow in rubber trees; increasing sugar production in sugarcane; synchronizing flowering in pineapple; and accelerating senescence of tobacco leaves. *See* PLANT GROWTH; PLANT HORMONES.

S. F. Yang

Bibliography. F. B. Abeles, *Ethylene in Plant Biology*, 2d ed., 1997; L. Albright, *Novel Production Methods for Ethylene, Light Hydrocarbons, and Aromatics*, 1991; M. Arshad and W. T. Frankenberg, Jr., *Ethylene: Agricultural Sources and Applications*, 2001; J. E. Brady and F. Senese, *Chemistry: Matter and Its Changes*, 4th ed., 2004.

Ethylene glycol

A colorless, nearly odorless, sweet-tasting, hygroscopic liquid, formula HOCH₂CH₂OH. It is relatively nonvolatile and viscous and is the simplest member of the glycol family. Ethylene glycol freezes at -3°C (8.6°F), boils at 197.6°C (387.7°F), and is completely soluble in water, common alcohols, and phenol. Low molecular weight, low volatility, water solubility, and low solvent action on automobile finishes make ethylene glycol ideal as a radiator antifreeze and coolant; water mixtures (58–80% glycol) freeze below -46°C (-50°F). *See* ANTIFREEZE MIXTURE; GLYCOL.

Worldwide, about half of all ethylene glycol is used in polyester resins, and about a third goes into antifreeze. Other uses for this commodity chemical are in explosives, brake and shock-absorber fluids, and alkyl-type resins. *See* POLYESTER RESINS.

The original commercial process for ethylene glycol involved hydrolysis of ethylene chlorohydrin derived from chlorine and ethylene. All commercial ethylene glycol is now produced by the vapor-phase oxidation of ethylene. A commercial process was developed to produce ethylene glycol directly from ethylene. This process uses a tellurium oxide/bromide ion catalyst to oxidize ethylene in acetic acid solution. Another patented process describes

production of ethylene glycol from synthesis gas (carbon monoxide and hydrogen).

Ethylene glycol undergoes the simple reactions of alcohols such as etherification, condensation, oxidation, and esterification. Mono- and dialkyl ethers of ethylene glycol are formed by its reaction with sodium hydroxide and dialkyl sulfates or alkyl halides. Dioxane, a cyclic diether, as well as diethylene glycol, can be prepared by acid-catalyzed dehydration of ethylene glycol. Di-, tri-, and polyethylene glycols are best produced by the acid- or base-catalyzed reaction of ethylene oxide with ethylene glycol. Commercially di-, tri-, and some tetraethylene glycols are produced as by-products of industrial ethylene glycol processes. Acid-catalyzed reactions of aldehydes and ketones with ethylene glycol produce five-membered cyclic acetals and ketals called 1,3-dioxolanes. Vapor-phase catalytic dehydrogenation (oxidation) of ethylene glycol can produce either 2-hydroxymethyl-1,3-dioxolane or glyoxal. Liquid-phase nitric acid or air oxidation produces oxalic, glycolic, and formic acids, formaldehyde, and glycol aldehyde and glyoxal. With nitric acid under dehydrating conditions, ethylene glycol forms the dinitrate ester, which is employed in conjunction with nitroglycerin to produce explosives that have low freezing points.

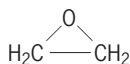
Dibasic acids or anhydrides react with ethylene glycol to form polyester condensation polymers. Reaction with terephthalic acid or its esters produces polyester resins which can be spun to fibers that find wide use in clothing and general fabrics applications. The polymer also has important film applications. Condensation of glycol with unsaturated diacids, for example maleic acid, followed by free-radical cross-linking reactions with polymerizable olefins such as styrene produces another commercially important polymer class identified as unsaturated polyester resins. Saturated and unsaturated polyester resins constitute the major use of ethylene glycol in the world today. See ETHER; ETHYLENE OXIDE; POLYETHYLENE GLYCOL.

Robert K. Barnes

Bibliography. J. E. Brady and F. Senese, *Chemistry: Matter and Its Changes*, 4th ed., 2004.

Ethylene oxide

The simplest cyclic ether or epoxide, with the formula C_2H_4O , and the structure



It is also called epoxyethane and oxirane. Ethylene oxide was discovered in 1859 by C. A. Wurtz. His preparation from ethylene chlorohydrin and aqueous base became the first commercial process, practiced during the early 1920s and until the mid-1940s. Thereafter the direct, silver-catalyzed vapor-phase oxidation of ethylene was the process of choice. Commercial processes use either air or oxygen to oxidize ethylene at low conversion and high selectiv-

ity to ethylene oxide. See ETHYLENE; HETEROCYCLIC COMPOUNDS.

Ethylene oxide is a colorless gas boiling at 10.4°C (50.7°F) and melting at -112°C (-170°F), with refractive index 1.3597 at 7°C (45°F), and density 0.8969 kg/liter at 0°C (32°F). Its vapors are flammable and explosive, and it is considered a relatively toxic liquid and gas. It is miscible in all proportions with water, alcohols, ethers, and other organic solvents. See ETHER.

Ethylene oxide reacts, usually by acid or base catalysis, with most active hydrogen compounds such as water, alcohols, acids, phenols, amides, hydrogen sulfide, and hydrogen cyanide; it also reacts with other organic compounds such as carbon dioxide, amines, ammonia, and Grignard reagents and with itself. Ethylene oxide's reactivity can be ascribed to the strained ring structure.

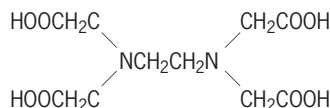
About 50% of the ethylene oxide produced is converted to ethylene glycol by reaction with water. Primary uses for ethylene glycol are in the manufacture of polyester resins and as an automobile antifreeze. About 30% of the ethylene oxide produced is used in the manufacture of nonionic surfactants, glycol ethers, and ethanalamines. Gaseous ethylene oxide in CO_2 or difluoromethane is used as a fumigant and a sterilizing agent for medical equipment. See ETHYLENE GLYCOL; POLYETHYLENE GLYCOL.

Robert K. Barnes

Bibliography. J. E. Brady and F. Senese, *Chemistry: Matter and Its Changes*, 4th ed., 2004; T. W. G. Solomons, *Organic Chemistry*, 7th ed., 1999.

Ethylenediaminetetraacetic acid

A chelating agent for metallic ions, abbreviated EDTA. The structural formula is



Tetrasodium EDTA is the most common form in commerce, but other metallic chelates are marketed, for example, iron, zinc, and calcium. Tetrasodium EDTA is a white solid, very soluble in water and forming a basic solution. Prepared from ethylenediamine, formaldehyde, and sodium cyanide in basic solution, or from ethylenediamine and sodium chloroacetate, EDTA is a strong complexing and chelating agent. It reacts with many metallic ions to form soluble chelates. As such, it is widely used in analysis to retain alkaline earths and heavy metals in solution. The iron chelate is useful in lawn management and gardening as a replacement for ferrous sulfate (copperas). Calcium EDTA is used to control the deterioration of natural seawater in aquariums. Calcium disodium salt of EDTA is used in pharmaceuticals to prevent calcium depletion of the body during therapy. See CHELATION; COORDINATION COMPLEXES

Frank Wagner
Bibliography. J. E. Brady and F. Senese, *Chemistry: Matter and Its Changes*, 4th ed., 2004.

Eucalyptus

A large and important genus of Australian forest trees; includes about 500 species in the family Myrtaceae. Only two species occur naturally outside Australia in the adjacent islands. Eucalypts occur throughout Australia except in coastal tropical and subtropical rainforests in Queensland and New South Wales and in temperate rainforests in Victoria and Tasmania. They are confined to water courses in the extensive arid zones of central and northwest Australia. Eucalypts grow from sea level to tree line (6600 ft or 2000 m).

Because of its large geographic range the genus exhibits many habits, from tall trees to multistemmed, shrubby species called mallees. The mountain ash (*Eucalyptus regnans*) of Victoria and Tasmania is the tallest hardwood in the world, reaching heights over 330 ft (100 m); only the coast redwoods of California are taller. Having epicormic buds in branches and stems, and lignotubers in the roots, many species are well adapted for surviving fire and drought. Some species have smooth bark, in which case they are referred to as gum-barked, while others have rough bark. Half-barked species have rough bark on the lower trunk and smooth bark on the upper stem and branch. Many species shed their bark in long ribbons.

Eucalyptus is an evergreen genus with four different leaf types—seedling, juvenile, intermediate, and adult—depending on plant maturity. Juvenile leaves of some species, particularly those that are silvery blue and oval, are extensively used for floral decorations. Most species have white or cream flowers; however, some species, particularly those from Western Australia, produce spectacular red (*E. ficifolia*) or yellow (*E. erythrocorys*) flowers and are planted widely as ornamentals. These ornamental species are small in habit and drought-resistant. Eucalypt seed cases are woody, and many species produce so-called gumnuts.

The major commercial species include mountain ash (*E. regnans*), alpine ash (*E. delegatensis*), messmate stringybark (*E. obliqua*), flooded gum (*E. grandis*), blackbutt (*E. pilularis*), spotted gum (*E. maculata*), jarrah (*E. marginata*), karri (*E. diversicolor*), tallowwood (*E. microcorys*), blue gum (*E. globulus*), and river red gum (*E. camaldulensis*). Heartwood is generally dense and varies in color from light brown, such as blackbutt, to dark red, such as jarrah.

Eucalyptus globulus (see **illus.**), with gray or brownish bark at the base and peeling bark above, is the common species throughout California. The red flowering eucalypt (*E. ficifolia*) is a popular ornamental in southern California. *Eucalyptus camaldulensis* has also been planted in Florida, and Hawaii has plantations of many species. Eucalypts have been planted widely for commercial use in Brazil and other South American countries, Africa, the Indian subcontinent, and the Middle East. They are used extensively for fuel and construction and are an important component of third world economies. Foliage of



Eucalyptus globulus, showing spray of mature foliage.

some species yields essential oils for medicines and perfumes. Tannins are extracted from the bark of certain species. See ESSENTIAL OILS; MYRTALES; TREE.

Robert L. Edmonds

Bibliography. D. J. Boland et al., *Forest Trees of Australia*, 1984; A. Keast (ed.), *Ecological Biogeography of Australia*, 1981; C. D. Pryor, *The Biology of Eucalypts*, 1976.

Eucarida

The largest and most highly evolved superorder of the crustacean class Malacostraca. It contains the orders Euphausiacea, Amphionidacea, and Decapoda. Despite the great morphological and ecological diversity of eucaridans, all share (1) a well-developed carapace that is fused to all the thoracic somites; (2) stalked and movable eyes, although in a few these have been secondarily reduced; (3) mandibles that lack a mobile element between the incisor and molar processes; (4) a telson without a caudal furca; (5) heart and gills that are thoracic in position; and (6) typically metamorphic larval development.

The Euphausiacea are pelagic, open-water crustaceans. The Amphionidacea, although also pelagic and broadly distributed, are represented by a single rather bizarre species. The Decapoda, including the shrimps, lobsters, crabs, and crablike crustaceans, abound from the intertidal to the abyssal seas but are found as well in terrestrial habitats and fresh-water rivers, streams, and lakes. Their morphological adaptations, although based on a common plan, are as ubiquitous as their habitats. See AMPHIONIDACEA; CRUSTACEA; DECAPODA (CRUSTACEA); EUPHAUSIACEA; MALACOSTRACA. Patsy A. McLaughlin

Bibliography. R. C. Brusca and G. J. Brusca, *Invertebrates*, 2d ed., 2002; P. A. McLaughlin, *Comparative Morphology of Recent Crustacea*, 1980; S. P. Parker (ed.), *Synopsis and Classification of Living Organisms*, 2 vols., 1982; F. R. Schram, *Crustacea*, 1986; F. R. Schram (ed.), *Crustacean Phylogeny, Crustacean Issues*, 1983.

Eucestoda

A subclass of tapeworms including most of the members of the class Cestoda, which is a taxon of the platyhelminthes. All species are endoparasites of vertebrates, living in the intestine or related ducts.

Morphology. Like other members of the class, the eucestodes have no digestive tract or mouth. Nutrition presumably occurs by absorption of food through the body surface. The body, or strobila, is usually very elongated and tapelike and frequently divided into segments, or proglottids, with replication of the hermaphroditic reproductive systems. In a few species there is duplication of both male and female organs within a single segment. The anterior end is usually modified into a holdfast organ, the scolex. Since a digestive tract is completely absent, the scolex is of solid construction, typically highly muscular with sucking depressions and hooks (Fig. 1). Behind the scolex there is usually a region of undifferentiated tissue, termed the neck, from which proliferation of new segments occurs. In some species a neck is not apparent, and segments begin immediately behind the scolex. The strobila, or chain of segments, typically arises by growth and differentiation, the smallest and least-differentiated segments being nearest to the neck. Young segments are always broader than long, but, as the segments grow in volume, they may become square or longer than broad. The body is clothed with what has been called a tegument, but it has been shown that this noncellular layer is relatively complex (Fig. 2). This might be expected since these animals actively absorb such nutrients as carbohydrates and amino acids.

Nervous system. The nervous system consists basically of a pair of longitudinal lateral nerves running through the length of the strobila. Additional nerve trunks are frequently present. The trunks are con-

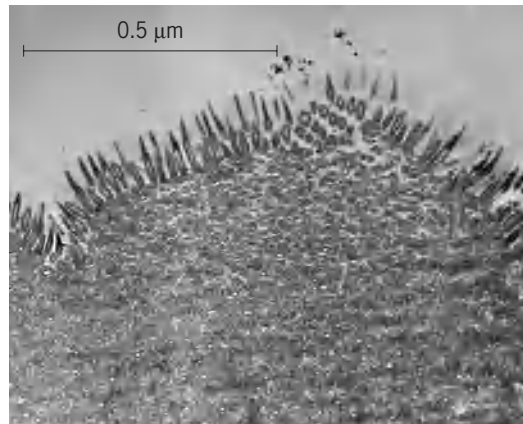


Fig. 2. The tegument of *Hymenolepis diminuta*. (Electron photomicrograph by A. H. Rothman)

nected by a ring-shaped commissure in each segment and by thickened commissures, a “brain,” in the scolex. Although the tapeworms lack special sensory structures, the body surface has many sensory nerve endings.

Excretory system. The excretory or osmoregulatory system consists of protonephridia, or flame bulbs, which connect ultimately through a system of fine ducts to a pair of longitudinal lateral canals on either side. These canals traverse the length of the body.

Reproductive system. The reproductive systems develop progressively along the strobila, the male system usually becoming mature before the female. In some forms there are two complete hermaphroditic systems. The male and female genital ducts commonly open into a common antrum, usually situated on the lateral margin of the segment. The testes are usually small rounded bodies, ranging in number from 1 to 1000 but most often about 100. Sperm are carried in a system of ducts to an eversible copulatory organ, the cirrus, which is frequently armed with spines or hooks. The ovary is single, commonly having two lobes. In the orders other than Cyclophyllida, yolk glands occur as numerous follicles. Yolk cells from these glands are gathered in a duct system which enters the oviduct from the ovary. The oviduct also receives the vagina and eventually enlarges to form the uterus. The uterus begins developing after the gonads and their associated ducts have matured and varies in shape from a simple sac to a ramified network.

With growth, the reproductive systems differentiate, and fertilization occurs in a modified region of the oviduct where the ovum is enclosed with yolk cells in a shell. These eggs develop into embryos, or oncospheres, within the uterus. Fertilization may occur within a single proglottid, between segments in the same strobila, or between proglottids of separate worms. Hypodermic impregnation occurs in a few cases. The posterior proglottids become distended with embryos. In some orders, such as Tetrathyrididae, the sexually mature proglottids are shed from the body before significant development of embryos, and the free proglottids undergo “ripening” in the host intestine. In most other orders the gravid,

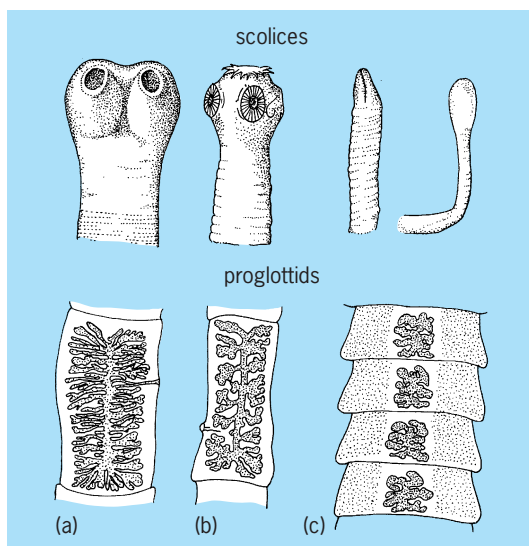


Fig. 1. Scolexes and mature proglottids of (a) *Taenia saginata*, (b) *T. solium*, and (c) *Dibothriocephalus latus*. (After T. I. Storer and R. L. Usinger, *General Zoology*, 3d ed., McGraw-Hill, 1957)

embryo-filled proglottids are shed from the strobila and pass out of the host, with or without degeneration of the segment tissue. In the Pseudophyllidea, gravid proglottids are retained and the shelled embryos are shed through a uterine pore. The growth rates of eucestodes show great variation. Different species of the same genus may show marked differences in growth rate and in the time during which a high growth rate, without senility, is maintained. The life of the sexually reproducing tapeworm varies from 4 or 5 days in some pseudophyllideans to 20 or more years in some cyclophyllideans. *See* CESTODA; TETRAPHYLLIDEA.

Larval development. The details of larval development are known for only a few tapeworms. Cleavage leads to the formation of a six-hooked embryo, the oncosphere, which is enclosed in a ciliated embryophore in the Pseudophyllidea. Further embryonic development occurs in some other host species. Except for Cyclophyllidea, two intermediate hosts are usually required for embryonic development. Except for some taeniaoid tapeworms, the first intermediate host is always an arthropod, while the second is often a fish, but may be an invertebrate. Invariably, the embryo attains entrance to the intermediate host by being eaten, and the vertebrate host is infected by eating the intermediate host.

The postoncosphere development of all eucestodes includes a proceroidlike stage. In most orders this is followed by the development of a plerocercoid larva in a second intermediate host. However, considerable modification of this developmental pattern occurs in the Cyclophyllidea, in which the postproceroid development continues in a single intermediate host and the plerocercoid is modified in one of several different ways. Other modifications of development in Cyclophyllidea include asexual reproduction of new individuals by budding. After differentiation and maturation of the scolex in the larval form, growth ceases until such time as the intermediate host is eaten by a vertebrate in which development of the strobila can occur. In general, the sexual strobilate phase of tapeworms shows considerable specificity with regard to the kind of vertebrate in which development occurs. It is rarely known whether this is due, in a given case, to a failure of contact of worm and host or to a physiological incompatibility. In many instances, such as the eucestodes of birds, tapeworms occur in only a single vertebrate host species. *See* PSEUDOPHYLLIDEA.

Nutrition. The nutritional relationships of adult tapeworms have been extensively studied. The worms require carbohydrate for growth and reproduction and this requirement is satisfied only from the host ingesta. On the other hand, nitrogenous nutrients and many micronutrients may be obtained from the body stores of the host; deleting such materials from the host's diet has no appreciable effect on the worms. In at least one case, *Dibothriocephalus*, it has been shown that the anemia appearing in some humans harboring the worm is attributable to the absorption of vitamin B₁₂ by the eucestode. The energy metabolism of tapeworms is primarily fermentative.

In air the rate of oxygen uptake is low and accounts for the oxidation of a small fraction of the carbohydrate metabolized.

Phylogeny. Most authorities agree that the tapeworms are ancient parasites, probably evolving as parasites of the earliest fishes. It seems probable that the tapeworms did not evolve from trematodes or other present-day groups of parasitic flatworms. The ancestry of the tapeworm may be directly derived from the acoele or rhabdocoele turbellarians and represents a line of evolution which is completely independent of other parasitic flatworms. These relationships remain obscure in the absence of any fossil records. *See* TURBELLARIA. Clark P. Read

Euclidean geometry

The word geometry is derived from two Greek words meaning earth measurement. It seems probable that many of the early discoveries in geometry were motivated by the need to make measurements of distances and areas on the Earth. However, euclidean geometry has a broader meaning. It is the chief subject matter of the monumental 13-volume work called *The Elements*, written about 300 B.C. by the Greek mathematician Euclid, who taught and founded a school of geometry at Alexandria. One of the milestones in the history of scientific thought, these books of Euclid still occupy an important position in mathematical instruction today.

Geometry, as developed by Euclid, was a systematic body of mathematical knowledge, built by deductive reasoning upon a foundation of three main pillars: (1) definitions of such things as points, lines, planes, angles, circles, and triangles; (2) the assumption of certain geometrical postulates regarded as true but perhaps not self-evident; and (3) the assumption of certain axioms or common notions which were taken to be self-evident truths. The body of Euclid's great work consists of a set of propositions or theorems, each derived systematically and logically from the definitions, axioms, and postulates of his foundation and from theorems already proved.

Definitions. Several of the more important definitions are given in this article. Other definitions and explanations are given in separate articles. *See* CIRCLE; POLYGON.

Point. A point, in the euclidean definition, is that which has no part. A point, according to most modern authors, is an undefined element in a geometry, such that to each pair of points corresponds a unique set of points called a line. In posteulidean metric geometry, each pair of points A and B determines a (real) number called the distance \overline{AB} , which is positive or zero according to whether A and B are distinct points or the same point.

Line. A line, according to the euclidean definition, is breadthless length, the extremities of a line are points, and a straight line is a line that lies evenly with the points of itself. Three points A , B , and C are said to be collinear if, and only if, they lie in

the same straight line. In modern usage the finite straight line of Euclid is called a line segment, and the single word line in geometry commonly refers to a euclidean straight line produced indefinitely in both directions. In metric geometry the point P is said to lie between A and B if $\overline{AP} + \overline{PB} = \overline{AB}$. Three points are collinear if one lies between the other two. A line segment $[AB]$ consists of two distinct points A and B and all the points between them. A line AB consists of two distinct points A and B and all the points collinear with them. A half line, or ray, $A(B)$ consists of the segment $[AB]$ and all points P such that B lies between A and P .

Any two lines (indefinitely extended) are congruent. Two lines having two distinct points in common are coincident; two lines having just one point in common are called intersecting lines. Any two rays (or half lines) are congruent. They are coincident if they have in common the vertex and one other point. Two line segments are congruent if and only if they have the same length.

Skew lines. Two lines that are not coplanar are called skew. Given any two skew lines in space, and a point P not on either, there is a unique line through P that intersects both skew lines. There is a unique plane that contains the first of two given skew lines and that is parallel to the second.

Parallels. Two lines (each extended indefinitely in both its directions) are called parallel if, and only if, they lie in the same plane but do not intersect. Two planes, or a line and a plane, are called parallel if they do not intersect. There is one and only one line parallel to a given line through a given point not on the line (Euclid's parallel postulate). Two distinct lines parallel to the same line (in space or in a plane) are parallel to each other. Two distinct planes parallel to the same plane are parallel to each other. If each of two parallel planes is intersected by a third plane, the lines of intersection are parallel (Fig. 1).

Perpendicular. When two straight lines intersect so that adjacent angles at their point of intersection are equal, these angles are called right angles, and the lines are said to be perpendicular, or normal, or orthogonal to each other (Fig. 2). Two skew lines are called perpendicular if one of them intersects at right angles a line parallel to the other. In general, a line perpendicular to one of two parallel lines is perpendicular to the other also. In a plane, but not in general in space, two lines each perpendicular to the same line are parallel. One, and only one, perpendicular can be drawn to a given line through a given point not on the line. The same is true in the plane, but not in space, if the given point is on the line. A line that intersects a plane in a point P is called perpendicular

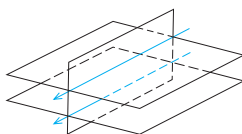


Fig. 1. Parallel lines and planes.

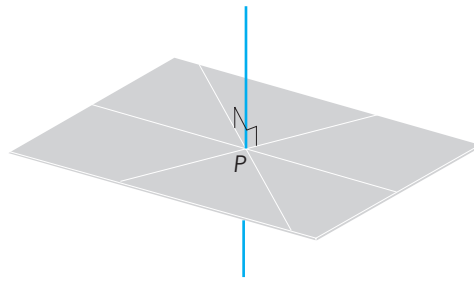


Fig. 2. Normal to a plane.

or normal to the plane if it is perpendicular to every line through P that lies in the plane.

The plane is then called perpendicular to the line. A line perpendicular to each of two intersecting lines is perpendicular to their plane and to every line in the plane. The distance from a point to a line or from a point to a plane is measured along the unique perpendicular from the point to the line or plane. A line perpendicular to a segment at its midpoint is called a perpendicular bisector of the segment. A perpendicular drawn from a vertex of a triangle to the opposite side is called an altitude of a triangle. The three altitudes of a triangle meet in a point H which is called its orthocenter. The three perpendicular bisectors of the sides of a triangle meet in a point P , called its circumcenter, equally distant from all three vertices.

Congruence. Two geometric figures are congruent to each other if the points in the one figure can be made to correspond in a one-to-one manner with the points of the other figure so that the distance between any two points in the one figure is equal to the distance between corresponding points in the other figure. (This is a posteulidean version of the definition.)

Postulates. Five postulates followed Euclid's list of definitions and may be paraphrased as:

1. A unique straight line can be drawn from any point to any other point.
2. A finite straight line can be extended continuously in either direction in a straight line.
3. A circle can be described with any given center and radius.
4. All right angles are equal.
5. If a straight line falling on two straight lines makes interior angles on the same side with a sum less than two right angles, the two straight lines, if produced indefinitely, meet on that side on which the angle sum is less than two right angles.

This fifth postulate is known as the parallel postulate and is essentially equivalent to the statement that "a unique line parallel to a given line can be constructed through any point not on the line." For 2000 years many unsuccessful attempts were made to prove that this postulate was a consequence of Euclid's other postulates, definitions, and axioms. Only as recently as the nineteenth century did N. I. Lobachevski (1793–1856), J. Bolyai (1802–1860), and G. F. B. Riemann (1826–1866) show that the parallel postulate was independent of the others,

by constructing so-called noneuclidean geometries in which the fifth postulate is not valid. See NONEUCLIDEAN GEOMETRY.

Axioms. Certain axioms were stated by Euclid and treated as self-evident truths.

1. Things which are equal to the same thing are equal to each other.

2. If equals are added to equals, the wholes are equal.

3. If equals are subtracted from equals, the remainders are equal.

4. Things which coincide with one another are equal to one another.

A fifth axiom, ascribed by Proclus to Euclid, is believed by Paul Tannery and T. L. Heath to have been added by later writers.

5. The whole is greater than the part.

This has also been replaced by the statement: The whole is equal to the sum of its parts.

Euclid's books of elements. Euclid's great work of systematic reasoning was divided into 13 books. Two other books, at one time ascribed to Euclid, are now believed to have been written by others.

The first six books of Euclid include most of the subject matter commonly taught in high school geometry; books 7, 8, 9, and 10 are concerned with number theory, square and cube roots, and incommensurable magnitudes; books 11 and 12 are concerned with solid geometry and mensuration; and book 13 is concerned with extreme and mean ratio and the five regular solids.

Book 1, after listing a large number of definitions and the postulates and axioms described above, contains 48 propositions, including many properties relating to perpendicular and parallel lines, angles, and congruent triangles, and the book culminates with the important pythagorean theorem.

Book 2 contains 14 propositions on geometrical algebra, establishing relationships between the areas of certain related squares and rectangles. For example, proposition 7 is equivalent to the identity

$$(a + b)^2 = a^2 + b^2 + 2ab$$

Book 3 deals with circles, chords, inscribed angles, and other figures related to the circle; book 4 with inscribed and circumscribed circles and polygons; book 5 with ratios; and book 6 with areas and similar triangles.

Book 7 is number-theoretical, not geometrical, and deals with units, primes, and squares and with the euclidean algorithm for least common divisor. Book 8 includes continued proportions and cube roots. Book 9 includes the factoring of numbers in continued proportion, a theorem on the infinitude of prime numbers, and a formula for even perfect numbers.

Book 10 includes a thorough treatment of quadratic irrationalities and related incommensurable magnitudes.

Book 11 introduces the fundamentals of three-dimensional geometry, including theorems on planes and lines, perpendiculars and parallels, solid angles,

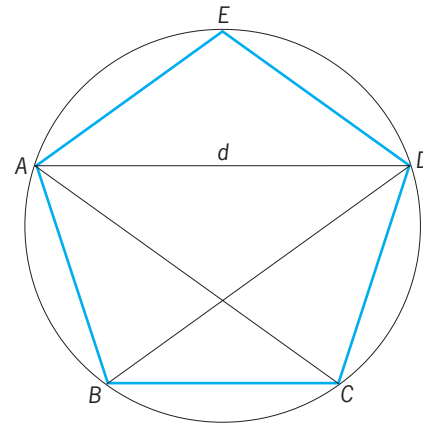


Fig. 3. Regular pentagon $ABCDE$ showing the diagonal d .

and the volumes of parallelepipedal solids. Book 12 deals with areas and volumes of pyramids and cones and spheres, deriving formulas for the latter by the method of exhaustion due to Eudoxus. Book 13 starts with the division of a segment in extreme and mean ratio, in such a manner that one part of the segment is to the second part as the second part is to the whole. This ratio is the so-called golden ratio of the sides of a regular pentagon to its diagonal d (Fig. 3). The ratio is equal to $(\sqrt{5} - 1)/2$, its reciprocal is $(\sqrt{5} + 1)/2$, and it plays a fundamental role in the study of two of the regular solids. The latter half of book 13 is concerned with a study of relations among the volumes and edges and diameters of the five regular solids, often called platonic solids. See DIFFERENTIAL GEOMETRY; POLYHEDRON; PROJECTIVE GEOMETRY; SOLID (GEOMETRY). J. Sutherland Frame

Bibliography. H. L. Baldwin, *Essentials of Geometry*, 1993; Euclid, *Elements*, ed. by T. L. Heath, 1926; R. D. Gustafson and P. D. Frisk, *Elementary Geometry*, 3d ed., 1991; P. G. O'Daffer and S. R. Clemens, *Geometry: An Investigative Approach*, 1992.

Eucoccida

The most important order of the protozoan subclass Coccidia, in the class Telosporae, subphylum Sporozoa. In this group there is alternation of asexual and sexual phases of the life cycle, and the parasitic stages occur within the cells of vertebrates or invertebrates. There are three suborders.

In the suborder Adeleina, the macrogamete and microgametocyte are joined during development (syzygy) and the microgametocytes produce only a few microgametes. Members of this suborder include parasites of invertebrates and lower vertebrates; a few even occur in higher vertebrates. In most of them, asexual development takes place in one host and sexual development in another, but the patterns of development and the hosts are too numerous to be described here.

In the suborder Eimeriina, there is no syzygy and the microgametocytes produce a large number of microgametes. This group contains several hundred species of coccidia, most of which occur in the

intestinal cells of vertebrates. They multiply asexually in their host cells, and some species cause coccidiosis, a disease that produces diarrhea, dysentery, and even death. They are especially important in domestic animals.

In the suborder Haemosporina, there is no syzygy and the microgametocytes produce a moderate number of microgametes. Asexual development takes place in the blood cells of vertebrates, and sexual development in blood-sucking insects or mites. This suborder contains *Plasmodium*, species of which cause malaria in humans, lower primates, birds, and other vertebrates. Other genera in this suborder are *Hepatocystis* (which is found especially in African monkeys), *Haemoproteus* (which occurs in birds; wild birds are often infected, but the parasites do not appear to be very pathogenic), and *Leucocytozoon* (which occurs in both wild and domestic birds and may cause a fatal disease in them). See COCCIDIA; MALARIA; PROTOZOA; SPOROZOA; TELOSPORA. Norman D. Levine

Eucommiales

An order of flowering plants, division Magnoliophyta (Angiospermae), in the subclass Hamamelidae of the class Magnoliopsida (dicotyledons). The order consists of a single family, genus, and species, *Eucommia ulmoides*, of China. It is a simple-leaved tree with reduced, solitary flowers that lack a perianth, and leaves without stipules, the paired basal appendages found on many leaves. The ovary has two unitegmic (with a single integument) ovules and ripens into a one-seeded, winged fruit. It is widely cultivated as an ornamental. See HAMAMELIDAE; MAGNOLIOPSIDA; PLANT KINGDOM. Arthur Cronquist

Eudicotyledons

One of the two major types of flowering plants (angiosperms), characterized by possession of three apertures in their pollen; the other major type is magnoliids. Although this difference in pollen development and form has been known for a long time, it has become clear, as a result of several studies of deoxyribonucleic acid (DNA) sequences, that this difference is very significant. The high degree of coincidence of the genetic data with this pollen distinction means that it is more important to recognize this distinction than the number of seed leaves, as previously thought. See DICOTYLEDONS; FLOWER; MAGNOLIALES; MONOCOTYLEDONS.

Mark W. Chase; Michael F. Fay

Euechinoidea

A subclass of Echinoidea (sea urchins, sand dollars, and heart urchins) in which the test plates are consistently arranged in double rows or columns.

Five double rows bear the tube feet (the ambulacral rows), and these alternate with five which do not (the interambulacral rows). Eighteen orders are included within the Euechinoidea. See ECHINACEA; ECHINODERMATA; ECHINOIDEA; NEOGNATHOSTOMATA. Andrew C. Campbell

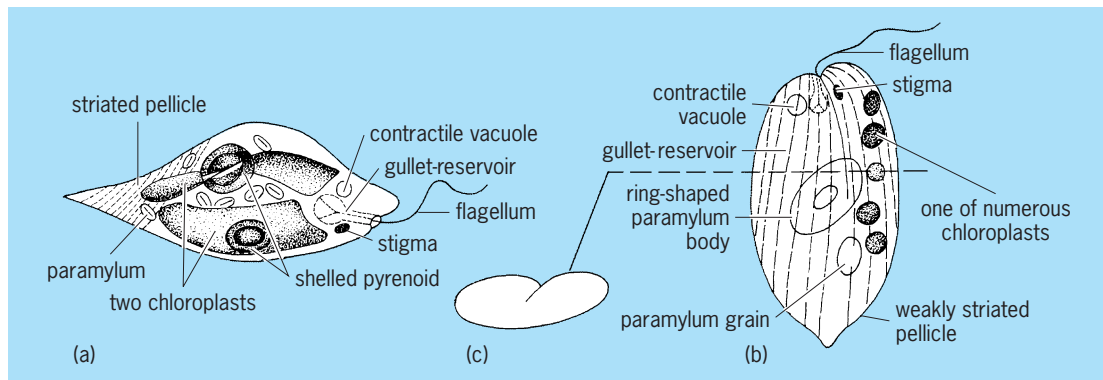
Euglenida

An order of the protozoan class Phytamastigophorea. This order of protozoa is also known as Euglenoidina, and botanists classify the group in the class Euglenophyceae. The euglenids include the largest green noncolonial flagellates, *Euglena ebrenbergii*, which are 400 micrometers long. Many of the colorless members are also large. They have one or two equal or subequal flagella. There are relatively few genera of Euglenida, even if doubtful ones such as *Clorachne* and *Otonia* are included. The freshwater genera *Euglena*, *Lepocinclis*, *Phacus*, and *Trachelomonas* have many species, as does the colorless genus *Petalomonas*. Others have few species, including the marine genus *Eutreptia*. *Calkinsia* is yellow, some euglenas are red at times, and all others are green or colorless, although *Trachelomonas* tests are frequently colored by iron. See EUGLENOPHYCEAE; PHYTAMASTIGOPHOREA.

Morphology. Euglenids generally show an anteroposterior elongation; however, *Trachelomonas volvocina* is practically spherical, as shown by young individuals, and *Euglena acus* is extremely needle-like. Most free-swimming members are round in cross section, but species of *Phacus* are flattened, as are *Petalomonas*, *Anisonema*, *Triangulomonas*, and others (see **illus.**). Many are enclosed in a thick pellicle sculptured into keels (*Phacus*, *Pleotia*, and *Tropidoscyphus*); others, such as *Euglena oxyuris* and *Phacus pyrum*, show strong torsion of the pellicle. This pellicle is sometimes ornamented with striae or small wartlike protuberances. It prevents distortion of the body according to its degree of thickening, but species which have a very thin pellicle (*Astasia*, *Peranema*, *Heteronema*, and some euglenas) shorten on their axis, then elongate. Such contortions are termed euglenoid or metabolic movements. There are sometimes fine openings in the pellicle through which mucous threads are extruded under certain conditions.

Colony formation is restricted to *Colacium*, in which a flagellated euglenoid cell usually attaches to a copepod by its anterior end and develops into an arboroid structure of several nonflagellated cells. Tests or shells are common; *Klebsiella* is a free-swimming marine form in an urn, *Ascoglena* a freshwater epiphyte in an urn. In *Trachelomonas* and *Strombomonas* the tests vary widely in shape and ornamentation and have a narrow pore through which the flagellum emerges.

Distinguishing features. Euglenids have two distinguishing features. They synthesize a starchlike polysaccharide, paramylum, of constant shape for, and often characteristic of, a given species. This



Euglenids. (a) *Euglena pisciformis*. (b) *Phacus strokesii*. (c) Cross section of cell.

assimilation product of both green and colorless forms is often large and lamellated.

The second characteristic is the gullet-reservoir system. An anterior mouth opens into a short gullet which widens into a reservoir. The flagella emerge from its floor, and at one side a contractile vacuole pulses, except in marine species. This system is present in some Chloromonadida (*Gonyostomum semen*) and is approximated by certain Cryptomonadida, but all euglenoids possess it. The mouth may serve to ingest food (*Peranema trichophorum*), but in most cases does not. Green forms possess a stigma near the reservoir, and a thickening on the longer of the two flagella is found in most of them, as well as in some colorless species.

Flagella and chromatophores. In most green species the flagellum is long, uniform, and movable throughout. *Euglena mutabilis* has an extremely short one, as does *E. spirogyra*. According to G. F. Leedale, there is no bifurcation of the locomotor flagellum; instead, a second short flagellum is found separate, but it does not extend beyond the mouth. *Calkinsia* has a long, uniform, anterior one, movable throughout, and a shorter, fine, trailing one. *Peranema* has a long tapering flagellum, vibratile at its tip during steady progression, and a very short, fine, posterior one. *Anisonema acinus* has a rather short, fine, anterior flagellum and a tapering, long, and very heavy posterior one. Chromatophores in the green members may be circular, stellate, or irregular disks, or variously shaped plates and short bands. The chromatophores are reticulate plastids containing bright green pigments closely related to, if not identical with, the chlorophylls of higher plants. The chromatophores of Euglenida never contain other masking pigments, though they are often associated with pyrenoids. In *Euglena pisciformis* both chromatophores have a central pyrenoid which has a shell of paramylum. Pyrenoids stain deeply with hematoxylin and may be proteinaceous. Besides paramylum, fats occur, and glycogen has been reported.

Nuclei. Nuclei are large, vesicular, and contain one or more endosomes which are deeply staining chromatoid bodies. These are surrounded by dispersed chromatin granules. In division, a specific number

of chromosomes is formed, but a spindle is lacking, and chromosomes are oriented around and parallel to the drawn-out endosomal mass. Cell division usually occurs in the active state, although *Euglena gracilis* divides while encysted. Palmella stages have been reported but are not common. Sexual processes are unknown. Daily isolations of several clones of *Entosiphon sulcatum* for many months showed no depression of the division rate, and the usual type of mitotic division was enough to maintain species vitality.

Nutrition. Nutrition varies greatly in Euglenida. Green members are holophytic or autotrophic, but there is evidence that some of them absorb soluble organic matter with no evidence of ingestion of solid matter. It seems questionable that more than a few green members are strict autotrophs which do not require organic carbon and nitrogen but form carbohydrates and proteins from carbon dioxide and inorganic salts. Certainly the colorless forms are heterotrophs (dependent on organic food), and some of the green ones readily assume a heterotrophic existence, if they ever were autotrophic. Thus *Euglena gracilis* and species of *Phacus* become colorless in an organic medium in the dark, whereas *Trachelomonas reticulata* is devoid of chlorophyll. *Peranema* ingests large food bodies; however, most of the colorless forms are saprozoic rather than holozoic and are readily cultivated in soluble organic media. They also vary widely in oxygen requirements; the same species occur in well-oxygenated or anaerobic situations. In a rich organic medium and under a film of oil, *Entosiphon* thrives, but often forms somatellalike masses, in which, however, cell division does not become complete. Parasitism is unknown in the class, although *Euglenomorpha* occurs in the intestine of tadpoles, where it is probably commensal.

Ecology. In habitat the Euglenida are widespread. Green members occur mostly in fresh water and frequently in such numbers as to form blooms. This is especially true for pools contaminated by cattle droppings. *Strombomonas* in the Ohio Valley of the United States is a potamophile form, *Trachelomonas* a lake or pond form. *Eutreptia* blooms densely in marine bays and estuaries, and the mud-water

interface in marine waters is a rich source of the colorless creeping members. Iron seepages and coalmine seepages containing sulfuric acid are often bright green with *Euglena mutabilis* even at pH values of 1.5, whereas cedar swamps have blooms of *Euglena polymorpha*. Unlike the Phytomonadida, the Euglenida appear in warmer waters but do not seem to occur in hot springs. They abound in citrus-waste lagoons and in later stages of sewage-treatment waters, an indication of their partially saprozoic nutrition.

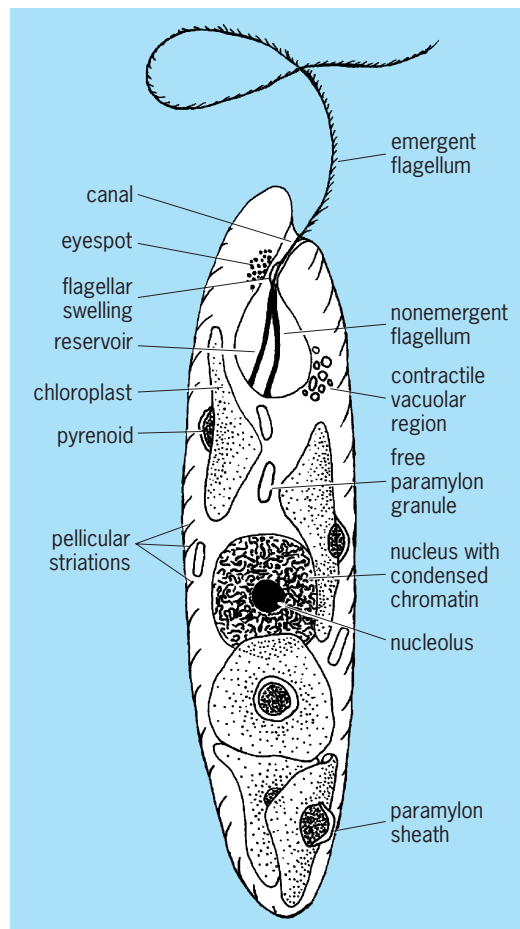
Relationships within the class are not well defined at present; each classification within recent years has departed from previous ones. Organisms are readily assigned to the Euglenida on a structural basis, because the class is quite specialized. However, intraclass relationships are difficult to assign. See CILIA AND FLAGELLA; EUGLENOPHYCEAE; MASTIGOPHORA; PROTOZOA. James B. Lackey

Euglenophyceae

A class coextensive with the division Euglenophycota, comprising unicellular colorless or photosynthetic flagellates with very distinctive cytological characters. In protozoological classification, these organisms constitute an order, Euglenida, of the class Phytomastigophora. Although photosynthetic euglenoids are like Chlorophycota in containing chlorophyll *a* and *b*, in most other respects they are so different from those algae as to suggest an independent phylogenetic origin of their pigments. About 1000 species have been described and classified into about 40 genera and 6 orders. See EUGLENIDA.

Most euglenoids are free-swimming and have two flagella, one of which may be nonemergent, arising from an anterior invagination known as a reservoir (see *illus.*). *Euglenamorphba* and *Hegneria*, both inhabitants of the rectum of tadpoles, have three to seven flagella. A few rhizopodial forms lack flagella. In *Colacium*, mucilaginous stalks join *Euglena*-like cells in dendroid colonies attached to small aquatic animals. In some species of *Euglena* and *Eutreptia*, cells cohere in palmelloid sheets. *Ascoglena*, although flagellate, is sessile.

Structure. Euglenoid cells have helical symmetry, usually with superimposed bilateral symmetry, and are 15–500 micrometers long. *Euglena gracilis*, the best-known species, is spindle-shaped, while other species are spherical, flattened, or twisted. All euglenoids have a characteristic proteinaceous pellicle composed of flexible or rigid, interlocking, spirally disposed strips lying beneath the cell membrane. In *Trachelomonas*, *Strombomonas*, and *Ascoglena*, an inorganic envelope (lorica) is deposited outside the cell membrane. In species with flexible pellicles, flagellar locomotion may be supplemented or supplanted by so-called euglenoid movement, a peristalsis involving the entire cell. The surface of all euglenoids is coated to a varying extent with mucilage secreted by muciferous bodies (mucocysts) that traverse the pellicle.



Euglena sp., showing overall morphology and intracellular structures.

Emergent flagella, which are covered with fine hairs, contain a paraflagellar rod in addition to the usual 9 + 2 array of axoneme microtubules. At the base of one flagellum is a swelling believed to be a photoreceptor that acts in concert with the eyespot in phototaxis. The eyespot, which is situated in the cytoplasm at one side of the reservoir, consists of lipid globules not surrounded by membranes. With the possible exception of a few marine and parasitic species, all euglenoids have a contractile vacuolar region adjoining the reservoir. A single nucleus is present in each cell. As in dinoflagellates, the chromatin is permanently condensed so that chromosomes are visible throughout the cell cycle. The nuclear membrane remains intact during mitosis. See DINOPHYCEAE.

Chloroplasts. Photosynthetic euglenoids contain one to many grass-green chloroplasts, which vary from minute disks to expanded plates or ribbons. Each chloroplast is surrounded by three membranes, the outer of which is discrete and not associated with endoplasmic reticulum. Photosynthetic lamellae are composed of three appressed thylakoids. They contain β -carotene and various xanthophylls in addition to chlorophyll *a* and *b*. An extraplastidial carotenoid causes a few species of *Euglena* to be red. Pyrenoids are present in many species. The reserve product of

photosynthesis is paramylon, a β -1,3 glucan stored outside the chloroplast in the form of rods or disks. Paramylon is also stored in colorless species. The chloroplast of some species can be experimentally eliminated or rendered nonfunctional by heat, ultraviolet light, or antibiotics.

Nutrition. Colorless euglenoids depend on osmotrophy or phagotrophy for nutrient assimilation. Some phagotrophs, such as *Peranema*, have a cytostome, an opening in the anterior of the cell through which particulate food is ingested. The cytostome, which is sometimes accompanied by other specialized ingestion organelles, is independent of the reservoir. Some photosynthetic euglenoids are facultative osmotrophs, but none are facultative phagotrophs.

Reproduction, use, and habitat. Euglenoids reproduce solely by longitudinal binary fission, sexuality being unknown. Thick-walled resting stages (cysts) are common in many genera. *Euglena gracilis*, which can be grown easily in bacteria-free culture, has been used extensively in the study of photosynthesis, nucleic acid synthesis, membrane generation, and other aspects of physiology and molecular biology. Its requirement for vitamin B₁₂, the antipernicious anemia factor in animals, makes it useful in medical bioassays.

Euglenoids are found most commonly in fresh water rich in organic matter, but they also occur in marine or brackish habitats, on mud or sand, and in ice or snow. A few species prefer very acidic water. There are three genera of endoparasites, inhabiting fresh-water flatworms, annelids, copepods, and tadpoles.

Paul C. Silva; Richard L. Moe

Bibliography. D. E. Buetow (ed.), *The Biology of Euglena*, vols. 1 and 2, 1968, vol. 3, 1982; G. F. Feedale, *Euglenoid Flagellates*, 1967; M. Gojdic, *The Genus Euglena*, 1953; S. P. Parker (ed.), *Synopsis and Classification of Living Organisms*, 2 vols., 1982.

Eukaryotae

The vast array of living and fossil organisms with complex cells in which the genetic material is organized into chromosomes (linear pieces of DNA bound with histone proteins and visible as compact structures during mitosis and meiosis) and contained within a membrane-bound nucleus. The Eukaryotae have a common origin from organisms with prokaryotic cells (bacteria and archaea) that do not have their DNA in chromosomes and lack membrane-bound organelles. The Eukaryotae thus include all multicellular plants, fungi, and animals, as well as a collection of the unicellular protists. See CELL NUCLEUS; CHROMOSOME.

Characteristics. All organisms in this group possess a eukaryote cell that is structurally more complex than the prokaryote cell of archaea and bacteria. In addition to a nucleus, the cell contains a number of other membrane-bound cytoplasmic organelles, such as mitochondria, Golgi bodies, lyso-

somes, peroxisomes, endoplasmic reticulum, and, in photosynthetic forms, plastids or chloroplasts. Rather than being locked within a rigid cell wall like prokaryotes, eukaryotes have evolved a complex cytoskeleton of microfilaments and microtubules that give structural support to the cell and allow the cell membrane to grow and change shape. Characteristically, centrioles, cilia, or flagella are also present, the latter locomotory organelles composed mainly of tubulin with microtubules generally in a 9 + 2 pattern. Therefore, there is a major discontinuity, at the cellular level alone, between the Prokaryotae and the Eukaryotae. See CELL (BIOLOGY); PROKARYOTIC CELL.

Biology. Organization at the organismal level ranges from solitary unicellular (as is common in many of the algal and protist groups) to colonies of cells, mycelial, syncytial (coenocytic), and truly multicellular with extensive tissue differentiation. Modes of nutrition run the gamut: absorptive, ingestive, photoautotrophic, plus combinations of these three. Life cycles vary tremendously and may include both asexual and sexual methods of reproduction. Asexual reproduction can be by cell division for unicellular forms, or by vegetative growth or budding in multicellular forms. Exchange of genetic material (sexual reproduction) takes place in the majority of eukaryotes by fusing haploid gametes such as egg and sperm or pollen to form a diploid embryo, but variations such as conjugation (where nuclei are exchanged) may also occur.

Aerobic metabolism is commonly exhibited, especially by aquatic and terrestrial forms; anaerobic mechanisms exist, however, for numerous species found in poorly oxygenated habitats, including various sites within bodies of host organisms.

Diversity. Size of species ranges considerably (from unicellular protists to large multicellular plants and animals). The size of the eukaryotic cell is usually 20–50 micrometers in diameter, but mammalian red blood cells are just 8 μ m in diameter while some single-celled protists, such as the ciliate *Spirostomum*, are over 1 mm (1000 μ m) in length. Habitats cover all possible ecological niches: aquatic, terrestrial, and aerial, for free-living forms; and in or on all kinds of hosts, for symbiotic or parasitic forms (internal habitats, including cells, tissues, organs, or various body cavities). Dormant stages include cysts, spores, and seeds; these are often involved in dispersion or propagation of the species.

Some species that cause fatal diseases in other eukaryotic organisms (including humans) are as important economically and medically as some of the parasitic members of the superkingdom Prokaryotae. Parasitic life cycles may involve multiple hosts, specialized vectors, and so on. See ANIMAL KINGDOM; FUNGI; PLANT KINGDOM; PROTOZOA.

Diana L. Lipscomb; John O. Corliss

Bibliography. L. Margulis and K. Schwartz, *Five Kingdoms: An Illustrated Guide to the Phyla of Life on Earth*, 3d ed., 1998; S. P. Parker (ed.), *Synopsis and Classification of Living Organisms*, 2 vols., 1982.

Eulamellibranchia

The largest subclass of the class Bivalvia in the phylum Mollusca. Eulamellibranchs are ciliary filter feeders with greatly enlarged ctenidia, in each of which the elongated gill filaments are held together in parallel series to form folded lamellae. The 30,000 living species constitute the most successful of molluscan groups, in terms both of individual numbers with their biomass and of ecological bioenergetics (calorie transfer in food chains).

Classification. There are several proposed systems of ordinal classification for the subclass Eulamellibranchia, but little agreement on the best structural bases. One system frequently employed comprises six orders: Anisomyaria, Taxodonta, Heterodonta, Schizodonta, Anomalodesmata, and Adapedonta. There is more general acceptance of the (possibly 18) superfamilies of eulamellibranchs, each with characteristic shell structure, hinge dentition, ligament, and degree of fusion of the mantle lobes involving the origin of the siphons. The more important superfamilies include Mytilacea (mussels), Ostreacea (true oysters), Cardiacea (true cockles), Unionacea and Corbiculacea (fresh-water clams), Myacea (including softshell clams), Mactracea (surf clams), Tellinacea, and Pholadacea (=Adesmacea). Within each larger superfamily there is a wide range of adaptive radiation, with different genera living as shallow or deep burrowers or as sessile attached forms. There can be considerable evolutionary convergence shown by genera from different superfamilies, and deep burrowers with massive fused siphons are found in five superfamilies. See BORING BIVALVES.

Filter feeding and digestion. The four folded lamellae of the paired gills are morphologically homologous (in blood vessel and ciliary arrangements) with the archetypic molluscan ctenidia and their functional orientation is also similar. The lateral cilia on the filaments produce a nearly continuous water current between adjacent filaments which flows from the outer (ventral) inhalant part of the mantle cavity to the inner (dorsal) exhalant part. All suspended material, including food organisms brought into the outer mantle cavity by the water current, is filtered and accumulates on the inhalant faces of the gill lamellae. Some of the filtration is accomplished by the frontal cilia, but most is due to the laterofrontals which are compound cilia with a finely pinnate structure. Accumulated material and food is then moved by the frontal cilia toward the ventral edges of the gills, and the fine food particles become concentrated in the food grooves. Coarser and heavier materials (such as sand grains) do not enter the food grooves but are rubbed or twitched off onto the mantle wall. Sorting is carried out in food grooves and on the labial palps on a size basis; palp cilia carry fine material into the mouth, from which ciliary tracks carry it slowly and continuously into the esophagus and to the stomach, where it undergoes further sorting.

In the stomach, extracellular digestion begins. In all eulamellibranchs, the peculiar secreted structure

termed the crystalline style is used as a cilia-driven gastric stirring rod and enzyme store, allowing a slow continuous release of amylase and glycogen-breaking enzymes. Food particles then pass into the tubules, the cells of which are phagocytic so that final intracellular digestion takes place. Indigestible material rejected from the stomach complex passes to the posterior intestine and rectum, where water reabsorption and a neutral pH consolidate the true feces into very hard pellets. This process avoids self-contamination of the feeding mechanism.

The coarser particulate material filtered off by the feeding structures but not entering the gut is termed pseudofeces. It is collected by the cleansing cilia of the inside of the mantle wall into ciliary vortices. The pseudofeces are expelled from the bivalve by spasmodic contractions of the adductor muscles which force water and accumulated pseudofeces out through the normally inhalant openings to the mantle cavity. The anus and the renal and genital openings are in the exhalant part of the mantle cavity (as in all mollusks); therefore, expulsion of wastes or reproductive products is not accomplished by these spasmodic cleansing movements but by the normal and continuous water flow of the feeding and respiratory current.

Feeding life-styles. A wide variety of eulamellibranch life-styles and corresponding shell forms can be based on this successful (and relatively uniform) filter-feeding mechanism. The majority of eulamellibranchs are suspension feeders filtering living phytoplankton and bacteria from the overlying water. A substantial minority are detritus feeders that siphon in and filter organic materials (living and dead) from the surface of the sea-bottom deposits. Different styles of suspension feeders include relatively active shallow burrowers with globular shells and short siphons such as cockles and quahogs, sedentary deep burrowers (and some borers) with long fused siphons such as softshell clams, attached forms such as mussels, and permanently cemented forms such as oysters. Detritus feeders (principally superfamily Tellinacea) are all active deep burrowers with long extensible and separate siphons, the inhalant siphon actively sweeping the sea bottom, and with particularly intensive sorting on gill lamellae and palps.

Ecology. The marine eulamellibranchs as suspension feeders consume a major part of the largest plant crop in the world, the marine phytoplankton. Certain eulamellibranch species are the most abundant marine benthic animals; beds of mussels and of oysters are found in satellite images of coasts and shallow seas, and less obvious burrowing clams can be as abundant. Computations have suggested that, in size of standing crop biomass for individual species (annual means), certain eulamellibranchs outrank all other animal species on the planet. In turn, eulamellibranchs are a major food of all bottom-living fishes (including cod, halibut, plaice, and haddock), and are second only to planktonic copepods (such as *Calanus*) in annual turnover of calories through animal tissues for all marine food webs. Birds, marine

mammals, and humans also feed directly on eulamellibranchs. See FOOD WEB; MARINE ECOLOGY.

Economic significance. A few eulamellibranch species have negative economic significance for humans. Teredinids or shipworms (of superfamily Pholadacea) bore destructively in marine wood structures such as dock pilings and boats. Certain small eulamellibranchs (*Corbicula* spp. and the zebra mussel, *Dreissena*) are recent introductions to fresh waters in North America and have quickly achieved pest status, fouling water intakes and clogging filters at power plants and municipal water supplies. See BIVALVIA; MOLLUSCA; SHIPWORM.

W. D. Russell-Hunter Bibliography. S. P. Parker (ed.), *Synopsis and Classification of Living Organisms*, 2 vols., 1982; W. D. Russell-Hunter, *A Life of Invertebrates*, 1979; K. M. Wilbur (ed.), *The Mollusca*, 11 vols., 1982-1989; C. M. Yonge and T. E. Thompson, *Living Marine Molluscs*, 1976.

Euler angles

Three angular parameters that specify the orientation of a body with respect to reference axes. They are used for describing rotating systems such as gyroscopes, tops, molecules, and nonspherical nuclei. These parameters are not symmetrical in the three angles but are simpler to use than other rotational parameters.

Unfortunately, different definitions of Euler's angles are used, and therefore it is confusing to compare equations in different references. The definition given here is the majority convention according to H. Margenau and G. Murphy.

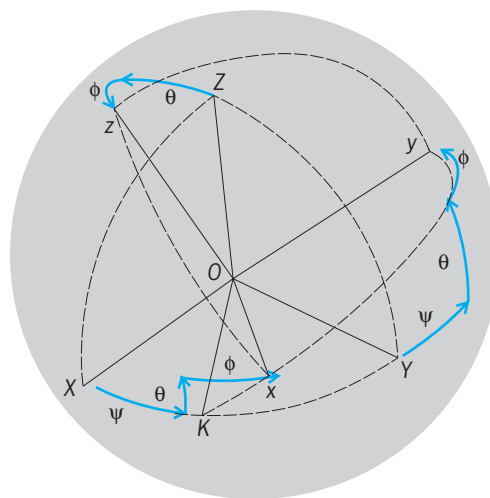
Let OXYZ be a right-handed cartesian (right-angled) set of fixed coordinate axes and Oxyz a set attached to the rotating body.

The orientation of Oxyz can be produced by three successive rotations about the fixed axes starting with Oxyz parallel to OXYZ. Rotate through (1) the angle ψ counterclockwise about OZ, (2) the angle θ counterclockwise about OX, and (3) the angle ϕ counterclockwise about OZ again. The line of intersection OK of the xy and XY planes is called the line of nodes.

Denote a rotation about OZ, for example, by Z (angle). Then the complete rotation is, symbolically, given by Eq. (1) where the rightmost operation is done first.

$$R(\psi, \theta, \phi) = Z(\phi)X(\theta)Z(\psi) \quad (1)$$

A point P will have coordinates (x, y, z) with respect to the body axes and (X, Y, Z) with respect



Euler angles. The successive movements of the axes on a unit sphere described in the text are shown by arrows. The complete rotation may also be obtained by a different sequence of rotations, namely, first through ϕ about OZ, then through θ about the displaced x axis (which is OK), then through ψ about OZ.

to the fixed axes. These are related by linear Eqs. (2)-(4), where (x, X) is the angle between the

$$x = X \cos(x, X) + Y \cos(x, Y) + Z \cos(x, Z) \quad (2)$$

$$y = X \cos(y, X) + Y \cos(y, Y) + Z \cos(y, Z) \quad (3)$$

$$z = X \cos(z, X) + Y \cos(z, Y) + Z \cos(z, Z) \quad (4)$$

axes Ox and OX, and so forth. The nine direction cosines are expressed in terms of the three Euler angles in the table.

Inspection of the illustration makes it apparent that no operation in $R(\psi, \theta, \phi)$ can be replaced by a combination of the other two. Therefore, three parameters are needed to specify the orientation, and the amounts of the angles are unique (barring additional 360° rotations). In dynamical problems of rotating bodies, ψ , θ , and ϕ can be used as independent angular coordinates.

Molecules and nuclei undergo oscillatory changes in shape while rotating. The body axes apply to the average shape. For another set of rotational parameters see CAYLEY-KLEIN PARAMETERS.

Bernard Goodman

Bibliography. A. P. Arya, *Introduction to Classical Mechanics*, 2d ed., 1997; R. Baierlein, *Newtonian Dynamics*, 1983; A. L. Fetter and J. D. Walecka, *Theoretical Mechanics of Particles and Continua*, 1980, reprint, 2003; H. Goldstein, C. P. Poole, and J. L. Safko, *Classical Mechanics*, 3d ed., 2002.

Direction cosines stated in terms of Euler angles			
	X	Y	Z
x	$\cos \psi \cos \phi - \sin \psi \sin \phi \cos \theta$	$\cos \psi \sin \phi + \sin \psi \cos \phi \cos \theta$	$\sin \psi \sin \theta$
y	$-\sin \psi \cos \phi - \cos \psi \sin \phi \cos \theta$	$-\sin \psi \sin \phi + \cos \psi \cos \phi \cos \theta$	$\cos \psi \sin \theta$
z	$\sin \psi \sin \theta$	$-\cos \phi \sin \theta$	$\cos \theta$

Euler's equations of motion

A set of three differential equations expressing relations between the force moments, angular velocities, and angular accelerations of a rotating rigid body. They are equations of motion in the usual dynamical sense, having the form of Eqs. (1)–(3). The formula-

$$I_1(d\omega_1/dt) + (I_3 - I_2)\omega_2\omega_3 = M_1 \quad (1)$$

$$I_2(d\omega_2/dt) + (I_1 - I_2)\omega_3\omega_1 = M_2 \quad (2)$$

$$I_3(d\omega_3/dt) + (I_2 - I_1)\omega_1\omega_2 = M_3 \quad (3)$$

tion employs as coordinate axes the three principal axes of rotational inertia of the body that can rotate about a body-fixed point, which is the center of mass if constraints are absent. These reference axes, which form a right-hand set, are indicated by subscripts 1, 2, and 3 in the equations, where I_1 , I_2 , and I_3 represent the principal moments of inertia; ω_1 , ω_2 , and ω_3 are the angular velocities about the axes; M_1 , M_2 , and M_3 are the corresponding force moments; and t is the time.

In the general case, these equations cannot be integrated, but solution is possible in special cases. Soluble problems of interest include that in which force moments are absent, the resulting complex behavior being called Poincaré motion, and that in which two of the principal moments of inertia are identical and only one force moment is present. The latter case includes spinning tops and gyroscopes. See RIGID-BODY DYNAMICS.

Russell A. Fisher

Bibliography. R. Baierlein, *Newtonian Dynamics*, 1983; H. Goldstein, C. P. Poole, and J. L. Safko, *Classical Mechanics*, 3d ed., 2002; T. W. B. Kibble and F. H. Berkshire, *Classical Mechanics*, 5th ed., 2004.

Eumagnoliids

An informal set of flowering plant orders that contains Laurales, Magnoliales, Piperales, and Winterales, plus the monocotyledons, and in total contains roughly one-third of angiosperm species. Formerly, many have equated this group (minus the monocotyledons) to the subclass Magnoliidae, but studies of sequences of deoxyribonucleic acid (DNA) have demonstrated that orders such as Nymphaeales and Ranunculales are distantly related to the others and that the monocotyledons are much closer. A formal taxonomy for this group is not yet feasible because until more data are collected the exact relationships of the eumagnoliids to the eudicots are not clear, and before formal names are proposed it is preferable to understand these taxa better. The eumagnoliids, including the monocotyledons, exhibit a set of characteristics that reflect their relatively close relationships; these include a tendency to trimerous flowers and a suite of chemical compounds restricted to the group. They also have a set of traits that are "primitive" by most estimations; these include vesselless wood, laminar stamens, and laminar placentation. It now seems quite clear that although

the eumagnoliids exhibit some primitive traits they are neither the most archaic nor the earliest of angiosperms, and they contain some of the most advanced taxa, such as many of the monocotyledons (grasses, palms, and orchids). See LAURALES; LILIOPSIDA; MAGNOLIALES; PIPERALES; PLANT TAXONOMY.

Mark Chase

Eumalacostraca

The major subdivision (subclass) of the crustacean class Malacostraca. It comprises three, four, or five Recent superorders: Syncarida, Peracarida, and Eucarida are most commonly included and all that are recognized here. However, some carcinologists still insist on ranking the Hoplocarida among eumalacostracans. Acceptance of a separate superorder, Pancarida for the Thermobaenacea, has not received general support, and thermobaenaceans are commonly now included with peracaridan taxa. See EUCARIDA; PERACARIDA; SYNCARIDA.

Formerly, the concept of the Eumalacostraca as a homogeneous evolutionary unit generated virtually no controversy. The subclass was defined by a series of characters referred to as the caridoid facies: carapace enclosing the thorax; stalked, movable eyes; biramous antennules; scalelike antennal exopod; thoracopods with natatory exopods; abdomen with well-developed, complex musculature designed for flexing it; telson and uropods forming a tailfan; biramous pleopods 1–5; and internal organs primarily excluded from the abdomen. Opponents of the caridoid theory point to the fact that several of these attributes are not restricted to eumalacostracans, and conversely, not all eumalacostracans possess all of them. The carapace, for example, is absent in syncarids and the peracaridan isopods and amphipods, but present in all other eumalacostracans. Eyes are sessile or lacking in bathynellaceans, thermobaenaceans, cumaceans, isopods, and amphipods, and either stalked or sessile in anaspidaceans. Similarly, the antennal scale is lacking in most Peracarida (except the Mysidacea).

Basic structural plan. The basic body plan common to all eumalacostracans includes a head of five somites, a thorax of eight somites, and an abdomen of six somites, plus a terminal telson. Zero to three pairs of thoracic somites may be fused to the head and their thoracopods and modified to form feeding appendages (maxillipeds). The thorax usually has five to eight pairs of ambulatory or specialized appendages (pereopods); the abdomen may or may not have paired appendages (pleopods) on the first five somites; and the sixth somite usually has a pair of well-developed appendages (uropods). Female gonopores (coxal or sternal) are present at the level of the sixth thoracic somite, while male gonopores are at the level of the eighth.

Carapace. The carapace, its origin and evolution, has provided the greatest controversy. Classically, the carapace, or cephalic shield, was believed to have had its origin as a fold of integument arising

from the posterior border of the maxillary somite of the cephalon. With expansion, anteriorly and posteriorly, cephalothoracic segments became dorsally fused to it. In one alternative hypothesis, the carapace is thought to originate as branchiostegal folds that grow outward and upward from one or more thoracic somites. A second suggests that the carapace is a dorsal expansion of the protocephalic area (antennal somite when differentiated), and that the cephalic and thoracic tergites do not enter into the process at all.

Optics. Three distinct types of optical mechanisms are present in eumalacostracan eyes: apposition optics, often with hexagonal facets; refracting superposition eyes with hexagonal facets; and superposition eyes with square facets. Apposition optics are common to decapod larvae and most adult brachyuran crabs. Refracting, superposition eyes with hexagonal facets occur in mysids and euphausiids, whereas similar eyes, but with square facets using mirror optics, are found in adult shrimps, lobsters, and dromiids crabs. Apposition optics are thought to reflect the primitive condition, with a trend to sessile eyes and eventually total eye loss associated with adaptation to an interstitial life. See EYE (INVERTEBRATE).

Locomotion. Two caridoid attributes shared by the majority of eumalacostracans are the strong abdominal musculature and well-developed tailfan. Together with accompanying neural elements, these appear to be designed as a highly efficient propulsion mechanism, that is, an escape reaction. Rapid flexion of the abdomen brings the tailfan forward, providing for the sudden propulsion of the animal backward. Studies suggest that the antennal scale also may play a role in this escape reaction.

Despite major differences in habitat and behavior, other aspects of ambulatory locomotion show considerable intraordinal resemblance. These similarities are attributed to the internal continuity in pereopodal skeletomusculature that exists among the various eumalacostracan orders. In thoracic appendages not modified as maxillipeds, the entire limb lies in a plane that remains unchanged with normal locomotion, but can bend to accommodate actions such as feeding, swimming, and burrowing. Walking is accomplished by extension and flexion within the limb plane, whereas "rowing" motions occur from the limb base. The presumably primitive pattern, found in syncarids and eucarids, involves the development of the pereopodal coxa as a gimbal. Dicondylic body-coxa articulation allows promotion-remotion, and dicondylic coxa-basis articulation allows abduction-adduction. In contrast, the body-coxa articulation in peracaridans, other than mysids, is either immobilized or capable of only limited abduction-adduction, whereas the coxa-basis articulation is monocondylic and capable of performing a complete suite of motions.

In swimming, which has been hypothesized to have evolutionarily followed ambulation, both the natatory exopods of the thoracic appendages and the pleopods work in concert. However, in interstitial eumalacostracans, some peracaridans, and many

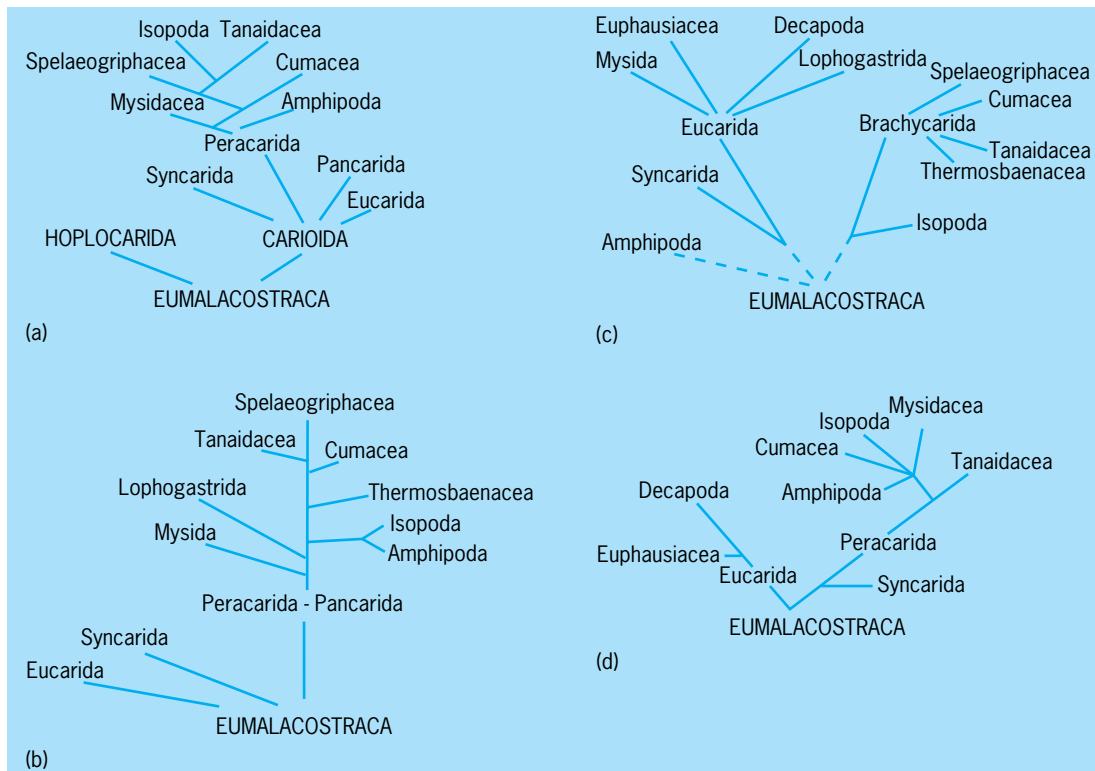
higher decapods, there is a tendency toward reduction of the pleopods. Swimmers among these groups rely almost exclusively on the thoracic exopods which are provided with long, marginal setae.

Respiration. Respiration in its simplest form is provided by the thoracic epipods. These vascularized, saclike outgrowths of the coxae form primitive gills; the blood vessels are marginal with transverse connectors. Epipodal gills are found in the Syncarida and several orders of the Peracarida. However, in isopods the pleopods have taken on the respiratory function. These gills may vary from simple branchial epithelium to tufts of branchial filaments, and in terrestrial species specialized air sacs (pseudotracheae) may develop from evaginations on the inner sides of the exopods. Gills of the Eucarida are more complex, each consisting of a central axis with numerous branches, and may also involve countercurrent circulation.

Circulation. The eumalacostracan circulatory system typically consists of an elongate dorsal blood vessel, part of which is expanded in the thoracic region to form a heart with one or more paired ostia, often with a descending (sternal) artery, and with anteriorly, posteriorly, and ventrally directed arteries. Exceptions are found in most peracaridans where a sternal artery is absent. In isopods where the heart is located in the abdomen, lateral arteries provide blood to the pereopods, pleopods, and pleotelson, and the anterior aorta supplies the head. The length of the heart is substantially reduced in euphausiids and decapods but remains positioned in the thorax.

Reproduction. Most eumalacostracans are dioecious; however, some isopods, tanaids, and a few decapods are hermaphroditic. Syncarids lay their eggs on the substrate after copulation. Peracarids brood their embryos in ventral (or in thermobae-naceans, dorsal) brood pouches. Female eucarids, for the most part, carry their eggs in egg sacs (euphausiids) or attached directly to the pleopods; however, penaeid shrimp and some euphausiids shed their eggs directly into the sea. Fertilization is by means of spermatophores and may be external or internal.

Syncarid sperm ultrastructure (based on *Anaspides*) is described as having a very elongate subacrosomal filament which is filiform and coiled and bypasses the nucleus. A nuclear membrane persists. In the Peracarida, spermatozoa consisting of two convergent linear components, the principal sperm body and a transversely striated, taillike, nonflagellate structure, are found in mysidaceans, amphipods, isopods, and cumaceans. Only in the Tanaidacea are the sperm rounded and lacking appendages, a condition similar to that of syncarids. Little is known of euphausiid spermatozoa except that they are ovoid and lack appendages. Both euphausiids and stenopodidean decapods appear to lack an acrosome. Sperm in dendrobranchiate and procarid shrimp have a single acrosomal spine, but rarely with arms corresponding to those characteristic of other decapods.



Phylogeny of the Eumalacostraca. (a) Based on the caridoid facies concept (after F. Schram, ed., *Crustacean Phylogeny, Crustacean Issues*, A. A. Balkema NE, vol. 1, 1983). (b) Based on general morphological characters (after F. R. Schram, *Crustacea*, Oxford University Press 1986). (c) Based on selected morphological characters (after F. Schram, ed., *Crustacean Phylogeny, Crustacean Issues*, A. A. Balkema NE, vol. 1, 1983). (d) Based on spermatozoa ultrastructure (after B. G. M. Jamieson, *Ultrastructure and phylogeny of crustacean spermatozoa*, *Mem. Queensland Mus.*, vol. 31, 1991).

Several spikes are characteristic of lobsters and crabs.

Systematics and phylogeny. A cursory review of eumalacostracan classifications reflects the current lack of consensus. Central to the differing views is disagreement on the origin and evolution of the carapace, and the identity of a eumalacostracan precursor. Proponents of retention of the Hoplocarida (Stomatopoda) in the Eumalacostraca have argued not only that the carapace has a maxillary origin but that the presence of the carapace in phyllocarids establishes its genesis well in advance of the advent of eumalacostracans. Additionally, certain caridoid attributes are primitive (plesiomorphic) or probably characterized the eumalacostracan ancestor. Advanced caridoid attributes (apomorphies) are, for the most part, related initially to adaptations to pelagic life and improved benthic ambulation.

Opponents point to the fact that if these advanced caridoid attributes were derived from an ancestral eumalacostracan, those eumalacostracans that lack such attributes, particularly the carapace, would have had to have arisen through degeneration of these attributes. Rather than the carapace being a primitive caridoid attribute, it is viewed then as an advanced character, evolved initially to serve a hydrodynamic function. Unity among the three eucaridan superorders relies on the presumed apomorphies of the escape reflex.

The possibility that the peracarids branched off

from a presyncarid ancestor before the eucarids also has been proposed. Attainment of the advanced caridoid attributes then arose independently within the Peracarida and Eucarida (see *illus.*). See CRUSTACEA; MALACOSTRACA. Patsy A. McLaughlin

Bibliography. D. E. Bliss et al. (eds.), *The Biology of Crustacea*, vols. 1, 2, and 5, 1982-1983; R. C. Brusca and G. J. Brusca, *Invertebrates*, 1990; A. A. Fincham, Eyes and classification of malacostracan crustaceans, *Nature*, 287:729-731, 1980; R. R. Hessler, The structural morphology of walking mechanisms in eumalacostracan crustaceans, *Phil. Trans. Roy. Soc. Lond.*, B296:245-298, 1982; P. A. McLaughlin, *Comparative Morphology of Recent Crustacea*, 1980; F. Schram (ed.), *Crustacean Phylogeny, Crustacean Issues*, 1983.

Eumycetozoida

An order of Mycetozoa. These are slime molds which form a plasmodium, a multinucleate stage that may in some species measure a foot or more across. The plasmodium, typically found on decaying plant material, is a migratory phagotroph preying on microorganisms. A mature plasmodium is a sheet of protoplasm containing a network of channels through which rapid endoplasmic streaming occurs. Direction of flow is reversible at intervals (shuttle-flow movement), but with a net unidirectional movement

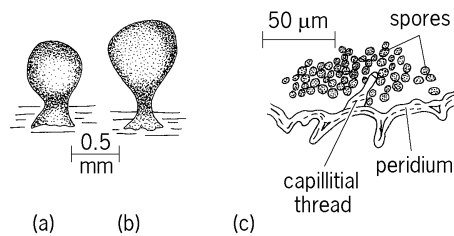


Fig. 1. *Arcyria cinerea*. (a, b) Development of sporangium. (c) Mature sporangium. (After R. P. Hall, *Protozoology*, Prentice-Hall, 1953)

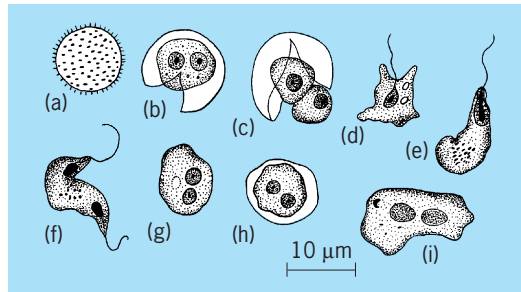


Fig. 2. Stages in the cycle of *Physarum polycephalum*. (a) Spore; (b, c) hatching of spore; (d, e) myxoflagellates; (f-h) zygote before nuclear fusion; (i) excysted zygote after first nuclear division. (After R. P. Hall, *Protozoology*, Prentice-Hall, 1953)

of the plasmodium. The cytoplasm often contains pigment (yellow, orange, green, or other colors).

Resistant cysts (spores), possibly haploid, are produced in sporangia (Fig. 1). Liberated spores (Fig. 2) hatch into myxoflagellates or myxamebae (which may change into flagellates). Either type may represent gametes, and mating types have been reported in certain species. In some species a large resistant stage may be formed by dehydration of the plasmodium and secretion of a membrane (sclerotium) around part or all of the plasmodium. Within the membrane, multinucleate cysts are produced. Such a mature sclerotium may remain viable for several years, and can be reactivated by rehydration. See PROTOZOA; RHIZOPODEA; SARCODINA; SARCOMASTIGOPHORA.

Richard P. Hall

Eumycota

True fungi, a group of heterotrophic organisms with absorptive nutrition, capable of utilizing insoluble food from outside the cell by secretion of digestive enzymes and absorption. Glycogen is the primary storage product of fungi. Most fungi have a well-defined cell wall that is composed of chitin and glucans. Spindle pole bodies, rather than centrioles, are associated with the nuclear envelope during cell division in most species. Typically, the fungal body (thallus) is haploid and consists of microscopic, branched, threadlike hyphae (collectively called the mycelium), which develop into radiating, macroscopic colonies within a substrate or host. The filamentous hypha may be divided by cross walls

(septa) into compartments. Hyphal growth is apical. Some species are coenocytic (without cross walls); others, including yeasts, are unicellular. Reproductive bodies are highly variable in morphology and size, and may be asexual or sexual.

Great changes in understanding the phylogenetic relationships of fungi have been brought about by the use of characters derived from deoxyribonucleic acid (DNA) sequences and the use of computer assisted phylogenetic analysis; these changes are reflected in current classification schemes. The Oomycota and Hyphochytriomycota, characterized by heterokont flagella and cellulose-containing cell walls, have been shown to occur in a clade (lineage) that includes organisms with heterokont flagella and chlorophylls *a* and *c*, including diatoms and brown algae. Myxomycota and other groups of slime molds also are excluded from the Eumycota. The Chytridiomycota has been recognized as the earliest diverging group of fungi, and many members are distinguished from all other fungi by their possession of a posterior whiplash flagellum, a character that signals the close relationship of fungi and animals. About 70,000 species of fungi have been described; however, some estimates suggest that there may be 1.5 million species. A modern classification of Eumycota follows:

Eumycota

Phylum: Chytridiomycota

Phylum: Zygomycota

Class: Zygomycetes

Trichomycetes

Phylum: Ascomycota

Class: Archiascomycetes

Hemiascomycetes

Euascomycetes

Phylum: Basidiomycota

Class: Hymenomycetes

Urediniomycetes

Ustilaginomycetes

Of cosmopolitan distribution, fungi are found in practically every type of habitat (soil, air, fresh or brackish water, ocean, and so forth). Most are strictly aerobic, although a few are anaerobes that live in the gut of herbivores. Some species are thermophilic. Many fungi form saprobic (including parasitic) relationships with animals and plants; the majority are saprobes. As now recognized, the Eumycota are a monophyletic group of the crown eukaryotes, presumed to have been present in the fossil record 900–570 million years ago. See FUNGI; MYXOMYCOTA; OOMYCOTA.

Meredith Blackwell

Bibliography. D. J. S. Barr, Evolution and kingdoms of organisms from the perspective of a mycologist, *Mycologia*, 84:1-11, 1992; D. Bhattacharya et al., Algae containing chlorophylls *a* + *c* are paraphyletic: Molecular evolutionary analysis of the Chromophyta, *Evolution*, 46:801-1817, 1992; N. J. Dix and J. W. Webster, *Fungal Ecology*, Chapman and Hall, London, 1995; D. Griffin, *Fungal Physiology*, 2d ed.,

Wiley-Liss, New York, 1993; D. D. Leipe et al., The straminopiles from a molecular perspective: 16S-like rRNA sequences from *Labyrinthula minuta* and *Cafeteria roenbergensis*, *Phycologia*, 33:369–377, 1994.

Euphausiacea

A well-defined marine order of the class Crustacea. These planktonic malacostracans are closely allied with, and sometimes included in, the Decapoda. The order contains two families, one with a single species, the other with about 84 species divided among 10 genera. Euphausiids are predominantly pelagic organisms of the open ocean, but a few species are neritic. All are shrimplike in appearance (see *illus.*).

Euphausiids form a significant proportion of the total planktonic biomass. This is especially true in higher latitudes and most pronounced in the Antarctic Ocean, where *Euphausia superba*, which attains a body length of some 2.4 in. (60 mm), occurs in vast numbers.

Distribution and ecology. All species are strictly marine and do not occur in fresh-water or brackish environments. Species are found in all oceans and seas, except the brackish Baltic and Black seas.

Known as krill to whalers, they constitute the diet of many whales, particularly the baleen whales. The main whale feeding grounds coincide with the areas of greatest concentration of euphausiids: areas of convergence, backwaters, vortices of mixed layers, centers of gyres, and fronts. They contribute to the diet of many other animals such as seals, herring, sardines, many birds, and even humans. Their concentration in areas of high productivity is correlated with their own diet, which consists mainly of phytoplankton, small crustaceans, and detrital matter. Their feeding appendages, which consist of not only

the mouthparts but also the anterior pairs of thoracic legs, are elaborate and can be used in several ways; this enables the animal to filter food particles from the water, collect organic material from the surface of the sediment, or catch live prey.

Most euphausiids live at considerable depth during daylight hours, and many undertake extensive diurnal vertical migrations to the surface layers at night. Their vertical position in the water column is thought to be photoregulated; during the day most species live at depth under conditions of blue-green light of low intensity. A few species live in the bathypelagic environment, having been captured at depths greater than 6000 ft (2000 m); these species have small eyes and attain considerable size, the maximum recorded body length of one species being 6 in. (150 mm).

Morphology. Euphausiids, with the exception of *Bentheuphausia amblyops*, possess photophores that emit a brilliant blue-green light.

The eyes are compound, bilobed in some genera, and contain three pigments: a carotenoid (astaxanthin), a melanoid, and a photolabile substance that is probably a visual pigment. The spectral sensitivity of the eye is greatest to blue-green light. See PHOTORECEPTION.

Respiration is by means of foliose, digitiform gills located at the bases of the second to eighth thoracic appendages.

The blood is a pale, leukocyte-bearing fluid with hemocyanin as the respiratory pigment. The heart is compact and has two pairs of ostia. See RESPIRATORY PIGMENTS (INVERTEBRATE).

Reproduction. The male copulatory organs are extremely complex and form the main criterion for specific identification. The phenomenon of swarming is common in many genera and is usually associated with reproduction. The male transfers a spermatophore to the spermatheca of the female at copulation. The early embryos usually live freely in the sea, but those of some 25 species are carried attached to the ventral regions of the thorax.

Life cycle. The larval stages are numerous. A nauplius hatches from the egg membranes and develops to a metanauplius. Three calyptopis stages follow during which the abdomen and stalked eyes develop. All species pass through these stages. The following developmental sequences, some 6–12 molts as furcillae, vary between species and within species in different localities and at different times. Sexual maturity is normally achieved at an age of 1 year, although some species such as *E. superba* and those living in the meso- and bathypelagic environments require 2 or more years. One species, *Thysanoessa longicaudata*, produces two generations each year in warmer parts of the North Atlantic.

Commercial exploitation. The possibilities of exploiting the Antarctic *E. superba* and other species in other regions have been investigated in some detail. Estimates of the possible yield of a commercial fishery for *E. superba* range 27–180 × 10⁶ ton (30–200 × 10⁶ metric tons) per year. This species forms

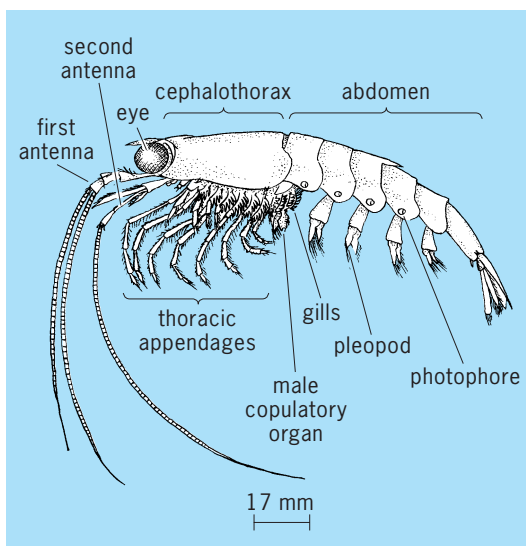


Diagram of a euphausiid crustacean.

large surface and subsurface aggregations that can be located visually or by echosounder. Ring and purse nets operated from small catchers or stern trawls operated from larger ships can be used to capture the swarms.

The euphausiids can be processed to yield large quantities of protein. Krill pastes and meal can be manufactured to feed domestic animals or used in therapeutic diets. Krill sausages, krill-stuffed eggs, and shrimp butter are already marketable products in Japan, Russia, and Germany.

In addition to *Euphausia superba*, some 30 species of euphausiids form surface or subsurface swarms or aggregations and are also therefore potentially exploitable resources. *Euphausia pacifica* is marketed for human consumption and also used, as is *Meganyctiphanes norvegica* in the North Atlantic, as food for fish in mariculture systems. See CRUSTACEA; DECAPODA (CRUSTACEA). John Mauchline

Bibliography. G. C. Eddie, *The Harvesting of Krill*, Food and Agricultural Organization, UN Development Program, Rome, Rep. GLO/SO/77/2, 1977; G. J. Grantham, *The Utilization of Krill*, FAO, UN Development Program, Rome, Rep. GLO/SO/77/3, 1977; J. W. Martin and G. E. Davis, *An Updated Classification of the Recent Crustacea*, Natural History Museum of Los Angeles, Science Ser. 39, 2001; J. Mauchline and L. R. Fisher, The biology of euphausiids, in *Advances in Marine Biology*, vol. 7, 1969; P. M. Mikkelsen, The Euphausiacea of eastern Florida (Crustacea: Malacostraca), *Proceedings of the Biological Society of Washington*, 100:275-295, 1987.

Euphorbiales

An order of flowering plants, division Magnoliophyta (Angiospermae), in the subclass Rosidae of the class Magnoliopsida (dicotyledons). The order consists of the large family Euphorbiaceae (about 7500 species) and 3 small satellite families which have fewer than 100 species among them. The Euphorbiales are a group of few-ovulate, mostly simple-leaved Rosidae in which the flowers have become unisexual and then undergone further reduction accompanied by aggregation. The trend toward reduction and aggregation culminates in the very large genus *Euphorbia*, with perhaps 1500 species, in which the pseudanthium has a pistillate flower consisting of a naked, tricarpellate pistil surrounded by several staminate flowers, each of which consists of a single stamen. The illusion that this cluster constitutes a single flower is fostered also by the cup-shaped involucre which often bears showy, petallike appendages. The Christmas poinsettia (*E. pulcherrima*) and the para rubber tree (*Hevea brasiliensis*) are well-known members of the Euphorbiaceae. Many African euphorbiads are spiny stem-succulents which resemble cacti in habit. Aside from the pronounced floral differences, the euphorbiads have a milky juice, which the cacti do not. See MAGNOLIOPSIDA; PLANT KINGDOM; ROSIDAE. Arthur Cronquist; T. M. Barkley

Europe

Although long called a continent, in many physical ways Europe is but a great western peninsula of the Eurasian landmass. Its eastern limits are arbitrary and are conventionally drawn along the water divide of the Ural Mountains, the Ural River, the Caspian Sea, and the Caucasus watershed to the Black Sea. On all other sides Europe is surrounded by salt water. The stretches of water between the Mediterranean and Black seas, however, are narrow, and it is possible to swim across the Bosphorus to Asia Minor. Not even the Mediterranean Sea is a real physiographic boundary inasmuch as the Atlas Mountains of northwestern Africa are structurally part of the great Tethys geosyncline. Of the oceanic islands of Franz Josef Land, Spitsbergen (Svalbard), Iceland, and the Azores, only Iceland is regarded, mainly on historical grounds, as an integral part of Europe; thus the northwestern boundary is drawn along the Danish Strait.

Maritime patterns. Europe has a large ratio of shoreline to land area reflecting a notable interfingering of land and sea. Excluding Iceland, the maximum north-south distance is 3529 mi (5680 km); and the greatest east-west extent is 3860 km. Of Europe's area of 3,880,000 mi² or 10,050,000 km² (that of the United States with Alaska and Hawaii is 3,620,000 mi² or 9,375,000 km²), 73% is mainland, 19% peninsulas, and 8% islands. Altogether, Europe has a shoreline length of 23,600 mi (37,900 km), that is, 0.205 mi/1000 ft² (3.55 km/1000 m²). Also, 51% of the land is less than 150 mi (250 km) from shores and another 23% lies closer than 300 mi (500 km). This situation is caused by the inland seas that enter, like arms of the ocean, deep into the northern and southern regions of Europe, which thus becomes a peninsula of peninsulas. The most notable of these branching arms of salt water are the White Sea, the North Sea, the Baltic Sea with the Gulf of Bothnia, the English Channel (La Manche), the Mediterranean Sea with its secondary branches—the Tyrrhenian, Adriatic, and Aegean Seas—and finally, the Black Sea. Even the Caspian Sea, presently the largest salt-water lake of the world, formed part of the southern seas before the folding of the Caucasus. Of these inland seas, only the Mediterranean and Black seas have depths greater than 400 fathoms (730 m). The penetration of the landmass by these seas brings marine influences deep into the continent and provides Europe with a balanced climate favorable for the evolution and settlement of humans. The favorable conditions are enhanced by the North Atlantic Drift (of Gulf Stream waters), which causes the northwestern coastal areas of Europe to be much warmer than the latitudes would suggest; for example, Labrador is much colder than the British Isles.

Land patterns. Europe has a unique diversity of landforms and natural resources. The relief, as varied as that of other continents, has an average elevation of 900 ft (300 m) as compared with North America's 1320 (440). However, 51% of the land lies below the 600-ft (200-m) contour lines, 27% is situated between 600 and 1500 ft (200 and 500 m), only



Fig. 1. Physiographic regions of Europe showing major units and subdivisions. A, Northern Uplands: 1, Kjölle Mountains; 2, Scottish Uplands; 3, Irish Uplands; 4, Pennines; 5, Wales; 6, Scandinavian Hill Land; 7, Southern Lake Region; 8, Iceland and Arctic Islands. B, Central Lowlands: 9, Arctic Lowlands; 10, Great Russian Plain; 11, Ukrainian Steppes; 12, Baltic Lowlands; 13, North Sea Plains; 14, Paris Basin; 15, Anglian Plains. C, Eastern Uplands: 16, Ural Mountains. D, Western Plateaus, Mountains, and Basins: 17, Iberian Meseta; 18, Massif Central; 19, Armorican Uplands; 20, Cornwall and South Ireland; 21, Ardennes, Shale Mountains, Vosges, and Black Forest; 22, Bohemian Massif; 23, Lysa-Gora; 24, Portuguese Lowlands and Andalusian Basin; 25, Aquitanian Basin; 26, Rhone Depression. E, Southern Highlands and Basins: 27, Baetic Cordillera; 28, Cantabrian Cordillera; 29, Pyrenees; 30, Alps; 31, Carpathians; 32, Balkan Ranges; 33, Caucasus and South Crimea; 34, Apennines; 35, Dinarides; 36, Pindos Mountains; 37, Peloponnesus; 38, Rhodope; 39, West Mediterranean Islands; 40, Aegean Islands; 41, Aragonian Basin; 42, Po Basin; 43, Carpathian Basin; 44, Wallachian Plain; 45, Thracian Basin. (1 mi = 1.6 km.)

5% is higher than 3000 (1000), and only 1% exceeds 6000-ft (2000-m) above mean sea level.

The shape and the overall physiographic aspect of the great peninsula are controlled by geologic structure that delimits the major regional units (Fig. 1). The triangular Central Lowlands (B on the map) with its acute apex in the Anglian Plains (B-15) and with younger strata covering the Precambrian Russian Table in the eastern parts (as in B-9, 10, 11) are bounded on each of the three sides by systems of mountain ranges with adjacent upland areas. On the northwestern side the erosional remnants of the once higher Caledonian ranges, folded during the early Paleozoic, are broken up into the Irish Uplands, Scottish Highlands and Uplands, and the Kjölle Mountains (A-3, 2, 1) bordering the predominantly granitic Baltic Shield of Precambrian age. In the east, the late Paleozoic Variscan mountain revolution folded up the Ural Mountains (C-16) and welded them to the Russian Table. The southern side of the triangle is more complex. Here the Paleo-Mesozoic Tethys geosyncline was folded twice into mountain chains; the Variscan revolution produced the

Hercynian Mountains and set the stage for the great Tertiary Alpine mountain revolution which elevated the chains of the Sierra Nevada, Pyrenees, Alps, Carpathians, Dinarides, Pindos, Balkan Ranges, and the Caucasus (E-27, 29, 30, 31, 35, 36, 32, and 33). These mountains form the Southern Highlands, whereas the eroded and block-faulted remnants of the Hercynian Mountains, from Iberia to the Lysa-Gora, complete the highly complex regional unit of Western Plateaus, Mountains, and Basins (D).

Climate. This aspect of Europe is determined by a number of factors. Probably the most important are a favorable location between 35° and 71° N latitudes on the western or more maritime side of the world's largest continental mass; the west-to-east trend (rather than north-south) of the lofty southern ranges and the Central Lowlands, as well as of the inland seas, which permit the prevailing westerly winds of these latitudes to carry marine influences deep into the continent; the beneficial influence of the North Atlantic Drift, which makes possible ice-free coasts far within the Arctic Circle; and the low elevation of the northwestern mountain ranges and the Urals,

which allows the free shifting of air masses over their crests.

Pressure, air mass, and moisture trends. The high atmospheric pressures built up in winter over central Asia obstruct the eastward passage of cyclonic storms, but the summer low of Asia, occurring simultaneously with the weakening of the Icelandic low, permits the prevailing westerlies to carry precipitation as far as the Urals and even beyond. Southeastward escape of these rain-bearing summer storms is barred by the Alps and Carpathians. Also, the seasonal shifts of air masses from and to the Icelandic and the subtropical Azores pressure centers produce precipitation maxima north of a line along the Pyrenees-Alps-Carpathian-Balkan Ranges and Caucasus Mountains during the warm months of the year, whereas south of this line the summers are more or less rainless. These main climatic trends are broken up into a mosaic of local climates, depending on the closeness to shore lines, local elevation, and slope direction, and influenced by human activities such as urbanization, industrialization, and agriculture.

Maritime climates. Among the marine climatic types, the Atlantic, the Mediterranean, and even the boreal polar climates are characteristic. True oceanic Atlantic climate occurs along the western fringes of Europe, from northern Portugal to central Norway and in the British Isles and Iceland. This west-coast marine climate is characterized by mild winters, with rarely a freezing period, and by somewhat cloudy, cool summers. The mean monthly range of temperature seldom exceeds 20°F (11°C). Precipitation, mostly in the form of drizzly rain, is regularly distributed over the year. Hence, runoff is small and soil erosion is at a minimum. Heather, wood anemone, small-leaved linden, English oak, and European beech are the representative plants of this climate. They are also prevalent in central Europe and along some of the shores of the western Mediterranean. In the Iberian and the Italian peninsulas, as well as in southern France, this climate is tempered by warm Atlantic influences reaching as far north as southern England and Denmark and eastward into Greece and the Aegean archipelago, with a plant community characterized by the ivy, English holly, box tree, and fig.

Mitigated by oceanic climate along northwest Norway's irregular shores, the tundra climate of the Barents Sea coast occupies but a narrow coastal strip. Average temperature of the warmest month is below 50°F (10°C) but above 32°F (0°C). Precipitation is scarce, soil is permanently moist, and true summers are absent. Tundra vegetation is only herbs, dwarf plants, and the characteristic tree of boreal regions, the dwarf birch.

On a relatively narrow strip along Europe's southern coasts the Mediterranean climate prevails. It influences to a certain extent some inland areas as far north as a line through the Pyrenees, Alps, Carpathian Basin, Crimea, and Caucasus. Mild winters with sunny weather and occasional rains change, through a delightful temperate spring, into dry and rainless summers with high temperatures and bright

blue skies. Some of the most characteristic plants of this climate are associated in the brushy macchia formation; others are mostly cultivated plants, such as the vine, maize (corn), chestnut, olive, and orange.

Continental climates. The counterparts to the marine climates are continental climates extending into Europe, across the Urals, from the great Asiatic landmass with its climatic belts running parallel to northern latitudes. Easterly air masses can move easily across the low, broadbacked Ural Mountains, and their impact reaches westward over Europe's central lowlands, where it mingles, especially in the summer, with the oceanic climate pushing eastward on the wings of the westerly winds. The northernmost zone of this continental influence has a subpolar continental climate, the southern boundary of which runs close to 60°N latitude. Brief summers are marked by temperatures above 50°F (10°C). The relatively meager precipitation covers the ground for almost one-half of the year in the form of snow, and late spring floods often occur as a result of the blocking effect of frozen estuaries of rivers flowing northward into Arctic waters. The vegetation consists mainly of coniferous forests known as taiga. Norway spruce and Scotch pine are the main species, intermixed northward with the characteristic boreal tree, the dwarf birch.

South of this zone, in the Great Russian Plains, the long cold-winter climate predominates. During 3 months of the winter period, temperatures drop below the freezing point and the snow is blown along land surfaces by icy continental winds. Temperatures averaging over 50°F (10°C) and the maximum rainfalls occur in the June-to-August growing season, but precipitation rarely exceeds 20 in. (50 cm) per year. Scotch pine on sandy soil and deciduous forests prevail in the northern parts of this zone, whereas oak forests alternate with the fertile cropland on chernozem soils to the south.

The southernmost continental climatic zone, covering roughly the Ukrainian Steppes, which are a western extension of the mid-latitude steppe and desert regions of Asia, is characterized by the semiarid steppe climate. With less than 5 in. (12.5 cm) of rainfall per year in the area of the lower reach of the Volga River, these eastern parts are a desert. Westward the 15–20 in. (37.5–50 cm) of yearly precipitation, with the maxima during the growing season, and mean annual temperatures just below 50°F (10°C) predominate in the natural grasslands. The most characteristic plant of these regions is the Stipa grass, which also appears to the westward in the Wallachian Plain and in the central portion of the Carpathian Basin.

Patterns of local climates. Among the more pronounced local climates of Europe, several are of importance (Fig. 2). An outstanding example is the semiarid climate of the Iberian Peninsula, partly due to the plateau character of the meseta upland and partly to warm Atlantic climatic influences. The peninsula has much scrubby vegetation and rocky soil covered with low grass.

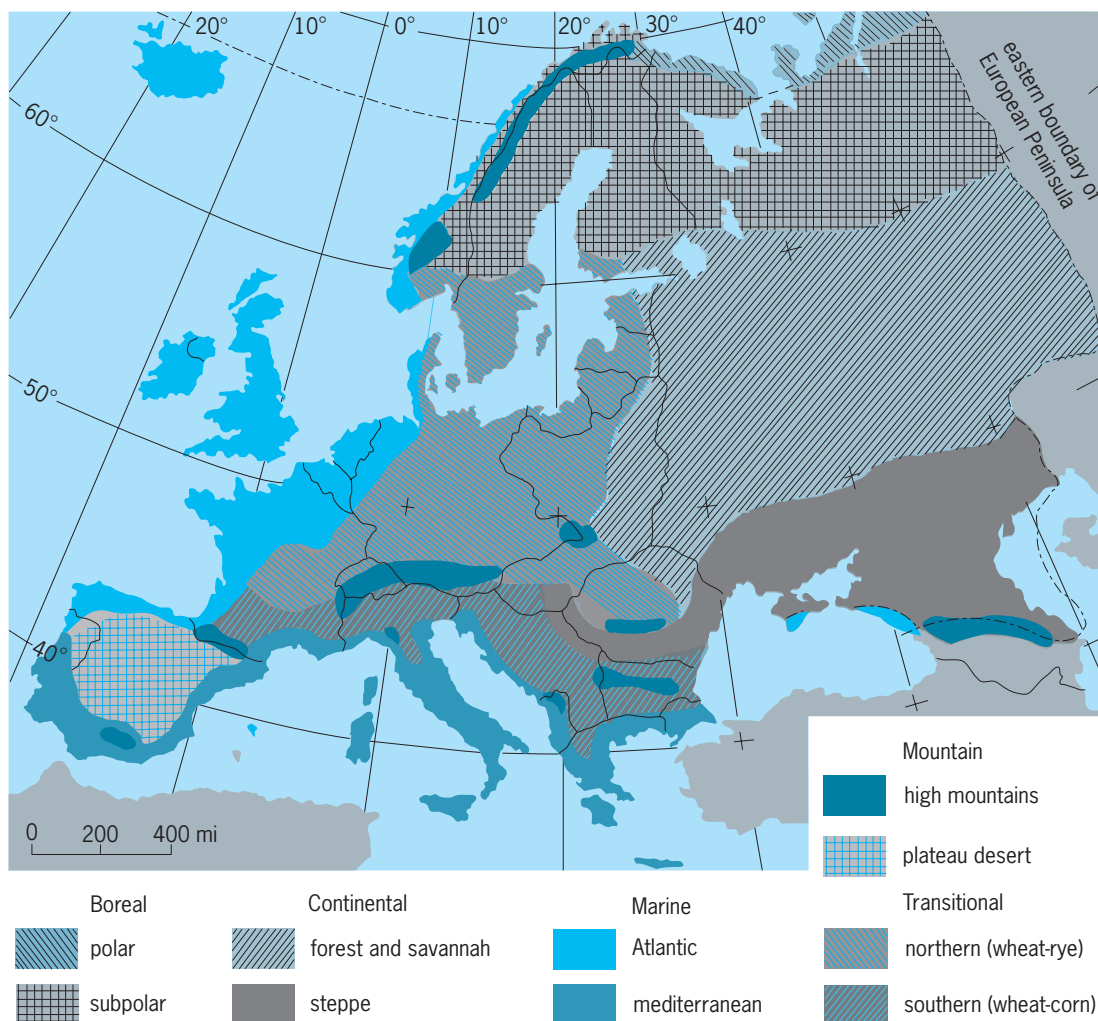


Fig. 2. Map of climatic regions of Europe. Regions and subdivisions are based on weather and vegetation data. 1 mi = 1.6 km.

Mountain climate is more a group of local climates than a single climatic type. Depending upon the elevation above sea level, the latitude in which they occur, the proximity to seas or landmasses, and the general slope directions, these climates show extreme variability. With glacial influences in the Alps, cordilleran character in the Pyrenees and in the southern Carpathians, subtropical influences in the western Caucasus, and boreal imprint in the Kjöllén Mountains, they do not have much in common. Only the occurrence of a fir line and a snow line, due to temperature decrease in the lofty sections of these mountain ranges, gives them a certain uniformity. Local winds in such areas have special names, such as the föhn of the Alps and the nemere of the eastern Carpathians. These climatic types intermix and compete in a nearly triangular area. This is comparable to Europe's structurally defined shape, with the apexes in the southern part of the Gulf of Bothnia, the western Pyrenees, and the western portion of the Black Sea. The center of this area is approximately near Wien, Austria. The northern transition zone is more tempered, the southern zone shows extremes of temperature, and toward the east both transitional climatic regions are apt to have variable extremes

of precipitation which quite often cause flood and drought years.

Drainage patterns. The intricate relief and the climates of Europe are well reflected in the drainage system. Extensive drainage basins with large slow-flowing rivers are developed only in the Central Lowlands, especially in the eastern part. Streams with the greatest discharge empty into the Black Sea and the North Sea, although Europe's longest river, the Volga, feeds the Caspian Sea. Second in dimension is the Danube, which crosses the Carpathian Basin and cuts its way twice through mountain ranges at the Gate of Bratislava and at the Iron Gate. The Rhine and Rhone are the two major Alpine rivers with headwater sources close to each other but feeding the North Sea and the Western Mediterranean Basin, respectively. Abundant precipitation throughout the year, as well as the permeable soils and the dense vegetation which temporarily store the water, provides the streams of Europe north of the Southern Highlands with ample water throughout the seasons. The combined effects of poor vegetation, rocky and desolate limestone karstlands, and slight annual precipitation result in intermittent flow of the rivers along the Mediterranean coast, especially on the eastern side

of peninsulas. Only the Alpine rivers carry enough water, and if it were not for the Danube and Rhone, both originating in regions north of the Alps, the only major river of the Mediterranean basin would be the Po. The main water divide between northern oceanic and southern Mediterranean drainage runs in a general east-northeast direction from Gibraltar along the Sierra Nevada, Cantabrian Mountains, Pyrenees,

Massif Central, Western Alps, Black Forest, Czech Forest, Northern Carpathians, Pripyat Marshes, and Valdai Hills to the northern central Urals. The Urals form the water divide between Europe and Asia.

Summary of physical patterns. Geologic structure, relief and drainage, climate, soils and vegetation, and human activities have shaped the physiographic aspects of Europe (Tables 1–4). Every regional

TABLE 1. Lakes

Name	Area, mi ² (km ²)	Countries, and area, mi ² (km ²)	Elevation, ft (m)	Approx. depth, ft (m)	Important tributaries	Important cities
Balaton	260 (660)	Hungary	348 (106)	36 (11)	Zala (outlet, Sio)	Siofok
Berre, Étangde	60 (156)	France	43 (13)	36 (11)		
Biesbosch, Lake	80 (200)	Netherlands	0	?		
Bodensee (Lake of Constance)	208 (539)	Germany, 127(328)	1,296 (395)	827 (252)	Rhine*	Konstanz, Bregenz
		Switzerland, 82 (211)				
Byeloseero, Lake	434 (1,125)	CIS	400 (122)	33 (10)		
Caspian Sea	168,000 (436,000)	CIS, Iran	-85 (-26)	3,100 (946)	Volga, Ural	Astrakhan, Baku
Chiemsee	30 (80)	Germany	1,700 (519)	240 (74)	Alz*	
Como, Lake	56 (145)	Italy	650 (198)	1,350 (410)	Adda*	Como
Corrib, Lake	73 (190)	Rep. of Ireland	26 (8)	144 (44)		
Femunnen, Lake	79 (205)	Norway	2,990 (910)	?		
Fertö, Neusiedler, Lake	120 (320)	Austria, 90 (232)	371 (113)	13 (4)		
		Hungary, 34 (88)				
Garda, Lake of	140 (370)	Italy	210 (65)	1,140 (346)	(outlet, Mincio)	
Geneva (Leman) Lake	225 (582)	France, 90 (234)	1,230 (375)	1,020 (310)	Rhone*	Geneva, Lausanne
		Switzerland, 134 (348)				
Ilmen Lake	354 (918)	CIS	59 (18)	10 (3)		
Inari, Lake	535 (1,385)	Finland	387 (118)	260 (80)	Ivalo	
Jalpug, Lake	85 (220)	CIS		?	Jalpug*	
Jannitza, Lake	40 (100)	Greece	1,345 (410)	72 (22)		
Ladoga, Lake	7,019 (18,180)	CIS 3,925 (10,166)	16 (5)	820 (250)	Svirj* (from Lake Onega)	
		Finland 3,094 (8,014)				
Lugano, Lake of	19 (49)	Italy, 7 (18)	899 (274)	915 (279)		Lugano
		Switzerland, 12 (31)				
Luzern, Lake of		Switzerland	1,433 (437)			Luzerne
Lomond, Loch	27 (71)	Scotland		?		
Maggiore, Lake	82 (212)	Italy, 73 (190)	636 (194)	367 (372)	Ticino*	Locarno
		Switzerland, 9 (22)				
Malaren, Lake	440 (1,140)	Sweden	3 (1)	210 (64)		
Mjosa, Lake	139 (359)	Norway	397 (121)	1,450 (443)		
Neagh, Lough	153 (396)	North Ireland	49 (15)	100 (31)		
Neuchâtel Lake of	83 (216)	Switzerland	1,400 (427)	510 (154)		Neuchâtel
Ohrida, Lake	142 (367)	Albania, 46 (119)	2,290 (698)	938 (286)	Drin*	
		Macedonia				
Onega, Lake	3,687 (9,550)	CIS	128 (39)	407 (124)		Petrozavodsk
Oulöjarvi, Lake	387 (1,002)	Finland	404 (123)	?		
Päijänne, Lake	503 (1,304)	Finland	26 (78)	305 (93)		
Peipus, Lake	1,383 (3,583)	CIS	28 (31)	59 (18)	Narva*	
Prespa, Lake	110 (286)	Greece, Albania, Macedonia	2,970 (906)	180 (54)		
Ree, Lake	64 (165)	Rep. of Ireland	130 (39)	120 (36)		
Saimaa, Lake	460 (1,200)	Finland	250 (76)	190 (58)		
Sasik, Lake	75 (195)	CIS	0	?		
Scutari (Shkoder), Lake	140 (370)	Albania, 57 (148)	18 (6)	145 (44)	Zeta (outlet, Bojana)	Shkoder (Scutari)
		Montenegro				
Segosevo, Lake	460 (1,200)	CIS	358 (109)	325 (99)		
Thingullatn, Lake	33 (85)	Iceland	2,050 (625)	?		
Trasimene, Lake	50 (129)	Italy	846 (258)	23 (7)		
Vänern, Lake	2,140 (5,550)	Sweden	145 (44)	320 (98)	Klar (outlet, Gotäelv)	Karlstad
Vättern, Lake	730 (1,900)	Sweden	290 (88)	390 (119)		Jönköping
Vierwaldstätter, Lake	44 (114)	Switzerland	1,430 (437)	702 (214)		
Zuider Zee (artificial)	?	Netherlands	0	?		Amsterdam
Zürich, Lake of	34 (89)	Switzerland	1,340 (409)	469 (143)		Zürich

*Tributaries which cross the lake.

TABLE 2. Selected peaks higher than 2800 m (9240 ft)

Name	Elevation, ft (m)	Location by mountain ranges, divisions and subdivisions, or island
Adamello	11,660 (3,554)	Prealps, Lombardian Alps, Alps of Chamonix
Adula	11,180 (3,406)	Western Alps, Lepontine Alps, Adula group
Aneto, Pico d'	11,170 (3,404)	Central Pyrenees, Maledetta group
Antelao	10,710 (3,263)	Eastern Alps, Dolomites, Antelao group
Argentea	10,820 (3,299)	Western Alps, Maritime Alps, Ligurian Alps
Bernina, Piz	16,290 (4,052)	Eastern Alps, Rhaetian Alps, Bernina group
Birkkar	9,700 (2,956)	Prealps, Bavarian Alps, Karwendel Mts.
Blanc, Mont	15,770 (4,807)	Western Alps, Pennine Alps, Mont Blanc group
Casse, Grand	10,110 (3,081)	Western Alps, Graian Alps, Massif de Vanoise
Coll' Alto	11,270 (3,435)	Eastern Alps, Noric Alps, Alps of Pusterthal
Corno, Monte	9,583 (2,921)	Appennines, Abruzzian Appennines, Gran Sasso
Cristallo, Monte	10,550 (3,216)	Eastern Alps, Dolomites, Tofano group
Dachstein	9,829 (2,996)	Eastern Alps, Salzkammergut
Dufour	15,300 (4,663)	Western Alps, Pennine Alps, Monte Rosa group
Écrins, Barre des	13,460 (4,103)	Prealps, Dauphinean Alps Pelvoux group
Elbrus	18,570 (5,660)	Western Caucasus
Etna*	10,760 (3,279)	Island of Sicily
Finsteraarhorn	13,970 (4,257)	Western Alps, Bernese Alps
Fluchthorn	11,165 (3,403)	Eastern Alps, Rhaetian Alps, Silvretta group
Gran Paradiso	13,320 (4,061)	Western Alps, Graian Alps
Grand Casse	10,110 (3,081)	Western Alps, Graian Alps, Massif de Vanoise
Grossglockner	12,460 (3,798)	Eastern Alps, Noric Alps Hohe (High) Tauern
Grossvenediger	12,010 (3,660)	Eastern Alps, Hohe Tauern
Hochfeiler	11,520 (3,510)	Eastern Alps, Noric Alps, Alps of Zillerthal
Hochgolling	9,393 (2,863)	Eastern Alps, Noric Alps Niedere (Low) Tauern
Hochkönig	9,639 (2,938)	Prealps, Alps of Salzburg
Jel Tepe	9,580 (2,920)	Rhodope, Pirin Planina
Jungfrau	13,670 (4,166)	Western Alps, eastern Bernese Alps
Kazbek	16,550 (5,045)	Eastern Caucasus
Leone, Monte	11,650 (3,552)	Western Alps, Lepontine Alps
Marmolata	10,970 (3,342)	Eastern Alps, Dolomites, Marmolata group
Matterhorn	14,780 (4,505)	Western Alps, Pennine Alps
Midi, Pic du	9,465 (2,885)	Central Pyrenees
Monte Rosa	15,190 (4,630)	Western Alps, Pennine Alps, Monte Rosa group
Mulhacen	11,420 (0.481)	Sierra Nevada
Mus Alla	9,596 (2,925)	Rhodope, Rila Planina
Olympos	9,573 (2,918)	Eastern Greek Mts., Chazia
Ortier	12,790 (3,899)	Eastern Alps, Rhaetian Alps, Ortier group
Paradiso, Gran	13,323 (4,061)	Western Alps, Graian Alps
Parseier	9,967 (3,038)	Prealps, Bavarian Alps, Alps of Lechtal
Rosa, Monte	15,190 (4,630)	Western Alps, Pennine Alps, Monte Rosa group
Sandspitze	9,393 (2,863)	Eastern Alps, Carnic Alps, Gailthal group
Scesaplana	9,734 (2,967)	Eastern Alps, Rhaetian Alps, Rhaetikon
Tabor, Mount	10,420 (3,177)	Western Alps, Cottian Alps (northern)
Todi	10,710 (3,623)	Western Alps, Glarnian Alps
Tofana	10,640 (3,243)	Eastern Alps, Dolomites
Triglav	9,393 (2,863)	Eastern Alps, Julian Alps
Vezzana, Cimi di	10,470 (3,191)	Eastern Alps, Dolomites, Pale group
Viso, Monte	12,600 (3,841)	Western Alps, Cottian Alps (southern)
Wildhorn	10,710 (3,264)	Western Alps, Bernese Alps (western)
Wildspitze	12,380 (3,774)	Eastern Alps, Rhaetian Alps, Alps of Oetzthal
Zuckerhütl	10,030 (3,057)	Eastern Alps, Rhaetian Alps, Stubai Alps
Zugspitze	9,724 (2,864)	Prealps, Bavarian Alps, Wetterstein Mts.

unit comprises regions with traits of similarity; but similarity does not mean uniformity. For example, the Alps have been sculptured by glaciers of the great Ice Age and their valleys are broad and long, whereas the Pyrenees and Caucasus, which bore few glaciers, have narrow, steep-sloped, V-shaped valleys, which give them a rougher aspect than that of the much higher Alps. Each region has its own characteristics of local relief, vegetation, animal life, and mineral resources. Some of the minerals occur in Precambrian and Paleozoic rocks of the ancient shields, blocks, and mountain chains,

for example, the iron ores of Scandinavia, England, Alsace-Lorraine, Spain, and Bohemia, or the coalfields of England, the Ruhr, Silesia, and the Donets Basin. Oilfields are mainly concentrated in the Carpathian region in Romania, Hungary, Poland, and Austria; a great variety of ores is bound to the metamorphic rocks and Tertiary volcanic areas affected by the Variscan and Alpine mountain revolutions; bauxite and manganese occur in pockets of the deeply weathered Mesozoic limestone areas, mainly in France, Hungary, and former Yugoslavia; and precious metals are most abundant in the Urals

TABLE 3. Major rivers

Name	Length, mi (km)	Drainage area, mi ² (km ²)	Discharges into	Main tributaries*	Source region and regions crossed
Bug (Polish)	484 (779)	28,370 (73,470)	Vistula	(R) Narew	Podolian Upland, Ukrainian and Baltic Lowlands
Bug (Ukrainian)	466 (750)	28,290 (73,280)	Black Sea		Podolian Upland, Ukrainian Steppes
Danube	1,800 (2,890)	312,192 (808,578)	Black Sea	(L) Vah, Tisza, Olt, Siret, Prut; (R) Inn, Drava, Sava, Morava	Black Forest, Bavarian Lowland, East Alpine Foreland, Carpathian Basin, Iron Gate, Wallachian Plain, Dobrudja, Moldavia boundary
Dnieper	1,340 (2,150)	204,000 (527,000)	Black Sea	(R) Pripjat, Beresina; (L) Seim	Great Russian Plain, Ukrainian Steppes
Dniester	850 (1,370)	29,700 (76,900)	Black Sea		Eastern Carpathians, Bessarabia (Ukrainian Steppes)
Don	1,160 (1,860)	166,000 (429,800)	Sea of Azov	(R) Donets	Great Russian Plain, Ukrainian Steppes
Donets	670 (1,078)	?	Don		Ukrainian Steppes
Duero	485 (780)	37,983 (98,375)	Atlantic Ocean	(R) Esla	Iberic Mts. (Cantabrian Cordillera), Iberian Plateau
Dvina	1,110 (1,780)	139,900 (362,300)	White Sea	Upper reach called Vicegda	Timan Mts., Arctic Lowlands
Ebro	580 (930)	33,000 (86,000)	Mediterranean Sea	(L) Segre	Cantabrian Cordillera, Aragonian Basin
Elbe	717 (1,154)	57,044 (147,744)	North Sea	(L) Saale	Sudeten Mts., Bohemian Basin, North Sea Plains
Garonne	400 (650)	32,741 (84,800)	Atlantic Ocean	(R) Lot, Dordogne;	Pyrenees, Aquitanian Basin
Guadalquivir	360 (579)	22,050 (57,120)	Atlantic Ocean	(L) Genil	Baetic Cordillera, Andalusian Basin
Guadiana	510 (820)	26,193 (67,840)	Atlantic Ocean		Sierra Morena, Iberian Plateau, Portuguese Lowland
Inn	320 (510)	10,000 (25,700)	Danube	(R) Salzach	Bernese Alps, Engadin, Bavaria Alps and Lowland
Kama	1,260 (2,030)	?	Volga	(L) Bjelaga; (R) Vjatka	Permian Hill Land, Western Ural foothills, Russian Plain
Kemijoki, Kemi	795 (494)	20,000 (50,000)	Gulf of Bothnia	(R) Ounas	Maan Hills, Northern Hill Land
Kuban	512 (824)	?	Black Sea		Caucasus, North Caucasian foothills
Loire	623 (1,002)	46,520 (120,500)	Atlantic Ocean	(L) Vienne	Plateau Central, South Paris Basin, Armorican Upland
Mernel	582 (936)	37,900 (98,100)	Baltic Sea		Pripjat Marshes, Baltic Lowlands
Meuse	575 (925)	13,000 (33,000)	North Sea		Plateau of Langres, North Sea Plains
Mosel	339 (545)	10,900 (28,200)	Rhine		Vosges, Hunsrück
Oder	561 (903)	48,136 (124,671)	Baltic Sea	(L) Neisse; (R) Warta	Sudeten Mts., Baltic Lowland
Pechora	932 (1,500)	123,700 (320,300)	Barents	(L) Ishma	Northern Ural, Arctic Lowlands
Po	418 (672)	28,950 (74,970)	Adriatic Sea	(L) Adda, Ticino	Cottian Alps, Po Basin
Rhine	824 (1,326)	86,640 (224,400)	North Sea	(L) Mosel; (R) Neckar, Main	Leptontine Alps, Vosges and Black Forest boundary, North Sea Plains
Rhone	505 (812)	38,000 (99,000)	Gulf of Lions	(R) Saone; (L) Durance	Leptontine Alps, Lake Geneva, Rhône Depression
Sava	564 (907)	38,800 (100,519)	Danube	(L) Drina	Julian Alps, Carpathian Basin
Seine	483 (776)	30,000 (77,800)	English Channel (La Manche)	(R) Marne, Oise	Plateau of Lengres, Paris Basin
Severn	210 (338)	8,100 (21,000)	Bristol Channel		Cambrian Mts., Wales
Tejo (Tagus)	628 (1,010)	31,250 (80,950)	Atlantic Ocean		Iberian Meseta, Portuguese Lowland
Thames	209 (336)	5,920 (15,340)	North Sea	Cotswold Hills, Anglican Plains	
Tiber (Tevere)	252 (405)	6,629 (17,169)	Tyrrhenian Sea		Tuscan Apennines, Apennine foothills (East)
Tisza	597 (961)	60,670 (157,135)	Danube	(L) Szamos, Maros, Koros (two sources, Black and White Tisza)	North-Eastern Carpathians, Carpathian Basin
Ural	1,480 (2,380)	84,900 (219,000)	Caspian Sea	(L) Ilek	Southern Ural, Kirghiz Steppes, Caspian Depression
Varda	200 (322)	?	Aegean Sea	(R) Crna	Albanian Alps, Thessalonian Plain
Vistula	862 (1,387)	76,560 (198,290)	Gulf of Gdynia	(R) Bug	Beskids (Western Carpathians), Baltic Lowlands
Volga	2,200 (3,570)	548,000 (1,420,000)	Caspian Sea	(R) Oxa; (L) Kama	Valdai Plateau, Great Russian Plain, Ukrainian Steppes, Caspian Depression
Warta	470 (760)	21,000 (54,000)	Oder		Lysa Gora foothills, Baltic Lowlands

* (L) left bank; (R) right bank.

TABLE 4. Volcanoes

Name (locate)	Last erupted	Elev. ft (m)	Rank*
Askja (Iceland)	1926	2,290 (698)	6
Etna (Sicily)	1979	10,760 (3,279)	1
Hekla (Iceland)	1948	5,108 (1,557)	3
Katla (Iceland)	1918	5,300 (1,600)	2
Santorin (Cyclades, Greek Aegean)	1928	1,920 (584)	7
Stromboli (Lipari Islands)	1934	3,038 (926)	5
Vesuvius (Apennine Peninsula)	1944	3,904 (1,190)	4
Vulcano (Lipari Islands)	1926	1,600 (500)	8

*By altitude.

and the Bihar Mountains. See AFRICA; ASIA; ATLANTIC OCEAN; BALTIC SEA; BLACK SEA; CONTINENT; EAST INDIES; MEDITERRANEAN SEA. Geza Teleki

Bibliography. E. A. Koster (ed.), *The Physical Geography of Western Europe*, 2005; M. Shahgedanova (ed.), *The Physical Geography of Northern Eurasia*, 2003; A. N. Strahler and A. H. Strahler, *Modern Physical Geography*, 4th ed., 1991; T. Unwin (ed.), *A European Geography*, 1998.

Europium

A chemical element, Eu, atomic number 63, atomic weight 151.96, a member of the rare-earth group. The stable isotopes, ^{151}Eu and ^{153}Eu , make up the naturally occurring element. The metal is the second most volatile of the rare earths and has a considerable vapor pressure at its melting point. It is very soft, is rapidly attacked by air, and really belongs more to the calcium-strontium-barium series than to the rare-earth series. See PERIODIC TABLE.

1																	18		
2																	2		
3	H																	He	
4	Li	Be											B	C	N	O	F	Ne	
5	Na	Mg	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	
19	K	Ca	Sc	Ti	V	Cr	Mn	Fe	Co	Ni	Cu	Zn	Ga	Ge	As	Se	Br	Kr	
37	Rb	Sr	Y	Zr	Nb	Mo	Tc	Ru	Rh	Pd	Ag	Cd	In	Sn	Sb	Te	I	Xe	
55	Cs	Ba	Lu	Hf	Ta	W	Re	Os	Ir	Pt	Au	Hg	Tl	Pb	Bi	Po	At	Rn	
87	Fr	Ra	Lr	Rf	Db	Sg	Bh	Hs	Mt	Ds	Rg								
lanthanide series			57	58	59	60	61	62	63	64	65	66	67	68	69	70			
actinide series			89	90	91	92	93	94	95	96	97	98	99	100	101	102			
			Ac	Th	Pa	U	Np	Pu	Am	Cm	Bk	Cf	Es	Fm	Md	No			

The element is attractive to the atomic industry, since the elements can be used in control rods and as nuclear poisons. These poisons are materials added to a nuclear reactor to balance the excess reactivity at start-up, and are so chosen that the poisons burn out at the same rate as the excess activity decreases. The television industry uses considerable quantities of phosphors, such as europium-activated yttrium

orthovana-dates, and other europium-activated yttrium phosphors have been patented. These phosphors give a brilliant red color and are used in the manufacture of television screens. See RARE-EARTH ELEMENTS.

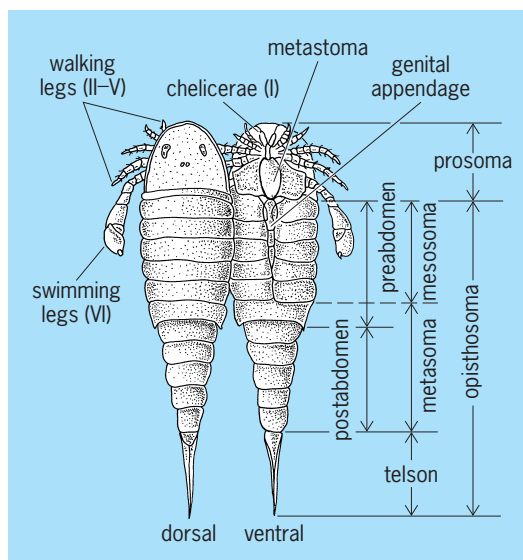
Frank H. Spedding

Bibliography. F. A. Cotton et al., *Advanced Inorganic Chemistry*, 6th ed., Wiley-Interscience, 1999; K. A. Gschneidner, Jr., J.-C. Bünzli, and V. K. Pecharsky (eds.), *Handbook on the Physics and Chemistry of Rare Earths*, 2005.

Eurypterida

An extinct Paleozoic group of aquatic arthropods, belonging to the subphylum Chelicerata and class Merostomata and thus related to the living marine xiphosurans (horseshoe crabs) and terrestrial arachnids (spiders, scorpions). Although most eurypterids were less than 10 in. (25 cm) in length, some members of the group were the largest arthropods of all time, reaching sizes up to 6 ft (2 m). See MEROSTOMATA.

The anterior portion of the eurypterid body is the prosoma. The prosoma is covered dorsally by the carapace, which bears a large pair of lateral compound eyes and a small pair of median simple eyes or ocelli (see *illus.*). Visible on the ventral surface of the prosoma are six pairs of appendages. The first pair, adjacent to the mouth, comprises small pincers known as the chelicerae. This feature unites the eurypterids with the xiphosurans and the arachnids. The chelicerae in some forms are enormous. The next four pairs of appendages were probably used in walking and food gathering. The distal portions of the sixth pair of appendages in many eurypterids are flattened and laterally expanded. They closely resemble the "swim paddles" of living blue crabs and were probably used to generate thrust and lift for



General morphology of eurypterids, in dorsal and ventral views.

swimming. The metastoma is a large plate which lies ventral to the bases of the sixth pair of appendages. It may represent a fused seventh pair of appendages and is unique to the eurypterids.

The posterior opisthosoma consists of 12 unfused segments and the terminal telson or tail spine. The opisthosoma can be divided into the anterior six-segmented mesosoma, which bears appendages, and the posterior six-segmented metasoma, which lacks appendages. The five pairs of platelike mesosomal appendages cover gills located on the ventral body wall. The first two mesosomal appendages carry the sexually dimorphic external genitalia. The opisthosoma can also be divided into a broad seven-segmented preabdomen and a narrower five-segmented postabdomen.

The telson of most eurypterids is styliform and closely resembles that of horseshoe crabs. In some forms the telson is distinctly curved, whereas in others it forms a broad flat plate, which may have functioned as a rudder during swimming.

Eurypterid fossils are found worldwide. The chitinous body of eurypterids is not mineralized, and so fossil remains are rare and generally restricted to a limited number of ancient aquatic environments. Almost all of the approximately 300 described species, belonging to about 65 genera, are from single localities or regions. The earliest definite eurypterids are found in the Ordovician, in both shallow and deep marine settings. The group diversified during the Silurian, peaked during the Late Silurian, and then declined during the Devonian. Silurian and Devonian eurypterids are found in both marine and nonmarine environments. Carboniferous (Mississippian Pennsylvanian), and Permian eurypterids are rare and are exclusively nonmarine (for example, coal swamps). The group became extinct sometime during the Permian. See DEVONIAN; MISSISSIPPIAN; ORDOVICIAN; PENNSYLVANIAN; PERMIAN; SILURIAN.

Eurypterids were almost exclusively aquatic, although some forms may have been amphibious. A number of forms were dominantly benthonic, but most were active and agile swimmers. Most eurypterids were probably carnivores, although they lacked the ability to crush heavily armored prey.

Although the eurypterids have no living descendants, they are very similar to (and are often found with) the earliest scorpions. The two groups probably share a common ancestor. See ARTHROPODA.

Roy Plotnick

Bibliography. R. Plotnick and T. Baumiller, The pterygotid telson as a biological rudder, *Lethaia*, 21:13–27, 1988; P. A. Selden, Autecology of Silurian eurypterids, *Spec. Pap. Palaeontol.*, 32:39–54, 1984; R. E. Snodgrass, *A Textbook of Arthropod Anatomy*, 1965.

Eustigmatophyceae

A small class of nonmotile, photosynthetic, unicellular algae, segregated from the class Xanthophyceae (Tribophyceae) on the basis of cytological, ultra-

structural, and biochemical features of the vegetative cell and zoospore. Because of their lack of chlorophyll *b*, these organisms may be placed in the division Chromophycota, but their lack of chlorophyll *c*, in contrast to other chromophytes, supports recognition at a higher level (such as the division Eustigmatophyta). Only a dozen species in three genera are known. These live chiefly in fresh water, but also in marine habitats and in soil.

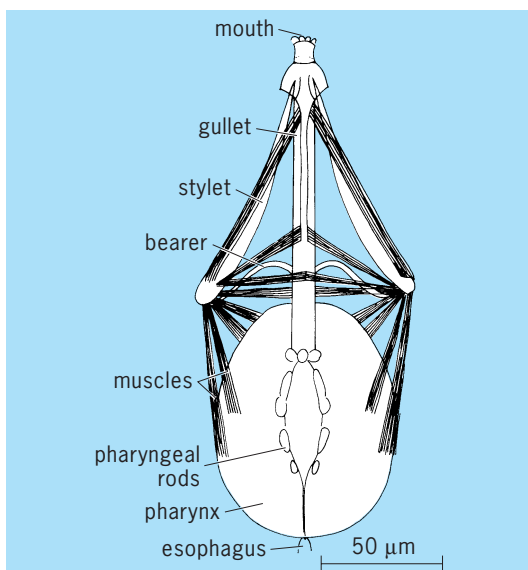
Vegetative cells usually are free-floating (attached in one species) and spherical, polyhedral, or elongate, with a maximum dimension of 5–60 micrometers. The cell wall is fibrillar-mucilaginous, of undetermined composition. Each cell contains one or more nuclei and usually a single chloroplast. The chloroplast is surrounded by three membranes, the outermost of which is continuous with endoplasmic reticulum that apparently is not confluent with the nuclear envelope. A distinctive pyrenoid—stalked and spherical or polyhedral—projects from the inner face of the chloroplast. Photosynthetic lamellae consist of three appressed thylakoids containing chlorophyll α , β -carotene, and violaxanthin. The presumed food storage product, of unknown composition, is deposited around the pyrenoid outside the chloroplast and appears lamellate in electron micrographs.

Sexual reproduction is unknown. Asexual reproduction is by the formation of two or four autospores or by zoospores. Zoospores are especially diagnostic of the class, although not known for all species. They are elongate, naked, and with two unequal flagella, only one of which is emergent in most species, the other being reduced to a basal body. The longer flagellum has two opposite rows of stiff hairs (tubular mastigonemes). Each zoospore has a single chloroplast without a pyrenoid. A large, reddish-orange eyespot not surrounded by membranes is situated in the cytoplasm near the base of the flagella. See ALGAE; XANTHOPHYCEAE. Paul C. Silva; Richard L. Moe

Bibliography. D. J. Hibberd, Notes on the taxonomy and nomenclature of the algal classes Eustigmatophyceae and Tribophyceae (synonym Xanthophyceae), *Bot. J. Linn. Soc.*, 82:93–119, 1981; D. J. Hibberd and G. F. Leedale, Observations on the cytology and ultrastructure of the new algal class, Eustigmatophyceae, *Ann. Bot.*, 36:49–71, 1972.

Eutardigrada

An order of tardigrades, lacking a cirrus lateralis, a sensory cephalic appendage, and a clava, or club-shaped appendage. Pharyngeal pockets are strengthened by separated rods or macroplacoids or are without thickenings. Claws are of different size, arranged in two pairs in which a larger and smaller claw are united. In *Milnesium* the claws are separated. *Haplomacrobotus* has two simple claws. Eight longitudinal muscles are associated with the midgut. At the beginning of the hindgut, there are three excretory glands, called the vasa malpighii, which are specifically stretched or trilobed. Each gland is composed of three cells. The gonoducts open into the rectum,



Buccal apparatus of Eutardigrada

resulting in a single opening known as the anogenital pore or cloaca.

Thermozodium esakii, from a warm spring in Japan, evidently represents a third order, the Mesotardigrada. This species combines such echiniscoidean features as the cirrus, clava, claws, and allusive plates with eutardigradan characters: pharyngeal rods (see **illus.**), vasa malpighii, and cloaca. See HETEROTARDIGRADA; TARDIGRADA.

Ernesto Marcus; Eveline Marcus

Eutectics

The microstructures that result when a solution of metal of eutectic composition solidifies. The eutectic reaction must be distinguished from eutectic microstructures.

Eutectic reaction. The eutectic reaction is a reversible transformation of a liquid solution to two or more solids, under constant pressure conditions, at a constant temperature denoted as the eutectic temperature T_E . The eutectic phase diagram is notable for displaying negative slopes of both pairs of solidus and liquidus lines, relative to the pure metal or compound terminal components A and B (**Fig. 1**). See PHASE EQUILIBRIUM.

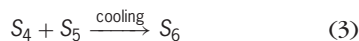
Eutectic reactions often constitute only a portion of a complex phase diagram, as in the Cu-Si system; other alloy systems, for example, the Cu-Mg system, may exhibit multiple eutectic reactions between intermediate phases. The transformation may occur between metals only, between metals and nonmetals (as in Ni-TaC), or between nonmetals only (KCl-NaCl).

In all cases the central liquid phase L touches the eutectic temperature line from above, as shown in **Fig. 1**; the eutectic alloy, at composition c_e , has the lowest melting point of any mixture of components A and B. Solidification of the eutectic composition

occurs with the simultaneous crystallization of the two solid phases α and β .

In terms of Gibbs' phase rule, at constant pressure, $F = C + 1 - P$, where F = number of degrees of freedom, C = number of components, and P = number of phases. Binary equilibrium involving three phases allows no degree of freedom: $F = 2 + 1 - 3 = 0$. Consequently, the eutectic reaction is known as an invariant reaction; the three phases must be at the equilibrium temperature and have fixed compositions.

A very similar invariant reaction involving liquid and solid phases is the peritectic reaction, in which two solids and one liquid are at equilibrium; in this case the central solid phase meets the invariant line from below (1). In the solid state the equivalent reactions with decreasing temperature are the eutectoid (2) and peritectoid (3), where S_i denotes solid phases.



Eutectic microstructures. Microstructures which are wholly eutectic in nature can occur only for a single, fixed composition c_e in each alloy system demonstrating the reaction (**Fig. 1**). However, any alloy whose composition passes through the eutectic invariant line (composition c_1 to c_2) undergoes the eutectic reaction as part of its solidification process. Initially, however, the liquid forms some primary crystals of the nearest end-point phase. Whatever liquid remains in the two-phase mixture at the equilibrium temperature then transforms by the eutectic reaction, and the primary crystals are left unchanged.

Microstructure classification. The great variety of patterns observed in eutectic microstructures, in some cases even in a single alloy system, has led to several attempts to develop classification systems. The most successful of these is due to E. Scheil, who classified eutectics on the basis of their mode of crystallization rather than on morphological characteristics, and this classification will be followed here.

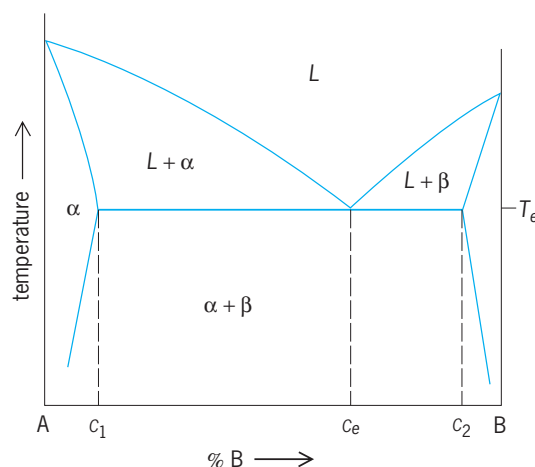


Fig. 1. Schematic eutectic phase diagram between components A and B.

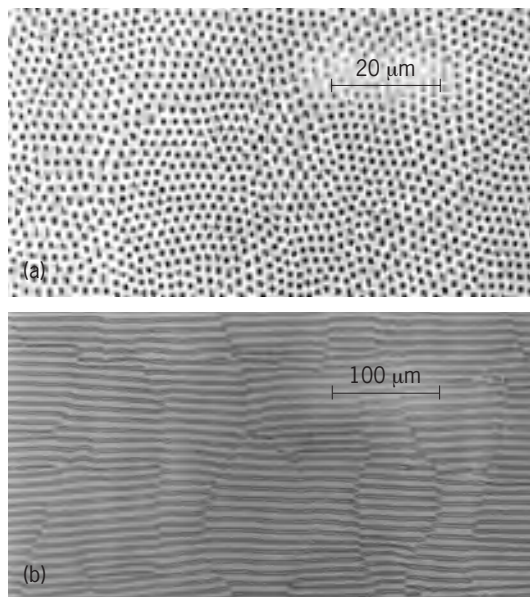


Fig. 2. Cross sections of aligned eutectic composites. (a) Fibers of Al_3Ni in Al matrix. (b) Alternating lamellae of Mg and $\text{Mg}_{17}\text{Al}_{12}$. Lighter etching phase is $\text{Mg}_{17}\text{Al}_{12}$.

Binary eutectic microstructures can be divided into normal and abnormal classes. Normal microstructures comprise primarily lamellar or fibrous types that are formed by simultaneous growth of the two solid phases in the form of parallel fibers in a continuous matrix (Fig. 2a) or as parallel lamellae (Fig. 2b). Normal growth occurs when the two solid phases grow locally at the same velocity, resulting in a planar solid-liquid interface. The fiber or lamellar axis tends to be perpendicular to the interface. When freezing begins at more than one center, each center gives rise to a eutectic grain, within which all lamellae or fibers have the same crystallographic orientation. Simple crystallographic relationships between coexisting eutectic phases have been noted; for example, in the lamellar Cd-Zn eutectic the following planes are parallel: $(0001)_{\text{Cd}}$ and $(0001)_{\text{Zn}}$; and the following directions are parallel: $(0110)_{\text{Cd}}$ and $(0110)_{\text{Zn}}$. Preferred growth directions have been reported for many systems, but there is no general agreement concerning their existence. See CRYSTAL STRUCTURE.

Abnormal eutectics are those in which the two solid phases do not grow at equal velocity. The faster-growing phase (usually that present in smaller volume) grows in a branching or dendriticlike pattern, and the other phase forms from the melt remaining between the branches.

Eutectic fusible alloys. Although technologically important alloy systems, particularly the Fe-C system (all cast irons used commercially pass through a eutectic reaction during solidification), exhibit at least one eutectic reaction, there has been little exploitation of wholly eutectic microstructures for structural purposes. A number of relatively low-melting eutectics are found in binary and higher-order mixtures of bismuth, lead, cadmium, tin, and indium. Table 1 shows the compositions and melting temperatures of eutectic fusible alloys which are used as solders, as

heat-transfer media, for punch and die mold and pattern applications, and as safety plugs. A silver-copper eutectic alloy is also used for high-temperature soldering applications. See SOLDERING.

Eutectic composites. During the early 1960s it was recognized that directional solidification of eutectic alloys so as to create a microstructure well aligned parallel to the growth direction could produce high-strength, multiphase composite materials. High-temperature gradients at the melt-solid interface, necessary to produce aligned alloys, are achieved by cooling the base of the crucible containing the liquid alloy as the crucible is withdrawn from the furnace. Excellent mechanical properties have been reported for metal-metal as well as nonmetal-reinforced metallic systems; examples of such alloys, all of which melt above 1832°F (1000°C), are given in Table 2. Note that lamellar eutectics are generally characterized by much higher-volume fractions V_f of reinforcing phase than fibrous alloys are. Many of the alloys listed demonstrate elevated temperature creep, stress-rupture, and fatigue properties which are superior to those of the best current conventional alloys. See COMPOSITE MATERIAL.

Several factors have been shown to exert major influences on the mechanical behavior of eutectic composites. These include: (1) Degree of alignment; the more perfect the alignment the stronger the alloy.

TABLE 1. Compositions and melting temperatures of eutectic fusible alloys

$T_E, ^\circ\text{F} (^{\circ}\text{C})$	Composition, percentage by weight				
	Bi	Pb	Sn	Cd	Other
116.2 (46.8)	44.70	22.60	8.30	5.30	19.10 In
136 (58)	49.00	18.00	12.00	—	21.00 In
158 (70)	50.00	26.70	13.30	10.00	—
196.7 (91.5)	51.60	40.20	—	8.20	—
203 (95)	52.50	32.00	15.50	—	—
216.5 (102.5)	54.00	—	26.00	20.00	—
255 (124)	55.50	44.50	—	—	—
281.3 (138.5)	58.00	—	42.00	—	—
288 (142)	—	30.60	51.20	18.20	—
291 (144)	60.00	—	—	40.00	—
351 (177)	—	—	67.75	32.25	—
361 (183)	—	38.14	61.86	—	—
390 (199)	—	—	91.00	—	9.00 Zn
430.3 (221.3)	—	—	96.50	—	3.50 Ag
457 (236)	—	79.7	—	17.7	2.60 Sb
477 (247)	—	87.0	—	—	13.00 Sb

TABLE 2. High-strength-high-temperature eutectic alloys

System	$T_E, ^\circ\text{F} (^{\circ}\text{C})$	Form	V_f
Ni-W	2732 (1500)	Fibrous	0.06
Ni,Cr,Al-TaC	2458 (1348)	Fibrous	0.05
Co-TaC	2556 (1402)	Fibrous	0.10
Ni-Ni ₃ Nb	2318 (1270)	Lamellar	0.26
Ni,Al,Cr-Ni ₃ Nb	2282 (1250)	Lamellar	0.33
Ni ₃ Al-Ni ₃ Ta	2480 (1360)	Fibrous	0.35
Ni ₃ Al-Ni ₃ Nb	2336 (1280)	Lamellar	0.44
Co-Cr ₇ C ₃	2377 (1303)	Fibrous	0.30
Co-NbC	2489 (1365)	Fibrous	0.12

(2) Growth speed; $\lambda^2 R = \text{constant}$, where $\lambda = \text{interphase separation}$ and $R = \text{growth rate}$; strength is inversely proportional to $\lambda^{-1/2}$. (3) Perfection of reinforcement; faults in the lamellae or fibers tend to reduce strength. (4) Ductility of coexisting phases; cracks nucleate readily in brittle phases, thereby limiting ductility of the composite. (5) Presence of intentional alloying additions; eutectics with up to nine components have been successfully aligned and reveal superior mechanical properties.

Among the major advantages of these alloys are extraordinary thermal stability of unstressed microstructures, retention of high strength to very close to the eutectic temperature of the respective alloys, and the ability to optimize strength by appropriate alloying additions to induce either solid-solution strengthening or intraphase precipitation of additional phases.

Future applications. The most likely applications for aligned eutectics are as gas turbine engine materials (turbine blades or stator vanes) or in nonstructural applications such as superconducting devices in which directionality of physical properties is an important requirement. See ALLOY; METAL, MECHANICAL PROPERTIES OF.

Norman S. Stoloff

Bibliography. K. Budinski and M. K. Budinski, *Engineering Materials: Properties and Selection*, 8th ed., 2004; J. C. M. Li (ed.), *Microstructure and Properties of Materials*, 2000.

Eutheria

A higher-level taxon that includes all mammals except monotremes and marsupials. Eutheria (Placentalia) is variously ranked as an infraclass or cohort within Mammalia. Eutheria includes over 4000 living species arranged in 18 orders; another 12 orders are known only from fossils. An ecologically diverse group, Eutheria includes primates, insectivores, bats, rodents, carnivores, elephants, ungulates, and whales. Like other mammals, eutherians are generally fur-covered and produce milk to nourish their young. In part because they can make their own body heat and regulate their body temperature, eutherians are widely distributed over most continents and occur in all oceans.

Reproduction. Eutherians, often called placental mammals, have a unique reproductive system in which unborn young are nourished for an extended period via a placenta. The placenta forms inside the uterus of the mother, and the developing fetus is attached to it by an umbilical cord. The placenta provides an interface between the bloodstream of the fetus and that of the mother, allowing the transport of oxygen and nutrients from the mother to the fetus and the transport of carbon dioxide and wastes from the fetus to the mother. This system permits retention of the young in the protective environment of the uterus during most of early development. Fetal survival rates are high under most conditions. Young are born in a relatively advanced state of development and are never sheltered in a pouch after

birth. Gestation time ranges from 20 days (for example, shrews and hamsters) to 22 months (elephants). Many eutherians have only one or two young per pregnancy, but as many as 20 offspring may be produced at a single birth in some species (such as some tenrecs).

Morphology and ecology. Eutherians range in size from insectivores and bats that weigh only a few grams to blue whales that can weigh over 190,000 kg (420,000 lb). All have a relatively large brain in which the two hemispheres are connected by a corpus callosum, which facilitates information transfer between the hemispheres. Eutherians exhibit more variation in ecology than any other group of vertebrates, and these differences are reflected in their morphological specializations. Most eutherians have a complex dentition, and structure of the teeth is indicative of diet. Plant-eating eutherians generally have teeth with grinding ridges, while meat-eating taxa have teeth for puncturing and shearing. Some groups (for example, rodents) have enlarged, ever-growing incisors that are used for gnawing. The sharp, bladelike incisor teeth of vampire bats are used to inflict tiny cuts in the skin of their prey, and the bats obtain a blood meal by lapping at the wound. Many ant- and termite-eating eutherians have a reduced dentition and use a long, sticky tongue to obtain their food. Long tongues are also used by other eutherians to feed on nectar in flowers (such as glossophagine bats). Some whales lack teeth and use flexible sheets of baleen to sieve small organisms from seawater. Baleen is made of keratin, the same material used to make hair, fingernails, claws, hooves, and the outer covering of horns.

Body form and limb structure reflect ecological habits. Eutheria includes the only mammals to evolve a fully aquatic life-style (for example, whales and sirenians), and these taxa have a fusiform body and paddlelike limbs for underwater locomotion. Whales have lost the hindlimbs and use a new structure, the tail fluke, for propulsion. Eutheria also includes bats, the only flying mammals, which have wings that consist of skin membranes supported by elongated forelimbs and fingers. Subterranean moles have forelimbs modified into shovellike structures for digging, arboreal sloths have hooklike hands for climbing and hanging under tree branches, ungulates have elongated limbs and hooves for running, and elephants have columnar limbs and modified foot pads to support great body weight. All of these modifications evolved from a primitive mammalian body plan that probably included small body size (less than 1 kg or 2 lb), an insectivorous diet, quadrupedal locomotion, and scansorial (scrambling) habits.

Sensory systems. Most eutherians have good vision, hearing, and a keen sense of smell. The relative importance of these sensory systems varies among taxa. Most diurnal eutherians (for example, ungulates, squirrels, and humans) have large eyes and well-developed vision. The eyes are reduced in many fossorial eutherians (such as moles). All bats retain functional eyes. Old World fruit bats (Megachiroptera) have excellent vision; echolocating bats

(Microchiroptera) rely on a sophisticated sonar system for navigation. High-frequency sounds produced in the larynx are emitted from the nose or mouth; the bats listen for and analyze returning echoes to obtain information about their environment. Some whales (odontocetes) also use echolocation.

Social behavior. All eutherians exhibit complex behaviors, and many live in social groups. Vocalizations, facial expressions, and scent marking play an important part in social interactions. Diurnal eutherians (for example, many primates and ungulates) often have colorful coat markings that may be used in social displays. Nocturnal forms often rely on olfactory rather than visual cues. The basic social bond is between mother and young. Apart from mother-offspring interactions, some species are solitary and interact with others only to mate. Even in social species, the sexes may remain separate except at mating season. However, many eutherians live in complex, relatively stable social groups (such as primates, elephants, and wolves). Monogamy is uncommon, occurring in fewer than 3% of eutherian species. Harem groups, in which large groups of females are guarded by a single male, are common, but many eutherians live in groups with multiple males and multiple females. One unusual social system occurs in naked mole rats (*Heterocephalus glaber*), which live in underground colonies organized like those of social insects, such as honeybees and ants.

Fossil record. The fossil record of Eutheria extends back at least into the Cretaceous Period. Several differences in the skull and dentition distinguish fossil eutherians from early members of other mammal lineages (for example, marsupials). The earliest eutherians were apparently small, nocturnal mammals that may have resembled some modern insectivores. Although Cretaceous eutherians are known from most continents, diversification of the modern orders apparently did not occur until the Paleocene and Eocene. Many species are known only from fossil teeth and skull fragments, but exceptionally well-preserved fossils of many groups have been found at unique sites such as Messel, Germany, and La Brea, California. See ARTIODACTYLA; CETACEA; CHIROPTERA; DENTITION; MAMMALIA; METATHERIA; THERIA.

N. B. Simmons

Bibliography. B. Grzimek, *Grzimek's Animal Life Encyclopedia: Mammals*, vols. 10–13, Van Nostrand, 1972; L. B. Halstead, *The Evolution of Mammals*, Eurobook, 1978; D. W. MacDonald, *The Encyclopedia of Mammals*, Facts on File, 1984; R. M. Nowak, *Walker's Mammals of the World*, 6th ed., Johns Hopkins University Press, 1999; W. Voelker, *The Natural History of Living Mammals*, Plexus Publishing, 1986.

Eutriconodonta (Triconodonta)

One of the taxonomically most diverse mammalian groups currently known from the Jurassic and Cretaceous (approximately 200 to 65 million years ago) when they were represented by 24 genera and 37

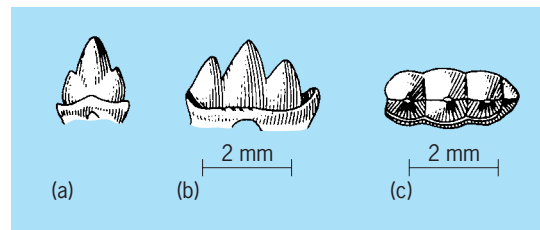


Fig. 1. Eutriconodontan molars had three main cusps. (a) External view of a lower molar of *Amphilestes*, an amphilestid eutriconodontan. (b) Internal view and (c) occlusal view of a lower molar of *Priacodon*, a triconodontid eutriconodontan.

named species (all extinct). The name Eutriconodonta refers to the characteristic shape of their posterior cheek teeth, which consist of three or sometimes four cusps aligned anteroposteriorly (Fig. 1). The lower molariform teeth in many eutriconodontans are linked by a tongue-in-groove articulation formed by a projection from the posterior end of one tooth that fits into a groove in the anterior end of the following tooth. See DENTITION; MAMMALIA.

Classification. In recent years, discovery and analysis of new fossils—some nearly complete skeletons—and revised interpretations of evolutionary relationships have greatly improved our understanding of the pattern of mammalian evolution during the Jurassic and Cretaceous. Previous interpretations grouped the mammals discussed here with morganucodontans and docodontans in the group Triconodonta. These three groups are now recognized as not being particularly closely related. For the sake of clarity, the core members of the Triconodonta have been united in a group designated Eutriconodonta, which includes three families. The Amphilestidae is made up of forms representing the primitive structural grade of the eutriconodontans. Although there is still uncertainty about their interrelationships, the other two families, Gobiconodontidae and Triconodontidae, appear to have evolved from an amphilestid ancestry. See ANIMAL EVOLUTION; DOCODONTA.

Fossil record. The oldest, but questionable, record of a eutriconodontan is from Early Jurassic deposits in India. Unquestionable records of the group are known from Middle Jurassic sites in England, China, and Mexico. The Late Jurassic and Early Cretaceous were the times of the eutriconodontans' greatest abundance, diversity, and biogeographic range, which included North America, western Europe, Asia, and Africa. Eutriconodontans survived into the latest Cretaceous of North America and possibly South America but disappeared from the fossil record before the mass extinctions marking the end of the Cretaceous. See EXTINCTION (BIOLOGY).

Characteristics. In comparison to modern mammals, the skulls of eutriconodontans were large relative to their bodies. Overall their skeletons were robustly built. Eutriconodontans exhibited a broad range of body sizes. Some were small, about the size

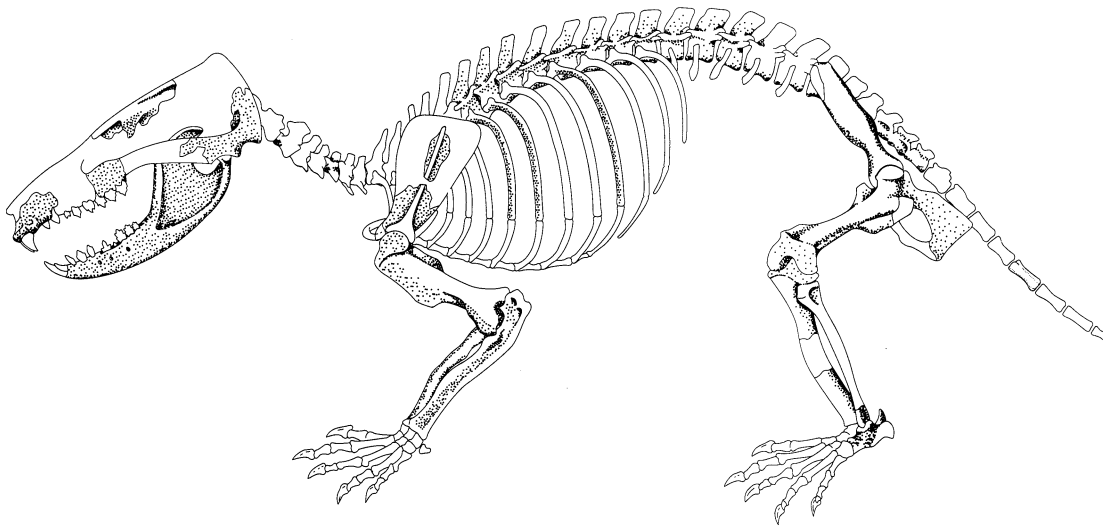


Fig. 2. Partial restoration of the skeleton of *Gobiconodon*, a gobiconodontid eutriconodontan. The length of the tail of *Gobiconodon* is unknown. In other gobiconodontids the tail contributed approximately 50% of the total length of the body. (After F. A. Jenkins, Jr., and C. Schaff, *The Early Cretaceous Mammal Gobiconodon (Mammalia, Triconodonta) from the Cloverly Formation of Montana*, *J. Vert. Paleontol.*, 8:1–24, 1988)

of modern mice. The largest eutriconodontans are found among the gobiconodontids. *Gobiconodon*, for example, grew to about the size of an adult Virginia opossum, *Didelphis virginiana* (approximately 5.5 kg or 12 lb) [Fig. 2]. The largest currently known Mesozoic mammals are species of the gobiconodontid *Repenomamus*, which are particularly well known from skeletons found in Early Cretaceous deposits in Liaoning Province, China. The largest species, *R. giganticus*, is estimated to have weighed 12–14 kg (26–30 lb). In any Late Jurassic or Early Cretaceous terrestrial fauna, particularly on continents in what is now the Northern Hemisphere, eutriconodontans usually were the largest mammals.

Evolution of an articulation between the dentary (the tooth-bearing bone of the lower jaw) and squamosal of the skull is a derived character of modern mammals. This joint replaced the primitive articular-quadrato articulation of their ancestors, the nonmammalian synapsids. A dentary-squamosal articulation was fully functional in all eutriconodontans. In contrast, in at least some eutriconodontans the primitive articular-quadrato articulation was retained. The articular remained connected to the dentary, while the quadrato contacted the squamosal. The articular and quadrato were not completely freed from their function of providing part of the articulation of the lower jaw to the skull. Possibly, in other eutriconodontans these bones were excluded from the jaw articulation and formed the malleus (articular) and incus (quadrato), which function as part of the sound-transmitting apparatus characteristic of the mammalian middle ear. See EAR (VERTEBRATE).

In other skeletal characteristics the eutriconodontans exhibit a mosaic of primitive and advanced conditions. For example, in the few eutriconodontan skeletons discovered so far, the posture of the fore-

limb with the elbow held relatively close to the body resembles the posture of most modern quadrupedal mammals. In contrast, the hindlimb preserved a primitive, sprawling posture. This and many other lines of evidence support the hypothesis that during the Jurassic and Cretaceous mammalian evolution was not characterized by simple, straight-line acquisition of new characters but involved a mosaic of changes in various lineages.

Probably many if not most eutriconodontans were carnivorous—predators, scavengers, or both. The morphology of their lower jaws, lacking an angular process but with a high coronoid process, indicates evolution of strong mandibular adductor musculature giving them a powerful bite. Evolution of a tongue-in-groove articulation linked the blades of the lower molars into an effective chopping or shearing mechanism. The canines and, in some species, the anterior incisors were large, stabbing teeth. Direct evidence of the diet of *R. giganticus* was provided by discovery of a skeleton with the remains of a juvenile dinosaur, *Psittacosaurus*, preserved in its stomach region.

Recent discoveries of new species of eutriconodontans have dispelled the view that during the Jurassic and Cretaceous mammals were uniformly small, insectivorous forms that scurried through the trees or hid on the ground to avoid the dominant dinosaurs. In addition to documenting the diversity of eutriconodontans, discoveries of aquatic otter- and beaverlike docodonts add to the evidence that mammals occupied a wide variety of niches in the ecosystems of the time.

William A. Clemens

Bibliography. T. S. Kemp, *The Origin and Evolution of Mammals*, Oxford University Press, 2005; Z. Kielan-Jaworowska, R. Cifelli, and Z.-X. Luo, *Mammals from the Age of Dinosaurs: Origins, Evolution, and Structure*, Columbia University Press, New York, 2004.

Eutrophication

The deterioration of the esthetic and life-supporting qualities of lakes and estuaries, caused by excessive fertilization from effluents high in phosphorus, nitrogen, and organic growth substances. Algae and aquatic plants become excessive, and when they decompose, a sequence of objectional features arises. Water for consumption from such lakes must be filtered and treated. Diversion of sewage, better utilization of manure, erosion control, improved sewage treatment and harvesting of the surplus aquatic crops alleviate the symptoms. Prompt public action is essential. See WATER CONSERVATION.

Extent of problem. In inland lakes this problem is due in large part to excessive but inadvertent introduction of domestic and industrial wastes, runoff from fertilized agricultural and urban areas, precipitation, and groundwaters. The interaction of the natural process with the artificial disturbance caused by human activities complicates the overall problem and leads to an accelerated rate of deterioration in lakes. Since population increase necessitates an expanded utilization of lakes and streams, cultural eutrophication has become one of the major water resource problems in the United States and throughout the world. A more thorough understanding must be obtained of the processes involved. Without this understanding and the subsequent development of methods of control, the possibility of losing many of the desirable qualities and beneficial properties of lakes and streams is great. See WATER POLLUTION.

Cultural eutrophication is reflected in changes in species composition, population sizes, and productivity in groups of organisms throughout the aquatic ecosystem. Thus the biological changes which are caused by excessive fertilization are of considerable interest from both the practical and academic viewpoints. See FRESH-WATER ECOSYSTEM; LAKE; LIMNOLOGY.

Phytoplankton. One of the primary responses to eutrophication is apparent in the phytoplankton, or suspended algae, in lakes. The nature of this response can be examined by comparing communities in disturbed and undisturbed lakes or by following changes in the community over a period of years during which nutrient input is increased. See ECOLOGICAL COMMUNITIES; PHYTOPLANKTON.

The former approach was utilized in studies at the University of Wisconsin, in which the overall structure of the phytoplankton communities of the eutrophic Lake Mendota, at Madison, and the oligotrophic Trout Lake, in northern Wisconsin, was analyzed. These investigations showed that in the eutrophic lake the population of species is slightly lower, although the average size of organisms is considerably larger, indicating higher levels of production than in the oligotrophic lake. When compared in terms of an index of species diversity, the community of the eutrophic lake displayed values lower than those observed in the oligotrophic lake. Seasonal changes and bathymetric differences in the

index of diversity were also more apparent in the eutrophic lake.

Often the low species diversity of the phytoplankton in eutrophic lakes is a result of high populations of blue-green algae, such as *Aphanizomenon flosaquae* and *Anabaena spiroides* in Lake Mendota. Frequently, however, species of diatoms such as *Fragillaria crotonensis* and *Stephanodiscus astrae* also attain high degrees of dominance in the community. Dense populations of the blue-green algae *Oscillatoria rubescens* were indicative of the deteriorating conditions in Lake Zurich (Switzerland) and Lake Washington (Seattle), but since diversion of sewage from these lakes, they are no longer a problem. The same species has been observed in several other lakes that have undergone varying degrees of cultural eutrophication. The relatively high nutrient concentrations in eutrophic waters appear to be capitalized on by one or two species that out-compete other species and periodically develop extremely high population levels. Because of the formation of gas vacuoles during metabolism, senescent forms of the blue-green algae rise to the surface of the lake, causing nuisance blooms.

In addition to nuisance scums in the pelagial, or open-water regions, the rooted aquatic plants and the attached algae of the littoral, or shoreward, region often prove to be equally troublesome in eutrophic lakes. Species of macrophytes, such as *Myriophyllum* and *Ceratophyllum*, and algal forms, such as *Cladophora*, frequently form dense mats of vegetation, making such areas unsuited for both practical and recreational uses.

Bottom fauna. Often in eutrophic lakes the bottom fauna display characteristics similar to those observed in the algal community. Changing environmental conditions appear to allow one or two species to attain high degrees of dominance in the community. Generally higher levels of production are associated with the change in structure of the community—the result being nuisance populations of organisms. In Lake Winnebago, in Wisconsin, for example, the lake fly or midge *Chironomus plumosus* develops extremely high populations, which as adults create an esthetic as well as an economic problem in nearby cities.

Great Lakes. It was generally thought that eutrophication would not be a major problem in large lakes because of the vast diluting effect of their size. However, there has been evidence that eutrophication is occurring in the lower Great Lakes. Furthermore, the undesirable changes in the biota appear to have been initiated in relatively recent years. Utilizing long-term records from Lake Erie, qualitative and quantitative changes have been observed in phytoplankton of that large body of water owing to cultural eutrophication. Total numbers of phytoplankton have increased more than threefold since 1920, while the dominant genera have changed from *Asterionella* and *Synedra* to *Melosira*, *Fragillaria*, and *Stephanodiscus*.

Other biological changes usually associated with the eutrophication process in small lakes have also

been observed in the Great Lakes. Of the five lakes, Lake Erie has undergone the most noticeable changes due to eutrophication. In terms of annual harvests, commercially valuable species of fish, such as the lake herring or cisco, sauger, walleye, and blue pike, were replaced by less desirable species, such as the freshwater drum or sheepshead, carp, and smelt. Similarly, in the organisms living in the bottom sediments of Lake Erie, drastic changes in species composition have been observed.

Oxygen demand. It is apparent that the increase in organic matter production by the algae and plants in a lake undergoing eutrophication has ramifications throughout the aquatic ecosystem. Greater demand is placed on the dissolved oxygen in the water as the organic matter decomposes at the termination of life cycles. Because of this process, the deeper waters in the lake may become entirely depleted of oxygen, thereby destroying fish habitats and leading to the elimination of desirable species. The settling of particulate organic matter from the upper, productive layers changes the character of the bottom muds, also leading to the replacement of certain species by less desirable organisms. Of great importance is the fact that nutrients inadvertently introduced to a lake are for the most part trapped there and recycled in accelerated biological processes. Consequently, the damage done to a lake in a relatively short time requires a many-fold increase in time for recovery of the lake.

Action programs. Lake eutrophication represents a complex interaction of biological, physical, and chemical processes. The problem, therefore, necessitates basic research in a wide variety of scientific disciplines. Moreover, to be profitable, such research must be coordinated into well-integrated team-research efforts requiring extensive monetary support.

Studies are needed in such areas as monitoring the amount of those nutrients that reach critical levels in lake waters and lead to the development of nuisance growths of plants and algae. Nitrogen and phosphorus are undoubtedly important.

Sewage effluent is the major contributor of nitrogen and phosphorus to lakes, followed by runoff from manured and fertilized land. Considering nitrogen alone, rain adds more than any single source. Its nitrogen content comes from combustion engines and smokestacks. Where sewage effluent and agricultural drainage have been diverted from lakes, an improvement in nuisance conditions occurs. Dead lakes become alive again; hence this treatment is the first step in alleviation.

The conditions observed in lakes and streams reflect not only the processes operating within the body of water but also the metabolism and dynamics of the entire watershed or drainage basin. After precise identification of the critical nutrient compounds, it is necessary to determine the nutrient budget of the whole drainage basin before acting to alleviate the undesirable fertilization of a lake or stream.

Methods of treating sewage plant effluent are

being explored, and further support for these efforts is justified. Agricultural practices such as low tillage reduces erosion, hence the amount of fertilizer in the runoff. See AGRICULTURAL SOIL AND CROP PRACTICES; EROSION; SEWAGE TREATMENT.

It is known that aquatic plants concentrate nutrients from the lake waters in their tissues. The removal or harvesting of aquatic plants in eutrophic lakes, consequently, is a good potential method for reducing nutrient levels in these lakes. Similarly, significant amounts of nitrogen and phosphorus are concentrated in fish flesh. Efficient methods of harvesting these organisms are important in impoverishing a well-fertilized lake.

Utilization or land disposal of farm manure is a major problem. Animal manures are largely unsewered, yet in the Midwest it is equivalent to the sewage of 350,000,000 people.

In addition to improvements in waste disposal, more research and development are needed on the profitable utilization of surplus algae, aquatic plants, fish, manure, and sewage. The overfertilized lake needs to be impoverished of its nutrients as well as protected from inflowing sources. Chemicals have been used to poison the plants and algae, but this is not a good conservation practice because the plants and algae rot and provide more nutrients. Moreover, eventual harm to other species has not been assessed.

It would seem desirable to set aside certain lake areas for research purposes. More information is needed to decide upon the best plans for allowing the domestic development of these areas with the least disturbance to the water resources. Steps will be necessary to devise optimum zoning laws and multiple-use programs in light of the intense economic and recreational uses made of water resources. Legislation and law enforcement in relation to public interactions undoubtedly will be a complex problem to overcome in this respect.

Cultural eutrophication is a paradoxical condition, since it is in large part due to human economic, agricultural, and recreational activities and at the same time eventually conflicts with these same activities of society in general.

Arthur D. Hasler

Bibliography. G. D. Cooke, *Restoration and Management of Lakes and Reservoirs*, 3d ed., 2005; F. de Jong, *Marine Eutrophication in Perspective: On the Relevance of Ecology for Environmental Policy*, 2006; A. J. Horne and C. R. Goldman, *Limnology*, 2d ed., 1994; M. C. T. Scholten et al., *Eutrophication Management and Ecotoxicology*, 2005.

Evaporation

The process by which a liquid is converted into a vapor. In the liquid phase, the substance is held together by intermolecular forces. As the temperature is raised, the molecules move more vigorously, and in increasingly high proportion have sufficient energy to escape from their neighbors. Evaporation is therefore slow at low temperatures but faster at higher

temperatures. In an open vessel, the molecules escape from the vicinity of the liquid, and there is a net migration from the liquid to the atmosphere. In a closed vessel, net evaporation continues until the number of molecules in the vapor has risen to the stage at which the rate of return from the vapor to the liquid is equal to the rate of evaporation. At this stage there is a dynamic equilibrium between the liquid and its vapor, with evaporation and its reverse, condensation, occurring at the same rate. The pressure of the vapor in the closed vessel is called the vapor pressure of the substance; its value depends on the temperature. Boiling occurs in an open vessel (but not in a closed vessel) when the vapor pressure is equal to the ambient pressure. *See* BOILING POINT; KINETIC THEORY OF MATTER; VAPOR PRESSURE.

Evaporation is an endothermic (heat-absorbing) process because molecules must be supplied with energy to overcome the intermolecular forces. The enthalpy of vaporization, $\Delta_{\text{vap}}H$ (formerly, the latent heat of vaporization), is the heat required at constant pressure per mole of substance for vaporization. The entropy of vaporization, $\Delta_{\text{vap}}S$, at the boiling point, T_b , is equal to $\Delta_{\text{vap}}H/T_b$. According to Trouton's rule, for many liquids the entropy of vaporization is close to $85 \text{ J/K} \cdot \text{mol}$. This value reflects the similar change in disorder that occurs when a liquid is converted into a gas. However, certain liquids (water and mercury among them) are more structured than others, and have a bigger entropy of vaporization than Trouton's rule suggests. *See* ENTHALPY; ENTROPY; THERMODYNAMIC PRINCIPLES.

Volatile liquids evaporate more rapidly than others at the same temperature. Such liquids have relatively weak intermolecular forces. In general, the rate of evaporation depends on the strengths of the intermolecular forces and the rate at which heat is supplied to the liquid. *See* INTERMOLECULAR FORCES; LIQUID.

P. W. Atkins

Bibliography. P. W. Atkins and J. de Paula, *Physical Chemistry*, 8th ed., 2006.

Evaporator

A device used to vaporize part or all of the solvent from a solution. The valuable product is usually either a solid or a concentrated solution of the solute. If a solid, the heat required for evaporation of the solvent must have been supplied to a suspension of the solid in the solution; otherwise the device would be classed as a drier. The vaporized solvent may be made up of several volatile components, but if any separation of these components is effected, the device is properly classed as a still or distillation column. When the valuable product is the vaporized solvent, an evaporator is sometimes mislabeled a still, such as water still, and sometimes is properly labeled, such as boiler-feedwater evaporator. In the great majority of evaporator installations, water is the solvent that is removed. *See* DRYING.

Uses. Evaporators are used primarily in the chemical industry. Common salt is made by boiling a satu-

rated brine in an evaporator. The salt precipitates as a solid in suspension in the brine. This slurry is pumped continuously to a filter, from which the solids are recovered and the liquid portion returned for further evaporation. In the manufacture of pulp and paper, the waste liquor from cooking the wood is a dilute solution of inorganic cooking chemicals and soluble, organic wood constituents. The liquor is disposed of by concentrating it in evaporators to a strength at which the liquor will burn in a boiler. This avoids a water-pollution problem, recovers the inorganic chemicals for reuse, and provides sufficient heat energy not only to operate the evaporator but also to supply other needs of the mill as well. In the alkali industry, salt brine passed through a diaphragm-type electrolytic cell yields a dilute solution of salt and caustic soda. Evaporation to a strength of about 50% NaOH purifies the caustic by causing precipitation of practically all of the salt. It also simplifies shipment by bringing about a more than fivefold reduction in volume. Further evaporation of the caustic solution to a final temperature of about 700°F (371°C) produces a practically anhydrous product that freezes at about 600°F (316°C) and is shipped as a solid. Evaporators are widely used in the food industry, usually as a means of reducing volume to permit easier storage and shipment. Evaporators are also the most commonly used means of producing potable water from seawater or other contaminated sources.

Classification. The vaporization of solvent requires large amounts of heat. Provisions for transferring this heat to the solution constitute the largest element of evaporator cost and the principal means of distinguishing between types of evaporators. Practically all evaporators fall into one of the following categories:

1. Those heated by a flame that burns below the liquid surface, and in which the hot combustion gases are bubbled through the liquid.
2. Those in which the flame and combustion gases are separated from the boiling liquid by a metal wall, or heating surface.
3. Those in which steam or other condensable vapor is the source of heat, and in which the steam condenses on one side of the heating surface and the heat is transmitted through the wall to the boiling liquid.

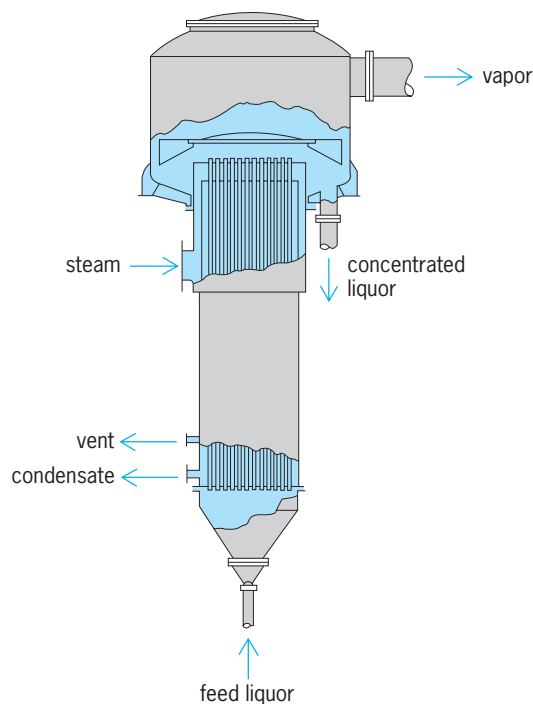
Submerged-combustion evaporators (type 1) are used primarily to concentrate solutions that would deposit a heat-insulating blanket of scale on the solid heating surfaces of other types of evaporators. Since the evolved vapor is mixed with the combustion gases, neither the solvent vapor nor its heat content can be recovered easily.

Direct-fired evaporators (type 2) are typified by the steam boiler and the old maple-syrup kettle. They are not commonly used for the concentration of solutions, primarily because local overheating can cause formation of insulating deposits on the heating surface, which then becomes overheated and may melt or burn through. Also, a large heating surface is needed to recover the heat in the combustion gases, and there is usually no cheap metal that will resist

attack by both the combustion gases and the boiling liquid.

Steam-heated evaporators (type 3) are by far the most common, primarily because condensing steam gives up its heat so readily. Condensing steam film coefficients are usually in excess of 1000 Btu/(h)(ft²)(°F) [5.7 kW/(h)(m²)(°C)]. Thus the design of the evaporator and materials of construction can be suited to the solution being concentrated instead of being dictated by the problem of getting heat to the heating surface. The heating surface is usually in the form of metal tubes, since this represents the most economical method of putting the largest heating surface in the smallest volume. The tubes may be vertical or horizontal, and the boiling liquid may be either inside or outside, depending on the characteristics of the solution, such as viscosity, ratio of feed to evaporation, and whether or not a salt is deposited. See DISTILLATION.

Operation. The vaporization capacity of an evaporator is determined by the usual rules of heat transfer, and is directly proportional to the area of heating surface, to the difference between condensing steam and boiling liquid temperatures, and to the coefficient of heat transfer. The heat-transfer coefficient is usually limited by conditions on the boiling liquid side, although the condensing steam film and the resistance of the metal wall have some influence. Various means are employed to increase the boiling-film coefficient and all involve movement of the liquid relative to the heating surface. A great many evaporators use only natural convection to accomplish this circulation. Typical is the long-tube vertical type shown in the **illustration**. Feed liquid enters the bottom of a nest of vertical tubes and begins to boil



Long-tube vertical evaporator.

as it passes up the tubes. The boiling causes a large increase in volume, which accelerates the liquid to high velocities and gives good heat-transfer performance. The vapor-liquid mixture is separated in the chamber at the top of the tubes. The liquor may all be discharged as product, or part may be recirculated to the feed inlet.

A pump or agitator may be used to force the liquid past the heating surface if even higher heat-transfer coefficients are needed, for example, when corrosive conditions dictate use of the smallest possible area of an expensive alloy. Forced circulation evaporators are also used for scaling liquids or those from which a salt is to be crystallized, since there is less tendency for solids to form on the heating surfaces. Such evaporators usually consist of a flash chamber or crystallizing chamber, a conventional shell-and-tube heat exchanger, and a pump to circulate fluid from the chamber to the exchanger and back to the flash chamber. Another type of forced-circulation evaporator for extremely viscous, heat-sensitive, or foamy materials either rotates the heating surfaces or employs wipers that sweep the material around the walls of a stationary surface.

The water vapor evolved in an evaporator is usually about the same in quality and quantity as the steam used to heat the evaporator. The only difference is that the vapor has a lower pressure and hence lower condensing temperature. It is possible to compress the vapor so that it can be used as the heating steam in the same evaporator. Such thermocompression evaporators require far less energy for the compressor than would be needed to generate fresh steam. However, to keep the compressor cost and power consumption within reason, it is necessary to use a narrow compression ratio and this requires a large and expensive evaporator.

Assuming a perfect compressor, the power requirement is given by the Carnot equation: $W = Q\Delta T/T$, where W and Q are work required and amount of heat pumped in the same units, ΔT is the difference between saturated steam temperatures at compressor discharge and suction pressures, and T is the absolute suction temperature. For pure water boiling at atmospheric pressure (212°F or 100°C), the latent heat is 970.3 Btu/lb (2257 kilojoules/kg). At a temperature difference of 10°F (5.6°C), the ideal work required is only $(970.3)(10)/(672) = 14.4$ Btu/lb (33.6 kilojoules/kg), or 35.3 kWh per 1000 gal of water evaporated (33.6 MJ per 1000 liters). This is only 1.5% of the heat energy required for a simple evaporator, but the energy is expensive mechanical energy rather than low-grade heat energy.

If pure water is being boiled, this 10°F (5.6°C) temperature difference established by the compressor is available as the driving force to transfer heat from the compressed steam through the heating surface to the boiling liquid. If the boiling liquid is an aqueous solution, it has a boiling temperature higher than that of pure water at the same pressure and the difference in these temperatures is called the boiling point elevation (bpe). The bpe cannot be utilized as a part of the

driving force. Thus, if the solution being evaporated had a bpe of 30°F (16.7°C), the compressor would have to work across a 40°F (22.2°C) compression range and would require four times the above power to establish the same 10°F (5.6°C) driving force in the evaporator.

An alternative method of reducing the energy requirement of an evaporator is to use the water vapor evolved in one part to heat another part in which liquid is boiled at a lower temperature. The water vapor evolved in this part, termed an effect, can then be used to heat another effect boiling at still lower temperature (and pressure), and so on. The ultimate limit is determined by the need to discharge to a heat sink the heat contained in the vapor from the last effect, which amounts to most of the heat supplied by steam to the first effect (the balance of the entering heat leaves with the condensate and concentrate, which are generally hotter than the feed solution). The heat sink is usually water from a river or other source, which limits the boiling temperature in the last effect to about 100°F (38°C).

If the condensate of last-effect vapors is valuable, the vapors are condensed in a shell-and-tube heat exchanger by the cooling water. If the solution is the desired product, the vapors are brought into direct contact with a shower of cooling water. This is usually done in a barometric condenser, which is elevated to such a height (about 34 ft or 10.3 m) that the high vacuum cannot prevent the water from draining out by gravity. To maintain the vacuum in the condenser, it also is necessary to remove noncondensable gases with a vacuum pump. These gases originate in the feed, as air dissolved in the condenser water and as air leakage into the evaporator. Only in the rarest of circumstances, as when a very low boiling temperature is needed to avoid degradation of the solution, is it feasible to compress the entire last-effect vapor for discharge at a higher temperature to the heat sink.

Such multiple-effect evaporators are more expensive than single-effect units because each effect can operate at only a fraction of the total difference, ΔT , between the temperature at which heat is accepted from the prime steam and the temperature at which it is rejected to the heat sink. For a single effect, the amount of heat transfer surface A required is given by the equation: $A = Q/U\Delta T$, where Q is the rate at which heat must be transferred to achieve the desired evaporation rate and U is the coefficient of heat transfer.

In a triple effect, each effect evaporates only about one-third of the water and must transfer only about one-third of the total heat, but the total temperature difference must also be divided between the three effects. Thus the heating surface in each effect may be written: $A_n = (Q/3)/U(\Delta T/3) = Q/U\Delta T$. Each effect must have substantially the same amount of surface as a single-effect evaporator. The steam required, however, goes down in inverse proportion to the number of effects. The choice of the proper number of effects involves an economic balance between the first cost of equipment and the continuing costs of steam and cooling water. The great majority of evap-

orator installations employ the multiple-effect principle and as many as a dozen effects have been used. Another means of achieving multiple-effect steam economy is the multistage flash cycle, which has been developed primarily for producing fresh water from seawater. See HEAT EXCHANGER; HEAT TRANSFER; SALINE WATER DESALINATION. Ferris C. Standiford Bibliography. J. A. Kent (ed.), *Kent and Riegel's Handbook of Industrial Chemistry and Biotechnology*, 11th ed., 2007; R. H. Perry and D. Green (eds.), *Perry's Chemical Engineer's Handbook*, 7th ed., 1997; R. Smith, *Chemical Process: Design and Integration*, 2d ed., 2005.

Evergreen plants

Plants that retain their green foliage throughout the year. Popularly, needle-leaved trees (pine, fir, juniper, spruce) and certain broad-leaved shrubs (rhododendron, laurel) are called evergreens. In warm regions many broad-leaved trees (magnolia, live oak) are evergreen, and in the tropics most trees are evergreen and nearly all have broad leaves. Many herbaceous biennials and perennials have basal rosettes with leaves close to the ground that remain green throughout the winter. See FOREST AND FORESTRY; LEAF; PLANT TAXONOMY. Nelle Ammons

Evolutionary computation

Evolutionary computation is a rapidly growing interdisciplinary science area that is concerned with modeling aspects of natural evolution in order to solve real-world problems. Living organisms, as well as those long extinct, demonstrate optimized complex behavior at all levels: cells, organs, individuals, and populations. Charles Darwin wrote of "organs of extreme perfection" when describing the ability of evolution to craft ingenious solutions to complex problems such as vision. Evolution is the great unifying principle of biology, but it extends beyond biology and can be used as an engineering principle where individuals in a population of candidate solutions to some particular problem undergo random variation (mutation and recombination) and face competition and selection based on their appropriateness for the task at hand.

Process. The common use of an evolutionary algorithm requires four elements: (1) an evaluation (fitness) function that describes the quality of any candidate solution in quantitative terms, (2) a representation or data structure that the computer uses to store solutions, (3) a random variation operator (or operators) that transform "parents" into "offspring," and (4) a means for selecting which solutions will survive to the next generation and which will be eliminated. In addition, the process must be initialized with a population of candidate solutions to the task at hand. This is often accomplished by seeding the first population with completely random solutions; however, if domain-specific knowledge is available

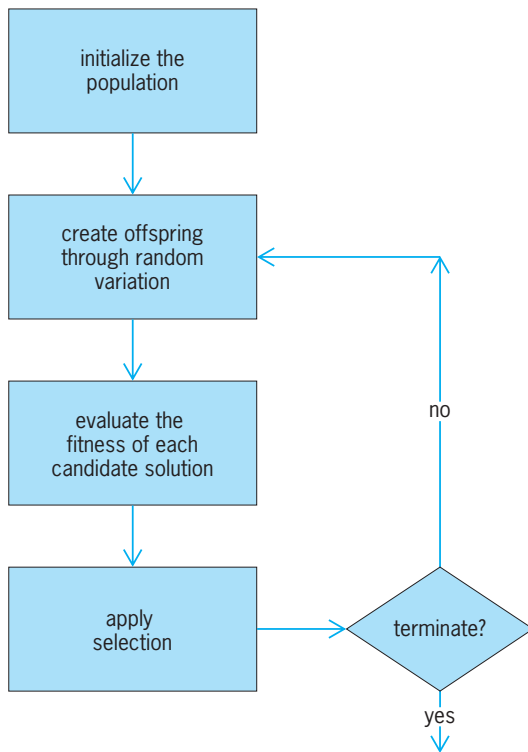


Fig. 1. Typical flowchart for an evolutionary algorithm.

regarding which solutions may be better than others, these hints can be used to bias the initial population and may accelerate the evolutionary optimization procedure.

An evolutionary algorithm (Fig. 1) is executed over a series of generations of random variation and selection. The variation to the existing solutions can come in the form of single-parent or multiparent operators. Alternative choices offer different sampling distributions from the space of all possible solutions. Each of the individuals in the population is scored with respect to how well the individual accomplishes the task at hand (fitness), and selection is used to eliminate some subset of the population or to amplify the percentage of above-average solutions. The algorithm terminates when some extrinsic criterion has been satisfied, such as prescribed maximum number of generations, or a suitable error tolerance. Over time, by eliminating poor solutions and extending the evolutionary search around those that appear better, the population can discover superior solutions to complex problems.

Origins. Although this idea of using evolutionary computation to solve problems is currently undergoing a renaissance that promises to have a significant impact on computing, its origins date back to the 1950s and 1960s. Some of the first experiments performed on John von Neumann's computer involved what would now be called artificial life. A program written for this machine by Nils Barricelli in 1953 simulated an environment, separated by cells in a grid. Numbers resided in each cell and migrated to neighboring cells based on a set of rules. When two numbers collided in the same cell, they competed

for survival. Even with very simple rules for propagating throughout the environment, certain numeric patterns would evolve and could persist only when other patterns were also present. These first experiments already demonstrated something akin to symbiosis.

The idea of simulating evolution on a computer arose as many as eight to ten times independently from 1953 to 1969 in diverse fields of study. Evolutionary algorithms were proposed in 1957 and again in the mid-1960s to simulate biological processes with the aim of understanding adaptation in nature. Another proposal (1960) was to use Darwin's ideas to generate artificial intelligence by evolving models that would predict future events in an environment. The use of evolutionary algorithms was also proposed (1962) for function optimization, and this led to work on the prospects of evolving neural networks. And in 1964, principles of evolution were used to search for improved designs for physical devices in fluid-mechanics problems (such as a supersonic flashing nozzle). The multiple simultaneous beginnings of evolutionary computation attest to the compelling nature of the idea of simulating the fundamental processes of evolution in fast time.

Putting evolution to work. As a stochastic search technique, an evolutionary algorithm is a means for generating useful solutions rather than perfect solutions. For example, a traveling salesperson problem comprising 1000 cities involved finding the minimum tour length required to cover each city once and only once and return to the starting position. The initial best solution required over 50,000 units to accomplish this tour, while the best found solution, after examination of 4×10^7 possible solutions, had a distance of just over 2000 units (Fig. 2a). This solution was not the best possible solution for this problem, but it was within about 5–7% of the expected minimum length for all 1000-city traveling salesperson problems. Moreover, the improvement in the quality of the solution was quite rapid in the initial phases of the evolutionary search (Fig. 2b). The size of the search space was over 10^{1200} (1 followed by 1200 zeros), indicating that the evolutionary algorithm had searched only a small fraction of all possible solutions while finding one of very high quality. Real-world applications of the traveling-salesperson problem include scheduling, the optimization of distribution networks, and computer chip layout and design.

Adaptability of technique. In comparison to traditional problem-solving techniques, evolutionary algorithms are often much faster and more adaptable to changes in the environment because whatever knowledge has been learned about how to solve a problem is contained in the collection of individual solutions that has survived up to that point. In contrast with, say, dynamic programming, an evolutionary algorithm need not be restarted when facing new data. This feature promises to have a significant potential in future problem solving. As the spread of data and processing speeds continue to increase,

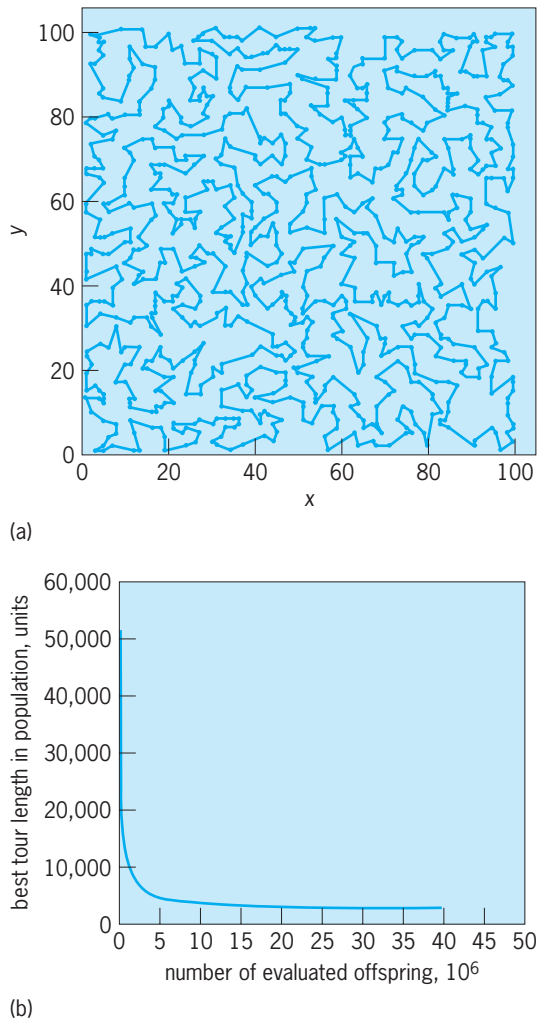


Fig. 2. Traveling salesperson problem comprising 1000 cities. (a) Best-evolved solution after examination of 4×10^7 possible solutions. (b) Rate of evolutionary optimization. (After D. B. Fogel, *Applying evolutionary programming to selected traveling salesman problems*, *Cybernet. Sys.*, 24:27–36, 1993)

it will be necessary to make decisions ever more quickly. The evolutionary approach of bootstrapping off the current basis set of knowledge may eventually become the standard approach to real-world, real-time problem solving.

Contrast with conventional approaches. The utility and applicability of evolutionary algorithms requires a shift in thinking that may require some time to overtake conventional engineering. The typical approach to problem solving is to think in terms of reducing complex situations into simpler elements (building blocks) and then working to optimize each of these components. The hope is that when brought back together, these optimized components will serve to make the best design. But it is well known that outstanding individuals do not necessarily make a great team. Evolution works in an entirely different manner. Only the entire cohesive functional unit is measured by selection. Owing to the complex effects of genetic interactions, where a single genetic change can affect myriad behavioral traits (a relationship

known as pleiotropy), evolution tinkers with the underlying genetics and simultaneously generates a diverse array of new behaviors. Selection then culls out the least appropriate designs.

Choice of evolutionary algorithms. Even though there have been some attempts to analyze evolutionary algorithms as a parallel to human engineering, focusing mainly on recombination as a mechanism for bringing together “good ideas,” there is no reason to expect that the evolutionary model of learning is the model that is used by engineers. Moreover, there is little reason to expect that it will be possible to improve evolutionary algorithms by capturing idealized models of genetic operators that occur in nature. The effects of those operators may be wholly different in the natural setting. A flapping wing will have different effects on a seagull in contrast to a jet aircraft. Airplanes can be designed to have increased biological fidelity by gluing feathers to their wings, but in the end all that will be accomplished is to increase their drag. What is termed the phenotypic effect, the manner of behavioral response to the environmental demands, is all that selection can assess. By consequence, in evolutionary algorithms, variation operators must be tailored to the task at hand in order to gain the most benefit and achieve the greatest efficiency.

This notion has been written down in formal mathematics in the “no free lunch theorem.” Within some overarching assumptions, all algorithms perform exactly the same on average when applied across all possible functions. The consequence of this result is that if an evolutionary algorithm is good for a certain class of problem, then it is not good (worse than a random search) on some other class. There is no single best construction for an evolutionary search (or any other search algorithm). Much of the early theory that supported using binary representations, one-point crossover, and proportional selection has now been shown to be either incorrect or of no practical value. What remains is a significant challenge for the future: Analyzing classes of optimization problems and determining useful means for identifying which types of evolutionary algorithms are appropriate in each case. This is a wide-open area of investigation and it promises to be quite difficult.

Prospects. Computing power has now caught up with the concept of evolving solutions to complex problems. Current real-world applications include supply chain optimization, screening for new drug leads, electronic circuit design, scheduling logistic and transportation services, and gaming.

The processing speed of modern desktop machines is rated at about 10^9 floating-point operations per second (flops). But evolution is an inherently parallel process in which potential solutions can often be evaluated simultaneously rather than sequentially. In these cases, a parallel distributed computing architecture can offer a tremendous speed-up, since many more trials can be conducted per unit of time. An evolutionary algorithm can be implemented with multiple populations on independent processors,

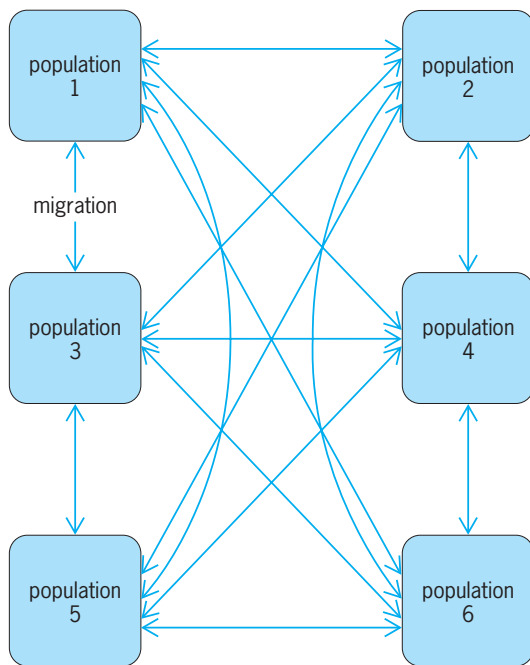


Fig. 3. Arrangement for distributing an evolutionary algorithm across multiple processors. Each processor holds its own population that undergoes random variation and selection. Occasionally, solutions will migrate from one population to another and begin competing for survival in their new environment.

with occasional migration being applied to introduce solutions from different populations (Fig. 3). In this manner, the parallel nature of evolution can be captured and many more trials can be conducted per unit of time. The best manner in which to set up such a distributed network of populations is problem-dependent and is an open question for research.

Experiments are being conducted with distributed architectures that comprise 1000 350-MHz computers. This computing power is roughly what can be expected for desktop machines around 2010. The products that will emerge from these massively parallel designs remain a matter of speculation; some foresee the use of evolutionary algorithms to design surrogate artificial brains that will supplant human cognition. Whether or not this conjecture becomes reality, with the price of desktop computers continuing to fall, the parallel approach to supercomputing is increasingly attractive and appears to be perfectly suited to evolutionary problem solving. See ALGORITHM; CONCURRENT PROCESSING; DISTRIBUTED SYSTEMS (COMPUTERS); GENETIC ALGORITHMS; MULTIPROCESSING; OPTIMIZATION; ORGANIC EVOLUTION; SUPERCOMPUTER. David B. Fogel

Bibliography. T. Bäck, *Evolutionary Algorithms in Theory and Practice*, Oxford University Press, New York, 1997; D. B. Fogel (ed.), *Evolutionary Computation: The Fossil Record*, IEEE Press, Piscataway, NJ, 1998; Z. Michalewicz and D. B. Fogel, *How To Solve It: Modern Heuristics*, Springer-Verlag, Berlin, 2000.

Exchange interaction

A consequence of the quantum mechanics of interacting identical particles that is responsible for magnetic interactions between electrons and atoms in solids, and magnetic order such as ferromagnetism and antiferromagnetism. The exchange interaction is actually not one of the fundamental interactions such as electromagnetism, but an effect that arises from the interplay of electromagnetism with quantum mechanics. See ELECTROMAGNETISM.

Quantum mechanics describes the state of a system of particles by a wave function that depends both on the positions of the particles and on any additional internal degrees of freedom needed to describe them (such as the spin of an electron). A fundamental axiom of quantum mechanics states that elementary particles fall into two classes, fermions and bosons. If two identical particles are interchanged, the wave function describing the state of the system is unchanged if they are bosons, and merely changes sign if (like electrons) they are fermions. See ELECTRON; ELEMENTARY PARTICLE; QUANTUM STATISTICS; SPIN (QUANTUM MECHANICS).

This antisymmetry of the wave function describing a many-electron state means that if the spins of two electrons are parallel the wave function must vanish when their positions coincide, making such electrons less likely to be found close together than a pair with antiparallel spins. This effect lowers the mutual electrostatic energy of electrons with parallel spins. Since electrons have a magnetic moment parallel to their spin, this effect helps explain the occurrence of magnetic moments in transition-metal and rare-earth atoms with a partially filled electronic shell (intraatomic exchange). This mechanism is used to explain the origin of ferromagnetism in metals such as iron. See ATOMIC STRUCTURE AND SPECTRA; ELECTRON CONFIGURATION; ELECTRON SPIN.

Similar considerations hold in the case of interatomic exchange of electrons between two such magnetic atoms in close contact, as in a magnetic solid. However, whether ferromagnetic (parallel) or antiferromagnetic (antiparallel) alignment of the two magnetic moments is favored now depends on the details of the most likely path by which exchange of electrons by quantum-mechanical tunneling between the atoms occurs. This path depends on the structure of the atoms and the solid, and may involve processes on nonmagnetic atoms in between the magnetic atoms. Both ferromagnetic and antiferromagnetic exchange interactions occur, though the latter are more common. Such interatomic exchange processes are used to explain the origin of magnetic order in nonmetallic magnetic materials. See TUNNELING IN SOLIDS.

Magnetic atoms also have a direct magnetic dipole interaction (like that between bar magnets) which can be understood purely in terms of classical electromagnetism, without invoking quantum mechanics. However, these interactions between neighboring atoms are much weaker than the exchange interaction and are not responsible for the

occurrence of magnetic order. Since the exchange interaction involves a quantum-mechanical tunneling process, its strength falls off rapidly (exponentially) with distance between atoms, and the long-range dipole interactions dominate at greater distances, controlling the large-scale magnetic-domain structure of ferromagnets. See ANTIFERROMAGNETISM; DIPOLE-DIPOLE INTERACTION; FERROMAGNETISM; QUANTUM MECHANICS.

F. Duncan; M. Haldane

Bibliography. D. C. Mattis, *Theory of Magnetism* vol. 1, 1981, reprint 1988; G. T. Rado and H. Suhl (eds.), *Magnetism*, vol. 1, 1963, vol. 4, 1966; R. M. White, *Quantum Theory of Magnetism*, 2d ed., 1983.

Excitation

Any of a number of different phenomena which (a) alter a system in some way or (b) result in some type of response. In electrical network theory, an excitation initiates a response or a response sequence. In other areas of technology, an excitation establishes an altered condition which causes an apparatus or system to exhibit useful response capabilities.

In electrical network theory, the term excitation designates a time-varying independent voltage or current in an n -port system or network. The independent voltage or current—which is also referred to as the excitation or the excitation signal—causes a network response which affects the voltage or current at the dependent port. The nature of the network response, and therefore of the dependent voltage or current, derives from the characteristics of the network. Thus an excitation signal at the input of a bistable multivibrator causes the output to shift from one of the circuit's stable states to the other; the next excitation signal causes the output to return to its former state. See MULTIVIBRATOR; NETWORK THEORY.

In atomic physics, excitation means the addition of energy to an atom at ground state to produce an excited state. See ATOMIC STRUCTURE AND SPECTRA; EXCITATION POTENTIAL.

In other contexts, excitation means the application of energy to one portion of a system or apparatus in a manner that enables another portion to perform a specialized function. Excitation energy may differ from the output energy in source, form, level, or location. That is, an excitation produces a primary effect that is linked, through an intermediate physical phenomenon, to a dependent secondary effect (see *illus.*).

Some examples will make this relationship apparent: A dynamic loudspeaker uses an excitation current in a field coil to generate a magnetic field; only then can a second magnetic field, generated by an audio signal, actuate the voice cone and produce sound waves. The sound-producing portion of a motion picture projector uses a lamp to provide excitation illumination; as the light passes through the sound-track portion of the moving film, its intensity is

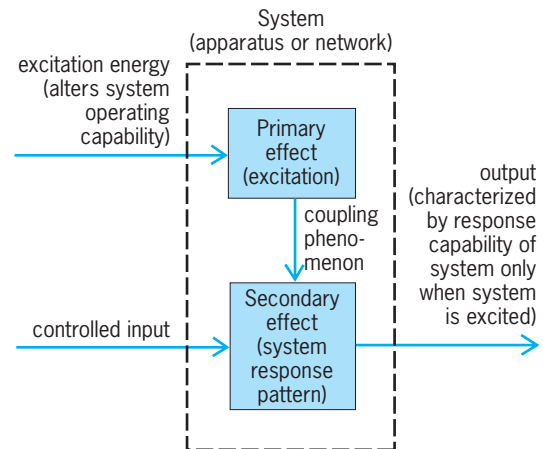


Diagram showing function of excitation.

modulated, producing audio-frequency fluctuations, which the photosensitive pickup converts into an electronic audio signal. See LOUDSPEAKER; OPTICAL RECORDING.

In some cases, the application of excitation energy is all that is necessary to generate the secondary response. A fluorescent lamp, for example, uses an excitatory current passing through an enclosed gas to generate ultraviolet radiation; when the phosphor coating inside the tube absorbs ultraviolet energy, it emits visible light. The human heart, to cite an example from biology, uses for excitation the decreasing electrical transmembrane potential of the pacemaker to generate an electrical impulse; when the surrounding striated muscle tissue receives this impulse, it transmits the electrical stimulus to other adjacent tissue and utilizes metabolic energy to contract, a response which produces the heartbeat. See CARDIOVASCULAR SYSTEM; FLUORESCENT LAMP.

William W. Snow

Excitation potential

The difference in potential between an excited atomic or molecular state and the ground state. The term is most generally used in connection with electron excitation, but it can be applied to excited molecular vibrational and rotational states.

A closely related term is excitation energy. If the unit of potential is taken as the volt and the unit of energy as the electronvolt, then the two are numerically equal. According to the Bohr theory, there is a relationship between the wavelength of the photon associated with the transition and the excitation energies of the two states. Thus the basic equation for the emission or absorption of energy is as shown below, where h is Planck's constant, c the

$$\frac{hc}{\lambda} = E_i - E_f$$

velocity of light, λ the wavelength of the photon, and E_f and E_i the energies of the final and initial states, respectively. If the final state is the ground state, then the difference in energies is just the excitation energy of the initial state. If neither of the two states

is the ground state, then the difference is numerically equal to the difference in excitation energies of the two states. This suggests that the excitation energy for many states may be determined spectroscopically. In fact, the careful measurement of the wavelengths associated with the transitions, along with the identification of the levels involved, permits the assignment of an excitation potential to each available energy state. This may be accomplished by either emission or absorption spectra. *See* ATOMIC STRUCTURE AND SPECTRA; ELECTRONVOLT.

Measurements have also been made by using electron impact means. Here electrons from a hot filament are accelerated by a grid and collected on an outer electrode after passing through an intervening space where they can interact with a gas of molecular or atomic species. The potentials are so adjusted that the electrons can reach the collector only if they lose no energy after passing the grid. As the accelerating grid potential is increased, a series of sharp drops is noted in the collector current. Each of these is interpreted as representing a case where the electrons have obtained just enough energy to produce excitation, and hence no longer have the energy needed to climb the potential barrier to the collector. The evidence thus obtained is very direct, but the accuracy is not as good as that for spectroscopic data. However, this is the method that was used in the historic Franck-Hertz experiment, which demonstrates the quantization of atomic energy levels. *See* GROUND STATE; IONIZATION POTENTIAL.

Glenn H. Miller

Bibliography. S. Gasiorowicz, *Structure of Matter: A Survey of Modern Physics*, 1979; I. I. Sobel'man, L. A. Vainshtein, and E. A. Yukov, *Excitation of Atoms and Broadening of Spectral Lines*, 2d ed., 1995.

Excited state

In quantum mechanics, a stationary state of higher energy than the lowest stationary state or ground state of a particle or a system of particles. Customarily, only bound stationary states, which generally are at most denumerably infinite in number, are spoken of as excited, although the formal quantum theory often treats the noncountable unbound stationary states on an equal footing with the bound states. Conventionally, the excited states are ranked in order of increasing energy; that is, the second excited state has higher energy than the first, which lies higher than the zeroth or ground state.

For actual physical systems it is usually, but by no means necessarily, the case that energy levels of high degeneracy (that is, energy levels associated with a comparatively large number of stationary states of the same energy) are encountered in excited states rather than in ground states. Illustrative systems possessing highly degenerate ground states include the 10-fold degenerate germanium-72 nucleus (32 protons, 41 neutrons) and the 13-fold degenerate uranium atom (92 electrons in the field of a massive positive point charge equal to 92 e , where e

is Sobel'man charge). *See* DEGENERACY (QUANTUM MECHANICS); GROUND STATE; METASTABLE STATE; NONRELATIVISTIC QUANTUM THEORY; QUANTUM MECHANICS.

Edward Gerjuoy

Bibliography. S. Gasiorowicz, *Quantum Physics*, 2d ed., 1995; R. B. Leighton, *Principles of Modern Physics*, 1959.

Exciton

An excited electronic state in a large molecule, or a semiconducting or insulating solid composed of a negatively charged electron and a positively charged hole bound together by their electrical attraction. Its creation through internal charge separation is most frequently caused by the absorption of light, and its demise is occasioned by the emission of light. The exciton is thus an important portal for the study of the properties of an insulator by optical means and for applications such as optoelectronics. Its external neutrality enables it to move through the system as an energy carrier, an important agent in energy transfer processes of special importance in biological functions such as photosynthesis. When the excited electron and hole are restricted to a plane, a wire, or a microscopic box in an artificially fabricated structure known, respectively, as a quantum well, quantum wire, or quantum dot, a new dimension is added to the ability to explore fundamental nonlinear and nonequilibrium physics and coherent quantum optics as well as quantum control of microscopic objects. *See* ELECTRON; HOLE STATES IN SOLIDS; PHOTOCHEMISTRY; PHOTOSYNTHESIS; QUANTIZED ELECTRONIC STRUCTURE (QUEST).

Creation and recombination. The many-electron ground state of the insulator is immune to excitation until the excitation energy reaches a threshold known as the energy gap. When external influences such as the electromagnetic field and lattice vibrations are ignored, the exciton may be viewed as a robust state of an excited electron and the hole which it has left behind in the sea of electrons. The hole acquires its positive charge from the loss of a charge from the sea whose total charge is neutralized by that of the ions in the molecule or solid. *See* BAND THEORY OF SOLIDS; ELECTRIC INSULATOR.

The photon-exciton interaction is responsible for the optical excitation (though not necessarily in the visible frequency range) of the exciton and for its spontaneous recombination emitting a photon (a quantum unit of light). The dipole matrix element responsible for the transition between the energy states is strong when the electron and hole wave functions overlap in space or match in wave vector. From Planck's law in quantum theory, the frequency of the emitting light is proportional to the energy loss in returning the exciton state back to the ground state. *See* ELECTRON-HOLE RECOMBINATION; NONRELATIVISTIC QUANTUM THEORY; QUANTUM MECHANICS.

Excitons in semiconductors. A semiconductor is distinguished from an insulator qualitatively by a smaller

energy gap, with the frequency of the emitting light from the exciton spanning the range from visible light to very far infrared. Because of the higher mobility of their excited electrons and holes, semiconductors are of supreme technological importance. For example, the electrons and holes can be electrically injected to form excitons whose recombination leading to light emission enables the conversion of electrical energy to light in such devices as light-emitting diodes and lasers. *See* LASER; LIGHT-EMITTING DIODE; SEMICONDUCTOR.

In addition to the direct “resonant” formation by the absorption of a photon (in which the photon energy equals the exciton formation energy), the exciton is commonly formed by the photoluminescence process during which the optically excited electron and hole with energy much higher than the exciton energy relax to the lowest states to form an exciton. The recombination of the exciton emitting a photon then completes the luminescence process. *See* LUMINESCENCE.

Indirect excitons. In an indirect-gap semiconductor such as silicon, where the momentum of the exciton does not match that of the photon, the excitons are generally formed after relaxation of optical excitations with initial energy much higher than the gap. The indirect exciton has a long lifetime because its recombination with the emission of a photon requires the conservation of momentum to be satisfied by the assistance of a lattice vibration or trapping by a defect. Consequently, the excitons have time to form a large pool known as an electron-hole drop.

In a direct-gap material, the spatial separation of the electron and hole can be enforced by housing them in two layers sufficiently close to maintain their electric attraction. The recombination of such indirect excitons may then be controlled by changing the separation with an electric field. An interesting phenomenon is the laser spot excitation of these indirect excitons, leading to their migration away from the laser spot as optically inactive excitons for long distances, some radiating in a bright spot when trapped by a defect, and others surviving the traps to form two concentric luminous circles centered at the laser spot. The outer one is a necklace of evenly spaced bright spots, thought to be a lattice of vortices. Since L. V. Keldysh predicted the condensation of excitons as bosons, the observation of exciton condensation has become a highly prized objective of research in the field. This system of excitons in a double quantum well is considered a good candidate for condensation at low temperature even though the necklace has a noncondensation origin. *See* BOSE-EINSTEIN CONDENSATION; TRAPS IN SOLIDS.

Frenkel and Wannier excitons. If the constituent electron and hole of the exciton are mostly localized at an ion, the exciton is localized, but with some probability to hop from site to site. Such a Frenkel exciton is common in molecules and molecular solids. At the other extreme, if the electron and hole wave functions are widespread as extended orbitals in a molecule or Bloch waves in a crystal, their bound

state as the exciton can have their center of mass moving through the system with ease. Such Wannier (or Wannier-Mott) excitons are most common in broad-band and small-gap semiconductors. They resemble the hydrogen atom or, more closely, the positronium system composed of an electron and a positron. Because of the dielectric screening of the electrical force in small-gap materials and sometimes the small effective mass of the electron, the Wannier exciton radius is from 10 to 100 times larger than the positronium radius of about 0.1 nm. The absorption spectrum of excitons also resembles that of a hydrogen atom, with a series of lines with rapidly decreasing spacing finally blurring into a continuum whose energy edge is theoretically the energy gap but practically difficult to discern. Such optical spectra are usually obtained at the temperature of liquid helium (around 2 to 4 K or -456 to -452°F) to avoid the broadening of spectral lines by lattice vibrations. *See* ATOMIC STRUCTURE AND SPECTRA; BLOCH THEORY; POSITRONIUM.

Exciton polaritons. Because of the interaction between electrons and light, not only can an exciton decay irreversibly into a photon or vice versa, but it can also exchange roles with the photon in a quantum-mechanically coherent fashion. Thus, the exciton may exist in the solid as a quantum superposition of an exciton and a photon, known as a polariton. Since the photon energy varies linearly with its momentum at the speed of light, the exciton at most values of its center-of-mass momentum has much too low an energy to mix with the photon. However, at low momentum the exciton and the photon can approximately match both their momentum and energy values, the coupling mixes the two states into two superpositions of photon and exciton with an energy splitting. Thus, the massless photon is slowed down by the massive exciton by the quantum-mechanical superposition—a version of slow light. *See* PHOTON.

Microcavities. Confinement of light in a microcavity enhances the coupling of the exciton with the confined light mode. It was discovered that two-dimensional exciton polaritons, when pumped resonantly and coherently, can be scattered into a low-energy mode where they accumulate to a macroscopic (Avogadro) number, causing a high-intensity luminescence. This phenomenon has been interpreted as an optically maintained boson condensation state. Experiments have shown that when quantum dots are planted in a nanocavity, their excitons couple strongly with cavity photons. This strong coupling is a promising basis for cavity quantum electrodynamics in a semiconductor for the proposed quantum control of excitons via photons. *See* NANOSTRUCTURE; QUANTUM ELECTRODYNAMICS.

Biexcitons. For the same reason that two hydrogen atoms can bind into a molecule, two excitons can bind into a biexciton. A biexciton takes two photons to create and can spontaneously emit two photons, usually one at a time. Nonlinear response of a semiconductor to three electric fields may be viewed as numerous replicas of the event with three photons,

some incoming and some outgoing, each supplied by an electric field. The biexciton created in the process plays a central role. Resonance of field frequencies with exciton or biexciton energy enhances the response. Not only has nonlinear optics flourished with the time-resolved measurement of dynamics, first at the picosecond and then at the femtosecond scale, but also the coherence of the laser light has been matched with the quantum-mechanical coherence of the exciton, giving rise to the study of the coherent optics of quantum properties in the solid-state system. The coherence of the laser light derives from the phase relation between different parts of the electromagnetic wave and the coherence of the exciton from the superposition of its constituent states. The delicate four-wave mixing has yielded information that extends beyond the structure of the excitations of the excitons to their dynamics far from thermal equilibrium. The exploration of the exciton dynamics of the first few femtoseconds after ultrafast intense laser excitation has contributed profoundly to our knowledge of nonequilibrium physics. *See* COHERENCE; NONLINEAR OPTICS.

Trions. Excitons can also form complexes. When an exciton is tethered to a positively charged impurity ion, the electron system in a donor is a bound state of two electrons and a hole, known as a negatively charged trion. Since the electron-hole pair in the trion may be coupled to light, light may be used to control the dynamics of an electron through the trion in an electronic Raman process in which no phonons need be involved. *See* DONOR ATOM; PHONON; RAMAN EFFECT.

Exciton spin and light polarization. The exciton possesses the quintessential quantum-mechanical property of an intrinsic angular momentum, known as spin, of one fundamental quantum unit \hbar (Planck's constant divided by 2π), which matches the intrinsic angular momentum of the photon perfectly. Now the plus-one- and minus-one-quantum-of-angular-momentum states of the photon are the building blocks, respectively, of right- and left-handed circularly polarized light. Angular momentum conservation in the photon-exciton conversion process means that it is possible to "polarize" excitons, all with their spins pointing along (or opposite) to the light propagation direction by excitation with right- (or left-) circularly polarized light. The electron constituent of the exciton has spin one-half quantum unit, which points in the opposite direction to the optically generated exciton. (Angular momentum conservation is maintained by the hole contribution of $3/2$ units from its orbital and spin angular momenta and no orbital angular momentum from the electron.) Thus, the optical excitation of excitons by circularly polarized light provides a good source of spin-polarized electrons (that is, with spin pointing mostly in one direction) when the electron and hole constituents of an exciton are separated by an electric field, for example, in a *pn* diode. Time-resolved photoluminescence by polarized light is an excellent probe of the spin dynamics of the excitons. *See* ANGULAR MOMENTUM; ELECTRON SPIN; POLAR-

IZATION OF WAVES; POLARIZED LIGHT; SPIN (QUANTUM MECHANICS).

Quantum wells. Forcing the electron and hole together by confinement increases their attraction and makes the exciton binding stronger. The exciton binding energy is the separation of its energy from the energy continuum of the unbound electron and hole pair. With the advent of molecular-beam epitaxy, it is possible to build a sandwich of semiconductors of different band gaps so that the electron and hole can be effectively confined in a quantum well. The exciton whose constituents are confined to move effectively in a plane can have as much as four times the binding energy of its counterpart in three dimensions. This increase is sufficient to make the exciton line visible in the optical spectrum at room temperature, whereas the exciton line in three dimensions would be indistinguishable from the continuum due to thermal broadening. For the same reason, the biexciton also has a larger binding energy in quantum wells. The study of exciton dynamics and spin dynamics flourishes in quantum wells and other confined structures. The exciton in the quantum well is also more susceptible to the influence of an external electric field. The modification of the exciton spectrum by an electric field has been applied to self-electrooptic-effect devices. *See* ARTIFICIALLY LAYERED STRUCTURES; ELECTROOPTICS; SEMICONDUCTOR HETEROSTRUCTURES.

Quantum dots. A semiconductor quantum dot of submicrometer dimensions serves as a trap of a single electron or an exciton. It behaves like a giant atom with the advantages of easier electrical and magnetic control because of the lower energy spacings and larger dipole transition matrix elements than real atoms. However, it also has the concomitant disadvantage of shorter radiative recombination time and decoherence time. It is possible to study the dots in ensemble or singly. In the method of microluminescence, by suitable masking to expose only a small area to light, it is possible to obtain the spectra of excitons, biexcitons, and trions in a single quantum dot from its luminescence. (The trion is generated by exciting an exciton in a dot doped with a single electron.) Similarly, coherent nonlinear spectroscopy can be performed by optically pumping and probing a single dot. This spectroscopy opens up the possibility of quantum control of a microscopic state such as the exciton. Just like coherent quantum optics in atoms, the coherent driving of the Rabi rotation between the ground state and the exciton has been demonstrated by several groups. The Rabi rotation is a change of the relative proportion of two states in a quantum superposition by a coherent electromagnetic pulse. This rotation stops when the pulse stops, a deterministic event to be distinguished from the probabilistic excitation from one state to another that is governed by the Fermi golden rule when the different frequency components of the electromagnetic field are not phase-coherent.

Quantum control. The demonstration of Rabi rotation is the first step toward operations on a quantum unit of information (qubit). The biexciton in a dot is

bound in the sense that it has a lower energy than the sum of energies of its constituent excitons when they have antiparallel spins. The biexciton can be coherently formed by driving the dot with two laser pulses with opposite circular polarizations. One might then imagine the absence and presence of each exciton as a qubit, thus forming a two-qubit system in the dot. The difference in energy between adding an exciton of spin up, say, to the ground state or to the single-exciton state of the opposite spin direction is key to the ability to drive the Rabi rotation between two double-digit states (0,0) and (1,0) or between (0,1) and (1,1) with light of the same polarization but a choice of different frequencies. Either rotation is a conditional dynamics (control of one exciton depending on the second exciton being absent or present). That is sufficient for the entanglement of the two excitons, which is the fundamental resource for quantum information processing. Such processing has indeed been demonstrated experimentally in a quantum dot at 4 K (−452°F). It is remarkable that it could be done in a dot that was housed in a solid-state system full of potentially disruptive decoherence sources, quite unlike the trapped-ion system in which such a feat had been demonstrated earlier. Housing two exciton qubits in a dot is not a model on which one can build a similarly controllable system of an arbitrary number of qubits. This two-qubit model is said to be “unscalable,” and thus not directly useful for the purpose of building a practical quantum computer. However, the demonstration of the optical control of an exciton in a single dot makes it possible to imagine the optical control (via Raman processes) of an electron in a dot whose two spin states constitute the qubit and of two electrons in separate dots, scalable to a large system. Thus, the exploration of quantum control of microscopic objects, now being electron spins residing in separate dots, can continue. *See* QUANTUM COMPUTATION.

Lu Jeu Sham

Bibliography. R. S. Knox, *Theory of Excitons (Solid State Physics, Advances in Research and Applications, Supplement 5)*, Academic Press, New York, 1963; E. I. Rashba and M. D. Sturge (eds.), *Excitons*, North-Holland, Amsterdam, 1982; W. Schäfer and M. Wegener, *Semiconductor Optics and Transport Phenomena*, Springer, Berlin, 2002; J. Shah, *Ultrafast Spectroscopy of Semiconductors and Semiconductor Nanostructures*, 2d ed., Springer, Berlin, 1999; P. Y. Yu and M. Cardona, *Fundamentals of Semiconductors*, 3d ed., Springer, Berlin, 2004.

Exclusion principle

No two electrons may simultaneously occupy the same quantum state. This principle, often called the Pauli principle, was first formulated by Wolfgang Pauli in 1925 and, for time-independent quantum states, it means that no two electrons may be described by state functions which are characterized by exactly the same quantum numbers. In addition to electrons, all known particles having half-integer

intrinsic angular momentum, or spin, obey the exclusion principle. It plays a central role in the understanding of many diverse phenomena, including the periodic table of the elements and their chemical activities, the electron contribution to the specific heat of metals, the shell structure in the atomic nucleus analogous to that of electrons in atoms, and certain symmetries in the scattering of identical particles. *See* ANGULAR MOMENTUM; NONRELATIVISTIC QUANTUM THEORY; QUANTUM NUMBERS; SPIN (QUANTUM MECHANICS).

Symmetry of wave functions. When a system of identical particles is described by a wave function, the indistinguishability of the particles implies that the wave function must have certain symmetry properties when the coordinates of any two of the particles are interchanged. Specifically, the wave function either remains unchanged or is changed only in sign when such a coordinate interchange is made. For the two-particle system, Eq. (1) holds, where \mathbf{x}_i denotes

$$\psi(\mathbf{x}_1\mathbf{x}_2) = \pm\psi(\mathbf{x}_2\mathbf{x}_1) \quad (1)$$

all the coordinates of the particle i , such as space coordinates, spin, and any others necessary, and ψ is the wave function.

When the plus sign applies, the wave function is said to be symmetric; the particles then are termed bosons and obey Bose-Einstein statistics. When the minus sign applies, the wave function is antisymmetric under particle exchange; the particles then are termed fermions and obey Fermi-Dirac statistics. Bosons are particles with integer spin, while fermions have half-integer spin, measured in units of Planck's constant. This connection between spin and statistical laws obeyed by the particles may be proved mathematically under appropriate assumptions. *See* BOSE-EINSTEIN STATISTICS; FERMI-DIRAC STATISTICS.

The Pauli principle follows from the symmetry of the many-fermion wave function in the special case when this wave function may be written as a product of single-particle wave functions. For example, consider two identical fermions and two single-particle states. Let $\psi_\alpha(1)$ denote particle 1 in state α , $\psi_\beta(2)$ particle 2 in state β , and so on, where α and β represent all quantum numbers labeling the states. Then a possible state of the two-particle system is shown by Eq. (2), which obviously changes sign if 1 and

$$\psi(1, 2) = \psi_\alpha(1)\psi_\beta(2) - \psi_\alpha(2)\psi_\beta(1) \quad (2)$$

2 are interchanged. If $\alpha = \beta$, then ψ is identically zero; that is, no such state exists which is exactly the statement of the Pauli principle. Note that ψ may be expressed as a determinant, as in Eq. (3).

$$\psi = \begin{vmatrix} \psi_\alpha(1) & \psi_\beta(1) \\ \psi_\alpha(2) & \psi_\beta(2) \end{vmatrix} \quad (3)$$

For more than two particles, one may still form product wave functions which are also determinants. The exchange of a pair of particles is equivalent to the interchange of two rows of the determinant which is known to change its sign. Furthermore,

if any two of the column labels, the α and β above, are the same, the determinant vanishes identically, which again is the Pauli principle. The most general wave function is a linear combination of all possible determinant wave functions, and the antisymmetry still holds, but the Pauli principle does not, since it applies only if there is a one-to-one correspondence between the number of single-particle states and the number of particles.

Spin. In the example above a single function was used to denote a single particle state. In many physical applications it is more convenient to express the wave function of a single particle as a product of a function which describes its spatial properties and one which describes its spin orientation. For spin s there are $2s + 1$ possible spin orientations ranging from $-s$ to s in integer steps. Electrons, protons, and neutrons all have $s = 1/2$. The wave function for a collection of identical particles may then also be expressed as the product of a spatial wave function for the collection times a spin wave function for the collection. For the case of two particles, the Pauli principle dictates that, if the spatial wave function is symmetric, the spin function must be antisymmetric and vice versa. For more than two identical fermions, mixed symmetries in each part may occur, but they must be complementary in the sense that the overall wave function obeys the Pauli principle.

Isotopic spin. It is often convenient to regard protons and neutrons as simply two different states of the same fermion, termed a nucleon. This approach takes into account that the proton is charged and the neutron is not, but ignores the small mass excess of the neutron. The additional quantum number needed to describe the two possible charge states is called isobaric spin. In this case the wave function of a single nucleon may be taken to be the product of a spatial part, an ordinary spin-orientation part, and a part to describe its charge state, that is, isotopic spin. The ideas about the symmetry apply similarly for a collection of nucleons. In this case one speaks of the generalized Pauli principle, in that the overall wave function is antisymmetric. For two nucleons this is attained by choosing the symmetries of the parts of the wave function as indicated in the **table**, where $+$ denotes symmetric and $-$ antisymmetric. Again, for more than two nucleons mixed symmetries occur, but the overall wave function obeys the generalized Pauli principle. See I-SPIN; SUPERMULTIPLY.

Atomic structure. One important consequence of the exclusion principle is that, if there are N different single-particle states corresponding to a given energy level, then, at most, N identical fermions may

have that same energy. For example, in the atomic case the single-electron states have allowed energies labeled by the quantum number $n = 1, 2, 3, \dots$ corresponding to the lowest energy, next lowest, and so on. For a given value of n , the orbital angular momentum l may have any one of the values $l = 0, 1, 2, \dots, n - 1$. For each l value there are $2l + 1$ states labeled by m_l , which range in integer steps from $-l$ to l , and specify the orbit orientation. In addition, there are the two possible spin orientations, $\pm 1/2$. For $n = 1$, then, there are two different single-particle states, while for $n = 2$, there are eight. In general, $2n^2$ is the maximum number of electrons that may have an energy specified by n .

Using the fact that a system will try to occupy the state of lowest possible energy, the electron configuration of atoms may be understood by simply filling the single-particle energy levels according to the Pauli principle. This is the basis of Niels Bohr's explanation of the periodic table. For example, the one electron of hydrogen goes into the $n = 1$ level, as do both of the electrons of helium. The $n = 1$ level is said to be full. The next element in the periodic table is lithium, and its third electron must go into the $n = 2$ level. For neon, of atomic number 10, two of its 10 electrons go into the $n = 1$ level and eight into the $n = 2$ level, so that both are full. Now, the chemical activity of an element is determined mainly by its outermost (largest n) electrons, and a filled level is essentially inert. One sees therefore why helium and neon are noble gases and have very low chemical activity, whereas lithium, which has one more electron than the very stable helium configuration, is very active chemically. By extending these ideas, such phenomena as binding energies and valences may be understood. Even if the Pauli principle applied only to this one area, it would still be a very important contribution to physics. See ATOMIC STRUCTURE AND SPECTRA; VALENCE. S. A. Williams

Bibliography. M. Chester, *Primer of Quantum Mechanics*, 1987, reprint 1992; R. Eisberg and R. Resnick, *Quantum Physics of Atoms, Molecules, Solids, Nuclei and Particles*, 2d ed., 1985; M. Jammer, *The Conceptual Development of Quantum Mechanics*, 1966, reprint 1989; D. A. Park, *Introduction to Quantum Theory*, 3d ed., 1992; R. L. Sproull and W. A. Phillips, *Modern Physics*, 3d ed., 1980, reprint 1989; R. F. Streater and A. S. Wightman, *PCT, Spin and Statistics and All That*, 1964, reprint 1980.

Excretion

The removal of excess material from a cell or a living animal. Such a definition is so broad that it is necessary to set some limits to the discussion that follows. The first limitation will be to the animal kingdom since excretion appears not to be a major problem in the plant kingdom—partly because of the difference in physiology.

Carbon dioxide is one of the chief products of the metabolism of cells. In the dark, plant cells

Choice of symmetries for two nucleons

Space	Spin	Isotopic spin
+	+	-
+	-	+
-	+	+
-	-	-

consume oxygen and produce carbon dioxide, but in the light, the actively synthesizing plant cells require a source of carbon dioxide and produce oxygen in excess. Gaseous diffusion over relatively short distances meets these respiratory problems in plants and very small animals, but in larger animals special mechanisms have evolved for the transport of carbon dioxide in blood. The excretion of carbon dioxide might properly be discussed here, but it has been found more satisfactory in the science of physiology to treat carbon dioxide and oxygen together under the heading of blood-gas transport. This case is cited to exemplify the assignment of excretion of heat to body-temperature regulation and of salts to osmoregulation. *See* PHOTOSYNTHESIS; RESPIRATION; THERMOREGULATION.

Physiological Problems

Though animals exhibit a wide variety of excretory structures, depending upon their individual needs, the excretory mechanisms are employed to overcome certain basic physiological problems.

Water and salts. One of the central problems of excretory physiology is the handling of water. To understand the role of water, it is important to speculate briefly about the very origins of life. Current views assume life to have begun in the form of self-reproducing molecules which grew by incorporating the already formed smaller molecules from the surrounding medium into their ordered structure. But, as essential ingredients in this medium diminished in quantity, self-reproduction became dependent upon a supply of the necessary material, obtained perhaps by modification of a closely related molecule. A primitive metabolism can thus be considered to have evolved. It would become essential to slow the diffusion processes which would inevitably carry away so valuable a material. This could be accomplished by the development of a relatively impermeable membrane from the many molecules which would be absorbed onto the surface of a complex molecule of this sort. But when a semipermeable membrane is placed around a quantity of organic matter osmotic pressure develops, and the primitive cell faces its first problem of excretion, namely, getting rid of the water which osmosis continuously brings into it. One answer to this problem was that developed by plant cells. By surrounding the protoplasm with an extremely strong wall of cellulose, osmotic pressures of many atmospheres could be allowed to develop. The cell walls prevent the cell rupture that would otherwise follow. Whole plants take advantage of differences in osmotic pressure between external cells and internal cells to move water along the roots and stems. The turgor of cells, petioles, and leaves which develops as a result of osmotic pressure is also used to close stomata or to alter the position of a leaf. Development of cellulose cell walls has therefore been made useful in manifold ways. *See* OSMOREGULATORY MECHANISMS; PREBIOTIC ORGANIC SYNTHESIS.

The planet provides an array of habitats—major ones being marine, brackish water, fresh water, and

terrestrial—and animals have established themselves successfully in each one. Particular osmotic problems are offered by each habitat, and the evolution of means of meeting these problems will be given brief consideration in phyletic sequence.

Nitrogen. Carbohydrate and fat are metabolized to carbon dioxide and water, and no other major excretory problems result from their use. Protein is in a different category. The amino acids used as an energy source must first be freed of nitrogen (deaminized), and the ammonia which is the first result of this process immediately constitutes a formidable excretory problem. Carnivorous animals of necessity have a diet high in protein, and omnivorous animals may prefer a high protein intake. One additional source of nitrogenous wastes, the nucleic acids, should be put in perspective. Each cell contains a small portion of this material in its nucleus so that purines may be present in the body in excess. An entire section below is devoted to the problem of nitrogen excretion. *See* NUCLEIC ACID; PROTEIN METABOLISM.

Other compounds. The kidneys of various animals have had to deal with compounds of considerable variety originating in plants taken as food. A surprising number of animal species employ all three of the mechanisms basic to urine production: filtration, reabsorption, and secretion.

Filtration of blood usually initiates the process and conserves blood cells, blood proteins, and fat droplets for the body. In the passage of the filtrate through tortuous channels, the substances most needed are recaptured by a secretory process termed reabsorption. Salts, amino acids, and carbohydrates fall into this category, and water may also be reabsorbed in quantity, thus concentrating the waste products left in the urine. The special compounds ingested but not used in the body may thus be excreted, but at this point the third mechanism, that of secretion of material from the blood into the urine, becomes pertinent. Some compounds reach a concentration in the urine far beyond that which would be possible by the combination of the first two processes.

Of major interest here is the ubiquity with which certain categories of substances are actively secreted. Phyla of animals separated by hundreds of millions of years of independent evolution have nevertheless retained the capacity for secreting the same compounds. The renal physiologist uses the phenolphthalein dye phenol red as a tool in exploring this field. This compound is actively secreted by almost every animal in which it has been tested. But it seems obvious that phenol red is not a compound found in nature. Its structure simply fits the transport units of an active pump retained to excrete other members of a large family of compounds. Hippuric acid, formed by the conjugation of the benzene rings of benzoic acid with glycine, might serve better as a model. This compound also fits the transport units which carry phenol red and is highly concentrated in the urine of many animals. The criticism of this compound, directed at some suggested

detoxification mechanisms, is that the products of the supposed detoxification are still somewhat toxic. The advantage may be simply the use of a compound whose concentration in the body is quickly lowered by the excretory pump. Benzoic acid, for example, is not as rapidly excreted as is hippuric acid.

Some of these compounds are excreted so efficiently that blood flowing through the kidneys is effectively cleared of the material. It appears that one mechanism may transport such diverse materials as phenol red, hippuric acid, penicillin, and diodrast. These compounds are but representatives of whole classes of compounds possessing the necessary structure to some degree.

Excretory Structures

Having listed the major excretory problems faced by animals, it will be profitable to proceed systematically through the animal kingdom to examine the structures that have developed to meet the problems and to follow their evolutionary sequence wherever possible.

Contractile vacuoles. In some plant gametes and in the protozoa and sponges, animals commonly placed at the bottom of the evolutionary array, there is a conspicuous mechanism for ridding the body of excess water and probably of wastes as well. These structures are called contractile vacuoles. Measurements show that some animals may have to put out a volume of water equal to the total body volume every 4 min. One view is that the water is secreted into channels in the protoplasm, through which it is gathered into the vacuole for extrusion. *See* PORIFERA; PROTOZOA.

It is possible that the contractile vacuole represents simply the interiorization of a phenomenon common to most cells. The reverse process, the taking up of water from a medium, has been watched and termed pinocytosis or cell-drinking. In such a case the surface membrane is very active. Drops of external fluid are trapped in infoldings of the membrane, which then leave the boundary of the cell to become vacuoles in the protoplasm. If the membrane is digested away, the droplet becomes a part of the protoplasm. It has been suggested that the secretion of water may be accomplished by a reverse of this process. The many internal membranes that are being revealed by the electron microscope surround droplets of water separated from the protoplasm. These vacuoles approach the surface membrane, incorporate into it, and discharge the droplet to the exterior. In the case of the contractile vacuole, the droplets discharge, instead, into the internal channels leading to the contractile vacuole because in these cellular forms it is important for the external pellicle to have considerable mechanical strength, whereas in the internal cells of multicellular animals the cell membrane can be a delicate and mobile one. *See* ENDOCYTOSIS.

Samples of fluid taken from the contractile vacuoles with micropipets show that, in the process of excreting water, salts are conserved by the cell; the

excreted fluid is hypoosmotic to the protoplasm. *See* VACUOLE.

Coelenterata. This important group of metazoan animals is usually placed just above the sponges and is known as the phylum Coelenterata or Cnidaria. This phylum warrants mention here because, while the cells do not possess contractile vacuoles, no specific mechanism for excretion or osmotic regulation has been found. Members of the group live in marine, brackish, and fresh-water environments, and it may safely be forecast that the cellular adaptations that make this possible will someday be identified and studied. *See* CNIDARIA.

Nephridia. A most significant development in the evolution of excretory processes occurs in the phyla beyond the Protozoa, Porifera, and Coelenterata. Structures termed nephridia are present in the Platyhelminthes, Rotifera, Nemertea, Acanthocephala, Priapulioidea, Entoprocta, Gastrotricha, Kinorhyncha, Cephalochorda, and some Archiannelida and Polychaeta, as well as in larval stages of the Polychaeta, Archiannelida, Echiuroidea, Mollusca, Phoronidea, and Cephalochorda. An eminent authority, E. Goodrich (1945), said that these organs (protonephridia) are in fact widely distributed and may be inferred to have been present in the common ancestor of all the metazoan Triploblastica.

Aside from the fact of its presence in all the animals at this intermediate level of evolution, the nephridium represents for the process of excretion the specialization of tissues that enable animals to become larger and more independent of their immediate environment. The cells of the body no longer live in antagonism to seawater or fresh water, but are bathed by an internal medium which is regulated to represent an optimal environment. Numerous tissues have as their special function the regulation of the composition of this medium. It is the part of the nephridium to clear the body fluids of wastes, without at the same time losing from the body the substances needed by the cells. As in the case of the contractile vacuole, it is not understood clearly how the nephridium accomplishes these functions, but assumptions can be made about several of the processes involved. *See* URINARY SYSTEM.

Two particular sorts of nephridia are worth describing briefly, the protonephridium and the metanephridium. The protonephridium is the primitive structure consisting of a single cell, shaped like a flask, within the neck of which there is found a cylindrical cavity, or lumen. In this lumen are found one or more threadlike structures, or cilia, which beat continuously, the appearance thus giving rise to the old term flame cell. The flame cell appears to filter through itself a part of the animal's circulating fluid, but its polarity is different from that of the other cells of the body, and a "urine" passes from the lumen to the exterior by way of a nephridiopore.

The metanephridium is a structure for which the functions can be more nearly assigned. The tubular structure opens into a body or coelomic cavity. There is always fluid in this cavity, resulting from filtration through the cells making up its lining. This fluid

contains not only wastes but also materials of value to the animal, and as it is propelled down the nephridium by a lining of cilia, it is subjected to changes as it passes the cells making up the tubule. Some salts and organic substances are reabsorbed before they leave the body. In the few cases which have been studied, it has been noted that phenol red is secreted into the urine by nephridial cells.

Higher invertebrate kidneys. The phylum Annelida has produced the peak development of small independent nephridia. Some individual animals may possess thousands, of which some are directed to the exterior and others to the digestive lumen. An evolutionary line from the group leads to the Mollusca and Arthropoda, and along this line it is found that the nephridial system is replaced by quite different structures. See ANNELIDA.

Mollusca. The trochophore larvae of the mollusks possess nephridia. As the larvae settle down and undergo metamorphosis, the nephridia are lost and replaced by kidneys proper. Filtration through the coelomic lining begins the process of urine formation in the animals possessing metanephridia. Now the kidney itself originates from the remnants of the coelom, which is very much reduced in modern mollusks. It is to be expected that filtration would be retained as a first step in urine formation, and this is found in all the mollusks so far studied. In the bivalve, or pelecypod, mollusks the filtration is actually through the walls of the heart; in the gastropods the filtration is probably from capillaries in the walls of the kidneys; and in cephalopods the filtrate appears to be formed by special organs, the branchial heart appendages. Filtration pressures are not high, but the mollusks generally have a low protein content in their blood, and so the reabsorptive forces which oppose filtration are not high either. As a result, surprisingly large volumes of filtrate have been obtained in the few studies made on the subject. From the pericardial space, or branchial heart appendage, the filtrate passes next through the renopericardial canal into the kidney proper, where two processes are simultaneously at work. Because the process of filtration cannot be finely discriminating, it merely separates water and solutes from blood proteins and cells. Many useful solutes remain in the filtrate and would be lost from the body, but they are absorbed by the cells of the kidney from the urine as it passes over the many folds into which the lining of the kidney is thrown. At the same time other cells of the kidney actively secrete materials into the urine. These three processes, filtration, reabsorption, and secretion, occur repeatedly in phylum after phylum, regardless of the origin of the kidney tissues.

Mollusks may be exposed to desiccation during tidal excursions or at certain seasons of the year. A low metabolic rate and the usual presence of an impermeable shell make them well suited for surviving prolonged drought. The kidneys may release no urine whatever during such a period, the chief nitrogenous excretion, uric acid, simply forming crystals within the kidney or in the kidney cells. At the

next wet period normal function is restored. See MOL-LUSCA.

Arthropoda. Of the very large phylum Arthropoda, only two major subgroups are treated here, the crustaceans and the insects.

1. *Crustacea.* The crustaceans which have been most studied originate their urine at a thin-walled sac held open by strands of connective tissue passing from the external surface of the sac or in a network of fine capillaries in its walls results in formation of a filtrate upon which other changes may later be made. The filtrate formed in the end sac runs next through a long canal lined with cuboidal or columnar epithelium before it is voided to the exterior. While the filtrate passes through the canal, salts and water may be reabsorbed to conserve them for the body, while other materials are secreted into the urine and thus excreted at a higher rate than they would have been by the process of filtration alone.

2. *Insecta.* Excretion in the insects follows an entirely different pattern from any described to this point. In his early work on the histology of the silkworm, M. Malpighi described the tubes which were later named the Malpighian tubes when their excretory function was learned. These tubes, 2-150 in number, lie in the abdominal cavity; when they are few in number they are long, and when they are more numerous they are short. The Malpighian tubes begin development as blind pouches, lie free in the body cavity, are convoluted or not depending upon their length, and finally open most frequently into the midgut but sometimes into the hindgut. There is no obvious mechanism for filtration of blood, and it is thought that the urine in all the members of this large group of animals is formed entirely by secretion of the cells lining the tubes. In fact the tubes sometimes secrete other than excretory substances; for example, silk for the cocoon is secreted by the Malpighian tubes of some of the Collembola and Neuroptera. In the absence of filtration, it might be surmised that there would be no need for reabsorption, but it has been shown that uric acid is secreted in the form of the relatively soluble sodium and potassium salts in the upper segment, and that the free uric acid is precipitated in the lower segment where the water and the liberated bases are reabsorbed and used again. In insects much reabsorption of water occurs in the hindgut. See INSECT PHYSIOLOGY; INSECTA.

Nephrons and the Vertebrata. The large number of phyla possessing nephridia has been emphasized above. One major evolutionary line whose forebears, though uncertain, were undoubtedly animals possessing nephridia, has progressed along a path that produced the echinoderms and led to the chordate group of animals. See CHORDATA; ECHINODERMATA.

The phylum Echinodermata, composed of relatively inactive marine animals, does not shed much light on the process since these animals do not have specialized excretory organs. Experimentally injected foreign substances are picked up by wandering amoeboid cells, the coelomocytes, and delivered to the digestive tract, from which they are excreted with the feces.

Among the chordates, the cephalochordates have nephridia comparable to those of the annelids. The vertebrates have a common origin with the other chordates, but the nephridia are called nephrons and are regarded as the unit structures of the kidney. With minor modifications, each class of vertebrates possesses kidneys made up of nephrons like those of its relatives, all the way from the cyclostomatous fishes up to the mammals. The small and primitive animal is characterized by few nephrons, and the large and highly evolved animal by many nephrons, so that a single kidney of the human may contain 1,250,000 of these units. *See* KIDNEY.

A nephron may be said to be composed first of a Malpighian body consisting of a tuft of capillaries, the glomerulus, surrounded by a close-fitting sac, Bowman's capsule. Blood pressure in the tuft of capillaries filters a large amount of fluid into the capsule, from which it passes into a much-convoluted tubule draining the fluid away through a narrow passage between closely packed kidney cells. These cells carry out the complicated processes of reabsorbing for the body the materials of usefulness and of adding to the urine, from the blood, materials of such a nature that special mechanisms could usefully be developed for their transportation across the kidney cells. Nevertheless, the utilization of the three processes of filtration, reabsorption, and secretion used throughout the rest of the animal kingdom can be noted. A few special features of excretion in the different vertebrate classes are worthy of comment.

Fishes. In the Insecta there is a specialization of structure such that secretion is the primary mechanism used by the Malpighian tube in fulfilling the function of excretion for the animal. A similar device has been developed by some marine fishes; the Malpighian body of the vertebrate nephron is gone, leaving only the tubule, closed at the end where the glomerulus should be found. Since the means of accomplishing filtration is gone, it is clear that secretion is the only remaining mechanism for accomplishing the excretory functions of the kidneys in these animals. These nephrons make up what is called an aglomerular kidney, and the aglomerular fishes represent a successful group of animals, being numerous and widely distributed. It is thought that the reduction or loss of the glomerulus represents a means of water conservation. Fish blood is osmotically inferior to the salt-water environment. The fishes consequently continuously lose water osmotically through the gills and pharyngeal membranes, and any means of reducing excretory water loss has survival value.

The general osmotic problem faced by marine fishes is of such overriding nature that an entirely extrarenal excretory mechanism has developed to meet it in the gills of fishes. The characteristic cation of living cells is potassium, whereas the abundant cation of body fluids and saline waters is sodium. Cell membranes are permeable to these small ions, so that a majority of all animal cells continuously pump sodium out of the cell and potassium into the

cell and must do so to continue normal life. Certain cells of the fish gill direct the pumping action not generally outward but polarize it, so that sodium entering from the blood and tissue fluid is pumped against a steep gradient into the seawater. In this way the kidney of the fish is spared the task of excreting the sodium ion, a task which the kidney of most terrestrial animals must perform. *See* RESPIRATORY SYSTEM.

Amphibia. Amphibians are characteristically bound to bodies of fresh water. Their osmotic problems are therefore quite different from those of the marine fishes just discussed. Their nephrons begin the process of urine formation with vigorous filtration. As the filtrate passes down the tubules, more salt than water is reabsorbed, water ordinarily being present in excess, and a hypotonic urine results. Even so, salt is continuously being lost in the urine and so is at a premium in these animals. The deficit is made up by the skin, which has the power of taking sodium and chloride ions from the environment into the blood against a steep gradient.

Reptiles. The discussion of excretory specialization in reptiles will be deferred until the section dealing with the excretion of nitrogenous wastes.

Birds. Marine birds face a situation similar to that of the marine fishes. Although they may not lose water by osmosis directly to the seawater, some of them remain in environments where fresh water is available. At the same time they may feed on marine invertebrates whose body cells are in osmotic equilibrium with the seawater, and so through the digestive tract they are exposed to salt loads that constitute a serious problem to the body. In these cases research shows that it is not the kidney which has assumed the excretion of salt as a higher rate but special glands of the head region, the nasal glands. These glands produce a salt solution very hypertonic to the blood, which is blown, drips, or is shaken from the beak. They are thus able to meet the imposed load easily. The function of these glands has been well integrated into the physiology of the animal and is controlled from the osmoreceptors of the body by way of the central nervous system. *See* SALT GLAND.

Mammals. Desert mammals meet somewhat the same excretory problem that marine fishes and some marine birds must meet, namely a limitation on the water available for excretion. It has been shown that their kidneys have developed the capacity to excrete salt far beyond that of the ordinary mammal. These animals can comfortably eat dry food and drink seawater and, because they excrete a urine hypertonic to seawater, maintain themselves in satisfactory water balance. Some marine mammals are not only limited in their fresh-water intake but eat marine invertebrates containing large amounts of salt. These animals show the same adaptation as the desert mammals: great ability to concentrate salt in the kidney tubules. In fact only in the mammals has the kidney been the organ of choice for the production of a hypertonic excretory fluid. The mechanism developed makes use of the ordinary structures of the kidney, but in an entirely new way.

Excretion of Nitrogenous Wastes

Nitrogenous wastes, the group of compounds which, it may be assumed, necessitated the development of excretory systems, are ammonia, urea, and the purine derivatives. From the standpoint of the enzymologist, these compounds make up a single family, all the members of which can be synthesized by well-known enzymes starting either from the ammonia generated from amino acids or with purines generated from nucleic acids. The intermediates are illustrated in **Fig. 1**. It has been estimated that 95% of the nitrogenous wastes of the average animal comes from proteins and only 5% from nucleic acids, but because of the metabolic interconnection, one cannot predict from these percentages what the form of excretion will be. Instead it turns out that the habitat and the need for conservation of water determine the form of excretion. In animals which spend their lives in an abundance of water—for example, amphibian larvae and many marine invertebrates—the ammonia split from protein is simply liberated to the water and carried away. The animals are termed ammonotelic. But these animals also produce nitrogenous wastes from nucleic acids. Instead of liberating adenine and guanine directly to the environment, the complement of enzymes shown in **Fig. 1** is invoked, the material is split all the way to ammonia, and the ammonia released therefore originates from protein and nucleic acid. This general principle of a common denominator of nitrogen excretion was enunciated by M. Florkin. The results of the principle, and some exceptions, are shown in **Fig. 2**.

Elaboration of urea. Animals of larger size, in particular those which might not have a steady supply of water, cannot tolerate the accumulation of ammonia which would result from the use of ammonia as the single nitrogenous excretory product, and it was necessary that a less toxic material be substituted. In retrospect it is clear that the next evolutionary step, the elaboration of urea, was both most important and very difficult. Early biochemists speculated that urea was formed from two molecules of ammonium carbonate in what appeared to be a relatively simple synthesis. Progress in working out the actual pathways was slow in spite of the fact that the action of the enzyme arginase in splitting urea from the amino acid arginine had long been known. There was simply not enough arginine available to account for the amounts of urea formed. It was through the genius of H. A. Krebs that the solution was found. Only catalytic amounts of arginine need to be present. Arginase splits off urea, leaving a molecule of ornithine. This is condensed with an ammonia to form arginine again. Ornithine and citrulline had long been known in nature as constituents of a few plants, but their possible roles in such a cycle had not been suggested until Krebs demonstrated that catalytic amounts of these compounds greatly accelerated the test-tube synthesis of urea by liver tissues. Since that time, much detailed information has been accumulated about this reaction, including the energy sources needed for the synthesis. Unlike the simple liberation of ammonia, the synthesis of urea demands the provision of

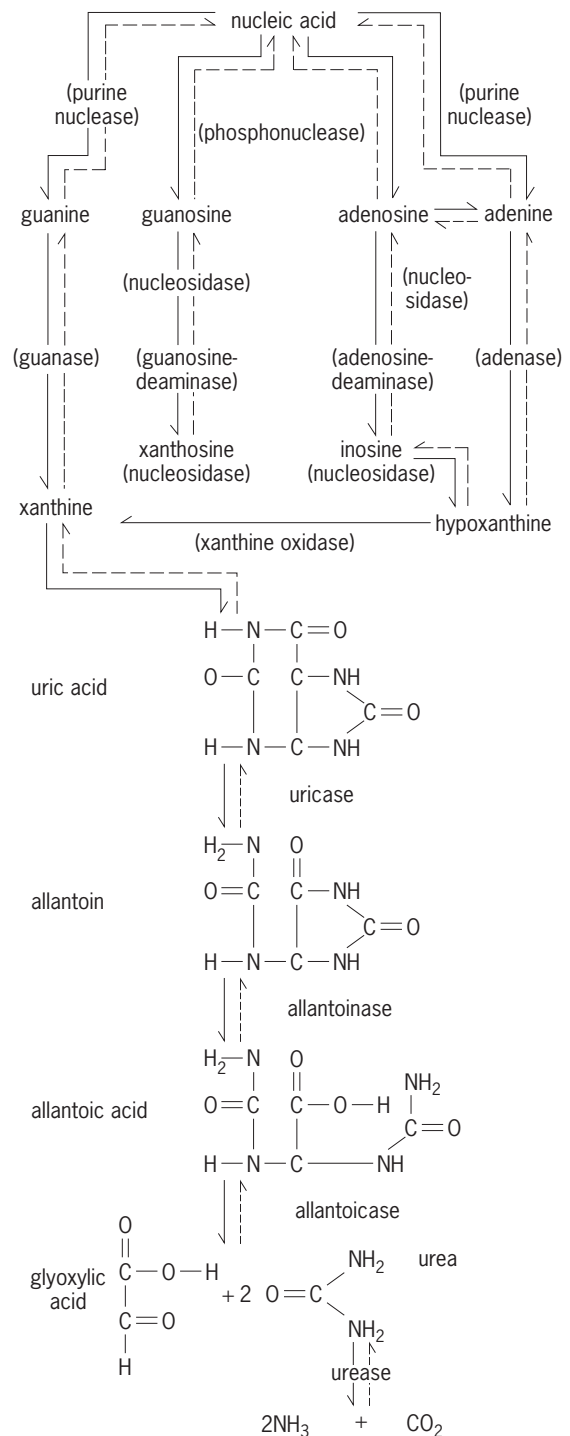


Fig. 1. Sequence of breakdown from nucleic acid to ammonia. Reverse pathways broken indicate either lack of knowledge of the mechanism or that other enzymes than those shown are involved. (After M. Florkin, *Biochemical Evolution*, Academic Press, 1949)

a special complement of enzymes and the utilization of some of the body's supply of energy. The result for the whole animal, however, is the substitution of urea, a relatively soluble and nontoxic material, for the highly toxic ammonia molecule. Animals which synthesize and excrete urea are spoken of as being ureotelic. Most fishes and many amphibians fall into this category. There is a decided tendency for the

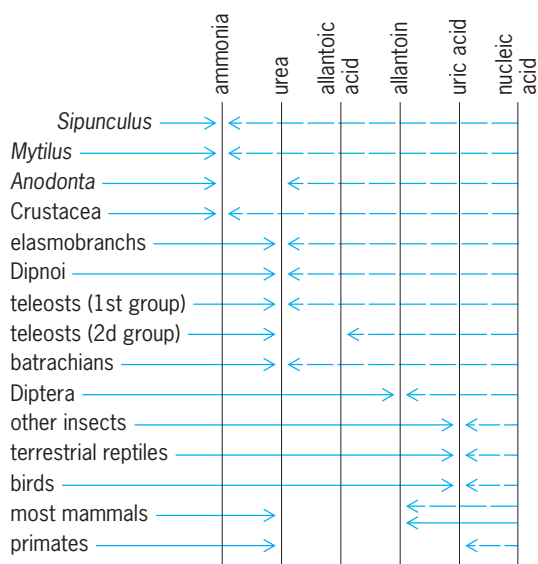


Fig. 2. Chart of the principle of excretion of a common nitrogenous waste, whether derived from protein or nucleic acid. (After M. Florkin and G. Duchateau, *Arch. Int. Physiol.*, 53:267-307, 1943)

purine bodies to be degraded to the level of urea and excreted in this form (Fig. 2). See NITROGEN EXCRETION.

Urea synthesis and excretion. The synthesis and excretion of urea made animals independent of a direct and abundant source of water. The solubility of urea, which appears actually to have been one of its virtues for the development of the primitive kidney, turns out to be a disadvantage to animals such as the birds and desert lizards which, for conservation of weight and water, need an excretory material which is both nontoxic and insoluble. At the same time the material must be capable of being secreted by kidney cells in solution in water and precipitated only in regions from which it can be expelled safely from the body. From the evolutionary point of view, such a material had been available all the time to the more primitive animals. It will be recalled that in the fish and amphibians (Fig. 1) the nitrogen from the purines was degraded by way of uric acid to urea. The birds and most reptiles arrest this degradation at the uric acid stage. By enzymatic processes not yet fully understood, they carry the synthesis of the protein ammonia past the stage of urea synthesis and all the way up to uric acid. The uric acid is precipitated in the ureters and cloaca, and the pasty material is eliminated through the rectum. In this way another great step is made in the conservation of water. These animals are called uricotelic. It may be noted from Fig. 2 that such a mechanism was independently developed by the insects, and for the same reasons. It should be pointed out that when such mechanisms develop there is not necessarily a clear-cut departure from the previous condition. Careful study shows therefore that various animals excrete varying mixtures of ammonia, urea, and uric acid, instead of only a single material. Arthur W. Martin

Bibliography. P. Dejours et al. (eds.), *Comparative Physiology: Life in Water and on Land*, 1987; N.

Heisler (ed.), *Acid-Base Regulation in Animals*, 1986; B. L. Muir, *Pathophysiology: An Introduction to the Mechanisms of Disease*, 2d ed., 1988.

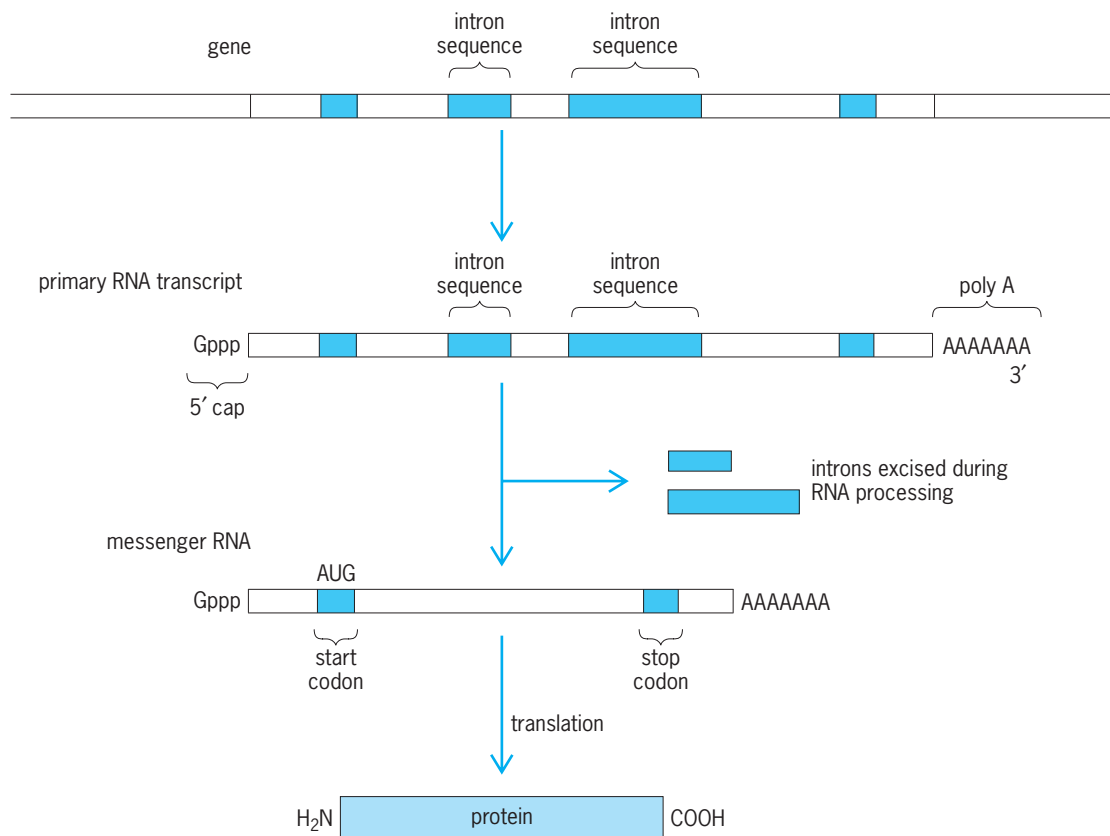
Exon

In split genes, a portion that is included in the ribonucleic acid (RNA) transcript of a gene and survives processing of the RNA in the cell nucleus to become part of a spliced messenger RNA (mRNA) or structural RNA in the cell cytoplasm. Split genes are those in which regions that are represented in mature mRNAs or structural RNAs (exons) are separated by regions that are transcribed along with exons in the primary RNA products of genes, but are removed from within the primary RNA molecule during RNA processing steps (introns). See INTRON; RIBONUCLEIC ACID (RNA).

Exons comprise three distinct regions of a protein-coding gene. The first is a portion that is not translated into protein, but contains the signal for the beginning of RNA synthesis, and sequences that direct the mRNA to ribosomes for protein synthesis. The second is a set of exons containing information that is translated into the amino acid sequence of a protein. The third region of a gene that becomes part of an mRNA is an untranslated end portion that contains signals for transcription termination and for the addition of a polyadenylate tract at the end of a transcript. See GENETIC CODE; PROTEIN; RIBOSOMES.

The evolutionary origin of the exon-intron organization of eukaryotic genes has not been clarified. However, it has been determined from comparisons of structural-functional domains in proteins with the exon domains in corresponding genes that there is often a correlation of the two. For example, the gene for one of the human cell-surface proteins that are important in the recognition of self and nonself by the immune system, the HLA-B7 gene, has nine exons. Exon 1 encodes the 5' (beginning) untranslated portion of the mRNA that is presumably responsible for steps in the initiation of protein synthesis from the mRNA. Exon 2 encodes the site of translation initiation and the sequence of amino acids (the signal peptide) that designates the protein being made as one that will either be a component of membrane or pass through a membrane to be secreted. Exons 3, 4, and 5 encode the three domains of the HLA-B7 protein that are modified by the addition of carbohydrates and are exposed to the cell exterior. Exon 6 encodes the transmembrane portion of the protein, being composed of codons for a series of amino acids that are able to interact with the lipids of the cell membrane. Exons 7 and 8 encode the cytoplasmic ("anchor") portion of the HLA-B7 protein, and exon 9 contains polyadenylation signals.

The mechanism by which the exons are joined in RNA copies of genes is called RNA splicing, and it is part of the maturation of mRNAs and some transfer and ribosomal RNAs (tRNAs and rRNAs) from primary transcripts of genes (see *illus.*). Three different RNA splicing processes have been identified. One involves mRNA precursors in nuclei, and



Map of a eukaryotic gene, primary RNA transcript, and messenger RNA for a protein. Exons, the unmarked segments of the map, are separated by intron sequences, which are spliced out during synthesis of messenger RNA. The appearance of the 5' cap, a modified guanine nucleotide, precedes or accompanies any splicing activity. The 3' end, always a series of adenine nucleotides, is added at the same time. The start codon, AUG, specifies methionine. The stop codon, which signals the termination of protein-chain translation, is a ribonucleotide triplet.

specific sequences at exon-intron junctions that are recognized by certain nuclear ribonucleoprotein particles that facilitate the cleavage and ligation of RNA. Another applies to nuclear precursors of tRNA, where splice sites are determined by structural features of the folded RNA molecules. The third form of splicing was discovered in studies of protozoan rRNA synthesis, and has also been shown to be a part of the maturation of both rRNA and mRNA in yeast mitochondria; it is an autocatalytic process that requires neither an enzyme nor added energy such as from adenosine triphosphate. In bacteria, it has been observed that certain mRNAs in bacteriophage-infected cells result from the same sort of self-splicing that has been described for protozoan rRNA. *See* GENE.

Peter M. M. Rae

Bibliography. J. Darnell et al., *Molecular Cell Biology*, 1986; J. D. Watson et al., *Molecular Biology of the Gene*, 5th ed., 2000.

Exopterygota

A division (also known as Hemimetabola) of the subclass Pterygota, including those insects that show relatively slight change in body form with growth. They develop through a series of immature forms (nymphs) from hatchling to adult, so that wings

grow as external pads and enlarge with each molt. The nymphs are often scaled-down copies of the adults, but in some cases, particularly in those orders with aquatic nymphs and aerial-terrestrial adults (such as in Ephemeroptera and Odonata), a considerable difference in body form exists between adults and their immature forms. This development, called incomplete metamorphosis, is characteristic of all Paleoptera and some members of the infraclass Neoptera (all ancestrally winged insects except Endopterygota). Most exopterygotan taxa have freely mobile nymphs with well-developed legs, although conspicuous exceptions occur in the form of the scale insects (Coccoidea) and some of their highly sedentary sap-feeding allies.

Exopterygota are a very diverse group, encompassing plant feeders, predators, and animal parasites, and living in nearly all habitats and areas where insects are found. Common examples of Exopterygota are mayflies, dragonflies, grasshoppers, termites, cockroaches, aphids, plant bugs, biting and sucking lice, and thrips. *See* ENDOPTERYGOTA; INSECTA; PTERYGOTA.

William L. Brown, Jr.

Bibliography. Commonwealth Scientific and Industrial Research Organization, *Insects of Australia: A Textbook for Students and Researchers*, 2d ed., 1991; O. W. Richards and R. G. Davies, *Imm's General Textbook of Entomology*, 2 vols., 10th ed., 1994.

Exotic nuclei

Nuclei with ratios of neutron number N to proton number Z much larger or much smaller than those of nuclei found in nature. Studies of nuclear matter under extreme conditions, in which the nuclei are quite different in some way from those found in nature, are at the forefront of nuclear research. Such extreme conditions include nuclei at high temperature and at high density (several times normal nuclear density), as well as those with larger or smaller N/Z ratios. The N/Z ratio depends on the nature of the attractive nuclear force that binds the protons and neutrons in the nucleus and its competition and complex interplay with the disruptive Coulomb or electrical force that pushes the positively charged protons apart.

A chart of the nuclides is shown in **Fig. 1**. The squares are the stable nuclei ($Z \leq 83$) and the very long-lived nuclei (with half-lives of the order of 10^9 years or more) found in nature. The first jagged lines to either side are the limits of the presently observed nuclei; very little is known about those at the edges. The light, stable nuclei (such as ${}^4\text{He}$) have $Z = N$, reflecting the preference of the nuclear forces for $N = Z$ symmetry, but as Z increases, the strength of the Coulomb force demands more neutrons than protons to make a particular element stable, and $N \approx 1.6Z$ for the heaviest long-lived nuclei (such as ${}^{238}\text{U}$). For $Z > 92$, the Coulomb force causes most nuclei to spontaneously fission. See NUCLEAR FISSION.

Stable nuclei lie in the so-called valley of beta stability. As the N/Z ratio decreases (proton-rich nuclei) or increases (neutron-rich nuclei) compared to that of the stable isotopes, there is, respectively, energy for a proton or neutron in the nucleus to undergo beta (β^+ , β^-) decay to move the nucleus back toward stability. Most knowledge of nuclear structure and decay has been gained from nuclei in or near the valley of beta stability. See RADIOACTIVITY.

Structure and decay. The spherical shell model was developed to explain the so-called spherical magic numbers for protons and neutrons, which give nuclei with these numbers a very stable structure and spherical shape. Spherical magic Z and N of 2, 8, 20, 28, 50, 82, and 126, and the weaker magic Z and N at 40, are shown in **Fig. 1** as horizontal and vertical lines. The nuclear shell model resembles the atomic shell model, where the noble gases (helium, neon, argon, and so forth) have filled shells and then there is a gap in the binding energy to the next electron shell (or orbit). A major question concerns what happens to these spherical magic proton and neutron gaps (orbits) in exotic nuclei. See NUCLEAR STRUCTURE.

Another major question concerns the decay modes whereby a nucleus rids itself of excess energy and returns to the stable forms of nuclear matter. As N/Z decreases or increases from the stable values, a point is reached for a given Z where, if one more neutron is pulled out, one proton becomes unbound (an isotope of that element cannot exist with that number of neutrons), or, if one more neutron

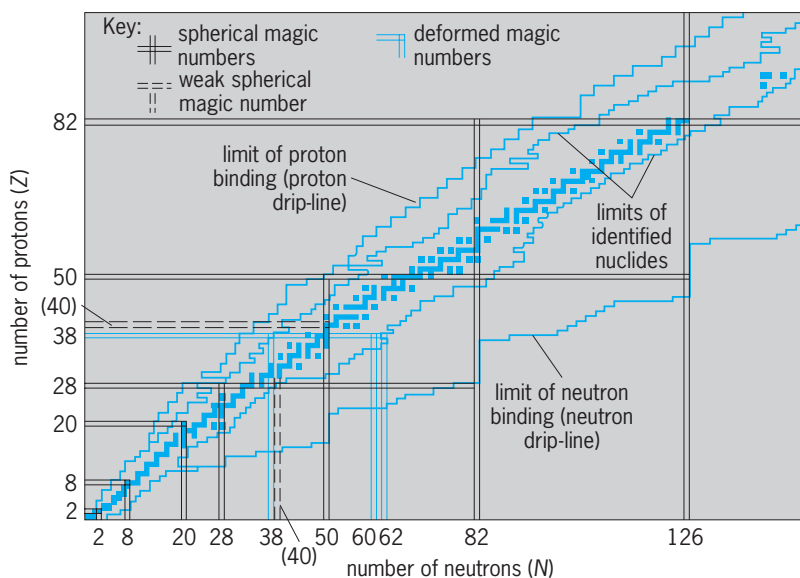


Fig. 1. Chart of the nuclides. The squares are the stable and very long-lived nuclei found in nature.

is added, that neutron is not bound to that nucleus. These limits define, respectively, the proton and neutron drip-lines (the jagged lines furthest from the valley of stability in **Fig. 1**).

Answers to the above questions give significant insights into the structure and decay modes of nuclear matter that make it possible to test and extend theoretical models of the nucleus and the understanding of the nature of both the strong and weak nuclear forces. These insights could not be gained by studies of nuclei in and near the valley of stability.

Shape coexistence. These insights include the discovery of nuclear shape coexistence, which was first observed in the very light mercury and selenium isotopes. Competing bands of levels occur in one nucleus, which overlap in energy but are quite separate in their decays because they are built on quite different coexisting nuclear shapes. Shape coexistence is now known to be important in many nuclei throughout the periodic table. Even multiple (four, and possibly more) different-shape coexisting structures are observed, for example, in ${}^{185}\text{Au}$ and ${}^{187}\text{Au}$.

Deformed magic numbers. A significant advance was made in the nuclear shell model with the discovery of new magic numbers associated with shell gaps in the energies of the proton and neutron orbitals. These new numbers may be called deformed magic numbers because they stabilize a nucleus in a deformed shape, just as the spherical magic numbers give stability to a spherical shape. The deformed magic numbers (shell gaps) identified so far include N and Z of 38 and N of 60 and 62. The deformed shell gaps at 38, 60, and 62 appear in **Fig. 2**, along with spherical shell gaps at 28 and 50, and other expected shell gaps are shown.

Shell gap reinforcement. These gaps led to a further advance in the shell model: the importance of proton and neutron shell gap reinforcement. Two new islands of the strongest ground-state deformations

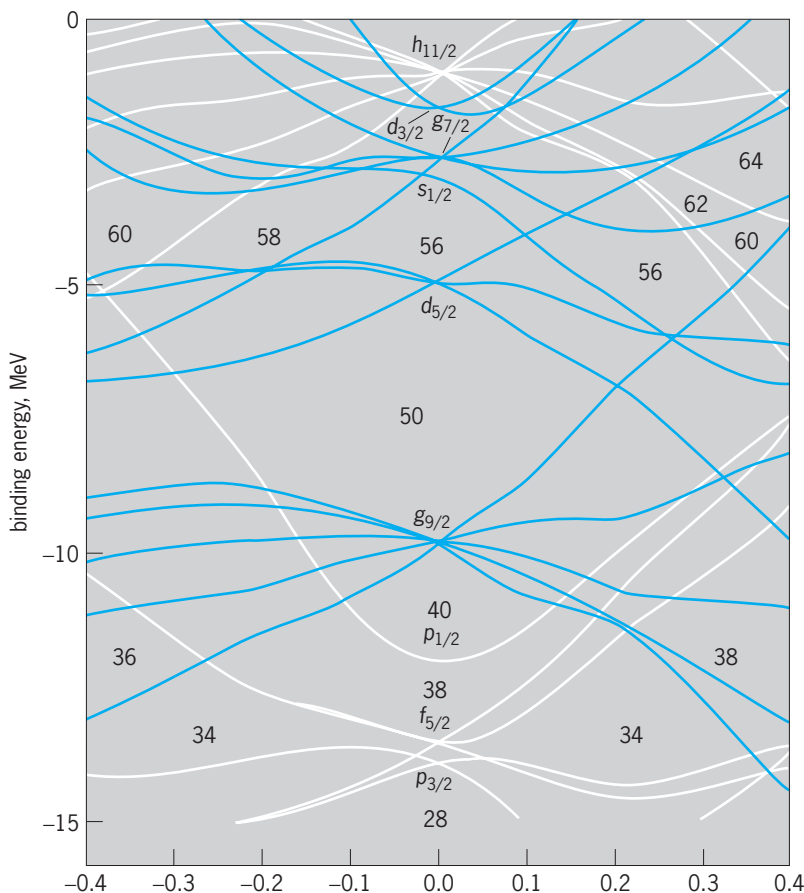


Fig. 2. Binding energies for neutrons and protons in nuclei in the mass region $A = 60\text{--}120$. These energies are plotted against the nuclear deformation parameter ϵ . The energies for a spherical shape occur in the center ($\epsilon = 0$), for prolate shapes to the right, and for oblate shapes to the left. A gap between levels occurs when a particular shell or orbital is filled. Numbers in gaps give total numbers of neutrons or protons in closed-shell or closed-subshell configurations. Subscripted letters give orbital and total angular momenta of single-particle states. (After R. Bengtsson et al., *Nuclear shapes and shape transitions*, *Phys. Scripta*, 29:402-430, 1984)

known above $A = 30$ (where $A = Z + N$ is the mass number) have been discovered centered on $N = Z = 38$ and $Z = 38, N = 60, 62$. In these cases, the protons and neutrons reinforce each other for the same strong deformed shape and override the weaker spherical shell gaps found at 38 and 40 in Fig. 2. On the other hand, the discovery of spherical ${}^{68}_{28}\text{Ni}_{40}$ shows that 40 is a good number for a spherical shell when it is reinforced by a proton partner (or neutron as in ${}^{90}_{40}\text{Zr}_{50}$) with a strong spherical shell gap like 28 (or 50). Shell gap reinforcement is important for both spherical and deformed shapes.

With all the shell gaps for different deformations in Fig. 2, the small number of known deformed magic numbers requires explanation. The 38, 60, and 62 deformed magic numbers are important only when the protons and neutrons reinforce each other for the same deformation. This reinforcement is likely to be necessary to observe the other gaps in Fig. 2. The competition between the different shell gaps also leads to shape coexistence as observed. It is very difficult to make exotic nuclei with Z and N shell gaps

at the same deformation, but the possibility of discovering additional deformed shell gaps is of much current interest.

Modes of decay. Exotic nuclei exhibit decay modes not seen near stability, such as proton radioactivity and beta delayed-particle emission. (After beta decay, the highly excited nucleus can emit one or more particles, such as one or two protons, an alpha particle, or one or two neutrons.) Joseph H. Hamilton

Neutron-rich nuclei. With calculations based on relatively simple models of the nucleus (for example, the liquid-drop model), it can be shown that nuclei with much greater neutron excess than those in, or near, the stability valley should remain stable with respect to one- and two-neutron emission, the limit that defines the neutron drip-line. This limit is approximately given by $N \approx 2.4Z$ for $Z > 10$. Between the stability valley, the neutron drip-line, and the limit imposed by spontaneous fission lie approximately 4000 neutron-rich nuclei, which should be stable with respect to the strong nuclear forces. They can, of course, decay by alpha and beta decay with lifetimes of $\gtrsim 1$ microsecond, which are very long times compared to the typical strong-force decay times of 10^{-23} s. Out of the 4000 neutron-rich nuclei, less than 1000, which lie close to the stability valley and can be reached by simple nuclear reactions, have been studied in great detail. Together with another ~ 700 nuclei on the proton-rich side of the valley, these constitute the sample on which almost all present knowledge of nuclear forces and nuclear structure is based. The remaining proton-rich (~ 1500) and neutron-rich (~ 3000) nuclei, which lie far from the stability valley and are difficult to reach, have earned the appellation exotic.

The exotic nuclei provide unusual opportunities for testing models of nuclear structure that have been proposed and optimized for nuclei in, or near, the stability valley. The neutron-rich exotics are especially important because of their astrophysical connection. The neutron-rich isotopes of the lightest elements, hydrogen and helium, provide models for neutron matter, and those of medium-weight elements play important roles in nucleosynthesis via the r -process. The heaviest neutron-rich nuclei, the superheavies, not only test the limits of nuclear stability but are of practical interest because of their expected unusual fission properties. See NEUTRON STAR; NUCLEOSYNTHESIS.

Techniques. A great variety of experimental techniques have been developed to create and study neutron-rich exotic nuclei. These include multi-nucleon transfer and charge-exchange reactions, fission, fusion, and projectile fragmentation and target fragmentation (spallation) reactions. The capabilities of these various techniques have been greatly enhanced by the use of medium- to high-energy beams of heavy ions, pions, and radioactive nuclei, together with stable and radioactive targets. The reaction products of interest are mechanically as well as electromagnetically transported and analyzed, and the stability, masses, half-lives, radii, spins, electromagnetic moments, and decay schemes have been

studied with ingenious techniques. See NUCLEAR REACTION; PARTICLE ACCELERATOR.

Limits of identified nuclei. The neutron drip-line appears to have been reached for elements with $Z \leq 7$; the heaviest particle-stable isotopes identified are ${}^3\text{H}_2$, ${}^8\text{He}_6$, ${}^{13}\text{Li}_8$, ${}^{14}\text{Be}_{10}$, ${}^{19}\text{B}_{14}$, ${}^{22}\text{C}_{16}$, and ${}^{23}\text{N}_{16}$. For oxygen the most neutron-rich isotope reached so far is ${}^{24}\text{O}_{16}$, although particle stability is expected to extend to ${}^{28}\text{O}_{20}$. For the alkali-metal sodium it has been possible to go as far as ${}^{34}\text{Na}_{23}$. For heavier elements the most neutron-rich isotopes reached so far seem to be at $N/Z \simeq 1.60 \rightarrow 1.70$. Examples are ${}^{50}\text{K}_{31}$, ${}^{99}\text{Rb}_{62}$, ${}^{148}\text{Cs}_{93}$, and ${}^{228}\text{Fr}_{141}$ from among the alkali metals, and ${}^{78}\text{Cu}_{49}$, ${}^{102}\text{Sr}_{64}$, ${}^{133}\text{Sn}_{83}$, ${}^{146}\text{Xe}_{92}$, ${}^{164}\text{Gd}_{100}$, ${}^{214}\text{Pb}_{132}$, ${}^{264}\text{Fm}_{164}$ from among the other elements. These N/Z ratios are still far from the neutron drip-line value of $N/Z \approx 2.4$; thus, much work remains to be done.

Nuclear structure. Several interesting nuclear structure results have emerged from the measurements of neutron-rich exotic nuclei.

1. Several nuclei that at one time or another were predicted to be unbound (mainly on the basis of assumed spherical shapes) have been found to be bound (for example, ${}^8\text{He}_6$, ${}^{13}\text{Li}_8$, and ${}^{14}\text{Be}_{10}$).

2. Mass extrapolations using relations based on systematics for the nuclei in, or near, the valley of stability, such as the Garvey-Kelson relations, have been found not to work well for the extremely neutron-rich nuclei.

3. An indirect measurement of the matter radius of ${}^{11}\text{Li}$ shows it to be very large, approximately the same as that of ${}^{28}\text{Si}$. It has been suggested that this reflects a unique phenomenon, a loosely bound neutron halo around a ${}^8\text{Li}$ or ${}^9\text{Li}$ core.

4. Particle transfer experiments with heavy ions have been used to infer that the nucleus ${}^6\text{H}$ is unbound with respect to one-neutron emission by only 2.7 MeV. If true, this would be of great interest. Unfortunately, a high-sensitivity pion double charge exchange experiment has failed to verify this claim.

5. Despite the fact that ${}^9\text{He}$ is found to be neutron-unstable by 1.1 MeV, the ground state and at least two excited states of ${}^9\text{He}$ are found to be ≤ 0.5 MeV wide. Surprisingly, the binding energies of all three states are well predicted by shell-model calculations whose parameters are optimized for the nuclei in the valley of stability.

6. Two-neutron separation energies, differences in charge radii, and in many cases location of the first 2^+ states have been studied for a large number of neutron-rich nuclei across the major shell-closing (magic) neutron numbers. Rapid unexpected onsets of large deformation have been found, for example, at $N \simeq 60$ (${}^{88}\text{Sr}_{60}$, ${}^{100}\text{Zr}_{60}$).

7. Elements up through $Z = 112$ have been successfully synthesized, as well as $Z = 114$, $Z = 116$, and $Z = 118$. See DARMSTADTIUM; ELEMENT 112; MEITNERIUM; ROENTGENIUM; TRANSURANIUM ELEMENTS.

Proton-rich nuclei. Highly unstable, proton-rich nuclei that define the limits of nuclear stability offer unique opportunities to study new or unusual ra-

dioactive decay phenomena. This boundary region, in which the proton binding energy goes to zero, has been delineated up to nuclides with mass $A \approx 60$ and has been established to a significant degree in the region $A = 100$ -185. It is relatively easier to reach this proton drip-line than the equivalent neutron drip-line, since, as successive neutrons are removed from a given stable isotope, the effects of the strong coulombic repulsion between the protons in the nucleus cause a rapid increase in the decay energy.

Apart from the observation of new radioactive decays (for example, proton radioactivity), studies of exotic proton-rich nuclei determine the nuclear mass surface under extreme conditions, test the effects of charge independence of nuclear forces, study the weak interaction in nuclei, and evaluate complex shell-model calculations of nuclear structure.

Experimental techniques. These must be capable of observing radioactive nuclei produced in low yield and with half-lives of 5-100 milliseconds to observe beta decays of proton-rich nuclides; further, at the drip-line, other decay modes such as proton decay may occur with considerably shorter half-lives than are predicted for beta decay.

Proton-rich nuclides can be produced by bombarding appropriate targets with either light-ion or heavy-ion beams. In principle, the observation techniques involve either prompt detection of the reaction products or on-line radioactivity measurements. An example of a prompt-detection technique is the use of one or more magnetic spectrometers followed by sophisticated detectors to study the fragmentation of heavy-ion projectiles with in-flight separation of the exotic species of interest. Examples of techniques involving radioactivity measurements are the helium-jet technique (in which nuclear reaction products are thermalized in helium, attach to an aerosol, and are transported rapidly to a detection area) and the use of on-line isotope separators.

Decay modes. The large beta-decay energies of these exotic nuclei permit the population of highly excited states in their daughters, so that beta-delayed proton and alpha-particle emission can be observed. More exotic decays are also possible.

Beta-delayed particle emission. Experimental information on how the many open beta-decay channels of a quite proton-rich nuclide manifest themselves can frequently be derived by measurements of delayed radiations (protons or alpha particles) from excited states of the daughter nuclide that are populated by beta (positron) decay. These decay modes are characteristic of nuclides with $N - Z \leq -3$. **Figure 3** illustrates this decay for the lightest aluminum isotope that exists, ${}^{22}\text{Al}$ ($N - Z = -4$), with a half-life of approximately 70 ms. This abbreviated decay scheme shows a favored positron decay of ${}^{22}\text{Al}$ (with 18.5-MeV available decay energy) to the analog state in its ${}^{22}\text{Mg}$ daughter, with subsequent emission of a proton (p) to either the ground state or first excited state of ${}^{21}\text{Na}$. (The analog state arises if nuclear forces are charge-independent, and it has the same nuclear structure as the parent state in ${}^{22}\text{Al}$.)

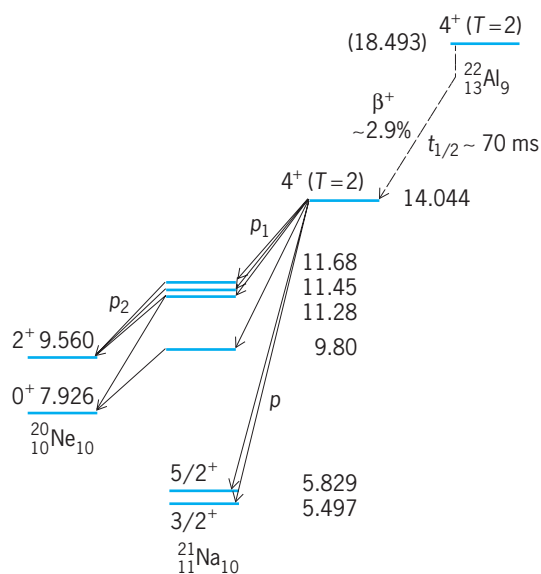


Fig. 3. A partial decay scheme for aluminum-22 (^{22}Al). Numbers above the bars at left give spins and parities of levels; numbers above the bars at right or to right of the bars give energies relative to the ground state of ^{22}Mg in MeV. T is the total isobaric spin. (After M. D. Cable et al., *Beta-delayed two-proton decays of ^{22}Al and ^{26}P* , *Phys. Rev.*, C30:1276–1285, 1984)

The high-energy proton decay from the analog state is detected (but not the preceding positron decay). See ANALOG STATES.

Beta-delayed two-proton emission. The large decay energy of ^{22}Al shown in Fig. 3 led to the nuclide being the first discovered (1983) in which the phenomenon of beta-delayed two-proton radioactivity could be observed. Interest centers on whether this additional branch in the analog-state decay is that of a proton pair being simultaneously emitted or of the sequential emission of the two protons. The mechanism is now known to be sequential emission of a proton (p_1) to a highly excited state in ^{21}Na followed by a second proton (p_2) to the final nuclide ^{20}Ne . This decay mode has also been established in the decays of six other light nuclei up to ^{39}Ti .

Proton radioactivity. This decay provides direct nuclear structure information and is the single-step emission of a proton from either the ground state or a long-lived isomeric state. As in alpha decay, the decay rate is determined by quantum-mechanical tunneling through the coulombic and angular-momentum barriers.

The first proton radioactive nucleus discovered (1970) was a 247-ms isomer of ^{53}Co , a nuclide that lies fairly close to the valley of stability. Ground state proton emission was first observed in 1981 in the decay of ^{151}Lu . Subsequent studies have observed this decay mode in 13 elements from antimony to bismuth.

Two-proton radioactivity. Direct two-proton radioactivity has been predicted since 1960 but has yet to be experimentally observed. This simultaneous emission of two protons from the nuclear ground state, or perhaps from a many-particle isomeric state,

is expected to occur in nuclei near the proton drip-line. Should this decay be observed, its mechanism would be of considerable interest, since this phenomenon could be compared with the tunneling of electrons between metals in superconducting and normal states, and thus could be said to be a nuclear Josephson effect. See JOSEPHSON EFFECT; SUPERCONDUCTIVITY.

Joseph Cerny
Bibliography. D. A. Bromley (ed.), *Treatise on Heavy Ion Science*, vol. 8, 1989; J. Cerny and A. M. Poskanzer, Exotic light nuclei, *Sci. Amer.*, 238(6):60–72, June 1978; J. H. Hamilton and J. Maruhn, Exotic atomic nuclei, *Sci. Amer.*, 255(1):80–89, July 1986; D. N. Poenaru (ed.), *Nuclear Decay Modes*, 1996; M. de Saint Simon and O. Sorlin (eds.), *International Conference on Exotic Nuclei and Atomic Masses*, 1995; B. M. Sherrill and D. J. Morrissey (eds.), *Enam 98: Exotic Nuclei and Atomic Masses*, AIP Conf. Proc. 455, 1999.

Exotic viral diseases

Viral diseases that occur only rarely in human populations of developed countries. However, many of these diseases cause significant human morbidity and mortality in underdeveloped areas, and have the proven capacity to be transported to population centers in developed countries and to cause explosive outbreaks or epidemics. Most of the exotic viruses are zoonotic, that is, they are transmitted to humans from an ongoing life cycle in animals or arthropods; the exception is smallpox.

Smallpox may be the most exotic pathogen on Earth, since it was successfully eradicated by a worldwide program launched in 1967, and now exists only in a few laboratories. Eradication is not possible for most other exotic viruses, however, since their life cycles include animal hosts. Control programs that may include vaccinations, insect control, or quarantine must be used in areas affected by outbreaks.

There are a number of important diseases caused by exotic viruses. (1) Yellow fever, caused by a mosquito-borne virus, particularly threatens population centers in tropical areas. (2) Venezuelan equine encephalitis is caused by a mosquito-borne virus that often infects humans. (3) Rift valley fever had a treacherous outbreak in Egypt in 1977. (4) Tick-borne encephalitis occurs in limited geographic areas where human habitats overlap with those of small mammals that normally serve as hosts to the causative virus. (5) Crimean hemorrhagic fever is caused by a tick-borne virus but is capable of transfer from human to human contact. (6) Rabies is found principally in dogs and many wild mammals. (7) Lassa fever is caused by a virus endemic to a single rodent species of west Africa. (8) Hemorrhagic fever with renal syndrome is caused by a rodent-borne virus and found in Russia and elsewhere. (9) Marburg and Ebola hemorrhagic fevers, usually lethal, are caused by two viruses of unknown natural history.

Control of these diseases is often very difficult because of the lack of detailed knowledge about the natural history of the viruses in their natural animal hosts, and because of the difficulty of controlling natural populations of alternative hosts such as insects or rodents. See ANIMAL VIRUS; ARBOVIRAL ENCEPHALITIDES; RABIES; SMALLPOX; YELLOW FEVER.

Frederick A. Murphy

Experiment

The test of a hypothesis under controlled conditions. The experiment is one of the distinctive tools of the scientist. It enables the scientist to put questions to nature and receive answers. These answers lead to new problems whose solutions require more complex experiments (seeking smaller differences), improved techniques, detailed plans, and better analysis of results.

Results of experiments are usually recorded numerically. Statistical methods are used to reduce the data to summary forms with estimates of the magnitude of such effects as averages, variances (dispersal of the data around the average), and the relationships between the variables being measured. Associated with these estimates are tests of significance that permit the researcher to go beyond the samples and make inferences about the characteristics (parameters) of the populations from which the samples are drawn. Statistics helps the researcher understand the magnitude of the imperfections of the experimental data. See BIOMETRICS; STATISTICS.

Experimental error. Results of experiments are affected not only by the treatments (experimental procedures whose effects are being measured and compared), but also by certain extraneous variations that tend to mask the effect of these treatments. Two main sources of experimental error which must be distinguished are (1) an inherent quality of the experimental material (the material on which the treatments are acting) and (2) lack of uniformity in the conduct of the experiment or failure to standardize techniques.

The presence and cause of experimental error need not concern the investigator, provided the results are sufficiently accurate to permit definite conclusions. Often, however, the results of experiments are greatly influenced by experimental errors, making decisions difficult.

Planning experiments. Inferences that can be made from results depend upon the way the experiment is conducted; thus planning of an experiment must include a detailed description of the proposal, together with such considerations as the following.

Objectives of the experiment. The purpose of the experiment must be clearly defined in terms of questions to be answered, hypotheses to be tested, specifications to be met, or effects to be estimated. The statement should indicate the extent to which inferences from the data will be applied.

Usually a sample is used. The size (that is, the number of units analyzed) of the sample needed depends

upon the accuracy desired for the results, the variability of the experimental material, the errors of measurements, the magnitude of the differences to be measured, and the time and money available for the study. The results of any piece of research are no better than the sample used. A sample is interesting only if it furnishes information about the population to which the conclusions are to be applied.

Description of the experiment. After the objectives have been determined, the researcher should describe the five following aspects of the experiment.

1. *Treatments.* The treatments whose effects are to be tested and measured are selected. It is important to understand how these treatments will help reach the objectives of the experiment. The choice of treatments may also have a substantial effect upon the precision of the experiment. Often it is desired to test many varieties, levels of chemicals, sources of vitamins, or kinds of spray material. These treatment combinations make up single-factor experiments. Striking gains in precision may be achieved in factorial experiments, where several factors are investigated simultaneously. Such experiments are useful in exploratory work to determine whether factors have any effect, to find interactions among factors, and to permit recommendations that apply to a wide variety of conditions.

2. *Experimental material.* The experimental material is chosen and its inherent variability is evaluated. It is useless to try to get the chemical determination to check within 2% when there is a 20% variation among samples of the source material. At the other extreme, it may be unwise to select uniform material because the responses obtained may not apply to the regular, unselected material.

The term experimental unit denotes the group of material to which a treatment is applied in a single trial of the experiment. The unit may be a patient in the hospital, an electronic tube, a group of pigs, a tree, or a missile. It is the characteristic of such units that they produce different results even when subjected to the same treatment.

3. *Size of the experiment.* This factor helps determine the sensitivity of the experiment, that is, the size of the detectable differences between the effects of treatments. One method of improving sensitivity is to increase the size of the experiment. This can be done by using more replications, that is, by increasing the number of times a complete set of treatments is run, or by using larger experimental units. Whatever the source of variation, replication of the experiment decreases the error associated with the difference between the average effects of two treatments, provided the precaution of randomization within replications has been used and inclusion of more experimental units or larger ones has not introduced additional variation. Statistics and the scientist's experience can assist in determining the efficient size for the experiment.

4. *Experimental techniques.* The principal objectives of good techniques are to secure uniformity in applying treatments, to impose sufficient control over external influences so that every treatment

produces its effect under comparable and desirable conditions, to provide unbiased measures of the effects of treatments, and to prevent gross errors. Refining the techniques helps to reduce the experimental variation.

5. *Related variables.* These variables, which predict to some extent the performance of the experimental units, if measured and used, will increase the efficiency of the experiments. Analysis of covariance enables the experimenter to estimate from the data the extent to which the results were influenced by variation in these supplementary variables.

Experimental designs. Experimental error often can be minimized by choosing an efficient design, a set of rules for allocating treatments to experimental units. Each restriction imposed by the experimental design has a definite purpose. Of the many types of designs available, only the basic ones and a few of the variations are to be mentioned.

Completely randomized. In this type the treatments are allocated to experimental units entirely by chance. These units should be handled in random order at all stages of the experiment at which order is likely to affect the results. This design is desirable for laboratory research, especially in physics, chemistry, bacteriology, or experimental cookery, in which mixing produces quantities of homogeneous material that can be tested under uniform conditions.

Randomized complete block. In this design the experimental units are in compact or homogeneous groups, each group containing enough units for all treatments. Usually several groups (replications) are needed to give an estimate of the effects being tested. During experimentation a uniform technique should be employed for all units in a group. The experimental error arises from variation within groups.

Latin square. The treatments are grouped into replications in two ways, thus providing an opportunity for two reductions in error by skillful planning. The location of leaves on a plant and the difference between plants may determine the basis for grouping.

Incomplete block. This design is used for experiments in which a large number of treatments is to be investigated and the number of homogeneous experimental units which can be grouped is small. These designs are composed of blocks or groups of experimental units smaller than a complete replication.

An example of how designs are selected to fit into experimental situations is given. In tests of mosquito repellents which involve exposure of treated arms to mosquitoes, the block consists of two arms of a subject at one time. To test six repellents the design in the **table** provides for each of five subjects to

submit his two arms to treatments three times.

Every possible pair of treatments occurs once on some individual on one day. Repellents 1 and 2 are assigned at random to the right and left arms of the individual A during the first test day. On the second day the arms are submitted to repellents 6 and 5. After 3 test days, the arms of each individual have been submitted to all six repellents.

Incomplete block plans, such as the one illustrated by the table, often are 20–100% more sensitive than randomized complete block designs; that is, smaller differences can be determined in an equivalent number of trials.

Split-plot. This design is used for experiments in which an extra factor is introduced by dividing each experimental unit into two or more parts. Any of the plans presented above could have each unit divided into parts permitting the use of an extra factor. In industrial experimentation frequently one series of treatments requires rather large bulks of experimental material, such as types of furnaces for the preparation of alloys. Another series, such as the molds into which the alloy is poured, can be compared through use of smaller units.

Factorial experiment. In factorial experiments more than one factor is included in the treatment combinations. The quantitative or qualitative factors may be either independent or not independent.

If factors are independent, these factorial experiments may be arranged in completely randomized, complete or incomplete block designs. As the number of factors or levels of each factor increase, the number of treatment combinations increases rapidly. When the number of homogeneous units in groups is limited, incomplete block designs are desired. For factorial experiments, a special series of incomplete block designs has been developed where information on certain treatments is sacrificed in order to secure better information on other treatments or where higher-order interactions are completely or partially confounded with blocks.

If quantitative factors are used regardless of independence, one thinks of the yield or response as a function of the level of the variables. Sequential experimentation is used to estimate the optimal point on a response surface and to explore the nature of this surface near this optimum.

The limitations of experimental materials, laboratory facilities, or time and convenience of investigators often dictate the specifications of the design, which is a logical plan. An adequate experimental design involves not only a satisfactory plan for conducting the experiment but also includes appropriate methods for evaluating results.

Analysis. The analysis of data calls for careful selection from the available efficient and flexible set of statistical methods. The use of these methods requires judgment and an understanding of the basic assumptions involved. It should be emphasized that statistical analysis cannot increase the validity of data. The accuracy of mathematical inferences from data must be limited by the precision, accuracy, and adequacy of the measurements of the observations. If

Incomplete block design					
Day	Individuals				
	A	B	C	D	E
1	1 2	1 3	6 2	3 6	4 5
2	6 5	4 6	4 1	5 1	2 3
3	3 4	5 2	3 5	2 4	6 1

the experiment has been conducted so that treatments are confounded with extraneous effects, statistics cannot give reliable estimates of treatment effects. For well-planned and well-executed experiments, the sample of observations will provide an unbiased estimate of treatment effects, with measures of the uncertainty of these estimates.

Interpretation. The statistician should help the experimenter plan his experiment so that the results will answer definite questions. The experimental results should be presented in concise and organized form so that others can draw accurate conclusions. Drawing conclusions without the help of statistics is similar to building a suspension bridge without knowledge of the tensile strength of cables. A perfectly safe bridge can be built without such knowledge, but it will be unnecessarily expensive.

Gertrude M. Cox

Bibliography. M. N. Das and N. C. Giri, *Design and Analysis of Experiments*, 2d ed., 1986; C. R. Hicks, *Fundamental Concepts in the Design of Experiments*, 5th ed., 1999; H. E. Klugh (ed.), *Statistics: Essentials of Research*, 3d ed., 1986; M. Lentner and T. Bishop, *Experimental Design and Analysis*, 1986.

Expert control system

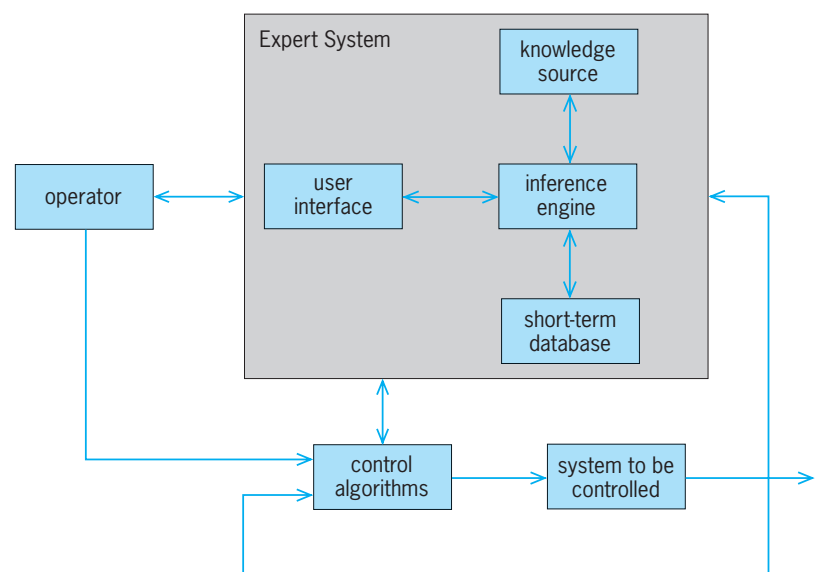
A type of intelligent control system which can emulate the reasoning procedures of a human expert in order to generate the necessary control action. Expert control systems seek to incorporate knowledge about control system design, practical operations, and abnormal system recovery plans to automate tasks normally performed by experienced control engineers (the experts). Techniques relating to the field of artificial intelligence are usually used for the purpose of acquiring and representing knowledge and for generating control decisions through an appropriate reasoning mechanism. As it operates essentially on a knowledge base, an expert control system is often referred to as a knowledge-based control system. One of the most important benefits associated with the use of an expert control system is the inherent capability of the system to deal with uncertainty in information. Information provided to these systems need not be complete or precise. It can be general, qualitative, or vague. These properties do not cause difficulty because, like humans, expert control systems possess functionalities to perceive, reason, infer, and deduce new information. They can learn, gain new knowledge, and improve their performance through experience.

Motivation. Conventional control systems operate, explicitly or implicitly, on the basis of a mathematical model of the system to be controlled. However, for a complex system, a mathematical model representative of all important aspects of the system either does not exist or provides at most an imperfect local approximation to the actual system. When the model is inaccurate, the performance of the control system may not be acceptable. Furthermore, conventional control systems often need a complete set of data to

be available, including sensory information and parameter values, in order to produce control actions. When these data are unavailable or unreliable, possibly due to failure of system components, the control system cannot function adequately. Therefore, in complex engineering systems, where great quantities of interacting numeric, symbolic, and quantitative information must be handled, and abnormal or emergency situations frequently arise, “hard” or “crisp” control with conventional control systems becomes inadequate. Human experts may be more suited to operate these complex systems, but unfortunately this is true only over a limited time span. People are quickly affected by fatigue, stress, emotion, and other environmental factors, which influence their performance to varying degrees. Therefore, for complex systems, intelligent control based on expert systems is useful to provide an efficient mimic of human expert judgment and control, yet with a robust level of proficiency and consistency unattainable with human experts.

Components. The structure of expert control systems may differ slightly due to the different application domains and the correspondingly different requirements. However, the principal components almost certain to be present in all expert control systems are the knowledge source, the database, the inference engine, the control algorithms, and the interface between the expert control system and humans (see *illus.*).

Central to an expert system is the knowledge source (or knowledge base), which contains knowledge in a specific domain (control in this context). This knowledge consists of domain-specific facts and heuristics, usually in the form of rules or frames, useful for solving problems in the domain. The knowledge source is implemented as an identifiable and separate part of the program. The integrity and correctness of the information within the knowledge source is the key to obtaining correct results or



Basic architecture of an expert control system.

solutions from the expert systems.

The database is a short-term memory component which contains the current problem status, inference states, and the history of solutions to date for reference purposes.

The inference engine operates on the information from the knowledge source, from the associated database, or from the user; guides the search process according to the programmed strategy and search algorithms; and uses the inferencing mechanisms, usually in the forms of rules of logic, to solve problems, arrive at conclusions, and activate the final control actions.

The control algorithms are the tools for the expert systems to perform the final control action. A rich library of control algorithms, with advanced features, is usually made available for the expert system to configure and select the appropriate one for commissioning. The algorithms may reside physically in the field controllers outside the expert system, and the control parameters may be downloaded dynamically from the expert system.

The remaining basic component is the user interface, which is provided for the user to interact with the overall system, browse through the knowledge source, edit rules, and perform many other interactive tasks.

Although not strictly a principal component, an expert system shell is an important tool in the development of expert systems. The shell is basically a computer program used to develop expert systems. Often a shell also has provisions for changing the reasoning processes of the inference engine to adapt to the specific problem. Expert system shells differ significantly from each other and offer the user a wide variety of capabilities.

A well-developed expert system will also have a degree of self-awareness or self-knowledge which allows reasoning and explanation of the decisions to the user. This ability is one of the most valuable features of expert systems. The user can request a complete trace for a consultation, and an explanation as to how the conclusion was inferred and why any particular information was needed. This feature allows the verification and validation of expert systems and improves their overall reliability and efficiency.

Applications. Expert systems have found application in control systems, design and engineering, computerized tasks, manufacturing, finance, science and medicine, geological work, and training. One prominent successful application of expert control systems is in process control. Most process control systems are typical large-scale complex engineering systems, where multiple control loops (there can be several thousand in an enterprise), with interacting and time-varying dynamics, are run concurrently. The processes are often subject to unpredictable material variations, extraneous disturbances, and inter-process loading effects. The control system must maintain an efficient level of performance in the face of these difficulties. This is one application domain naturally akin to the functions of expert systems. In process control, the expert control system

typically performs the following functions:

1. *Process operations and monitoring:* Supervise the normal operation of the system (process) and field controller, monitor the relevant signals, and take corrective actions when necessary.

2. *Diagnostics:* Examine possible failure or fault of the system components, replace these faulty components, or revise control algorithms to maintain the necessary performance of the system. See FAULT ANALYSIS.

3. *Control configuration:* Select suitable control strategies for specific situations to adapt to the variation of the system parameters and the environment.

4. *Scheduling and planning:* Coordinate the capabilities or individual controllers within the system to optimize production or increase efficiency. See ARTIFICIAL INTELLIGENCE; CONTROL SYSTEMS; EXPERT SYSTEMS; PROCESS CONTROL.

T. H. Lee; K. K. Tan; C. C. Hang; K. C. Tan; L. C. Woon
Bibliography. K. J. Åström, J. J. Anton, and K.-E. Årzén, Expert control, *Automatica*, 22:276-286, 1986; C. W. de Silva and T. H. Lee, Knowledge-based intelligent control, *Meas. Control*, 164:102-113, 1994; T. H. Lee, Q. G. Wang, and K. K. Tan, A knowledge-based approach to dead-time estimation for process control, *Int. J. Control*, 61:1045-1072, 1995.

Expert systems

Methods and techniques for constructing human-machine systems with specialized problem-solving expertise. The pursuit of this area of artificial intelligence research has emphasized the knowledge that underlies human expertise and has simultaneously decreased the apparent significance of domain-independent problem-solving theory. In fact, a new set of principles, tools, and techniques have emerged that form the basis of knowledge engineering.

Expertise consists of knowledge about a particular domain, understanding of domain problems, and skill at solving some of these problems. Knowledge in any specialty is of two types, public and private. Public knowledge includes the published definitions, facts, and theories that are contained in textbooks and references in the domain of study. But expertise usually requires more than just public knowledge. Human experts generally possess private knowledge that has not found its way into the published literature. This private knowledge consists largely of rules of thumb or heuristics. Heuristics enable the human expert to make educated guesses when necessary, to recognize promising approaches to problems, and to deal effectively with erroneous or incomplete data. The elucidation and reproduction of such knowledge are the central problems of expert systems.

Importance of expert knowledge. Researchers in this field suggest several reasons for their emphasis on knowledge-based methods rather than formal representations and associated analytic methods. First, most of the difficult and interesting problems do not have tractable algorithmic solutions. This is

reflected in the fact that many important tasks, such as planning, legal reasoning, medical diagnosis, geological exploration, and military situation analysis, originate in complex social or physical contexts, and generally resist precise description and rigorous analysis. Also, contemporary methods of symbolic and mathematical reasoning have limited applicability to the expert system area; that is, they do not provide the means for representing knowledge, describing problems at multiple levels of abstraction, allocating problem-solving resources, controlling cooperative processes, and integrating diverse sources of knowledge in inference. These functions depend primarily on the capacity to manipulate problem descriptions and apply relevant pieces of knowledge selectively. Current mathematics offers little help in these tasks.

The second reason for emphasizing knowledge is pragmatic: human experts achieve outstanding performance because they are knowledgeable. If computer programs embody and use their knowledge, they too attain high levels of performance. This has been proved repeatedly in the short history of expert systems. Systems have attained expert levels in several tasks: mineral prospecting, product configuration, chemical structure elucidation, symbolic mathematics, chess, medical diagnosis and therapy, electronics analysis, manufacturing control, and credit analysis, among others.

The third motivation for focusing on knowledge is the recognition of its intrinsic value. Knowledge is a scarce resource whose refinement and reproduction creates wealth. Traditionally, the transmission of knowledge from human expert to trainee has required education and internship periods ranging from 3 to as much as 20 years. By extracting knowledge from human beings and transferring it to computable forms, the costs of knowledge reproduction and exploitation can be greatly reduced. At the same time, the process of knowledge refinement can be accelerated by making the previously private knowledge available for public test and evaluation.

In short, expert performance depends critically on expert knowledge. Because knowledge provides the key ingredient for solving important tasks, it reflects many features characteristic of a rare element: it justifies possibly expensive mining operations; it requires efficient and effective technologies for fashioning it into products; and a means of reproducing it synthetically would be "a dream come true." See DATA MINING.

Distinguishing characteristics. Expert systems differ in important ways from both conventional data-processing systems and systems developed by workers in other branches of artificial intelligence. In contrast to traditional data-processing systems, artificial intelligence applications exhibit several distinguishing features, including symbolic representations, symbolic inference, and heuristic search. In fact, each of these characteristics corresponds to a well-studied core topic within artificial intelligence. Often a simple artificial intelligence task may yield to one of the formal approaches developed for these core problems. Expert systems differ from the broad

class of artificial intelligence tasks in several regards. First, they perform difficult tasks at expert levels of performance. Second, they emphasize domain-specific problem-solving strategies over the more general, "weak" methods of artificial intelligence. Third, they employ self-knowledge to reason about their own inference processes and to provide explanations of justifications for the conclusions they reach. As a result of these distinctions, expert systems represent an area of artificial intelligence research with specialized paradigms, tools, and system-development strategies.

Accomplishments. The first laboratory examples of expert systems appeared during the latter half of the 1970s, the first commercial applications appeared in 1981-1983, and by 1988 the number of applications in daily use worldwide was estimated to exceed a thousand. Since 1988, the number of applications has grown considerably. During this period, expert systems demonstrated relevance to a broad class of industrial and commercial problems, particularly order entry and configuration, field service, eligibility analysis and candidate screening, and process control. In addition, many companies claimed to have developed both cost-saving and strategic types of applications.

The following were among the first applications to demonstrate the practical utility of expert systems: PROSPECTOR discovered a molybdenum deposit whose ultimate value will probably exceed \$100 million; XCON (originally R1) configured customer requests for computer systems; DENDRAL supported hundreds of international users daily in chemical structure elucidation; CADUCEUS embodied more knowledge of internal medicine than any human (approximately 80% more) and could correctly diagnose complex test cases that baffled experts; and PUFF integrated knowledge of pulmonary function disease with a previously developed domain-independent expert system for diagnostic consultation and routinely provided expert analyses at a California medical center.

Among later applications, the following illustrate the diversity and importance of commercial expert systems: ESPM remotely monitored the error logs of thousands of computers for signs of impending disk failures and scheduled appropriate maintenance preventively; ACE performed a function similar to ESPM for scheduling maintenance of telephone cables; BRUSH accepted product performance requirements from an engineer and automatically designed, drew, and documented the appropriate electric motor parts; MIX used expert rules of thumb to blend ingredients into finished products that satisfied customer requirements and production constraints; and CHARLIE collected physical vibration data from rotating equipment and diagnosed all types of mechanical problems. In each of these applications, the expert systems met the high performance requirement of equaling or surpassing the capabilities of highly trained human experts.

Undoubtedly, most applications that developers label expert systems would not generally meet this

Generic categories of knowledge engineering applications	
Category	Problem addressed
Interpretation	Inferring situation descriptions from sensor data
Prediction	Inferring likely consequences of given situations
Diagnosis	Inferring system malfunctions from observables
Design	Configuring objects under constraints
Planning	Designing actions
Monitoring	Comparing observations to plan vulnerabilities
Debugging	Prescribing remedies for malfunctions
Repair	Executing a plan to administer a prescribed remedy
Instruction	Diagnosing, debugging, and repairing students' knowledge weaknesses
Control	Interpreting, predicting, repairing, and monitoring system behaviors

same high performance standard. Instead, companies have developed many so-called competence systems that automate more-routine decision-making. The appeal of the expert systems approach for these applications derives from three factors. First, the rule-based formalism of expert system shells makes it easy to express simple judgmental systems. Second, many knowledgeable people want to automate and distribute their particular know-how. Third, the computing equipment required to run the development tools and applications has become nearly ubiquitous.

Types of systems. Most of the knowledge-engineering applications fall into a few distinct types (see table).

Interpretation systems infer situation descriptions from observables. This category includes surveillance, speech understanding, image analysis, chemical structure elucidation, signal interpretation, and many kinds of intelligence analysis. An interpretation system explains observed data by assigning symbolic meanings to them that describe the situation or system state accounting for the data. *See* CHARACTER RECOGNITION; COMPUTER VISION.

Prediction systems present the likely consequences of a given situation. This category includes weather forecasting, demographic predictions, traffic predictions, crop estimations, and military forecasting. A prediction system typically employs a parametrized dynamic model with parameter values fitted to the given situation. Consequences which can be inferred from the model form the basis of the predictions. *See* WEATHER FORECASTING AND PREDICTION.

Diagnosis systems predict system malfunctions from observables. This category includes medical, electronic, mechanical, and software diagnosis. Diagnosis systems typically relate observed behavioral irregularities with underlying causes by using one of two techniques. One method essentially uses a table of associations between behaviors and diagnoses, and the other method combines knowledge of system design with knowledge of potential flaws in design, implementation, or components to gener-

ate candidate malfunctions consistent with observations.

Design systems develop configurations of objects that satisfy the constraints of the design problem. Such problems include circuit layout, building design, and budgeting. Design systems construct descriptions of objects in various relationships with one another and verify that these configurations conform to stated constraints. In addition, many design systems attempt to minimize an objective function that measures costs and other undesirable properties of potential designs. This view of the design problem can subsume goal-seeking behavior as well, with the objective function incorporating measures of goal attainment. *See* COMPUTER-AIDED DESIGN AND MANUFACTURING.

Planning systems design actions. These systems specialize in problems with the design of objects that perform functions. They include automatic programming, robot, project, route, communication, experiment, and military planning problems. Planning systems employ models of agent behavior to infer the effects of the planned agent activities. *See* ROBOTICS.

Monitoring systems compare observations of system behavior to features that seem crucial to successful plan outcomes. These crucial features, or vulnerabilities, correspond to potential flaws in the plan. Generally, monitoring systems identify vulnerabilities in two ways. One type of vulnerability corresponds to an assumed condition whose violation would nullify the plan's rationale. Another kind of vulnerability arises when some potential effect of the plan violates a planning constraint. These correspond to malfunctions in predicted states. Many computer-aided monitoring systems exist for nuclear power plant, air traffic, disease, regulatory, and fiscal management tasks, but few autonomous expert systems have been reported publicly. Many companies are using applications of this sort either in pilot applications or in advisory capacities.

Debugging systems prescribe remedies for malfunctions. These systems rely on planning, design, and prediction capabilities to create specifications or recommendations for correcting a diagnosed problem. Computer-aided debugging systems exist for computer programming in the form of intelligent knowledge base and text editors, but none of them qualify as expert systems.

Repair systems develop and execute plans to administer a remedy for some diagnosed problem. Such systems incorporate debugging, planning, and execution capabilities. Expert systems participate with human technicians and more traditional computing applications in the domains of automotive, network, avionic, and computer maintenance, as well as others.

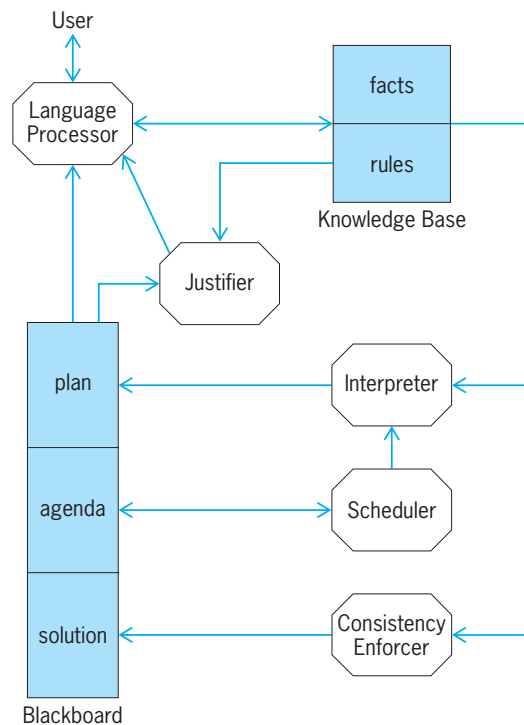
Instructional systems incorporate diagnosis and debugging subsystems that specifically address student behavior. These systems construct a hypothetical description of the student's knowledge. From this model they diagnose weaknesses in the student's knowledge and identify an appropriate remedy. Finally, they select or plan a tutorial interaction

intended to convey the remedial knowledge to the student. Such systems, called intelligent tutoring systems or knowledge-based tutors, typically include additional expert systems that allow the instructional system to solve problems assigned to the student. This domain expertise capability allows the tutor to demonstrate problem-solving skills, answer questions, and compare expert knowledge and problem solving to the inferred model of the student's knowledge and problem solving.

An expert control system adaptively governs the overall behavior of a system. To do this, the control system must repeatedly interpret the current situation, predict the future, diagnose the causes of anticipated problems, formulate a remedial plan, and monitor its execution to ensure success. Problems addressed by control systems include air traffic control, business management, battle management, mission control, and others. Prototype systems developed as part of the U.S. Strategic Computing Initiative have illustrated these capabilities in two different kinds of applications. In the autonomous land vehicle, vision, planning, and real-time control functions were combined to drive a test vehicle at slow speed on both paved streets and open terrain. Although this system achieved unprecedented results in the field of mobile robotics using artificial intelligence techniques, none of the individual components achieved a very high level of expertise, measured by human standards. In the second program, a collection of independent but cooperating expert systems combined to perform the function of a pilot's associate, an artificial "backseater" who would replace the second human in future one-seater aircraft. This kind of application combines independent expert systems addressing separate areas of expertise, including situation assessment, mission planning, tactical planning, system status, pilot-vehicle interface, and overall mission management. Collectively, the pilot's associate subsystems incorporate most of the generic categories of knowledge engineering applications. *See* CONTROL SYSTEMS; INTELLIGENT MACHINE.

Systems components. The ideal expert system (see **illus.**) contains a language processor for problem-oriented communications between the user and the expert system; a blackboard for recording intermediate results; a knowledge base comprising facts plus heuristic planning and problem-solving rules; an interpreter that applies these rules; a scheduler to control the order of rule processing; a consistency enforcer that adjusts previous conclusions when new data or knowledge alters their bases of support; and a justifier that rationalizes and explains the system's behavior.

The user interacts with the expert system in problem-oriented languages, usually some restricted variant of English, and in some cases through means of a graphics or structure editor. The language processor mediates information exchanges between the expert system and the human user. Typically, the language processor dissects, or parses, and interprets user questions, commands, and volunteered infor-



Anatomy of an ideal expert system. No existing system contains all the components shown, but one or more components occur in every system.

mation. Conversely, the language processor formats information generated by the system, including answers to questions, explanations and justifications for its behavior, and requests for data.

Existing expert systems use a wide variety of language processors, depending on the type of user. In many commercial applications, the ideal format for dialog is the business form, depicted graphically on a cathode-ray tube. In military applications, both natural language and annotated maps provide more appropriate languages for system-user dialog. In other applications, the most common languages of dialog are graphic illustrations, such as computer-aided design (CAD) drawings, or semiformal languages, such as part-assembly specifications. These language processors increasingly employ stand-alone interface packages developed specifically for reuse by multiple applications with similar dialog requirements. *See* COMPUTER GRAPHICS; PROGRAMMING LANGUAGES.

The blackboard is a global database that records decisions, which the expert system manipulates. Every expert system uses some type of intermediate decision representation, but only a few explicitly employ a blackboard for the various types of ideal expert system decisions. Three types of decisions that are recorded on the blackboard are plan, agenda, and solution elements. Plan elements describe the overall or general attack the system will use to solve the current problem, including current plans, goals, problem states, and contexts. For example, a plan may recommend processing all low-level sensor data first; formulating a small number of most-promising hypotheses; refining and elaborating each of these hypotheses until the best one emerges; and, finally,

focusing exclusively on that candidate until the complete solution is found. This kind of plan has been incorporated in several expert systems. The agenda elements record the potential actions awaiting execution. These generally correspond to knowledge-base rules that seem relevant to some decision placed on the blackboard previously. The solution elements represent the candidate hypotheses and decisions that the system has generated thus far, along with the dependencies, called links, that relate decisions to one another.

The scheduler maintains control of the agenda and determines which pending action should be executed next. Schedulers may utilize considerable abstract knowledge, such as “do the most profitable thing next” and “avoid redundant effort.” To apply such knowledge, the scheduler needs to prioritize each agenda item according to its relationship to the plan and to other solution elements. To do this, the scheduler generally needs to estimate the effects of applying the potential rule.

The interpreter executes the chosen agenda item by applying the corresponding knowledge-base rule. Generally, the interpreter validates the relevance conditions of the rule, binds variables in these conditions to particular solution blackboard elements, and then makes rule-prescribed changes to the blackboard. Interpreters of this sort are generally written in LISP because of its facilities for manipulating and evaluating programs, but other languages are also suitable. *See* PROGRAMMING LANGUAGES.

The consistency enforcer attempts to maintain a consistent representation of the emerging solution. This may take the form of likelihood revisions when the solution elements represent changing hypothetical diagnoses and when some new data are introduced. Alternatively, the enforcer might implement truth maintenance procedures when the solution elements represent changing logical deductions and their truth-value relationships. Most expert systems use some kind of numerical adjustment scheme to determine the degree of belief in each potential decision. This scheme attempts to ensure that plausible conclusions are reached and inconsistent ones avoided.

The justifier explains the actions of the system to the user. In general, it answers questions about why some conclusion was reached or why some alternative was rejected. To do this, the justifier uses a few general types of question-answering plans. These typically require the justifier to trace backward along blackboard solution elements from the questioned conclusion to the intermediate hypotheses or data that support it. Each step backward corresponds to the inference of one knowledge-base rule. The justifier collects these intermediate inferences and translates them to English for presentation to the user. To answer “why not” questions, the system uses a heuristic variant of this technique. It can identify a possible chain of rules that would reach the questioned conclusion but that did not apply because the relevance condition of some rule failed. The jus-

tifier explains the system’s decision to reject a possible conclusion by claiming that such failed conditions impede all reasoning chains that can support the conclusion.

Finally, the knowledge base records rules, facts, and information about the current problem that may be useful in formulating a solution. While the rules of the knowledge base have procedural interpretations, the facts play only a passive role. The language PROLOG is the most widely available general-purpose programming language that explicitly supports this fact-plus-rule view of programs. However, relatively few expert systems are developed directly in PROLOG. Instead, most developers adopt a higher-level knowledge engineering tool or expert system shell as their programming language. These shells offer additional features for enhanced representations, inference, and control. Some of the most popular enhancements include expressive notations and built-in mechanisms for uncertainty; procedural blocks for more explicit control; classes and type hierarchies for organizing semantic relationships among categories, specializations, and instances; and change-driven control mechanisms to support modeling, reactive systems, and message-passing styles of programming. The portrait of the ideal expert system ignores these embellishments which, although occasionally important, often obscure the common essence of the expert systems technology.

Agents. Agents can vary from simple programs, lacking knowledge, to sophisticated autonomous real-time control systems that provide both deliberative reasoning (for example, including sophisticated planning and scheduling) and fast real-time response. The vast majority of simple utilitarian agents, such as off-line browsers that automatically download web sites, or agents that monitor web pages for changes, are not knowledge-based. Most knowledge-based agents are still experimental.

Agents are variously defined as autonomous programs; programs that communicate with an agent communication language (for example, KQML); mobile programs that can travel from one computer to another; or programs that collaborate with humans in solving tasks, offloading tedious tasks, and acting as monitors and liaisons with other agents. Agents can also be viewed as higher-level architectural components, providing distributed open architectures in which the agents can be registered, requests can be brokered, and different organizational structures can be used to solve higher-level tasks. In this architectural view, they are distinguished from objects by their persistence, autonomy, communications capabilities, and behavior. Another key aspect of most agent definitions is the agent’s dynamic and uncertain environment. Agents also may or may not have graphical representations, commonly called avatars or personas.

Internet agents may make use of information retrieval techniques such as term frequency and inverse document frequency analyses. In these approaches, the frequency of terms in documents

relative to their overall frequency in a document corpus is compared to the same measure for the terms in a query (for retrieval) or in another document (to gauge similarity). Most Internet agents rely more on information retrieval techniques such as these or machine learning techniques (such as collaborative filtering, a variant of the nearest-neighbor algorithm) than the reasoning and representation techniques used in most knowledge-based applications. These Internet agents can provide recommendations, provide personalized searches and web page recommendations, monitor multiple sources of information, negotiate, and so forth, using the Web as the medium they are situated in for perception and action. Stand-alone robotic or software agents designed for entertainment may have a different architecture comprising reactive structures or simple learning architectures (for example, neural networks). See INTERNET; NEURAL NETWORK; WORLD WIDE WEB.

Knowledge-based agents overlap with expert systems: When a knowledge-based program is a real-time autonomous program or if its software architecture comprises a multiple-agent system, these knowledge-based systems are agents or, conversely, the agents are also knowledge-based systems.

An example of a knowledge-based agent is Guardian, an on-line agent designed to operate in an intensive-care unit to monitor and control a ventilator unit. It has sensors and effectors, and operates in a dynamic, uncertain environment. It also has models of human biological responses, disease processes, and intervention strategies to handle crises that arise. It can be viewed as an agent that incorporates multiple expert systems to handle different tasks (diagnosis, monitoring, and repair). Guardian has been tested only in a simulated environment.

The New Millennium Remote Agent (NMRA) is an autonomous expert system and agent that operates in deep space. Developed by the National Aeronautics and Space Administration, it controls the functions of an entire spacecraft. It includes planning and monitoring functions and operates entirely autonomously. Like Guardian, it has both reactive and deliberative capabilities. It has a procedural executive for fast real-time response coupled with a model-based diagnosis system and planner that can be invoked by the executive for recovery from plan execution errors. The model-based diagnosis system is used to infer the state of the satellite components and to provide input to the procedural executive.

Falling between heavy-weight knowledge-based agents such as Guardian and NMRA and simple utilitarian agents that do not rely on knowledge-based techniques, there are agents that rely on a changing knowledge base to provide useful assistance, but do not constitute expert systems in the traditional sense. Rather than expert knowledge, they provide informed guidance. Examples are guides to web sites and virtual guides to software and other repositories.

Another area of overlap between expert systems and agents is the use of pedagogical agents in intelligent tutoring systems. Such systems operate in a

dynamic, uncertain environment, and are, by definition, knowledge-based, having knowledge of the subject matter and the student, and some qualitative model of the teaching process. Many intelligent tutoring systems use graphical personas as pedagogical agents. For example, one program uses a talking character to teach botany to middle-school students, and another uses a virtual-reality instructor to teach troubleshooting and maintenance of Navy propulsion systems. In both cases the tutor agents are knowledge-based: in the first case, the tutor knows about botany, and in the second case, the tutor knows about equipment operation and the current state of its controls and internal parts.

Expert systems and agents can also overlap architecturally: agents can be used as higher-level building blocks to build knowledge-based systems. ARCHON is one example of a real-time expert system that is architecturally composed of agents. It provides a methodology for building a multiagent community that incorporates existing control systems. ARCHON has been operationally deployed in an electrical transmission system in northern Spain. It wraps existing on-line control and expert systems as agents so they can communicate, share knowledge, and act in concert where this was not previously possible.

Similarly, expert systems can be used as building blocks of agents. The NMRA agent can be viewed as a collection of expert systems to handle different spacecraft systems. Intelligent tutoring systems and pedagogical agents typically have one expert system to provide domain expertise, and may incorporate additional expert systems for student modeling and pedagogical control.

Frederick Hayes-Roth; William R. Murray

Evaluation. Expert systems can be evaluated from three perspectives: subjectively by evaluating the system's stated requirements and level of usability, technically by evaluating how well the knowledge base was built, and empirically by evaluating the system's effect on user performance.

Subjective evaluations. The goal of subjective evaluations is to assess the system from the perspective of potential users and sponsors. To accomplish this, users should be involved in specifying the system's performance requirements. Users are increasingly involved in requirements validation, that is, evaluating whether the system is being built to achieve the right performance criteria.

Users are also involved in system evaluation through judging the adequacy of system performance and usability. User involvement may take the form of seeing a general system demonstration, systematically walking through each system feature, or actually using the expert system to solve decision problems representative of their actual job.

Technical evaluations. Static testing focuses on detecting logical problems in the knowledge base without actually executing it by using test cases. Static testing may be performed manually or with automated static testors. The more prevalent logic problems that can be detected include redundant rules,

subsumed rules, conflicting rules, circular rules, unnecessary “if” conditions, and unreferenced or illegal attribute values.

Dynamic testing uses experts who have not participated in development and test cases to evaluate the knowledge base’s functional completeness and predictive accuracy. Ideally, there will be correct answers for the test cases; otherwise, it is necessary to rely on the judgment of the evaluation expert or, preferably, the consensus judgments of a group of experts. If the knowledge base’s predictive accuracy is acceptable, the evaluation experts will be unable to tell whether they or the system generated the best conclusions. This approach is referred to as a Turing test.

Empirical evaluations. Experiments help to generalize from a test sample to the larger population. There are two kinds of experiments. The first tests the system against objective benchmarks representing performance constraints. If the system passes, it proceeds further; if it fails, it undergoes further development or is set aside.

The second kind of experiment is a factorial design where one or more factors are systematically varied as the independent variables, and system performance is measured by one or more dependent variables. Dependent variables include objective measures (such as performance and speed), observational measures (such as regarding how the system is used), and subjective measures (such as user confidence in the solution). The goal is to show that people using the expert system perform better than those who have access only to their typical tools.

Ideally, field experimentation would be used to assess if the expert system significantly improves performance in the targeted setting. When the sample size and randomization requirements of true experiments are not possible, quasi-experimental designs can be used. These include time-series designs, where the performance of an organizational unit is measured before and after receiving the system; multiple time series designs, which use a control group that does not receive the system; and nonequivalent (and nonrandomized) control-group designs, which rely on statistical techniques to compare performance improvement in expert-system and control groups. See ARTIFICIAL INTELLIGENCE.

Leonard Adelman

Bibliography. L. Adelman, *Evaluating Decision Support and Expert Systems*, 1992; J. C. Giarratano and G. Riley, *Expert Systems: Principles and Programming*, 3d ed., 1998; U. Gupta (ed.), *Validating and Verifying Knowledge-Based Systems*, 1991; B. Hayes-Roth, An architecture for adaptive intelligent systems, *Artif. Intell.*, 75:195–240, 1995; F. Hayes-Roth and N. Jacobstein, The state of knowledge-based systems, *Commun. ACM*, 37(3):27–39, 1994; W. L. Johnson (ed.), *Proceedings of the 1st International Conference on Autonomous Agents*, ACM Press, 1997; S. C. Shapiro (ed.), *Encyclopedia of Artificial Intelligence*, 2 vols., 2d ed., 1992; D. A. Waterman, *A Guide to Expert Systems*, 1986.

Explosive

A substance containing a large amount of stored energy that can be released suddenly, thereby converting the substance into compressed gases or numerous fragments that expand with great force or velocity. An explosion is a sudden expansion of matter into a much larger volume than it formerly occupied, or a sudden increase in the pressure exerted by confined matter.

Explosions

Examples of explosions are expansion of highly compressed gas and throwing of fragments following brittle rupture of a pressurized container; rapid combustion of a fuel-air mixture in an internal combustion engine; deflagration of smokeless powder in a gun followed by sudden expulsion of the projectile and gases from the muzzle; rapid deflagration of a gas-producing pyrotechnic mixture such as black powder; detonation of a charge of chemical explosive such as TNT; release of nuclear energy in a mass of nuclear explosive such as uranium-235 with sufficient rapidity to suddenly vaporize it and any adjacent material; explosive eruption of a volcano; and explosion of a star to produce a supernova. See INTERNAL COMBUSTION ENGINE; NUCLEAR EXPLOSION; SUPERNOVA; VOLCANO.

Blast and shock. Divergent motion, in the form of a pressure wave of high amplitude above ambient pressure (that is, high overpressure), propagates outward from the source of an explosion into the surrounding medium, typically air, water, or rock. In those parts of the wave where the overpressure is low, the velocity of wave propagation is not much higher than the velocity of sound in the unperturbed medium. Here, the velocity of the medium itself (the particle velocity) is relatively low. In those parts of the wave where the overpressure is high, the sudden compression has increased the temperature and the sound velocity in the medium, and here the particle velocity is also relatively high. The outward velocity of a local perturbation in the wave is the sum of the sound velocity and the particle velocity. Therefore, perturbations and the energy they carry move outward faster in the high-pressure parts of the wave, thereby changing the shape of the wave. Consequently, if the amplitude of the wave is sufficiently high, it finally evolves into one or more shock waves of characteristic shape (**Fig. 1a** and **b**) regardless of the initial wave form.

A shock wave has a shock front with a typical thickness of about 3×10^{-7} ft (1×10^{-7} m) which travels at supersonic velocity and at which the pressure jumps almost instantaneously from ambient pressure to a value high above ambient. Immediately ahead of the shock front, the medium is usually at rest, while immediately behind it the medium is moving forward with a high particle velocity. The kinetic energy of this motion, per unit volume of the medium, is called the dynamic pressure. The overpressure of a shock wave decays with the distance traveled, not only due to divergence, which causes its energy to be spread

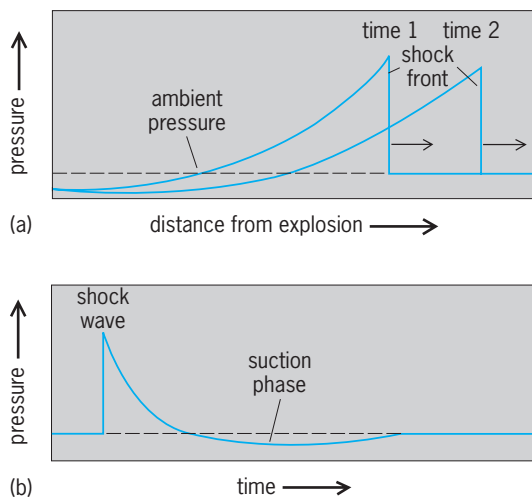


Fig. 1. Pressure waves from explosions. (a) Pressure versus distance near the shock front at two times. (b) Pressure versus time in air at one distance.

over a surface of increasing size, but also due to an irreversible degradation of its energy to heat resulting from the discontinuous jump in pressure at the shock front. Its velocity of propagation approaches that of sound (about 1050 ft/s or 320 m/s in air, and about 5000 ft/s or 1520 m/s in water) as the overpressure approaches zero. See SHOCK WAVE.

In air, the particle velocity of the pressure wave is sensed as a sudden blast of wind, so the wave is sometimes called a blast wave whether or not it is also a shock wave.

Effects of scale. The smaller the amount of energy released, the shorter must be the time and the smaller the volume in which it is released in order to obtain the transient effects of an explosion. Conversely, as the amount of energy increases, explosion effects can be produced with larger durations and larger volumes for the energy release. The release of energy may produce a pressure wave that never develops into a shock wave because the wave has too slow a rise time, too low an overpressure, or is in a medium that rapidly degrades much of the energy of deformation into heat. Such an energy release may nevertheless produce motion that is sufficiently sudden to produce damage or other mechanical effects associated with explosions, and may therefore be classified as an explosion.

Under certain conditions, it is possible to use relatively simple relationships to describe the behavior of explosions over a very large range of energies. Necessary conditions are that the effects of radiation have negligible effect on the explosion phenomenon of interest and that the material in which the explosion occurs not be changed. Then the desired simple relationships are obtained when the independent variable is the scaled time $tW^{1/3}$ or the scaled distance $DW^{1/3}$, and the dependent variable is a gravity-independent variable, where t is the time after explosion, D the distance from the explosion, and W the explosion energy. W is usually expressed as the mass of TNT having the same energy yield as the explosion.

When gravity effects cannot be neglected, as for explosion phenomena resulting from heaving rock, soil, or water upward against gravity or against centrifugal force, similar results will also be obtained for explosions at equal values of the scaled time $tW^{-1/3}$ and scaled distance $DW^{-1/3}$, provided that equal values of tg and Dg are also maintained, g being the gravitational or centrifugal acceleration of the explosion environment.

Such scaling laws permit the effects of very large explosions to be predicted from the behavior of very small explosions. In some cases, these may be made in a centrifuge to obtain a high value of g commensurate with the small values of t and D .

Damage to structures. The damage done to a structure by an explosion or a rapid series of them can be the result of pressure exerted against the outside or inside of the structure, impulse delivered to it, or vibrational energy imparted to it.

When a shock wave in air strikes a solid surface head-on, part of the wave is reflected. The change in momentum associated with the resulting impact of the moving air on the surface can generate an overpressure that is more than twice that of the incident wave, and for very strong shocks and very rigid surfaces can be almost eight times as great. If the natural period of vibrations of a structural member is very short compared to the duration of the positive phase, the damage-causing potential of the explosion depends on the peak pressure or the reflected peak pressure appropriate for the angle of incidence of the wave on the member. This is commonly the case for very large explosions, such as nuclear explosions, or for very thin, light members, such as windowpanes. If the natural period of vibration of a member is long compared to the duration of the positive phase, the damage-causing potential of the explosion tends to depend on the positive-phase overpressure impulse per unit area.

A pressure wave can be described as a sum of sinusoidal waves of various frequencies and amplitudes. If the arriving wave is not a single pulse but a wave train with multiple peaks resulting from a series of explosions or the effects of the intervening media that transmitted it, then energy carried by the wave train when plotted as a function of frequency will have many high peaks and deep valleys. If one or more peaks of sufficient height in this power spectrum contain frequencies that match natural frequencies of vibration of the structure, the wave train may excite vibrations in the structure of sufficient amplitude to cause damage. The approximate vibration levels that buildings of various types can tolerate without damage are known, and industrial blasts near buildings are designed so as to stay well below those levels. In rock blasting, which is commonly done by detonating a number of charges with short delays between detonations, the delays are chosen to maximize rock breakage and minimize air blast and ground vibration.

Even a mild explosion that produces no easily perceptible wave in the open air can be very damaging if it occurs inside an enclosed structure, such as a

building, where rapid venting cannot occur. Most structures of this type will fail at internal pressures easily generated by the burning of a combustible dust or vapor in air.

Combustion and detonation. Chemical explosives are metastable compositions that can react rapidly within themselves to produce large quantities of hot gas. The reaction is generally initiated by heat or by mechanical deformation of the explosive that produces heat. Once initiated, the reaction may proceed through the mass of explosive at a relatively slow (subsonic) rate, called combustion, burning, or deflagration, or at a relatively rapid (supersonic) rate, called detonation. In combustion, the rate at which the reaction front advances into the unreacted explosive is governed by the conduction and convection of heat from the hot products of reaction to the unreacted explosive.

In detonation, the rate of advance of the reaction front is the supersonic velocity of a shock wave, called the detonation wave, driven by the high-pressure products of reaction. The chemical reaction is initiated in fresh material by the mechanical deformation produced by this shock wave. The velocity of a detonation wave is determined by the energy released in the chemical reaction, the equation of state of the reaction products, and the constraints of conservation of mass, momentum, and energy. A detonating explosive is classified as a primary explosive if heat alone will cause a small unconfined volume of it to detonate, and as a secondary explosive if not. Primary explosives can be detonated by an influence that will produce even a minuscule hot spot, such as contact with a small flame or hot particle, a static electrical discharge, or a relatively low level of friction or impact.

Combustion of an explosive can accelerate and transform into a detonation, and vice versa. Many of the damage effects of explosions can be produced by sufficiently rapid combustion as well as by detonation, but detonations can produce much higher peak pressures with a more sudden onset than combustion, and damage effects characteristic of such suddenly applied high pressure. *See* COMBUSTION; FLAME.

The propellants used in guns and rockets are explosives whose combustion under the conditions of use occurs in a reproducible and controlled manner without converting to a detonation, which would burst the chamber. Explosives such as propellants that combust in normal use are called low explosives; those that detonate are called high explosives. But given suitable conditions, low explosives will detonate, and high explosives will combust, even though these conditions may not occur in normal practice.

Double-base powder, the propellant most commonly used in artillery and firearms, is made from nitrocellulose with various amounts of glycerine trinitrate (nitroglycerine) and the nitrates of mono-, di-, and triethyleneglycol, plus other ingredients that modify the burning rate, prolong shelf life, reduce barrel erosion, and reduce muzzle flash and smoke. When a substantial amount of nitroguani-

dine is added to such propellants in order to reduce muzzle flash and barrel temperatures, they are called triple-base powders. The term powder is often used for explosives and propellants, whether or not their physical form is powder.

Modern solid rocket propellants are also often of the double-base type, but those with the highest performance may be made from high-explosive compounds such as RDX, sometimes with an added high-energy fuel (usually powdered aluminum), an oxidizer (usually ammonium perchlorate), an uncross-linked monomer, a cross-linker, and other minor additional additives. After the viscous mixture is poured into a mold or rocket casing, the monomer cross-links to form a solid elastomer, resulting in a body of high-performance solid propellant that will yield rather than fracture under the transient stresses imposed during initiation and combustion. A fracture could produce a disastrous increase in the area of the burning surface.

Although solid propellants are inherently detonable, the shock or chamber pressures required to initiate detonation in them are so high that in practice they deflagrate rather than detonate. *See* PROPELLANT.

No sharp distinction exists between energy-releasing compositions that are explosives or can produce an explosion, and those that are not and cannot. The term pyrotechnics covers a continuum of such compositions, ranging from high explosives and low explosives, through fuel/oxidant gaseous mixtures, incendiaries that burn with a flame but do not explode, to mixtures of solids that react to produce high temperatures but little or no gas or flame and no explosive effects. *See* PYROTECHNICS.

Combustion and detonation fronts move through a mass of explosive at rates typically within the ranges 3×10^{-2} to 3×10^2 ft/s (10^{-2} to 10^2 m/s) and 3×10^3 to 3×10^4 ft/s (10^3 to 10^4 m/s), respectively. For both modes of reaction, the velocity of propagation tends to increase with the energy of the explosive, the amount of confinement provided by the surroundings, and the diameter of the explosive mass normal to the direction of propagation. For combustion, the velocity also tends to increase with the surface area of the mass of explosive and its permeability to gas. For detonation, it increases with the density of the mass of explosive if the charge has a sufficiently large diameter.

The velocity s with which the combustion front moves into a single grain of a typical propellant follows a relationship of the form $s = AP^n$, where P is the gas pressure and A and n are constants characteristic of the propellant. The propellant is generally formulated to have a pressure index n less than 1 to prevent the burning rate from accelerating to a level that could burst the chamber. To control the burning rate, propellant grains are also shaped to maintain a nearly constant surface area as they burn (**Fig. 2**).

Gas and dust explosions. Explosions that result from the ignition of a mixture of air with fuel in a confined space are a frequent source of industrial accidents. It is possible for such explosions to be

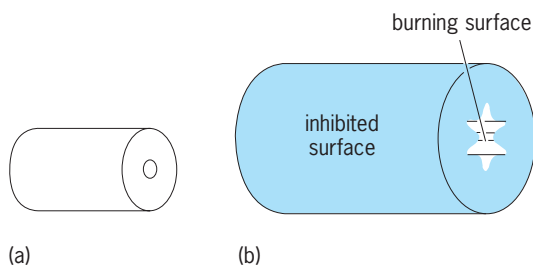


Fig. 2. Propellant grains. (a) A single powder grain in the form of a perforated cylinder, which may have a wall thickness (or web) from less than 0.04 in. (1 mm) to 0.25 in. (6 mm). (b) A rocket grain, which may be several feet or more in diameter, depending on the size of the rocket.

detonations (the knocking of an internal combustion engine being an example), but they are usually much slower, burning reactions. The damaging explosion pressure results from the confinement of the gases as they are heated by the reaction.

Environments that are potentially explosive include storage bins, tanks, exhaust ducts, and underground spaces containing combustible vapors or dusts. Such environments are often encountered in grain elevators, ships, manufacturing spaces, and coal mines, where flammable gases or powders are handled or created. Examples of gases that frequently create such a potential hazard are methane and hydrogen; examples of dusts are starch, synthetic-resin dust, coal dust, and powdered aluminum, magnesium, and zirconium. Dust and dusty air can generate static voltages and electrical discharge when they flow, particularly at low humidity. Such discharges have been sources of ignition for many dust explosions.

The maximum pressure that can be produced is limited by the amount of heat and the resulting pressure that can be produced by the combustion of all the available oxygen in the air. This is about 10 bars (1 megapascal) for conventional fuels in air initially at 1 bar (10^5 pascals). The volume in which the explosion occurs and the surfaces against which the pressure acts are usually very large, so the total energy and force mobilized in such an explosion can be very damaging. Ease of ignition and burning rate of dusts increase with decreasing particle size and vary greatly with their composition. Similarly, the ease of ignition of gas mixtures varies greatly with the nature of the gases and their proportions, ignition being possible at all only within certain limits of composition. The most sensitive methane-air mixture requires about 20 times the energy for ignition as the most sensitive hydrogen-air mixture. An air-methane explosion in a gassy part of a coal mine can produce a pressure wave and flame front that first disperses coal dust settled on beams, ledges, and floors and then ignites it, producing an explosion that can propagate for long distances through the mine. The ignitability of coal dust is reduced by diluting it with powdered limestone spread where the dust collects. The ignitability of methane is reduced by maintaining good ventilation. Other practices that reduce the possibility of a mine explosion include use of ex-

plosives formulated to be able to detonate in an air-methane mixture without igniting it, use of water-filled bags rather than coal dust to stem the charges, and use of electrical devices designed to function without producing sparks or flame. See MINING.

Types of Explosives

For reasons of safety as well as economics, large explosive charges are usually made of explosives that are relatively difficult to initiate. The amount of very sensitive explosive that is easy to initiate is kept to a minimum by using only enough of it to reliably initiate a larger charge of a less sensitive explosive, and using that to initiate a still larger charge of a still less sensitive explosive.

Either failure to detonate or premature explosion can be disastrous. Therefore, it is important that explosives be formulated, manufactured, and used so as never to explode accidentally and to detonate with high reliability when desired. Achieving this balance is an overriding consideration in all work with explosives, both industrial and military.

Industrial explosives. Black powder (a mixture of charcoal, sulfur, and potassium nitrate) was the only industrial explosive until Alfred Nobel's introduction of liquid nitroglycerine for rock blasting in 1863. Black powder is still used extensively in pyrotechnic devices. The modern era of industrial explosives began with Nobel's invention of the detonator in 1865 and of the first dynamites and blasting gelatin in the period 1867-1875. These inventions were a great and lasting contribution to the world's industrial development because dynamite had much greater safety and shattering power than black powder, and much lower cost.

Uses of industrial explosives include blasting ore, coal, and rock in mining and construction, generating vibrations in seismic prospecting for oil and gas, stimulating and perforating gas and oil wells, bonding sheets of dissimilar metals to each other, and synthesizing industrial diamonds.

Rock, ore, and coal are blasted by drilling them with a pattern of holes, filling most of each hole with explosive containing a detonator, plugging the exit of the hole with a material such as drill cuttings or bags of water (stemming), and detonating the charges in a rapid and carefully timed sequence that results in efficient rock breakage and low vibration of nearby structures.

In contrast to military explosives, industrial explosives tend to have larger critical diameters, lower densities, lower detonation velocities, and lower explosion pressures; and to have more complex compositions and lower cost. Most military explosives are rigid solids, whereas most industrial explosives are formulated to be plastic, pumpable, or free-flowing to permit filling the cross section of deep, rough, and irregular holes in rock.

Industrial explosives usually contain separate fuel and oxidizer ingredients in intimate combination. Usually, they also contain a sensitizer to aid in the initiation and propagation of detonation. Other ingredients may be used to increase or decrease density, to

increase explosion energy, to prevent the detonation from igniting methane in coal mines, to provide plasticity, pumpability, or flowability, to prevent setting or stiffening during storage or at low temperature, to prevent separation of ingredients and chemical instability during storage, and to provide resistance to desensitization by water, low temperature, hydrostatic pressure, or transient pressure from explosions in nearby holes. Industrial explosives are usually packaged in bags or cartridges of polymer film or paper, or are carried in bulk form to the blasting site, where they are blown or pumped into the holes through hoses. Sometimes they are mixed at the blasting site.

To avoid toxic fumes underground or fumes that can give secondary explosions when mixed with air, industrial explosives are usually oxygen-balanced to minimize carbon monoxide (CO), nitrogen oxides (NO_x), and hydrogen (H₂) in the reaction products. In a typical modern dynamite, nitroglycerine containing a nitrocellulose thickener is adsorbed on a mixture of ammonium nitrate and cellulosic fuel to form a cohesive, plastic mass that is packaged in paper cartridges. The composition, including the cartridge, is approximately oxygen balanced; that is, the oxygen balance is about zero.

In an oxygen-balanced composition, the amount of oxygen present is just sufficient to oxidize all the carbon to CO₂ and all the hydrogen to water (H₂O) and any metals present to their oxides. If the oxygen is insufficient to do this, the oxygen balance is said to be negative, and if it is more than sufficient, it is said to be positive. Explosives containing only carbon (C), H, nitrogen (N), and oxygen (O) exhibit their maximum energy at a slightly negative oxygen

balance and a minimum production of toxic NO_x and CO at about zero oxygen balance.

Most industrial blasting is done with ANFO, a free-flowing mixture of ammonium nitrate (AN) in the form of small spheres (prills) and 6% No. 2 diesel fuel oil, sometimes with aluminum powder to increase the energy. Ammonium nitrate prills are made by the solidification of droplets of molten ammonium nitrate during free fall in a shot tower. ANFO is a very low-cost, oxygen-balanced explosive with adequate sensitivity and energy for most industrial blasting in dry holes of relatively large diameter. Prilled ammonium nitrate alone is a relatively weak explosive having a positive oxygen balance. It is too difficult to initiate and has too large a critical diameter to be a useful explosive in itself. See PRILLING.

Aqueous explosives are used for much blasting where the holes have small diameter or contain water, or have too thick a burden of rock for ANFO to break. These compositions contain an aqueous phase saturated with inorganic nitrates or perchlorates, one or more fuels such as ground coal, starch, oil, wax, aluminum, or sugar, and a sensitizer such as trinitrotoluene (TNT), monomethylammonium nitrate, pigmentary aluminum, microscopic hollow glass spheres, or bubbles. This fluid mixture is thickened, gelled, or emulsified through the use of other ingredients. The resulting semisolid consistency prevents separation of the ingredients, provides water resistance, and prevents loss of explosive through fissures in the rock. Aluminum (Al) fuel, when used in sufficient quantities to leave little or no water in the explosion products, can be particularly effective in giving aqueous compositions a very high energy

Combustion and detonation parameters for various explosives

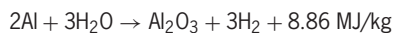
Composition	Charge density, Mg/m ³	Detonation velocity, ft/s*	Chapman-Jouguet pressure, kPa	Approximate critical diameter, in. (mm)
2H ₂ + O ₂ (moles of gas at 1 bar and 20°C)		9.28 × 10 ³	1.80 × 10 ³	
2H ₂ + O ₂ + 5N ₂ (moles of gas at 1 bar and 20°C)		5.97 × 10 ³	1.44 × 10 ³	
Nitroglycerine (liquid at 20°C)	1.60	25.3 × 10 ³	2.53 × 10 ^{7†} , 2.51 × 10 ^{7‡}	0.08 (2)
TNT (liquid at 100°C)	1.44	21.4 × 10 ³		1.6–2.4 (40–60)
TNT (particulate)	0.95	15.9 × 10 ³		0.08–0.9 (2–23)
TNT (cast)	1.62	22.1 × 10 ³	2.10 × 10 ^{7†} , 2.23 × 10 ^{7‡}	0.47–1.4 (12–35)
RDX	1.77	28.8 × 10 ³	3.38 × 10 ^{7†} , 3.48 × 10 ^{7‡}	0.2 (5)
PETN	1.76	27.1 × 10 ³	3.35 × 10 ^{7†} , 3.32 × 10 ^{7‡}	
Tetryl	1.71	25.8 × 10 ³	2.60 × 10 ^{7‡}	0.43 (11)
Lead azide	3.8	18 × 10 ³		0.02 (0.5)
Black powder	1.0	4.43 × 10 ³		
40% special gelatin dynamite	1.59	16.4 × 10 ³	8.10 × 10 ⁶	0.55 (14)
ANFO made with prilled AN	0.85	15.6 × 10 ³	4.82 × 10 ⁶	4 (100)
Amatol 80/20	1.6	17 × 10 ³		3 (80)
Cyclotol 75/25	1.76	27.2 × 10 ³		0.16 (4)
Pentolite 50/50	1.70	24.7 × 10 ³	2.55 × 10 ^{7‡}	0.3 (7)
Composition B-3	1.72	25.9 × 10 ³	2.87 × 10 ^{7†}	0.16 (4)
HBX-3	1.84	23.4 × 10 ³		
PBX-9407	1.60	26.0 × 10 ³	2.87 × 10 ^{7†} , 3.00 × 10 ^{7‡}	

* 1 ft/s = 0.3048 m/s.

† Measured.

‡ Calculated.

from the reaction below, with all carbon appearing



in the products as CO. See EMULSION; GEL.

Military explosives. Most military explosives are simple compositions formulated for high energy density, loading in munitions plant, and long storage life. Most of them are based on explosive chemical compounds that incorporate both oxidizer and fuel components in the same molecule. Many of these compounds are used for special purposes in industrial explosives also (see table).

TNT, which was once used extensively as a military explosive, either alone or in Amatol (80% AN, 20% TNT), is a relatively insensitive and inexpensive explosive of moderate energy. TNT melts at 177°F (81°C), which allows it and compositions containing it to be easily melt-cast. It is mainly used in admixture with more powerful explosives such as RDX and PETN. Examples are Cyclotol 75/25 (75% RDX, 25% TNT), composition B-3 (60% RDX, 40% TNT), HBX-3 (31% RDX, 29% TNT, 35% Al, 5% wax), and Pentolite 50/50 (50% PETN, 50% TNT; Fig. 3).

RDX (hexagen) is a fairly insensitive explosive of high density and melting point with the highest detonation pressure of any compound in common use. It is often combined with an elastomeric or polymeric binder (Fig. 3). An example is PBX-9407 (94% RDX, 6% vinyl chloride-chlorotrifluoroethane copolymer). HMX is an analog of RDX, having an eight-membered rather than a six-membered ring, and has a higher density and therefore a higher detonation pressure.

PETN has moderate sensitivity, a small critical diameter, and a high detonation pressure, which makes it useful in detonating cord (Fig. 3). Detonating cord is flexible line having an explosive core covered sometimes with lead but usually with a layer of fibers and plastic or wax. It is used to transfer a detonation from one point to another, both in military devices and in commercial blasting.

Tetryl, like PETN, is a relatively sensitive explosive, but it is easier to press into pellets for use as booster charges. Booster charges are detonated with a still smaller charge of an even more sensitive explosive, such as lead azide, and in turn they detonate the main charge consisting of a relatively insensitive explosive.

HNS (Fig. 3) is a secondary explosive, having unusually high stability at elevated temperatures and the unusually small critical diameter of 0.020 in. (0.5 mm). TATB (1,3,5-triamino-2,4,6-trinitrobenzene) is another secondary explosive. It is used in munitions because it has an unusually high resistance to transition from deflagration to detonation in a fire.

Mercury fulminate (Fig. 3) is the primary explosive that was used in the first dynamite detonators (blasting caps). It has since been replaced by explosives such as lead azide, lead styphnate (2,4,6-trinitroresorcinate), and lead mononitroresorcinate, which have greater stability at high temperatures or are less subject to failure when pressed to high density.

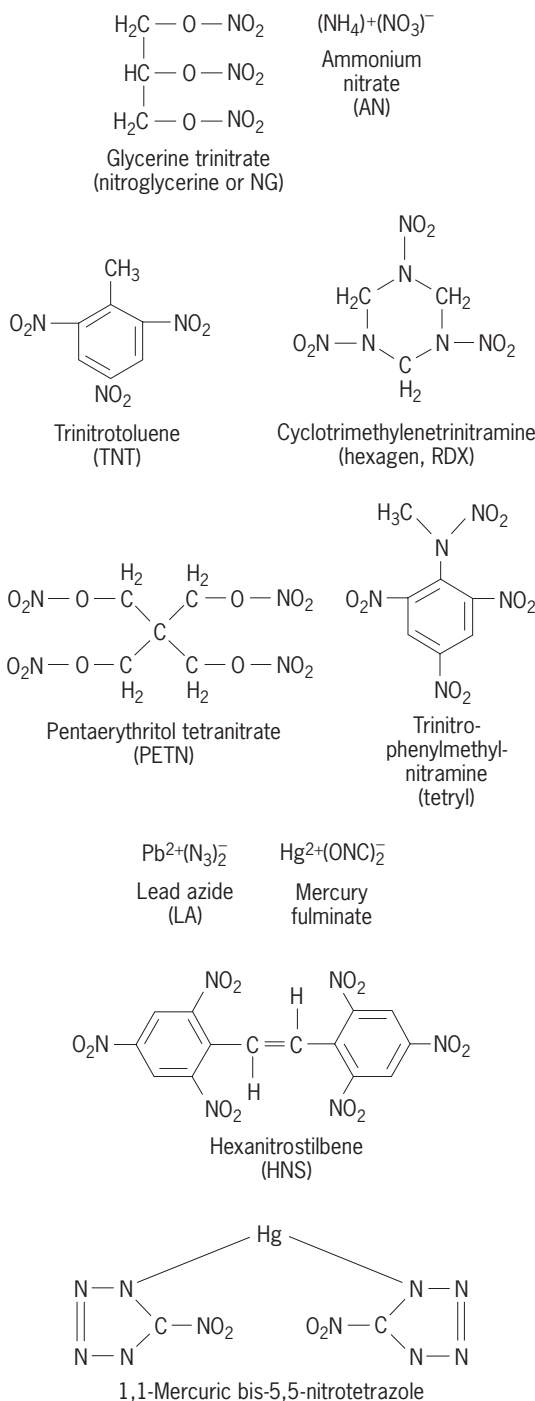


Fig. 3. Structural formulas of some common explosives.

Lead azide (Fig. 3) is a primary explosive much used in detonators. It has good thermal stability and detonates promptly without undergoing a perceptible stage of deflagration, but it is subject to hydrolytic instability. When stored in contact with copper or brass, it can react to form dangerously unstable copper azide.

1,1'-Mercuric bis-5,5'-nitrotetrazole (DXN-1; Fig. 3) has high thermal stability and detonates promptly upon application of heat without exhibiting a well-defined burning zone as do mercury fulminate and lead styphnate.

Detonators. Large charges are usually made of explosives that are relatively insensitive. The amount of very sensitive, more dangerous, and more costly primary explosive required for initiating detonation is kept to a minimum. Only enough is used in an assembly known as a detonator to reliably initiate a larger charge of a less sensitive explosive. This in turn initiates a still larger charge of a still less sensitive explosive.

Most detonators employ increasing violence to produce an explosion accompanied by a shock wave that initiates explosion of the small primer charge of a pyrotechnic composition. The latter's flame contacts a charge of a very sensitive primary explosive, such as lead azide, in contact with another small base charge of a powerful secondary explosive, such as PETN or RDX.

The base charge of the detonator is placed in or against the receptor charge. Upon being heated or shocked, for example, by a fuse or an electrically ignited match-head, the charge of primary explosive detonates. This in turn detonates the base charge. The resulting shock wave or impact of missiles detonates the receptor charge.

In other detonators, the charge of primary explosive is detonated by impact, friction, an intense flash of light, or an intense surge of electric current through a bridgewire or across a spark gap. In still others, the primary explosive is omitted and the secondary explosive is detonated by igniting it while it is strongly confined in a thick-walled metal tube, or by the impact of a thin metal plate driven by an electrically ignited charge of secondary explosive.

Detonators are frequently provided with one or more safety devices, such as an electrical shunt across the lead wires or a temporary barrier between the detonator and the receptor charge. These safety devices are removed only when preparing to fire the charge.

In mining and civil engineering, blasts usually use many separate charges, each in a separate borehole and each with a separate detonator cap. These charges are detonated in a timed sequence chosen to reduce the vibration of nearby buildings, increase rock breakage, and control the throw of the broken rock. The timing of the detonations may be controlled by timing the pulse of electricity used to initiate each detonator, or by providing a pyrotechnic delay train or an internal electronic time delay circuit for each detonator. The delay train may be in the detonator or in a line of detonating cord or shock tube, if that is used to initiate the detonator.

Detonators can be accidentally detonated by rough handling, electric current, or heat, and, though deceptively small, have sufficient energy to maim. Therefore, extensive precautions are required in their transportation, storage, and use.

Primers, boosters, igniters, and squibs. The charge that is initiated by a detonator (if used to detonate a main charge that is still larger and more insensitive) is called a booster or primer. Bags of a water gel explosive, cartridges of dynamite, or cast cylinders of pentolite are used as boosters in industrial blasting.

Pressed pellets of tetryl are often used as boosters in military ordnance.

Sometimes the term primer is also used loosely to mean detonator. The term so used is not to be confused with the same term more correctly used to mean a flame-producing pyrotechnic charge or an assembly used to ignite a fuse or a larger mass of gun propellant, or to provide the heat to detonate a small charge of primary explosive in a detonator. A primer may be initiated by percussion, stab action, friction, or heat.

The terms igniter and squib are variously used to mean heat or flame-producing devices for igniting rocket propellants and other deflagrating pyrotechnics.

The amount of primer used is usually about 0.0035 oz (0.1 g) or less. A common example is the tiny globule that may be seen on the ends of the fuse wire post in a photographic flashbulb; this material serves to ignite the aluminum foil.

Although primer compositions deflagrate rapidly, they do not detonate. In fact, a detonating action would be undesirable because it would tend to blow away rather than to ignite the charge. The total burning time of a primer is about 500 microseconds.

Primer mixtures sensitive primarily to friction may consist simply of potassium chlorate as the oxidizer and antimony sulfide as the fuel. Such a mixture produces hot particles as well as gases. Ground glass or silicon carbide may be added to increase the frictional effect, and sulfur may be added to make the mixture more sensitive. Also, very fine grained black powder (meal powder) is sometimes added. The ingredients are made into a paste with a small amount of gum arabic that binds them together when the paste dries.

Percussion primer for small-arms ammunition usually contain a primary explosive such as mercury fulminate in addition to the materials mentioned above. However, primary explosives are not essential; they may be replaced by other materials such as lead thiocyanate. Lead compounds seem to be preferred, probably because of the hot particles of lead oxide that they produce. Lead styphnate (2,4,6-trinitroresorcinat) is used as a percussion-cap component.

In detonators, the primer is often of a match-head composition. In fact, the formulation of primers is quite similar to that of matches and pyrotechnics; and ingredients such as potassium chlorate, sulfur, and antimony sulfide are common in all applications.

Fuse. This is conduit that leads a hot or high-pressure chemical reaction front from one place to another. The word is not to be confused with "fuze" (sometimes spelled fuse), which is a mechanism used to fire an item of explosive or pyrotechnic ordnance after it has been launched. A fuze is usually a mechanical or electrical device but may contain a fuse.

A typical deflagrating fuse is miner's safety fuse, used in commercial blasting and comprising a train of black powder sheathed in woven and waxed fabric and made in types burning at 120 and 90 s/yard.

(131 and 98 s/m). Firecracker fuse is a train of black powder wrapped in tissue paper. Gasless delay fuse comprises a core of reactants [for example, lead oxide (Pb_3O_4) plus boron (B)] that react exothermally to produce mainly solid products, encased in a short tube (often of lead), made so as to give relatively accurate overall burning times ranging between 10^{-3} and 10^{+2} ; it is used in delay detonators for commercial blasting with multiple charges, and also is used in fuzes.

A detonating fuse contains a train of detonating explosive. The most common type, ordinary detonating cord, comprises a train of PETN or RDX sheathed with wax- or asphalt-impregnated fiber and detonating at 17,000–26,000 ft/s (5200–7900 m/s), and is used for initiating multiple charges in commercial blasting. Shock tube is a flexible, plastic tube on the inner wall of which is an adherent coating of very fine HMX and aluminum powders. Initiation with a shotgun shell primer results in a shock wave supported by a dust explosion, which travels through the tube at 1950 m/s (6400 ft/s), leaving the tube unbroken and producing no noise. It is used in commercial blasting for detonating multiple charges, particularly when airborne noise is to be minimized.

D. Linn Coursen

Bibliography. A. Bailey and S. G. Murray, *Explosives, Propellants, and Pyrotechnics*, 2d ed., 2000; International Society of Explosives Engineers, *The Blasters Library*, 1971–1993; *Kirk-Othmer Encyclopedia of Chemical Technology*, 4th ed., vol. 9, 1998; B. Lewis and G. von Elbe, *Combustion, Flames, and Explosions of Gases*, 3d ed., 1987; P. A. Persson, R. Holmberg, and J. Lee, *Rock Blasting and Explosive Engineering*, 1993; R. A. Sickler, *Explosive Principles: An Essential Guide to Understanding Explosives and Detonations*, 1992; U.S. Department of Defense and Department of Energy, *The Effects of Nuclear Weapons*, 3d ed., 1977, reprint 1983.

Explosive forming

The shaping or modifying of metals by means of explosions. The explosives may be of either the detonating or deflagrating type. Explosive gas mixtures or stored gas at high pressure may also provide the motive power (Fig. 1). See EXPLOSIVE.

Most types of explosive forming involve a mold into which a flat sheet of metal is pressed by the explosion. The metal is stretched uniformly by the explosive impulse. Even welds in the original blank survive the deformation without damage. The advantages of explosive forming over conventional forging methods accrue especially in the case of intricate shapes of which only a few items are required. Tooling costs are low.

Cold welds can be made between dissimilar metals by driving the two parts together under explosive impact. In other applications of explosive-forming methods, powders are pressed into solid billets.

In a quite different application, high explosives are

used to cut large blocks of metal and even to split thin sheets into two layers of exactly one-half the original thickness. When an explosive charge is detonated in contact with the metal, it produces a compression wave. On reflection from a free surface, the compression wave turns upside down to produce a tension wave. Along the planes in the metal where the tension in the wave front exceeds the strength of the metal, rupture occurs.

Explosives can also be employed to extrude metal shapes and to punch hard metals with the aid of dies. Shapes produced explosively are very exact and free from the fine cracks that sometimes result when pressure is slowly applied (Fig. 2). Forces far exceeding those of the largest hydraulic presses can be applied by explosives.

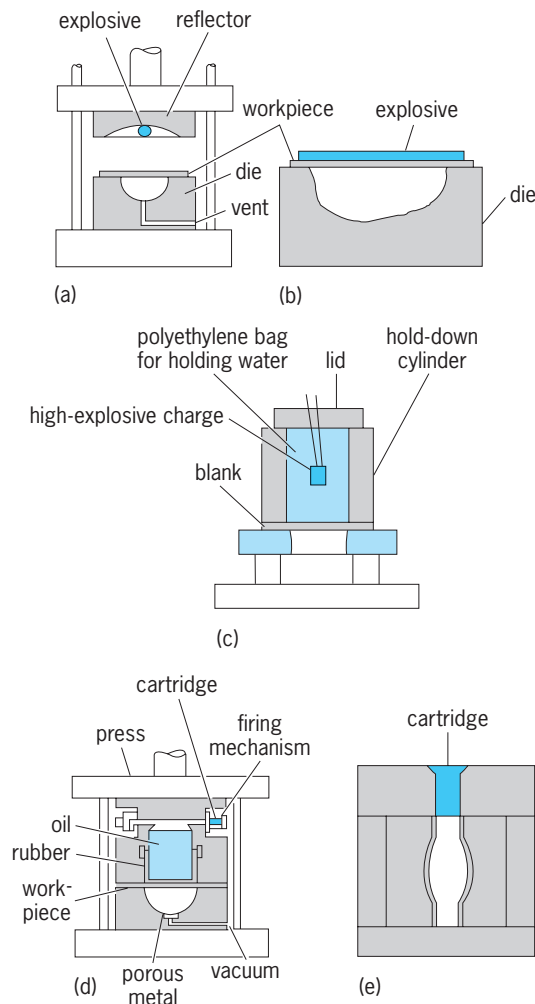


Fig. 1. Five methods of explosive forming. (a) Shaped charge and parabolic reflector. Shock and pressure are motive power. Reflector and charge must be shaped to direct explosive force toward workpiece. (b) Flat high explosive is simply placed on top of sheet which lies on female die. Forming over a punch produces wrinkles. (c) The hold-down cylinder is often a carton of water or a bolted ring on lower die. (d) Gunpowder cartridges are usually applied in a press or enclosed die. This one uses oil and rubber to distribute force. Vacuum prevents air pocket which would retard action. (e) Bulge forming with cartridge power is comparatively simple and safe. Such devices properly vented are less noisy than many standard forming operations. (After *Explosives form space age shapes*, *Steel*, August 25, 1958)



Fig. 2. Typical metal shapes produced by explosive forming. Note weld in bowl. (National Northern Corp.)

Metals can also be hardened under explosive impact with results that compare favorably with those from slow cold-working methods. See METAL FORMING. William E. Gordon

Exponent

In mathematics, a symbol or number written to the right of and above another symbol or number to denote how many times the latter is to be multiplied by itself. For example, $7^3 = 7 \cdot 7 \cdot 7$, and $a^3 = a \cdot a \cdot a$. By use of this convenient mathematical device, the number 420,000,000,000,000 can be expressed unmistakably and more compactly as 4.2×10^{14} , which is read as, "four and two-tenths times ten to the fourteenth power." This abbreviated notation is particularly valuable in expressing the extremely large or extremely small numbers encountered in modern scientific work. For example, 0.000000143 may be expressed as

$$\frac{143}{1,000,000,000} \text{ or } \frac{143}{10^9} \text{ or } 1.43 \times 10^{-7}$$

Operations with exponents are governed by the following rules:

1. $x^a x^b = x^{a+b}$.
Example, $x^2 x^3 = xx \cdot xxx = x^5$.
2. $(x^a)^b = x^{ab}$.
Example, $(x^2)^3 = (x^2)(x^2)(x^2) = x^6$.
3. $(xy)^a = x^a y^a$.
Example, $(xy)^3 = xyxyxy = xxxxyyy = x^3 y^3$.
4. If $y \neq 0$, then $\left(\frac{x}{y}\right)^a = \frac{x^a}{y^a}$.
Example, $\left(\frac{x}{y}\right)^3 = \frac{x}{y} \cdot \frac{x}{y} \cdot \frac{x}{y} = \frac{x^3}{y^3}$.
5. If $x \neq 0$, then $\frac{x^a}{y^a} = x^{a-b}$.
Example, $\frac{x^5}{x^2} = \frac{xxxxx}{xx} = xxx = x^3$.
6. If $x \neq 0$, and $-n$ is a negative integer, then $x^{-n} = \frac{1}{x^n}$.
Example, $\frac{x^2}{x^5} = \frac{1}{x^3} = x^{-3}$.
7. If $x \neq 0$, then $x^0 = 1$.

8. If p/q is any rational number expressed with positive denominator q , then $x^{p/q} = (\sqrt[q]{x})^p$.

Example, $x^{-2/3} = \sqrt[3]{x^{-2}} = (\sqrt[3]{x})^{-2} = \frac{1}{(\sqrt[3]{x})^2}$.

9. If a is any irrational number, then x^a is equal to the limit of the sequence, $x^{a_1}, x^{a_2}, x^{a_3}, \dots, x^{a_n}$, where $a_1, a_2, a_3, \dots, a_n$, etc., are the one-place, two-place, ..., n -place approximations of a .

Example, x^π is equal to the limit of the sequence, $x^3, x^{3.1}, x^{3.14}, x^{3.141}, x^{3.1416}, \dots$

Any number which can be written in the form $a + ib$, where a and b are real numbers and $i = \sqrt{-1}$, is called a complex number. Exponents which are complex numbers are frequently encountered. For example, any number may be expressed as a power of $e = 2.718\dots$, the base of the system of natural logarithms. Many quantities in natural phenomena are so expressed. Where e has a complex exponent, e^{a+ib} , by rule 1, $e^{a+ib} = e^a \cdot e^{ib}$. The value of e^{ib} may be found by substituting ib for x in the expansion

$$e^x = 1 + \frac{x}{1} + \frac{x^2}{1 \cdot 2} + \dots + \frac{x^n}{1 \cdot 2 \cdot 3 \cdot \dots \cdot n} + \dots$$

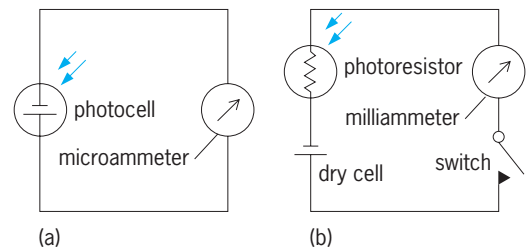
By separating the real and imaginary parts, the Euler equation is obtained, $e^{ib} = \cos b + i \sin b$. Thus, $e^{a+ib} = e^a(\cos b + i \sin b)$. See ALGEBRA.

Arnold N. Lowan; Salomon Bochner

Exposure meter

An indicating instrument used in photography to determine lens aperture and shutter speed. An exposure meter may be used either in the darkroom to determine approximate printing time for a contact print or an enlargement or, more usually, with a camera to determine exposure of film.

Exposure meters (see **illus.**) are of two basic types, photovoltaic and photoconductive. In photovoltaic meters, a selenium cell converts photons to electrons, producing a current directly proportional to received light. A sensitive microammeter indicates this current to the user. The photocell has rapid response, spectral sensitivity closely approximating that of panchromatic film, and indefinitely long life. Because the power to deflect the microammeter comes entirely from the received light, design of the exposure meter is a compromise between a



Basic exposure meter circuits. (a) Photovoltaic type circuit. (b) Photoconductive type circuit.

cell of large area and a meter movement of delicate construction. See PHOTOVOLTAIC CELL.

In the photoconductive meter, a cadmium sulfide cell changes conductivity in proportion to the received light. A battery (usually a mercury cell) supplies power through the cadmium sulfide cell to the meter movement. Because the received light serves to control the power from the battery, this type of instrument is about three orders of magnitude more power-sensitive than the photovoltaic type. As a consequence, cell area can be small and meter movement can be somewhat more sturdy. Spectral response is greatest in the green and the yellow portion of the spectrum and least in the blue, so a filter is built into the meter. The photoconductive cell recovers slowly from bright light; therefore, it is preferably covered when not in use. Also, a switch opens the electric circuit when the meter is not in use to conserve battery life, nominally about a year. See PHOTOCONDUCTIVE CELL.

For the photovoltaic meter, a honeycomb (or fly-eye) lens with acceptance angle comparable to that of a camera lens of normal focal length gathers light for the cell. A multi-iris cover may serve to reduce sensitivity for use in bright light, with a mechanically linked change in meter scale. For the photoconductive meter, a small lens controls acceptance angle; or the cell may be mounted inside the camera and thus can use the camera lens, which determines acceptance angle. As an accessory meter separate from the camera, the exposure meter may be fitted with a lens having an acceptance angle of about 3° and a corresponding viewfinder, so that the meter can be used to measure from a distance the brightness of small portions of the subject.

Used at the camera, either type of meter serves as a brightness meter to measure reflected light from the subject. Used at the subject, with a diffuser head over the exposure meter lens, the instrument serves as an illumination meter to indicate incident light. See PHOTOMETRY.

To aid in interpreting the indicated light as exposure conditions, the meter may carry a circular slide rule. If built into the camera, the meter movement may either operate an indicator in the field of the viewfinder, or it may control the camera iris to produce an average exposure. In any case, the user presets the rated speed of the emulsion in use into the meter. See CAMERA; PHOTOGRAPHY. Frank H. Rockett

Bibliography. J. Neubart, *The Photographer's Guide to Exposure*, 1988; L. Stroebel and R. D. Zakia (eds.), *The Focal Encyclopedia of Photography*, 3d ed., 1996.

Extended x-ray absorption fine structure (EXAFS)

The structured absorption on the high-energy side of an x-ray absorption edge. The absorption edges for an element are abrupt increases in x-ray absorption that occur when the energy of the incident x-ray matches the binding energy of a core electron (typically a 1s

or a 2p electron). For x-ray energies above the edge energy, a core electron is ejected from the atom. The ejected core electron can be thought of as a spherical wave propagating outward from the absorbing atom. The photoelectron wavelength is determined by its kinetic energy, which is in turn determined by the difference between the incident x-ray energy and the core-electron binding energy. As the x-ray energy increases, the kinetic energy of the photoelectron increases, and thus its wavelength decreases. See ABSORPTION OF ELECTROMAGNETIC RADIATION; LIGHT; PHOTOEMISSION; QUANTUM MECHANICS.

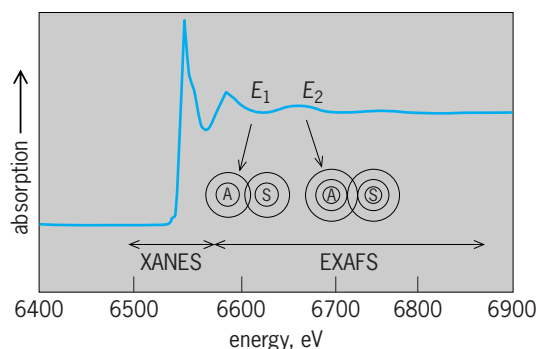
The x-ray-excited photoelectron will be scattered by the neighboring atoms surrounding the absorbing atom. The portion of the photoelectron wave that is scattered back in the direction of the absorbing atom is responsible for the EXAFS oscillations. If the outgoing and backscattered photoelectron waves are out of phase and thus interfere destructively, there is a local minimum in the x-ray absorption cross section. At a higher x-ray energy (shorter photoelectron wavelength), constructive interference leads to a local maximum in x-ray absorption (see *illus.*). EXAFS thus arises from photoelectron scattering, making it a spectroscopically detected scattering method. See INTERFERENCE OF WAVES; SCATTERING OF ELECTROMAGNETIC RADIATION.

EXAFS typically refers to structured absorption from approximately 50 to 1000 eV or more above the absorption edge. X-ray absorption near edge structure (XANES) is often used to refer to the structure in the near (around 50 eV) region of the edge. X-ray absorption fine structure (XAFS) has gained some currency as a reference to the entire structured absorption region (XANES+EXAFS).

Form of oscillations. The energy can be given in units proportional to the inverse photoelectron wavelength [the photoelectron wavevector, or k , defined as in Eq. (1), where m_e is the electron mass, E is

$$k = \sqrt{2m_e(E - E_0)/\hbar^2} \quad (1)$$

the incident electron energy, E_0 is the core-electron binding energy, and \hbar is Planck's constant divided by 2π]. If this is done, the EXAFS oscillations take



X-ray absorption spectrum for manganese, showing XANES and EXAFS regions. As the x-ray energy increases from E_1 to E_2 , the interference of the outgoing and backscattered photoelectron wave (shown schematically by concentric circles around the absorbing, A, and scattering, S, atoms) changes from destructive to constructive.

on a particularly simple form. For a single absorber-scatterer pair (for example, a diatomic gas) the alternating constructive and destructive interference gives rise to regular, sinusoidal oscillations in the absorption coefficient. The fractional modulation in the x-ray absorption coefficient, $\chi(k)$, can then be described by an equation similar to (2), where the

$$\chi(k) = \sum_s \frac{N_s A_s(k)}{k R_{as}^2} \exp(-2k^2 \sigma_{as}^2) \cdot \sin[2k R_{as} + \phi_{as}(k)] \quad (2)$$

right-hand side is a sum over the various types of scattering atoms.

In Eq. (2), the parameters of greatest interest are the number of scattering atoms of type s , N_s , the absorber-scatterer distance, R_{as} , and the mean-square deviation in R_{as} , σ_{as}^2 . To a first approximation, these can be related to the depth of modulation (that is, the EXAFS amplitude), the frequency of the modulations, and the shape of the amplitude envelope describing the oscillations, respectively. In order to determine these structural parameters, the intrinsic scattering amplitude, $A_s(k)$, and the scattering phase shift, $\phi_{as}(k)$, which also appear in Eq. (2), need to be known. These can be obtained by fitting the EXAFS data for structurally characterized models (models where N_s , R_{as} , and σ_{as}^2 are known), or they can be calculated from first principles, typically with empirical calibration. Both $A_s(k)$ and $\phi_{as}(k)$ depend on the chemical identity of the scatterer and thus provide the information necessary to identify the scattering atom. However, both $A_s(k)$ and $\phi_{as}(k)$ depend only weakly on scatterer identity, and thus it is difficult to identify the scatterer with precision. This typically limits scatterer identity to the nearest row of the periodic table (that is, sulfur and oxygen can be distinguished, while nitrogen and oxygen cannot).

All absorber-scatterer pairs contribute to the observed oscillations, as reflected by the sum in Eq. (2). In practice, however, it is not usually realistic to refine all of the different absorber-scatterer interactions. Consequently, it is necessary to group absorber-scatterer interactions into shells. A shell is a group of similar scatterers at approximately the same distance from the absorber.

Fourier transform. Most commonly, EXAFS data are not presented in the form shown in the illustration, but as the Fourier transform of the data. The Fourier transformation separates the frequency-space signal $\chi(k)$ into its different constituent frequencies, with the peaks in the Fourier transform representing the different shells of scatterers. The Fourier transform thus resembles a radial distribution function of electron density around the absorbing atom. However, the Fourier transform peaks are shifted to lower R by about 0.5 Å (50 picometers) due to the phase shift $\phi_{as}(k)$, and the amplitude of the peaks cannot be related directly to electron density due to the effect of the amplitude terms $A_{as}(k)$. See FOURIER SERIES AND TRANSFORMS.

Structural information. EXAFS spectra contain structural information comparable to that obtained

from single-crystal x-ray diffraction. The principal advantage of EXAFS in comparison with crystallography is that EXAFS is a local structure probe and does not require the presence of long-range order. This means that EXAFS can be used to determine the local structure in noncrystalline samples. In ideal circumstances, EXAFS data can be analyzed to determine the absorber-scatterer distance with an accuracy of around 0.02 Å (2 pm) and a precision that is substantially better (0.004 Å or 0.4 pm, and perhaps even better under carefully controlled conditions). Coordination numbers can be determined with an accuracy of around 25%. See X-RAY CRYSTALLOGRAPHY.

The information available from EXAFS is relatively limited. Scattering atoms can be resolved into different shells only if they are separated by around 0.15 Å (15 pm) or more and only scatterers within about 4–5 Å (0.4–0.5 nanometer) of the absorbing atom contribute to the signal. Equation (2), which describes single scattering, does not contain any angular information. In some cases, multiple scattering, in which the photoelectron is scattered off more than one atom, can give EXAFS sensitivity to geometry. However, multiple scattering is a large effect only if the structure contains a nearly linear arrangement of atoms. Consequently, EXAFS does not generally provide significant angular detail. Despite these limitations, EXAFS often provides the only available structural information for noncrystalline systems. This has made EXAFS one of the preferred structural probes in many areas, particularly for studies in materials science, catalysis, and biological systems.

James Penner-Hahn

Bibliography. J. Goulon, C. Goulon-Ginet, and N. B. Brookes (eds.), *X-ray Absorption Fine Structure*, *J. Phys. IV*, 1997; D. C. Koningsberger and R. Prins (eds.), *X-ray Absorption: Principles, Applications, Techniques of EXAFS, SEXAFS, and XANES*, 1988; P. A. Lee et al., *Extended x-ray absorption fine structure: Its strengths and weaknesses as a structural tool*, *Rev. Mod. Phys.*, 53:769–806, 1981; J. Penner-Hahn, *X-ray absorption spectroscopy in coordination chemistry*, *Coord. Chem. Rev.*, 190–192:1101–1123, 1999; B. K. Teo, *EXAFS: Basic Principles and DATA Analysis*, 1986.

Extinction (biology)

The death and disappearance of a species. The fossil record shows that extinctions have been frequent in the history of life. Mass extinctions refer to the loss of a large number of species in a relatively short period of time. Episodes of mass extinction occur at times of rapid global environmental change; five such events are known from the fossil record of the past 600 million years. Human activity is causing extinctions on a scale comparable to the mass extinctions in the fossil record.

Nature and record. An extinction may be of two types; phyletic or terminal. Phyletic extinction occurs when one species evolves into another with time; in this case, the ancestral species can be called

extinct. However, because the evolutionary lineage has continued, such extinctions are really pseudoextinctions. In contrast, terminal extinction marks the end of an evolutionary lineage, termination of a species without any descendants. Most extinctions recorded in the fossil record and those occurring today are terminal. *See* ORGANIC EVOLUTION.

It has been estimated that 99% of all species that have ever lived are now extinct. Because of problems of fossil preservation and sampling, the measurement of extinction in the fossil record is not simple. Except in a few circumstances, the fossil record of individual species is not used. Rather, a higher taxonomic level (typically the family) is used as a crude measure of species-level extinction. This is done because of the potential of sampling error in the geologic record. The absence of a species from some geologic stratum may be the result of poor preservation of fossils, inappropriate rock types, or inadequate sampling rather than true extinction. The absence of all species within a family may be a more reliable indicator of actual extinction. Because the number of species varies from family to family, the number of family extinctions may be an inaccurate measure of the number of species extinctions. Yet, where extinctions have been studied in detail, at species, genus, and family levels, high rates of extinction at the higher taxonomic level result from high rates at the lower level. For example, the record of family extinctions since the Permian reflects the level of genus extinctions.

An additional difficulty in measuring extinctions in the fossil record is estimating their time of occurrence. Although radiocarbon dating gives precise estimates for extinctions of the past 30,000 years, the method cannot be used for older rocks and fossils. The extinction of a species (or genus or family) is marked by the horizon where it last appeared in the rock record. Such horizons can be dated with a precision of a few million years at best. Generally speaking, the farther back in the rock record, the more difficult it is to pinpoint the time of extinction.

The fossil record is best known for marine organisms. The extinctions of marine organisms over the past 570 million years were relatively minor, ranging from one to seven families per million years. The extinction rates never dropped to zero, indicating that extinction is a normal part of the history of life. Such low levels are sometimes referred to as background extinctions, those that are always occurring as a consequence of normal environmental changes or local catastrophes. Background rates of family extinction appear to have declined through geologic time. Although this may reflect an increase in average fitness through evolution, the decline in family extinction rate may also be an artifact; there has been an increase in the number of species per family since the Cambrian. Because a family becomes extinct only when all its component species become extinct, families with many species are less likely to become extinct than families with few species.

The mass extinctions of the marine fossil record

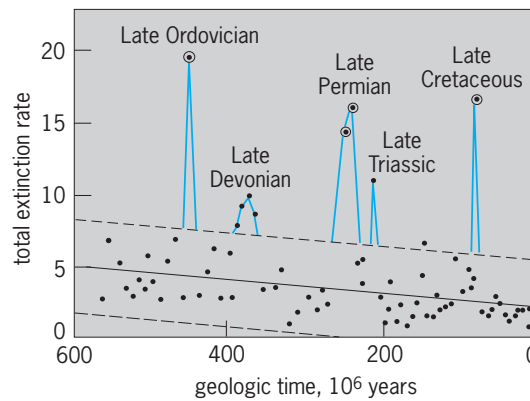


Fig. 1. Record of extinction of families of marine organisms. Spikes show times of mass extinctions. (After D. M. Raup and J. J. Sepkoski, Jr., *Mass extinction in the marine fossil record*, *Science*, 215:1501–1503, 1982)

are from two to five times as severe as background extinctions. The mass extinctions are indicated by the five spikes in **Fig. 1**: the Late Ordovician, Late Devonian, Late Permian, Late Triassic, and Late Cretaceous. These mass extinctions affected a variety of organisms in many different ecological settings.

Figure 2 shows the effect of these mass extinctions on the diversity of marine life since the latest Precambrian. Each sharp drop in diversity is associated with a mass extinction. The most severe drop is in the Late Permian, when 52% of marine families became extinct. As many as 96% of all species may have become extinct at that time.

A 26-million-year periodicity has been detected in the marine extinctions of the past 250 million years; that is, rates of family and genus extinction seem to increase, sometimes to mass extinction levels, every 26 million years. The most recent episode was about 13 million years ago in the Middle Miocene. Although there is debate regarding the quality of both the data and the analytical methods, the implications are profound. Periodic extinctions suggest the periodic causes.

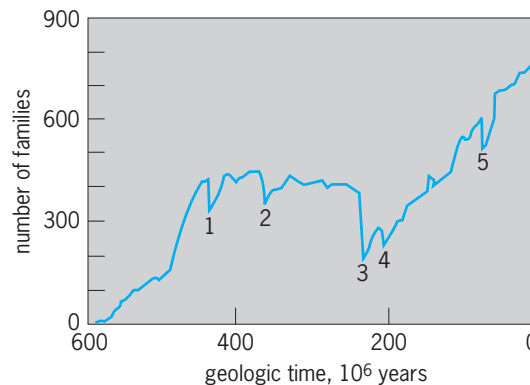


Fig. 2. Diversity of marine families in the fossil record. Sharp drops (1–5) indicate times of mass extinctions. 1 = Late Ordovician (–12%). 2 = Late Devonian (–14%). 3 = Late Permian (–52%). 4 = Late Triassic (–12%). 5 = Late Cretaceous (–11%). (After D. M. Raup and J. J. Sepkoski, Jr., *Mass extinction in the marine fossil record*, *Science*, 215:1501–1503, 1982)

The fossil record of terrestrial animals and plants is less well known than for marine organisms. Terrestrial animals and plants are not as likely to be preserved as fossils. As a result, it is more difficult to estimate the time and magnitude of the extinction of terrestrial organisms.

Terrestrial and marine mass extinctions seem to occur at about the same time. The Late Permian, Late Triassic, and Late Cretaceous are also times of extinction for terrestrial vertebrates; the most dramatic extinction of terrestrial vertebrates took place at the end of the Cretaceous, when the last dinosaurs died off. Correlation between the terrestrial and marine record is difficult, but it appears that the terrestrial and marine extinctions at the end of the Cretaceous occurred at the same time. *See* DINOSAURIA.

The best record of terrestrial vertebrate extinction is that of the Pleistocene. Late Pleistocene extinctions in North America are especially well known—33 genera of mammals vanished during the last 100,000 years. These extinctions were concentrated among the large mammals—those over 100 lb (44 kg) in weight—and most occurred during a short time interval approximately 11,000 years ago.

Causes. Ever since the work of Georges Cuvier, the French naturalist who demonstrated the reality of extinction, explanations have fascinated both scientists and the general public. Cuvier invoked sudden catastrophic events, whereas his contemporaries favored more gradual processes. These two themes, catastrophism and gradualism, are still debated.

In 1980 high concentrations of iridium were reported precisely at the Cretaceous-Tertiary boundary. Iridium is rare in most rocks but more abundant in meteorites. It was proposed, therefore, that an asteroid struck the Earth 65 million years ago. The

impact darkened the atmosphere with dust, caused a catastrophic short-term cooling of the climate, and thus led to the extinction of dinosaurs and many other Cretaceous species. The iridium-rich layer at the boundary marks this terminal Cretaceous event.

Subsequent research has strengthened the case. Mineral grains with textures characteristic of high shock pressures have been found in the boundary layer. Soot, possible evidence of fires ignited by the impact, has also been recovered at the boundary.

Astronomical theories have been put forward to explain the Late Cretaceous extinctions as well as the 26-million-year periodicity. In one theory, the Sun has a distant companion star that would pass in orbit near the solar system's cloud of comets every 26 million years. This might perturb many comets, sending a few into the Earth. A comet would produce the same effects as an asteroid. Although this companion star has yet to be discovered, it has been named Nemesis, and a search for it has been undertaken.

Opponents of the impact theory contend that many features of the Late Cretaceous rock and fossil record are better explained by terrestrial causes such as volcanism, lowered sea level, and climatic change. They argue that the fossil record indicates that Late Cretaceous extinctions occurred over several millions of years and were not catastrophic.

Although the impact theory has been offered as an explanation for almost all the extinctions in the fossil record, well-confirmed discoveries of iridium-rich sediments in association with extinctions have been made only at the Cretaceous-Tertiary boundary and in the late Eocene.

Other explanations for mass extinctions include lowered sea level, climatic cooling, and changes in oceanic circulation. It has been pointed out that

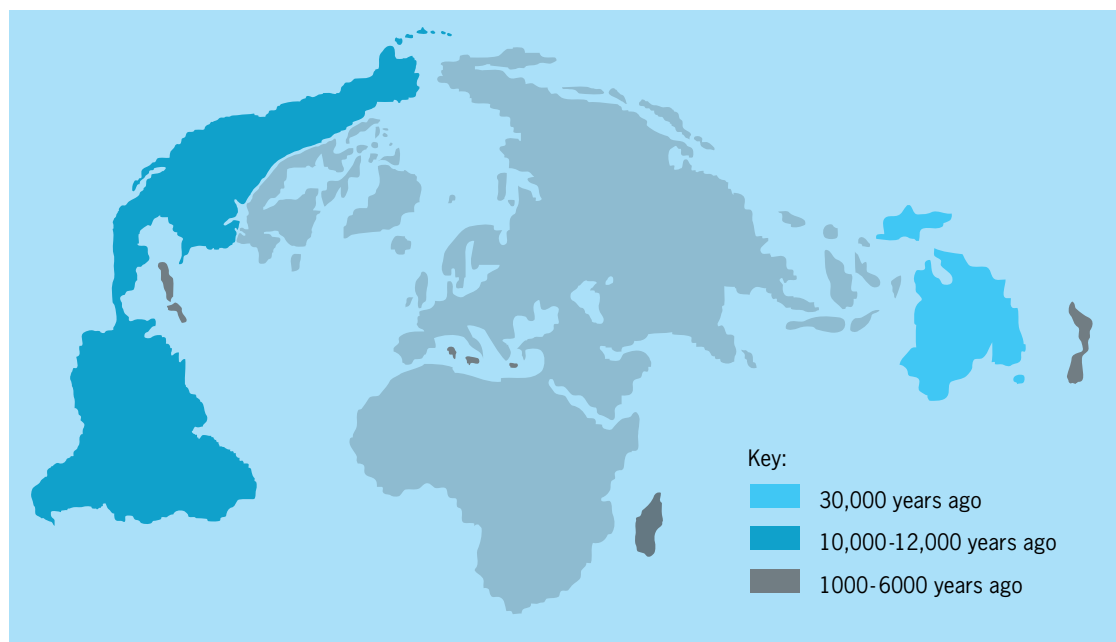


Fig. 3. Times and places of extinctions of large mammals during the late Pleistocene. (After P. S. Martin and R. G. Klein, eds., *Quaternary Extinctions*, 1984)

the Late Ordovician extinctions coincided with an episode of global cooling and sea-level lowering prompted by the growth of the polar ice caps. Explanations for the Late Devonian, Late Permian, and Late Triassic mass extinctions are more elusive. Some evidence for climatic cooling is present in the Late Devonian, though other evidence suggests more sudden causes. The Late Permian marks a time in the Earth's history when sea level was exceptionally low. Habitat for marine organisms may have been severely restricted, and global climates may have been strongly seasonal. The Late Triassic extinctions are also associated with a lowered sea level.

Biotic processes such as disease, predation, and competition may also cause the extinction of species but are difficult to prove from the fossil record because they leave little evidence. Thus there is no evidence that mammals preyed upon dinosaur eggs. In any event, biotic factors usually affect only one or a few interdependent species. Disease, predation, and competition probably played only minor roles in mass extinctions.

Predation has been implicated in the extinction of Pleistocene mammals, which may have resulted from human hunting. Their large size and their disappearance at a time of human immigration support this possibility. **Figure 3** shows the times of large-mammal extinctions during the past 30,000 years. Pleistocene extinctions have also been attributed to rapid climatic change.

Predation and competition are important causes of more recent extinctions, which continue today. Human activities such as hunting and fishing (predation), habitat alteration (competition for space), and pollution have probably destroyed thousands of species. These activities, together with continued tropical deforestation and resulting changes in climate, are likely to cause extinctions that will be comparable to the mass extinctions seen in the fossil record. See FOSSIL.

Karl W. Flessa

Bibliography. J. C. Balovet, *Extinct Species of the World: 40,000 Years of Conflict*, 1990; P. Ehrlich and A. Ehrlich, *Extinction*, 1981; D. K. Elliott (ed.), *Dynamics of Extinction*, 1986; P. S. Martin and R. G. Klein (eds.), *Quaternary Extinctions*, 1984; M. H. Nitecki (ed.), *Extinctions*, 1984; M. J. Novacek and Q. Wheeler (eds.), *Extinction and Phylogeny*, 1992; D. M. Raup, *The Nemesis Affair*, 1986, revised 1999; S. M. Stanley, *Extinction*, 1987; E. O. Wilson (ed.), *Biodiversity*, 1988.

Extraction

A method of separating the constituents of a mixture utilizing preferential solubility of one or more components in a second phase. Commonly, this added phase is a liquid, while the mixture to be separated may be either solid or liquid. As a mundane example, the preparation of tea or coffee is a process of liquid/solid extraction whereby water selectively dissolves certain components of the mixture, leaving behind the insoluble residue (as tea leaves or coffee

grounds). If the starting mixture is a liquid, then the added solvent must be immiscible or only partially miscible with the original and of such a nature that the components to be separated have different relative solubilities in the two liquid phases.

Principles. The ratio of the concentrations of a particular dissolved substance (solute) in two coexisting liquid phases at equilibrium is shown in Eq. (1),

$$D_A = \frac{x_A}{y_A} \quad (1)$$

where D is the distribution coefficient, and x_A and y_A are the concentrations of A in the two phases. The ease of separation of two components in a mixture is conveniently measured by the separation factor α , which is the ratio of the distribution coefficients of the two components between the two solvents (the equivalent of relative volatility in distillation), as in Eq. (2).

$$\alpha_{AB} = \frac{D_A}{D_B} \quad (2)$$

Although the basic concepts of equilibrium outlined above define the potential for separation, they give no indication of the rate of the process. When two phases that are not at equilibrium are contacted together, the rate of transfer of solute between them depends on the extent to which the concentrations of the solute in the two phases differ from the equilibrium value as fixed by the distribution coefficient. In the classical two-film theory, it is assumed that the two phases are actually in equilibrium at the interface and that the resistance to mass transfer is concentrated in thin films on either side of the interface. Mass transfer through these films takes place by molecular diffusion. Thus, Eq. (3) applies, where

$$\text{Rate} = k_x I(x_A - x'_A) = k_y I(y'_A - y_A) \quad (3)$$

k is the film mass-transfer coefficient for a particular phase; I is the interfacial area over which mass transfer takes place; x'_A and y'_A are the solute concentrations in the bulk of a phase; and x_A and y_A are the concentrations in a phase adjacent to the interface.

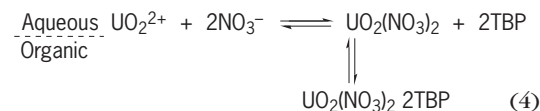
Important variables affecting rate of mass transfer are temperature and agitation. When two immiscible liquids are mixed together, one will break into droplets, forming the dispersed phase while the other remains coherent as the continuous phase. Increasing the degree of agitation gives smaller droplets, hence greater total surface area, enhancing the mass-transfer rate.

Having achieved dispersion and mass transfer, it is necessary to segregate the phases for separate removal. Coalescence of droplets then becomes important, and various chemical and mechanical aids are available for speeding this step.

Chemistry of extraction. Many applications of extraction, particularly those concerned with the separation of metals, depend upon chemical interaction between solute and solvent. In nonreacting systems, phase behavior is determined by physical differences

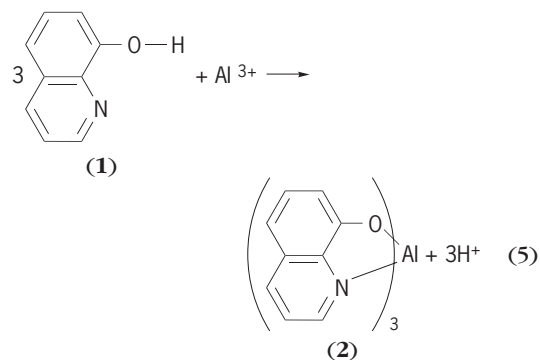
such as polarity and can be analyzed by using thermodynamics. See CHEMICAL THERMODYNAMICS.

Ion association systems. Ion association systems involve the pairing of oppositely charged ions under the influence of electrostatic attraction. Metals may be incorporated as either the cationic or anionic species in an aqueous phase; by association with an oppositely charged ion, an uncharged molecule is formed which can be extracted into the organic phase. Examples of metals in the anionic form include halides and thiocyanates, while the organic solvent is typically an oxygenated compound (for example, ethers, alcohols, ketones). An example is the extraction of uranium by tributyl phosphate (TBP), shown in reaction (4). The position of equi-



librium in metals extraction reactions of this type depends upon the pH of the system. Control of the pH therefore allows sharp separations between different metals.

Chelate systems. Metal chelates are cyclic coordination compounds containing a metal atom in the ring. A reagent that has found widespread application is 8-quinolinol (**I**) which contains two atoms (N and O) that will coordinate with metal ions to form a five-membered ring as shown in reaction (5). Thus



aluminum 8-quinolinate (**II**) is an uncharged chelate that resembles many other unionized molecules in its extremely low water solubility and relatively high solubility in organic solvents. Metal separations may be achieved by the control of pH and oxidation state. Other chelating agents which have proven to be useful in extraction and dithiozone (diphenylthiocarbazon) and acetylacetone. See CHELATION.

Nonreacting systems. Separation of a liquid mixture into two liquid phases in thermodynamic equilibrium can be understood via energy considerations. The total energy of the system is described by the Gibbs free energy G , defined in Eq. (6), where H

$$G = H - TS \quad (6)$$

is the enthalpy, T is temperature, and S is entropy. When a mixture is prepared from pure components, the change in Gibbs free energy, or energy of mixing,

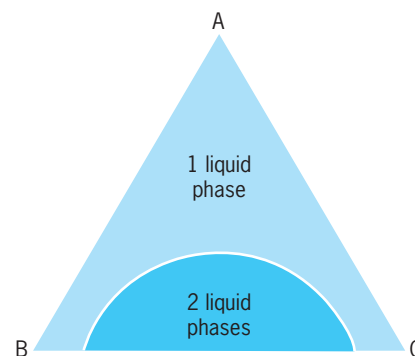


Fig. 1. Phase diagram of a three-component system.

g_m , can be expressed as in Eq. (7), where g_i is the

$$g_m = g_i + g_e \quad (7)$$

energy of mixing for an ideal component and g_e is the "excess" Gibbs free energy. When g_e is zero, an ideal system exists and no phase separation occurs. Highly nonideal systems with significant g_e contributions exhibit liquid-phase separation over certain ranges in composition.

The phase behavior of three-component systems is commonly represented by using a ternary diagram (Fig. 1) which forms a useful tool for quantitative analysis and prediction. A, B, and C define the three pure components; binary mixtures correspond to the sides of the figure, while any point within the triangle is a ternary composition, which may be miscible (one liquid phase) or form two liquid phases.

Contacting. The three basic methods for contacting solvent (S) with the feed mixture (F) are shown in Fig. 2.

Countercurrent extraction is the most efficient and is the choice for commercial operation whenever possible. In this latter approach, raffinate and extract phases flow countercurrently and emerge at opposite ends of the contactor. Many stages can be used to give a high degree of separation while maintaining a modest solvent requirement. See COUNTERCURRENT TRANSFER OPERATIONS.

Liquid/solid extraction. Liquid/solid extraction may be considered as the dissolving of one or more

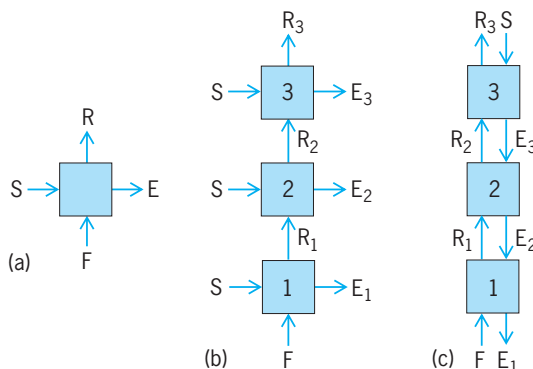


Fig. 2. Contacting methods used in extraction: (a) single stage; (b) crosscurrent; and (c) countercurrent.

components in a solid matrix by simple solution, or by the formation of a soluble form by chemical reaction. The largest use of liquid/solid extraction is in the extractive metallurgical, vegetable oil, and sugar industries. The field may be subdivided into the following categories: leaching, washing extraction, and diffusional extraction. Leaching involves the contacting of a liquid and a solid (usually an ore) and the imposing of a chemical reaction upon one or more substances in the solid matrix so as to render them soluble. In washing extraction the solid is crushed, permitting the valuable soluble product to be washed from the matrix. Sugar recovery from cane is a prime example of this method. In diffusional extraction the soluble product diffuses across the denatured cell walls (no crushing involved) and is washed out of the solid. The recovery of beet sugar is an excellent case in point.

Particle size is significant in all cases since it is a direct function of the total surface area that will be available for either reactions or diffusion. In ores, porosity and pore-size distribution greatly affect the rate of extraction, because the solvent must flow or diffuse in and out of the pores; and, in many cases, the solute moves through the pores to the particle surface by diffusion.

Liquid/liquid extraction. In crosscurrent extraction, fresh solvent is added at each stage; enhanced separation is obtained, but solvent requirement is high.

Liquid/liquid extraction has found application for many years in the coal tar industry. The recovery of tar acids from crude tar oil by washing with an aqueous solution of alkali is an example where chemical interaction between solute and solvent determines differential solubility.

On a smaller scale, extraction is a key process in the pharmaceutical industry for recovery of antibiotics from fermentation broths. Penicillin is obtained by extraction into solvents such as amyl acetate at relatively low pH values (2 to 2.5) and is then stripped from the organic phase by treatment with a buffered aqueous solution at about pH 7 to 7.5. Other examples in this field are the recovery and separation of vitamins and the production of alkaloids from natural products.

Liquid/liquid extraction separates the components of a homogeneous liquid mixture on the basis of differing solubility in another liquid phase. Because it depends on differences in chemical potential, liquid/liquid extraction is more sensitive to chemical type than to molecular size. This makes it complementary to distillation as a separation technique. One of the first large-scale uses was in the petroleum industry for the separation of aromatic from aliphatic compounds. The original process employed liquid sulfur dioxide as solvent. More recently, sulfolane (thio-cyclopentane-1,1-dioxide) has replaced sulfur dioxide for extraction of lighter aromatics due to its greater selectivity and ease of recovery. For the selective separation of higher-molecular-weight aromatics and aliphatics as in lubricating oil manufacture, phenol and furfural are the most widely used solvents.

Supercritical fluid extraction. With a conventional liquid solvent and a given system, temperature is the only variable which can be manipulated to influence solvency. However, if the solvent is a fluid near its critical temperature, changes in pressure significantly affect the fluid's density, hence solvency for a given solute. One advantage of this approach is that the solvent can be separated again merely by reducing the pressure, offering energy savings. A commercial application of the technique is the use of supercritical carbon dioxide for selective removal of caffeine from coffee extract.

Equipment. On a laboratory scale, the separatory funnel is the simplest device for achieving mixing and subsequent phase separation. To obtain continuous extraction of a small feed sample, the Soxhlet extractor is employed.

If the starting mixture of the sample is a solid, it must first be ground down to the required particle size. On a commercial scale, the equipment which is used for achieving efficient contacting with the liquid solvent tends to be specific to the type of industry and may employ screw-type conveyors, successive mixing/settling tanks, or a variety of proprietary designs.

Liquid/liquid extraction using countercurrent processing is preferably carried out in a vertical column. The internals are designed to disperse one phase into droplets, and to accomplish mass transfer between this phase and the continuous phase in a series of stages. These stages may be physical partitions comprising plates or trays, or theoretical concepts in the case of differential contactors such as a packed column.

A novel means of supplying mechanical energy to enhance mixing in an extractor is by pulsing the inlet flow. This concept of pulsed extraction has been applied to both packed and perforated-plate columns, but the mechanical problems associated with pulsing in large-capacity units have limited the commercial adoption of this technique. See CHEMICAL EQUILIBRIUM; CHEMICAL SEPARATION TECHNIQUES; COUNTERCURRENT TRANSFER OPERATIONS; SOLVENT EXTRACTION; TRANSPORT PROCESSES.

Bruce M. Sankey

Bibliography. J. Rydberg, *Solvent Extraction Principles and Practice*, 2d ed., 2004; L. T. Taylor, *Supercritical Fluid Extraction*, 1996; J. Wisniak and A. Tamir, *Liquid-Liquid Equilibrium and Extraction: A Literature Source Book*, 1987.

Extrapolation

A process in mathematics used to find the value of a function outside its tabulated values. This is done as in interpolation by assuming that over a small range of x the function may be closely approximated by a polynomial or some other readily computed function. See INTERPOLATION.

Formulas. Any of the interpolation formulas can be used, therefore, and the desired value of x substituted in them. Thus, for example, if $y = f(x)$ has been

Difference table				
x	$y = \log x$	δy	$\delta^2 y$	$\delta^3 y$
1.00	0.0000 000			
1.01	0.0043 214	43 214		
1.02	0.0086 002	42 788	-426	8
1.03	0.0128 372	42 370	-418	9
1.04	0.0170 333	41 961	-409	8
1.05	0.0211 893	41 560	-401	7
1.06	0.0253 059	41 166	-394	$7 + \epsilon$
1.07	$0.0293\ 838 + \epsilon \times 10^{-7}$	$40\ 779 + \epsilon$	$-387 + \epsilon$	$7 + \epsilon$
1.08	$0.0334\ 237 + 4\epsilon \times 10^{-7}$	$40\ 399 + 3\epsilon$	$-380 + 2\epsilon$	$7 + \epsilon$
1.09	$0.0374\ 263 + 10\epsilon \times 10^{-7}$	$40\ 026 + 6\epsilon$	$-373 + 3\epsilon$	

tabulated at $x = x_{-N}, x_{-N+1}, \dots, x_{-1}, x_0$, the Gregory-Newton interpolation formula (1) may be used to

$$y = y_0 + u \delta y_{-1/2} + \frac{1}{2!} u(u+1) \delta^2 y_{-1} + \frac{1}{3!} u(u+1)(u+2) \delta^3 y_{-3/2} + \dots + \frac{1}{m!} u(u+1) \dots (u+m-1) \delta^m y_{-m/2} \quad (1)$$

determine a polynomial equation passing through the $m + 1$ ordinates $y_{-m}, y_{-m+1}, \dots, y_{-1}, y_0$. These differences give expression (2), and the differences

$$u + \frac{x - x_0}{h} \quad (2)$$

$\delta^k y_{-k/2}, k = 1, 2, 3, \dots$, are the same as those used in interpolation.

If $-1 < u < 0$, then $x_{-1} < x < x_0$ and the formula is used to interpolate. On the other hand, substitution of positive values of u permits its use for extrapolation for y beyond y_0 , the last value tabulated.

If the function $y = f(x)$ is known, the error introduced by using a polynomial to extrapolate for the value of the function can be expressed by adding to Eq. (1) a remainder term, shown as notation (3), where now ξ is any value of x lying be-

$$\frac{1}{(m+1)!} u(u+1) \dots (u+m) f^{(m+1)}(\xi) \quad (3)$$

tween the smallest and the largest of the numbers $x_{-m}, x_{-m+1}, \dots, x_{-1}, x_0$, and x . This term will be larger for extrapolation, $u < 0$, than for interpolation, $u > 0$, for two reasons. First, since u is positive for extrapolation, the coefficient of $f^{(m+1)}(\xi)$ will be larger. Second, since the range of values permitted ξ is larger, $|f^{(m+1)}(\xi)|$ must be assumed to be larger. It is necessary, in calculating the error, to take the largest absolute value of this $(m + 1)$ th derivative of $f(x)$ in the above range of ξ . If there is a singularity of $f(x)$ or of its derivatives near the value of x required in the extrapolation, the remainder term

in notation (3) would indicate that the extrapolation could involve a large error.

If it is necessary to extrapolate a distance greater than b , the interval of the table, beyond the limits of a table, it may be helpful to proceed by first extending the entries in the table by extrapolation. This may be done by assuming, for instance, that some order of difference remains constant, or by letting u take on positive integral values in Eq. (1). An estimate of the errors introduced by extending the table in this way can be made by attempting to extrapolate for the last few entries in the table from the earlier entries. For values not near a singularity of the function tabulated, these errors should be of about the same size as those introduced in extending the table. Of course, the same number of entries should be added in the two cases.

Having extended the entries in the table by extrapolation, one can look upon the problem of finding y for an intermediate value of x as just the problem of interpolation. If the degree of the interpolating polynomial is large enough, the error of this interpolation can be ignored in comparison with the error in extrapolating for the additional entries.

Example. From the portion of the difference table lying above the line, find by extrapolation the logarithms of 1.07, 1.08, and 1.09 and determine a probable limit for the error.

If it is assumed the third difference stays constant at seven units in the seventh decimal place, the difference table can be extended, as shown below the line, by working from right to left. For this purpose the ϵ 's added to the number are ignored.

To determine the maximum error in the extrapolated values it must be recognized that the true third differences would not necessarily be equal to 7×10^{-7} . The maximum error would occur if all these third differences assumed were too high or too low. Therefore the true values of logarithms, rounded off to seven decimal places, will differ from those computed above by less than the ϵ s attached to the numbers in the table.

Since a reasonable value for ϵ is ± 1 reasonable upper limits for the errors in the extrapolation for the logarithms of 1.07, 1.08, and 1.09 are, respectively, 1, 4, and 10 units in the last decimal place. Comparison of these values with the true values reveals that the errors are actually 0, -1 , and -2 units in the last decimal place. It should be clear from this example, however, that the error in extrapolating beyond the limits of a table can be expected to grow rapidly for each new entry added. See GRAPHIC METHODS. Kaiser S. Kunz

Extrasolar planets

Planets that orbit stars other than the Sun. Since 1995, the number of such planets known has increased from zero to more than 100. Prior to the discovery of extrasolar planets, theoretical studies suggested that most planetary systems would look like the solar system, with small rocky planets orbiting near the star and the larger gas giants farther out, all in nested concentric quasicircular orbits. These theories have been disproven, and in their place a much richer understanding of planet formation and evolution has emerged.

More than 200 years ago, the observation that the orbits of the major solar system planets are quasicircular, coplanar, and aligned with the Sun's equator led Pierre-Simon Laplace and Immanuel Kant to the conjecture that the planets formed from a common disk of material that must have circled the early Sun. This theory has been greatly strengthened by observations over the past 20 years that more than half of all nearby young stars are swaddled in disks of gas and dust. See PROTOSTAR; SOLAR SYSTEM.

It is extremely challenging to detect planets orbiting Sun-like stars. Planets shine only by the reflected light of their host stars. The largest (and therefore easiest to detect) planet in the solar system, Jupiter, has 300 times the mass and 11 times the diameter of the Earth. It is nonetheless 10^9 times fainter than the Sun. Taking a picture of a Jupiter-like planet orbiting another star would be similar to snapping a photo from New York of a firefly buzzing around a megawatt searchlight in Los Angeles. Though this is beyond current technology, a number of strategies are being studied that might be capable of taking direct images of Jupiter-like, or even Earth-like, planets within 20 to 50 years.

In the intervening time we are left to deduce the presence of extrasolar planets by the small effects a planet has on the host star. A planet will alter the position and velocity of its host star by gravitationally tugging on it. For example, the Sun and Jupiter jointly orbit a common center of mass, which lies on the line connecting the Sun and Jupiter, just outside the surface of the Sun. A hypothetical alien astronomer could detect the presence of Jupiter by noting that the position of the Sun is periodically wobbling against the background stars (astrometry) or by measuring the periodic velocity variation of the Sun (Doppler spectroscopy; Fig. 1). Such mea-

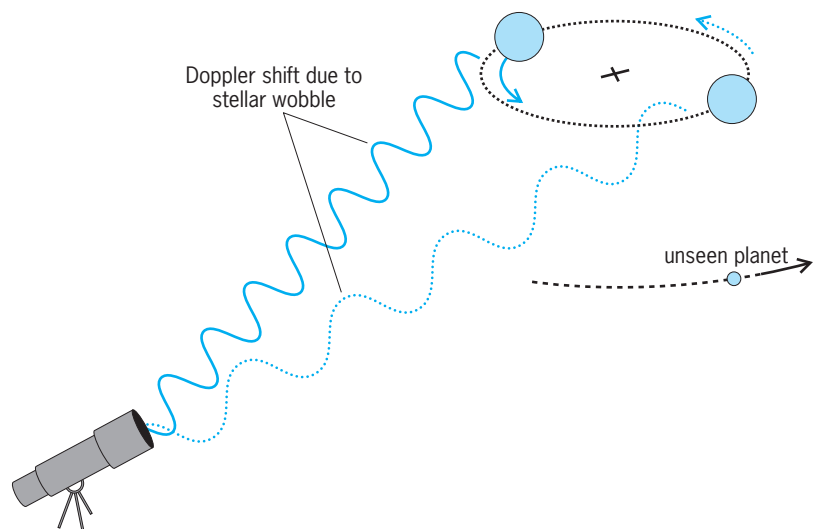


Fig. 1. Stellar wobble that results when an unseen planet orbiting a star gravitationally tugs the star in a small counterorbit. The planet can be detected by observing the star's periodic wobble relative to background stars (astrometry) or by observing the periodic change in the star's line-of-sight velocity (Doppler velocity). (Courtesy of California & Carnegie Planet Search Program)

surements would reveal the orbital period and the magnitude of the wobble. From these data the hypothetical alien astronomer would be able to calculate the orbital radius and mass of the unseen Jupiter from Kepler's third law of planetary motion (1619) and Newton's principle of momentum conservation. See ASTROMETRY; CELESTIAL MECHANICS; DOPPLER EFFECT; KEPLER'S LAWS.

If the Sun-Jupiter orbital plane were aligned edge-on to the alien's line-of-sight, Jupiter would block 1% of the Sun's light as it passed directly in front of (transited) the Sun once per orbit. From this fact, the alien could deduce that Jupiter has 1% the surface area, and hence 10% the diameter, of the Sun.

The planetary zoo. The Doppler velocity method (Fig. 1) is responsible for the discovery of essentially all of the known extrasolar planets, although new space-based telescopes should begin making astrometric detections by about 2015. The Doppler velocity technique can detect only massive planets, and is biased toward detecting planets in small orbits. This bias is due to two effects. More than half an orbit (and preferably more than two orbits) needs be observed to detect a planet. This duration favors the detection of planets in short-period orbits. In addition, the closer a planet is to its star, the stronger the mutual gravitational attraction and hence the larger the velocity that the planet imposes on the star.

In place of the outdated notion that all planetary systems would be similar to the solar system, recent discoveries reveal three primary types of planets. Among the systems of multiple planets that have been discovered to date, some have more than one type of planet. These varying types of planets and systems provide valuable clues about the formation and evolution of planetary systems.

Solar system analogs. Of the 135 known extrasolar planets (as of 2004), only a couple bear even passing resemblance to the giant planets of the solar

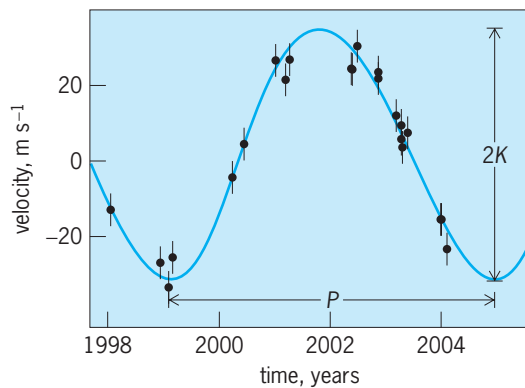


Fig. 2. Doppler velocity curve of the solar system analog HD 70642. The orbital period P of the planet can be directly read off the figure, while the mass of the planet can be calculated from velocity amplitude K . Here the period is 6 years. The semiamplitude (K) is 30 m/s (100 ft/s) and the full amplitude ($2K$) is 60 m/s (20 ft/s). The planet thus orbits at 3.3 AU (compared to Jupiter at 5.2 AU), with a mass about twice that of Jupiter. The shape of the velocity curve reveals that the orbit is approximately circular. (*The Anglo-Australian Planet Search Program*)

system with quasicircular orbits at several astronomical units (1 AU equals the Earth-Sun distance, about 93 million miles or 150 million kilometers). The planet orbiting HD 70642 is the system that most closely resembles Jupiter in the solar system, with an orbital distance of 3.3 AU, twice the mass of Jupiter, and a circular orbit. **Figure 2** shows the discovery data for this planet. The orbital period, 6 years, is immediately obvious from the Doppler velocity curve, while the velocity amplitude yields the planet mass.

Detection of true solar system analogs, Jupiter- and Saturn-like planets, requires 10 to 30 years of monitoring. Doppler velocity programs with sufficient sensitivity to detect these planets have been in existence for only a few years. The first crop of long-period planets should emerge by about 2010.

Eccentric planets. The vast majority of the planets found to date have markedly elliptical (noncircular) orbits. The planet-star distances for these systems vary by 10 to 70% over the course of the orbit. The surface temperatures of these planets have huge seasonal variations as the heating from the host star changes by hundreds of degrees between “winter” and “summer.”

The chance of a life-supporting planet in a system with a giant eccentric planet is very small. Many of the eccentric planets have “Earth-crossing” orbits that would spell immediate doom for small Earth-like planets. Even for the eccentric planets that are not “Earth-crossing,” the interaction with a small Earth-like planet would gravitationally destabilize the orbit of the small planet.

The unexpected discovery of eccentric planets has driven new theoretical understanding of the evolution of planetary systems. Planets forming in protoplanetary disks are initially in circular orbits. The gravitational interaction between a planet and the disk can cause the planet’s orbit to migrate, moving toward or away from the central star, and driving orbital eccentricity.

For planets that survive the disk phase of forma-

tion and evolution, another danger awaits: planet-planet interactions. A system with multiple giant planets might stably orbit for millions or even billions of years, then interact chaotically with some planets thrown into smaller orbits and others thrown farther out or even ejected from their system. Such a fate would have befallen the solar system if Saturn had ended up about twice as massive as it did. The Jupiter-Saturn interaction, like billiard balls on a gravitationally warped pool table, would have been devastating to the orbital stability of the “little” planets of the inner solar system, including Earth. *See CHAOS.*

Hot Jupiters. About 20 planets have been found orbiting their host stars with periods of 3 to 5 days. Due to the biases of the Doppler velocity method, these are the easiest type of planets to detect. The orbital distance of these planets is typically 0.04 AU, only about 10 stellar radii from the central star itself. At this distance, these planets are beyond blazing hot, about 2500°F (1400°C). Even at these blast furnace temperatures, hydrogen gas is gravitationally bound to Jupiter- and Saturn-mass planets. The proximity of the star raises enormous tides on the planets, and theorists agree that such tides will circularize a planet’s orbits typically in a million years or less. As a result, all the “hot Jupiter” planets are in circular orbits.

The chance that a given planet will transit (pass directly in front of) its host star is inversely proportional to the orbital distance. A hot Jupiter has about a 10% chance of transiting, while the odds for a planet at the Earth’s orbital distance are only one in 200. Of the 20 known hot Jupiters found from precision Doppler surveys, three are known to transit their host stars. Another eight transiting hot Jupiters have been found from photometric surveys. The prototype, HD 209458, passes directly in front of its star for about 2 hours every orbital period. During the transit the planet blocks about 1.4% of the host star’s light. The combined information from Doppler velocities and the transit yield the exact mass (63% of Jupiter) and size (40% larger than Jupiter) of the planet. The density of the planet, about half that of Saturn, implies that the planet is composed primarily of hydrogen and helium, similar to Jupiter and Saturn. This planet is “puffed up” relative to Jupiter and Saturn due to the intense heat of the nearby host star, like a balloon that expands when heated. *See TRANSIT (ASTRONOMY).*

The hot Jupiter planets could not have initially formed so close to their host star. The enormous heat and tidal forces of the star would prevent the planet from gathering the necessary amount of hydrogen and helium. The hot Jupiters must have formed farther out and then migrated into their roasting hot orbits.

Multiple-planet systems. The solar system consists of multiple planets (eight or nine depending on the current status of Pluto). While multiple-planet systems are probably common, about 80% of the known planets are single. This is almost certainly due to the biases of the Doppler detection technique, which cannot detect small earth-mass planets, and are most sensitive to giant planets in small orbits. Even with

these limitations, several multiple-planet systems have already emerged. These include systems with just eccentric planets (HD 38529, HD 217107) and systems with both a hot Jupiter and eccentric planets (upsilon Andromedae, 55 Cancri). Systems with multiple circular planets, such as the solar system, have yet to emerge.

It is expected that, as the time baseline of precision Doppler surveys increases, many of the single-planet systems will be revealed as multiple systems. When next-generation techniques capable of detecting small Earth-like planets come on line, virtually all planetary systems may be shown to be multiple.

Mass function. Since the Doppler technique is sensitive only to motion perpendicular to the “plane-of-the-sky,” Doppler velocities produce only a lower limit on a planet’s mass. For the case of an orbit that is edge-on as seen from Earth, the Doppler technique yields the full mass of the planet, while it is insensitive to orbits that are face-on or in the “plane-of-the-sky.” Without knowledge of the orbital inclination (the angle of the orbit with respect to the “plane-of-the-sky”), Doppler velocities give the product of the planet’s mass M times the sine of the inclination angle i (so-called $M \sin i$).

For studies of planet mass, the uncertainty of the inclination angle is not a major hurdle. Assuming random orbital inclinations as seen from Earth (it would be beyond extraordinary if the orbital axes of planetary systems were systematically pointed either toward or perpendicular to the Earth), two-thirds of all planets will have a true mass within a factor of 2 of the observed $M \sin i$, and only one in a thousand planets would have a true mass ten times larger than the derived $M \sin i$. The typical planet will have a true mass about 1.5 times larger than the minimum $M \sin i$ mass.

Figure 3 shows the observed planet mass function, the number of known planets as a function of the minimum $M \sin i$ mass of the planet. As this diagram shows, there are very few planets with a minimum mass greater than 5 Jupiter masses. There are many more low-mass planets relative to high-mass planets. The real planet mass function must increase

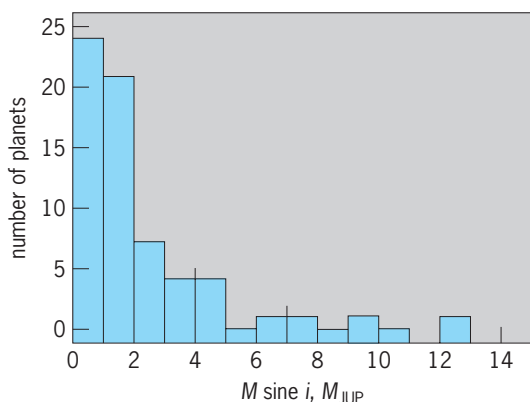


Fig. 3. Planet mass histogram, showing the number of known planets as a function of the minimum $M \sin i$ planet mass, measured in units of Jupiter’s mass, M_{JUP} . (California & Carnegie Planet Search Program)

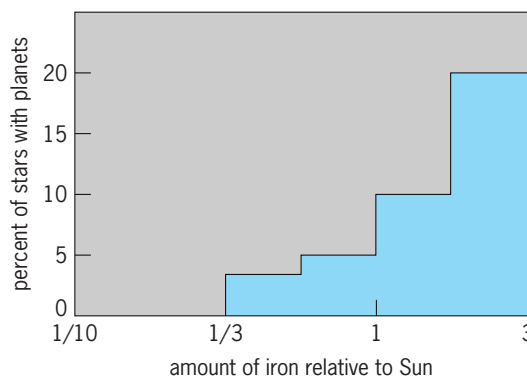


Fig. 4. Percent of nearby stars with known planets as a function of stellar metallicity. The plot shows that stars rich in elements heavier than hydrogen and helium (called “metals” by astronomers) are much more likely to have detectable planets, while stars with less than one-third the Sun’s heavy-element abundance have yet to yield any planets. (Courtesy of Debra Fischer and Jeff Valenti)

even more rapidly, since massive planets are the easiest to detect with the Doppler velocity technique, while the less massive planets are likely to be overlooked. Planets less than a Jupiter mass are more likely to be missed by the Doppler velocity technique, and below a Neptune mass (15 Earth masses) the Doppler technique is currently insensitive.

While the observed planet mass function reveals that there are roughly 20 planets of 1 Jupiter mass for every planet of 5 Jupiter masses, it would be wrong to conclude that there should be many tens or hundreds of Earth-mass planets for every Jupiter-mass planet. The solar system provides the only example of a system where both large and small planets can be detected. The inventory of the solar system consists of two massive planets (Jupiter and Saturn), two intermediate-mass planets (Uranus and Neptune), and four small planets (Mercury, Venus, Earth, Mars).

Metallicity. The two most common elements in the universe, accounting for about 99% of all matter, are the two simplest, hydrogen and helium. Not surprisingly, stars and giant planets also consist primarily of hydrogen and helium. The abundance of heavier elements in a star can be calculated from the spectrum of the star. Stars with greater than average abundances of heavier elements are much more likely to yield detectable planets (**Fig. 4**). This suggests that planets initially form from an accumulation of dust, rocks, and ices in the protoplanetary disk. Theoretical models show that protoplanets that reach 10 Earth masses can then gravitationally grab hydrogen and helium gas from the protoplanetary disk and rapidly grow to 100 Earth masses or more. See STELLAR POPULATION.

Prospects. Precision Doppler programs will continue to be the principal means of detecting planets orbiting nearby stars, at least until about 2015. All 2,000 Sun-like stars (stars with masses between 0.5 and 1.1 solar masses) within 160 light-years (50 parsecs) have been placed under survey, primarily by the two groups that have discovered about 90% of the known planets. By about 2015, these surveys will be sensitive to Jupiter- and Saturn-mass planets

orbiting beyond 4 AU. The central looming question is whether these planets will commonly be found in circular orbits, or alternatively if the architecture of the solar system is rare.

Of the greatest anthropocentric interest are Earth-like planets in quasicircular orbits capable of supporting life. The National Aeronautics and Space Administration (NASA) and the European Space Agency (ESA) have plans for new space-based telescopes to detect possible terrestrial analogs, including the transit missions *COROT* and *Kepler*, which are scheduled to launch in 2006 and 2007, respectively. These may be sensitive to Earth-size planets at 1 AU, providing valuable information on the incidence of such planets. As Earth-size planets block only 1 part in 10,000 of the light from the host star, these missions may be subject to astrophysical and instrumental effects that make interpretation of potential transit signals ambiguous. Transit missions cannot determine orbital eccentricity, and thus cannot determine if planets are in circular, solar system-like, orbits. These space-based transit missions are targeting stars at several hundred light-years, making follow-up by other techniques difficult.

The NASA Space Interferometry Mission (SIM) is extended to beyond 2015. A primary goal of SIM is the astrometric detection of planets as small as 3 Earth masses in Earth-like orbits around the nearest stars. SIM has the promise of measuring both the masses and orbital eccentricities of these planets, thus determining if they are solar system analogs.

NASA and ESA have long-range plans to build giant space-based telescopes capable of directly detecting (imaging) extrasolar planets around the nearest stars. Two techniques currently being researched, coronagraphs and nulling interferometers, block or “null” the light from the central star before it reaches the detector, allowing the faint light from the planets to be seen. The resulting images would be unresolved (a planet would appear as a single dot), but it would be possible to take spectra of the planets for the first time. Spectra could reveal water, oxygen, and other possible signposts of life. Though many technological hurdles remain to be overcome, these systems may optimistically be in place in 20 to 50 years. *See* ASTRONOMICAL SPECTROSCOPY; CORONAGRAPH; INTERFEROMETRY; PLANET.

R. Paul Butler

Bibliography. D. Goldsmith, *Worlds Unnumbered*, University Science Books, Sausalito, CA, 1997; M. D. Lemonick, Searching for other worlds, *Time*, 147(6): 52–57, Feb. 5, 1996; G. Marcy and P. Butler, Hunting planets beyond, *Astronomy*, 28(3):42–47, March 2000; D. Shiga, Imaging exoplanets, *Sky Telesc.*, 97(4):44–52, April 2004.

Extraterrestrial intelligence

Postulated entities beyond Earth with a level of intelligence and comprehension at least equal to that of humans at the present time. While extraterrestrial intelligence is usually envisioned as an advanced civilization, populated by creatures that have evolved

via Darwinian evolution on a planet vaguely similar to Earth, it could conceivably be artificial intelligence initially created by biological beings. Extraterrestrial intelligence is a subset of astrobiology, which encompasses all aspects of the existence of, and search for, extraterrestrial life. Astrobiology is sometimes referred to as exobiology or bioastronomy. *See* ASTROBIOLOGY.

Scientific rationale. The expectation that extraterrestrial intelligence exists derives from two facts and one assumption: (1) The universe is vast, with approximately 10^{11} galaxies (a total of about 10^{22} stars) within the reach of telescopes. This number is so large that even if the emergence of intelligence is improbable, such intelligence could still have arisen frequently. (2) The physics and chemistry of the universe are everywhere the same. This is known from astronomical observation. (3) Habitable, Earth-like planets of the type that might spawn intelligence, with thick atmospheres and liquid water on their surface, are not extraordinarily rare. This is a hypothesis, sometimes called the principle of mediocrity. According to this principle, the Earth is not extraordinary in any of its important properties. The principle dates, in its modern form, to Nicolaus Copernicus (1473–1543), who dethroned the Aristotelian idea of an Earth-centered cosmos. *See* UNIVERSE.

In addition to these general arguments, research gives support to the idea that extraterrestrial biology (not necessarily intelligent) might be plentiful. Since 1995, astronomers have detected planets around other, Sun-like stars. At minimum, 10–20% of such stars are now thought to have a solar system consisting of at least one orbiting body. The discovery technique is to measure the small motions of the star induced by the planet, a scheme that is most sensitive to massive worlds in tight orbits. It is still unknown what fraction of stars have small, rocky planets similar to Earth, but the *Kepler* and *Darwin* space-based telescopes, now being planned, should decide this question. *See* EXOPLANETS.

The possibility that some other worlds in the solar system could have spawned life has also increased. There is good photographic evidence that Mars once had lakes and possibly oceans, and may still harbor liquid water hundreds of meters below its surface. Direct evidence for ancient Martian life, claimed to have been found in a meteorite that is known to have come from the Red Planet, is highly controversial. A surprising discovery has been the growing indication for enormous oceans beneath the surface crust of several of Jupiter's satellites (Europa in particular, but also Callisto and Ganymede). Even Titan, a large satellite of Saturn swathed in a thick, hydrocarbon-laced atmosphere, is considered a possible (albeit unlikely) habitat for simple life. Whereas Earth was once thought to be the only solar system body that could support life, there are now several other candidates. If any of these has spawned indigenous life, it would demonstrate that biology is a commonplace occurrence. *See* JUPITER; MARS; SATURN.

But even if life is widespread, can intelligent life be expected to frequently evolve? This question will

probably remain unanswered until we either detect extraterrestrial intelligence or learn what drove the emergence of intelligence on Earth.

Search schemes. Given the enormous distances between the stars (the nearest is 4.2 light-years distant, or 4.0×10^{13} km or 2.5×10^{13} mi), it is beyond human capability to send robotic probes to directly search for intelligence elsewhere. A more promising approach is to look for signals that are either deliberately or inadvertently transmitted from their world to Earth. While there are many possible ways to signal, the most promising is electromagnetic radiation, and more specifically, light and radio waves. In the twentieth century, radio was recognized as an effective way to send information across space; and as early as 1900, Nikola Tesla mistakenly thought he had picked up transmissions from Mars. *See* ELECTROMAGNETIC RADIATION.

Radio waves travel at the speed of light which, according to current understanding of physics, is the fastest possible. Stars produce relatively little radio emission, and the universe is very “quiet” at radio frequencies, particularly in what is called the microwave part of the spectrum (approximately 1000–100,000 MHz), thereby making communication easier. Microwave signals also pass unperturbed through the gas and dust clouds that float between the stars. In 1959, Philip Morrison and Giuseppe Cocconi made calculations regarding the equipment and power necessary to signal over interstellar distances. It turned out that the requirements were not much beyond the type of installation that humans could build now. Consequently, the two physicists urged that a search be made for signals broadcast by other societies that had reached or surpassed the human level of technology. They also noted the advantages of the microwave band, and suggested that 1420 MHz (21 cm wavelength), which happens to be the frequency at which interstellar hydrogen naturally radiates, was the best part of this band to monitor. Since hydrogen is the most abundant element in the universe, this frequency will be known to all technologically sophisticated civilizations. *See* MICROWAVE; RADIO ASTRONOMY.

Frank Drake had independently reached the same conclusions, and in April 1960 he searched for artificial radio emissions from two nearby Sun-like stars, Tau Ceti and Epsilon Eridani. Drake used an antenna at the National Radio Astronomy Observatory in Green Bank, West Virginia, that was 26 m (85 ft) in diameter, and tuned his receiver near the hydrogen frequency. His search, whimsically named Project Ozma, became the prototype for today’s more comprehensive experiments, known as SETI (Search for Extraterrestrial Intelligence).

In the 1970s, the National Aeronautics and Space Administration (NASA) began a modest SETI program to build equipment and develop search strategies. In late 1992, the NASA program initiated experiments using the 305-m (1000-ft) Arecibo Radio Telescope in Puerto Rico and the 34-m (111-ft) Goldstone antenna in California. However, action by the U.S. Congress stopped the program within a year. Today, radio

SETI bears a strong likeness to the halted NASA effort, although advances in digital technology have improved the equipment. The two largest contemporary experiments are Project Phoenix, which used the Arecibo antenna and other large radio telescopes in a sensitive hunt for signals from nearby Sun-like stars, and SERENDIP, which sweeps large sections of the sky, using the same telescope, in a survey for powerful transmissions. Some of the data from SERENDIP are freely distributed for analysis by home computers, a project known as SETI@home. There are also SETI experiments in Australia, Italy, and Argentina. *See* RADIO TELESCOPE.

So far, no confirmed extraterrestrial signals have been found. However, as technology improves, radio SETI will speed and extend its search for artificially produced transmissions. A major development is the construction of the Allen Telescope Array, which can be used nearly full-time for SETI. It will eventually consist of hundreds of 6-m-diameter (20-ft) antennas located at the Hat Creek Observatory in northern California.

A second search method is to monitor star systems for brief flashes of light, signals that might be deliberately beamed in the direction of Earth with powerful lasers. By using mirrors to focus the lasers, it is straightforward to produce flashes lasting a nanosecond (10^{-9} s) or less that can greatly outshine the light from the transmitting civilization’s home star. Optical SETI, as searches for such light pulses are called, has already examined several thousand star systems. The research is being conducted at observatories in California and at Harvard and Princeton universities. An optical SETI telescope that uses a sky-scan strategy, rather than concentrating on individual star systems, has been built at Harvard.

While searching for signals is generally regarded as the most promising scheme for proving the existence of extraterrestrial intelligence, other approaches have been suggested. One might look for examples of astroengineering by very advanced societies, or possibly the radiation produced by high-powered, interstellar spacecraft. In much of the public’s mind, the many thousand sightings of unidentified flying objects (UFOs) each year are proof that some extraterrestrials are actually visiting Earth. Nearly all scientists are skeptical of such claims, and one of their principal objections is the lack of convincing physical evidence, despite more than a half-century of publicized sightings.

In addition to the expected improvements in search technology, advances in astronomy will aid SETI efforts. In particular, it will soon be possible to construct space-based telescopes that can directly image planets around other stars. By analyzing the light reflected from these planets’ atmospheres, it should be possible to find biomarkers such as oxygen and methane that would indicate the presence of biology. Future SETI experiments will benefit from being able to concentrate their efforts on star systems where worlds with life are known to exist.

Drake equation. In 1961, Drake devised a simple formula to calculate N , the number of broadcasting

societies in the Milky Way Galaxy. The computation multiplies the rate at which such societies arise, R , by the average length of time they survive, L . This is akin to computing the number of students at a university as the product of the number admitted each year times the number of years (4) before the average student graduates. The rate R is composed of terms that estimate the birth rate of stars, the fraction that have habitable planets, the chances that biology will appear, and the fraction of the worlds with life that develop intelligent, technically competent beings. While many of the terms in the Drake equation remain very speculative (particularly L , which depends on sociological factors), the formula remains a highly useful framework for discussing the subject of extraterrestrial life.

Fermi paradox. Enrico Fermi is said to have posed the question “Where is everybody?” in 1950. His remark was intended to point out that, while humans are in no position to colonize other star systems, advanced extraterrestrials—if they exist—might do so. Even if they require thousands of years to travel from one star to the next, an ambitious society could spread itself throughout the entire Milky Way Galaxy in only a few tens of millions of years. Since this is far less than the age of the Galaxy (about 12 billion years), it suggests that, if sophisticated and ambitious societies arose in the past, evidence of their presence should now be everywhere. Since that evidence is lacking, the implication is that the Galaxy is inhabited solely by humans.

There are many suggestions of how this paradox might be resolved; in other words, how the failure to see local evidence of alien activity could be reconciled with a Galaxy that we think might house many sophisticated societies. For example, it could be that interstellar colonization is so daunting that no one ever undertakes it. Perhaps humans are incapable of recognizing the widespread presence of intelligence. Or the extraterrestrials could know about humans, but have arranged for a “one-way mirror,” whereby they can watch, but humans cannot detect them (the “zoo hypothesis”). The Fermi paradox, while intriguing, remains a point of discussion rather than a key to new knowledge.

Consequences of discovery. If a signal or other verifiable proof of extraterrestrial intelligence is discovered, what would be the effects on human society? Clearly, this depends on how distant the intelligence is, and whether (in the case of a signal) an embedded message can be deciphered. It is possible that any intelligence that is found will be hundreds or even thousands of light-years distant, in which case two-way conversation or direct, physical contact would be exceedingly difficult.

In the case of a SETI detection, it is quite likely that the transmitting intelligence will be far more advanced than our own. This is because there is a good chance of detection only if L is large (at least thousands of years), and it is unlikely that the first signal received would be from a society that is only a few hundred years beyond us. (Of course, signals will not be found from societies that are not at least

as accomplished as ours.) If such advanced beings wish to transmit useful information, human society would be greatly altered. On the other hand, it may be that humans would never be able to decode the transmissions from such a sophisticated civilization. Even in that case, merely learning that intelligence has developed elsewhere in the cosmos would be a profound event.

Seth Shostak

Bibliography. I. Shklovskii and C. Sagan, *Intelligent Life in the Universe*, Holden-Day, New York, 1966, reissue edition, Emerson-Adams Press, 1998; S. Shostak, *Sharing the Universe*, Berkeley Hills Books, Berkeley, CA, 1998; S. Shostak and A. Barnett, *Cosmic Company*, Cambridge University Press, 2003; S. Webb, *If the Universe Is Teeming with Aliens, Where Is Everybody? Fifty Solutions to Fermi's Paradox and the Problem of Extraterrestrial Life*, Copernicus Books, New York, 2002.

Extrusion

The forcing of solid metal through a suitably shaped orifice under compressive forces. Extrusion is somewhat analogous to squeezing toothpaste through a tube, although some cold extrusion processes more nearly resemble forging, which also deforms metals by application of compressive forces. Most metals can be extruded, although the process may not be economically feasible for high-strength alloys.

Hot extrusion. The most widely used method for producing extruded shapes is the direct, hot extrusion process. In this process, a heated billet of metal is placed in a cylindrical chamber and then compressed by a hydraulically operated ram (Fig. 1). The opposite end of the cylinder contains a die having an orifice of the desired shape; as this die opening is the path of least resistance for the billet under pressure, the metal “squirts” out of the opening as a continuous bar having the same cross-sectional shape as the die opening. By using two sets of dies, stepped extrusions can be made. The small section is extruded to the desired length, the small split die is replaced by the large die, and the large section is then extruded.

Extrusion pressures and speed vary considerably, depending upon the size and shape of the section

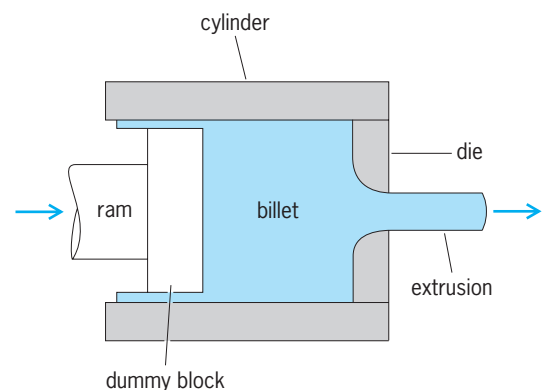


Fig. 1. Schematic representation of the direct extrusion process (hot).

and the mechanical properties of the metal. Some metals, such as magnesium and some aluminum alloys, require slow speeds of a few feet per minute, while others, such as some copper alloys, lead, and steel, are extruded at speeds of over 1000 ft/min (300 m/min). The extrusion speed is also somewhat dependent upon the temperature of the alloy. Considerable heat is generated by the process; if the extrusion speed is high, this heat cannot be dissipated, resulting in a rise in temperature. In some instances, the rise in temperature may be sufficient to melt or at least weaken the metal to the point at which the frictional stresses at the surface cause cracking.

The flow of metal is not uniform during extrusion of the billet. Because of the restraining effect of the die face and the frictional effects at the cylinder walls, the outer zones of the billet resist deformation and the flow occurs most rapidly at the center of the billet. Eventually, as the billet shortens, the different rates of flow at the center and surface result in the billet's becoming hollow. If the extrusion ram travels further, defects appear at the center of the extruded section. Therefore, a portion of the billet may be left in the cylinder and discarded. To prevent the oxidized surface layers of the billet from getting into the extruded product, the dummy block (or follower plate) in front of the ram is of slightly smaller diameter than that of the cylinder; thus, a thin sleeve of the billet is extruded out of the ram end of the cylinder and is subsequently discarded.

Lubricants are used to minimize friction and protect the die surfaces. Graphite is a common lubricant for nonferrous alloys, whereas for hot extrusion of steel, glass is an excellent lubricant.

Indirect, or inverted, extrusion was developed to overcome such difficulties as surface friction and entrainment of surface oxide of direct extrusion. In the indirect process, the ram is hollow, the die opening is in the dummy block, and the opposite end of the cylinder is closed. As the ram advances, the billet does not move as in the case of direct extrusion, and the metal is extruded backward through the die and the hollow ram. However, the process is not very popular because the hollow ram is weaker, resulting in lower machine capacity; trouble-free operation requires that the extruded product be straight and not hit the inside of the ram.

Tubular shapes are produced by a mandrel of the desired inside shape of the product. The mandrel may be fastened to the dummy block in a direct extrusion press if a pierced billet is used; or, in the case of a solid billet, a mandrel must first pierce the billet, after which the main ram advances.

In the production of lead cable sheathing or cored solder wire, the core material (cable or rosin flux) passes through a core tube (or die block) around which is heated lead under pressure. The core material passes from the core tube into a die cavity and the lead is extruded out with it, forming a casing around the cable or rosin core. The process is semicontinuous, molten lead being added to a vertical cylinder and pressure applied periodically. The lead is solid when it reaches the core material.

Cold extrusion. The extrusion of cold metal is variously termed cold pressing, cold forging, cold extrusion forging, extrusion pressing, and impact extrusion. The term cold extrusion has become popular in the steel fabrication industry, while impact extrusion is more widely used in the nonferrous field.

The original process (identified as impact extrusion) consists of a punch (generally moving at high velocity) striking a blank (or slug) of the metal to be extruded, which has been placed in the cavity of a die. Clearance is left between the punch and die walls; as the punch comes in contact with the blank, the metal has nowhere to go except through the annular opening between punch and die. The punch moves a distance that is controlled by a press setting. This distance determines the base thickness of the finished part. The process is particularly adaptable to the production of thin-walled, tubular-shaped parts having thick bottoms, such as toothpaste tubes.

A process requiring less pressure than backward extrusion is the forward-extrusion process, originally called the Hooker process (Fig. 2). A formed blank (usually a thick-walled cup) is placed in a die cavity and struck by a punch having a shoulder or enlarged section a short distance from the end. Upon contact with the blank, the nose or end of the punch starts to push the center of the blank through the die cavity, in a manner similar to the action occurring in deep drawing of sheet metal. After the punch has advanced a short distance, the shoulder comes in contact with the top of the thick wall of the blank. The punch shoulder then extrudes the metal through the annular space between the die and the end of the

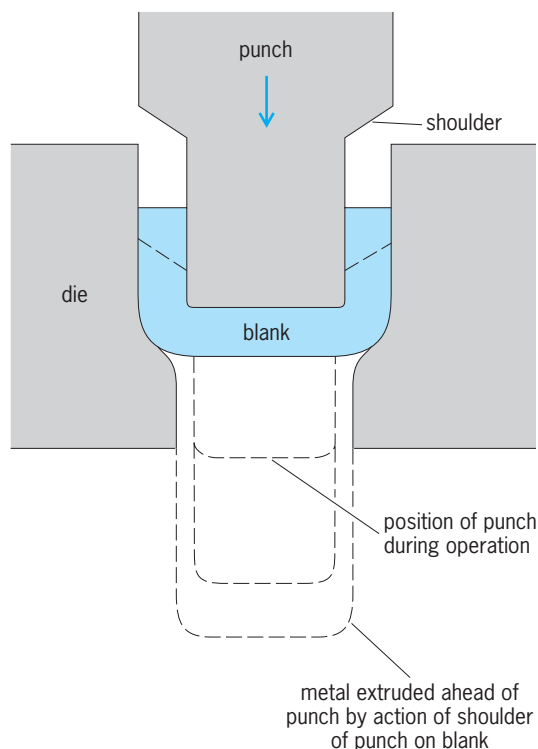


Fig. 2. Schematic of the forward cold-extrusion process, originally called the Hooker process.

punch. Thus, in forward extrusion, the metal moves in the same direction as the punch, whereas in backward extrusion the metal moves in the opposite direction. See SHEET-METAL FORMING.

The application of the cold-extrusion process to steel was developed in Germany about 1930 but was not released until after World War II. The success of cold extrusion of steel hinged on the discovery of a suitable lubricant and surface treatment. The steel is given a phosphate surface coating, which absorbs and holds the lubricant (soap emulsion, vegetable oils, or dry metal stearates) and prevents seizure between the metal and tools.

Advantages of cold extrusion are higher strength because of severe strain-hardening, good finish and dimensional accuracy, and economy due to fewer operations and minimum of machining required. See METAL FORMING. Ralph L. Freeman

Bibliography. M. Bauser, G. Sauer, and K. Siegert (eds.), *Extrusion*, 2d ed., 2006; P. Saha, *Aluminum Extrusion Technology*, 2000; S. L. Semiatin (ed.), *ASM Handbook Volume 14A: Metalworking: Bulk Forming*, 10th ed., 2005.

Eye (invertebrate)

An aggregation of photoreceptor cells together with any associated optical structures. Eyes occur almost universally among animals, and are possessed by some species of virtually every major animal phylum. However, the complexity of eyes varies greatly, and this sense organ undoubtedly evolved independently a number of times within the animal kingdom.

Structure and occurrence. The simplest invertebrate organs that might be considered to be eyes are clusters of photoreceptor cells located on the surface of the body. Pigment cells are usually interspersed among the photoreceptors, giving the eye a red or black color. Accessory structures, such as the lens and cornea, are usually absent. Simple eyes of this type, called pigment spot ocelli, are found in such invertebrates as jellyfish, flatworms, and sea stars (Fig. 1).

The most basic image-forming type of invertebrate eye probably arose from such patches of photoreceptor cells by an in-sinking of the sensory epithel-

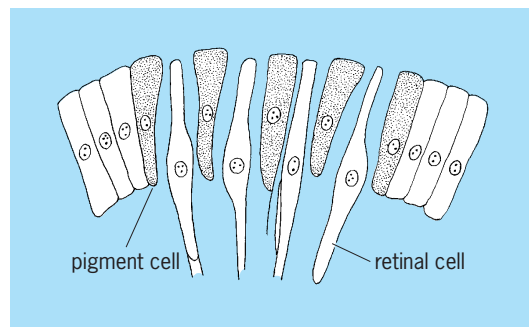


Fig. 1. Pigment spot ocellus, the simplest invertebrate eye. (After R. D. Barnes, *Invertebrate Zoology*, 2d ed., W. B. Saunders, 1968)

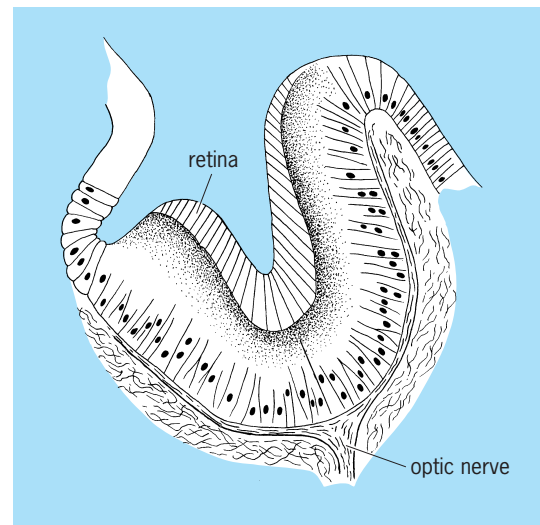


Fig. 2. Simple open eye of the marine limpet *Patella*. (After R. D. Barnes, *Invertebrate Zoology*, 2d ed., W. B. Saunders, 1968)

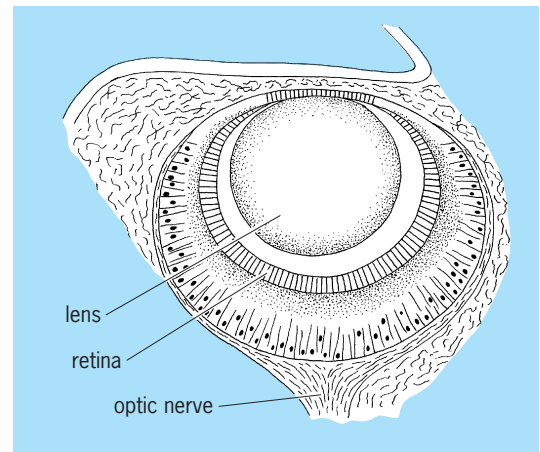


Fig. 3. Large, complex eye of the marine snail *Murex*. (After R. D. Barnes, *Invertebrate Zoology*, 2d ed., W. B. Saunders, 1968)

ium to form a cup. The cup, which represents the retina, may have remained as a pit open to the surface, may have sunk beneath the surface, or may have become closed in conjunction with the evolution of a cornea and lens. Such an evolutionary history is clearly suggested by the embryology and comparative anatomy of many invertebrates. In gastropod mollusks (snails), for example, the eyes vary from simple pigment cup ocelli to large complex eyes (Figs. 2 and 3). In embryonic development, the cups form as an in-sinking of the surface ectoderm. The eyes of free-living flatworms and marine annelids indicate a similar history (Figs. 4–7).

The photoreceptor cells in the wall of the cup are generally cylindrical or rounded with an axon fiber directed toward the brain. The photoreceptor cells may be directed toward, or away from, the light source, and the eye is said to be direct or indirect respectively. The distribution of direct and indirect eyes is variable. The eyes of mollusks are direct; the

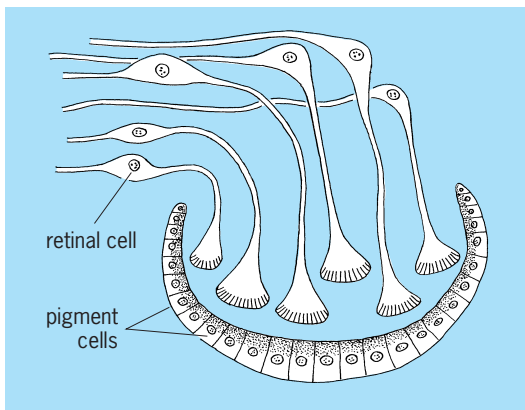


Fig. 4. Simple pigment cup ocellus of flatworm. (After R. D. Barnes, *Invertebrate Zoology*, 2d ed., W. B. Saunders, 1968)

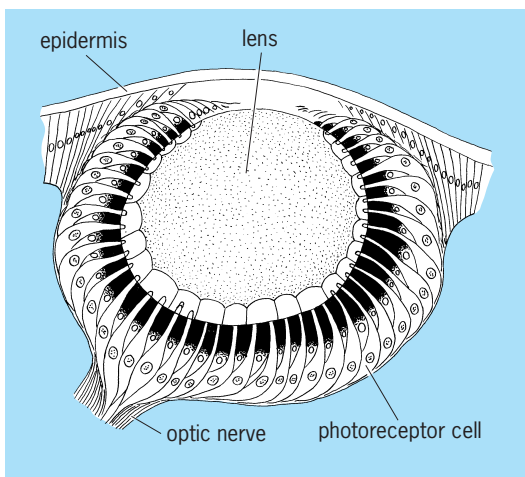


Fig. 5. Large, complex eye of the clam worm *Nereis*. (After R. D. Barnes, *Invertebrate Zoology*, 2d ed., W. B. Saunders, 1968)

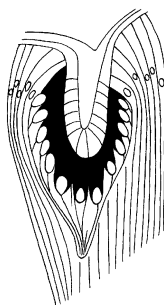


Fig. 6. Simple eye of *Ranzania*, a marine annelid. (After R. D. Barnes, *Invertebrate Zoology*, 2d ed., W. B. Saunders, 1968)

eyes of many flatworms, like those of vertebrates, are indirect. Most spiders possess 8 eyes, of which 2 are direct and 6 are indirect.

In bilateral cephalic invertebrates, the eyes are typically paired and located at the anterior end of the body. Although one pair is usual, as in mollusks and many arthropods, multiple pairs are not uncommon. Some polychaete annelids have 4 eyes, and scorpions may have as many as 12. The greatest number of

eyes is found in marine flatworms, where there may be over 100 ocelli scattered over the dorsal anterior surface and along the sides of the body.

The occurrence of eyes on parts of the body other than the head is usually correlated with radial symmetry or unusual modes of existence. Thus eyes are found around the bell margin of jellyfish and at the tips of the arms of some sea stars. Scallops possess many eyes along the margin of the mantle (Fig. 8). Some tube-dwelling fanworms have eyes on the plumelike head appendages, which project from the opening of the tube and are used in filter feeding. See FEEDING MECHANISMS (INVERTEBRATE).

Function. The primitive function of animal eyes was merely to provide information regarding the intensity, direction, and duration of environmental light. Many invertebrate eyes, even those with lens and cornea, retain those primitive functions. Cornea and lens merely provide the eye with a protective covering and concentrate the light on the retinal surface. The perception of objects is dependent upon several other factors, namely, the number of photoreceptors in the retina, the quality of the optics, and central processing of visual information.

Image formation. This has evolved as an additional capacity of the eyes of some invertebrates. The number of photoreceptor cells composing the retinal surface is of primary importance, since each photoreceptor cell or group of cells acts as the detector for one point of light. An image is formed by the retina through the association of points of light of varying intensity, much as an image is produced by an array of pixels on a computer monitor. The ability of an eye to form an image and the coarseness or fineness of the image are, therefore, dependent upon the number of points of light that are distinguished which, in turn, is dependent upon the number of photoreceptor cells composing the retina. A large number of photoreceptor cells must be present to produce even a coarse image. The great majority of invertebrate eyes cannot form a detailed image because they do not possess a sufficient number of photoreceptor cells. The number of photoreceptor cells might be sufficient to detect movement of an object, but is inadequate to provide much information about the object's form. Clearly, the question as to whether an animal can "see" depends on the meaning of the word. See PHOTORECEPTION; TELEVISION.

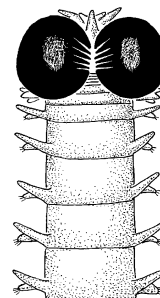


Fig. 7. Highly developed eyes at anterior end of pelagic marine annelid, Alciopidae. (After R. D. Barnes, *Invertebrate Zoology*, 2d ed., W. B. Saunders, 1968)

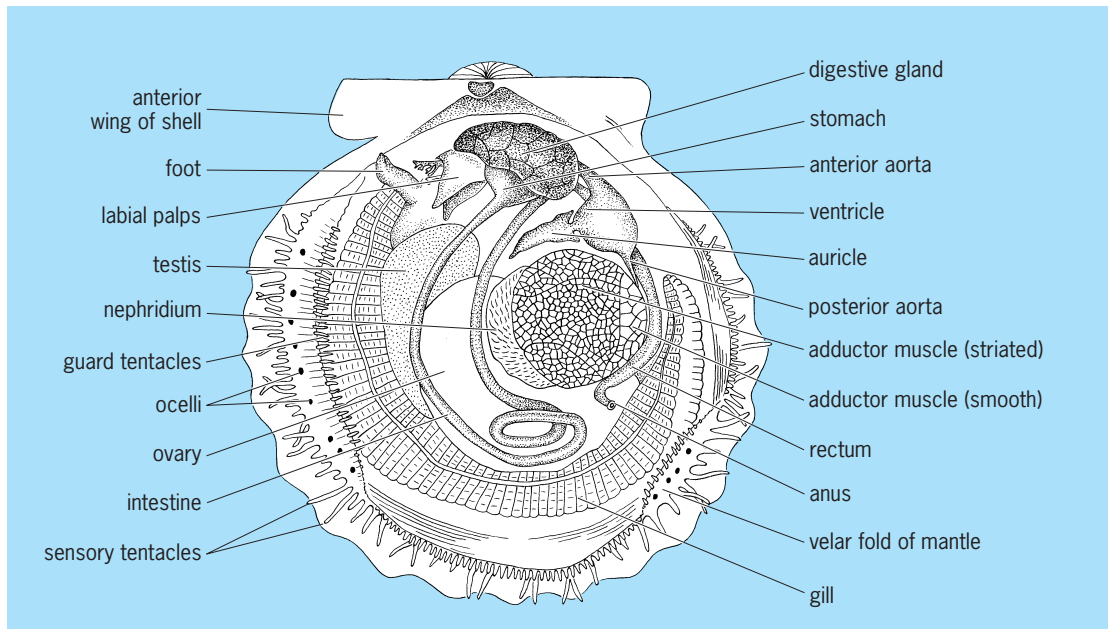


Fig. 8. Scallop *Argopecten* with one valve removed showing ocelli along the margin of the mantle. (After R. D. Barnes, *Invertebrate Zoology*, 2d ed., W. B. Saunders, 1968)

Focus. The focusing mechanisms of invertebrate eyes vary considerably. The focus of arthropod eyes tends to be fixed, that is, the distance between the optical apparatus and the retina cannot be changed. Thus objects are in focus only at a certain distance from the eye, determined by the distance between the lens and the retina. For example, jumping spiders, which are able to detect an object 30 cm away (a considerable distance for so small an animal), have tubelike anterior median eyes (Fig. 9). The tubular form increases the distance between the lens and the retina, enlarging the image for distant vision.

The oceanic family of swimming polychaete worms, Alciopidae, have the most highly developed

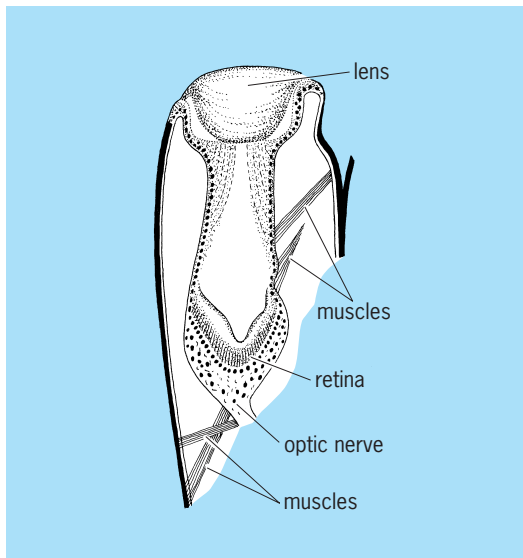


Fig. 9. Longitudinal section through the tubelike anterior median eye of jumping spider *Salticus*. (After R. D. Barnes, *Invertebrate Zoology*, 2d ed., W. B. Saunders, 1968)

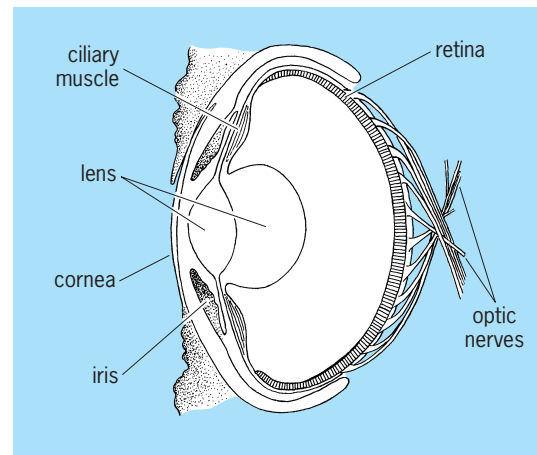


Fig. 10. Octopus eye. (After R. D. Barnes, *Invertebrate Zoology*, 2d ed., W. B. Saunders, 1968)

eyes of any of the segmented worms (Fig. 7). The eyes of these worms are focused hydrostatically. A bulb to one side of the eye injects fluid into the space between the retina and the lens, forcing the lens outward. Another mechanism is employed in octopods whereby lens movement is brought about by a ciliary muscle attached to the lens (as in aquatic vertebrates, like fish) [Fig. 10].

Interpretation. Despite the emphasis placed upon eye structure, an image is ultimately an interpretation by the central nervous system of signals received from the photoreceptors. Visual interpretation is still very poorly understood among invertebrates. The little information available is based upon a few experimental animals, notably the octopus and the honeybee.

Color discrimination has been demonstrated in crustaceans and also in insects, such as bees, flies, and butterflies. Invertebrate color vision can be

extremely complex, involving up to eight primary color classes in some crustaceans.

Visual object perception in invertebrates is correlated with a variety of behavioral adaptations. In fact, behavior frequently indicates that the eye is capable of some degree of image formation. Invertebrates which commonly possess complex eyes, which are probably capable of detecting and recognizing objects, include predators such as cursorial hunting spiders, squids and octopods, certain insects, some crabs and other crustaceans; those with highly controlled flight, such as bees, wasps, flies, dragonflies, and butterflies; those displaying courtship behavior, such as jumping spiders and fiddler crabs; or those whose orientation is by celestial clues, such as bees and some crustaceans.

Compound eyes. The compound eye of crustaceans, insects, centipedes, and horseshoe crabs has a sufficiently different construction from that of other invertebrates to warrant separate discussion. The structural unit of the compound eye is called an ommatidium (Fig. 11). The outer end of the ommatidium is composed of a cornea, which appears on the surface of the eye as a facet. Beneath the cornea is an elongated, tapered crystalline cone; in many compound eyes the cornea and cone together function as a lens. The receptor element at the inner end of the ommatidium is composed of one or more central translucent cylinders (rhabdome), around which are located several photoreceptor cells (typically 7 or 8).

The rhabdome is the initial photoreceptive element, and it in turn stimulates the adjacent photoreceptor cells to depolarize. Note that the photoreceptor element of each ommatidium functions as a unit and can respond only to one point of light. Thus

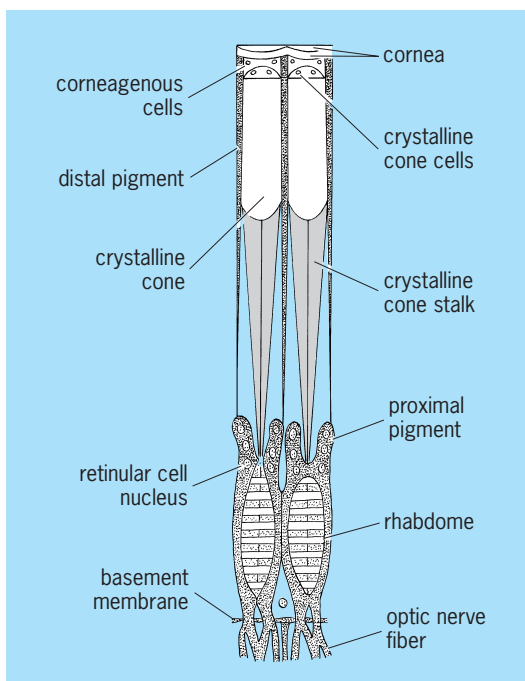


Fig. 11. Two ommatidia from compound eye of crayfish *Astacus*. (After R. D. Barnes, *Invertebrate Zoology*, 2d ed., W. B. Saunders, 1968)

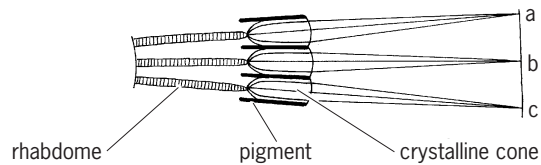


Fig. 12. Compound eye adapted for opposition image formation. Light rays from points a, b, and c. (After R. D. Barnes, *Invertebrate Zoology*, 2d ed., W. B. Saunders, 1968)

image formation is dependent upon the number of photoreceptor units present. The number of ommatidia composing a compound eye varies greatly. Some wood lice (crustaceans) possess no more than 25 ommatidia per eye, whereas the lobster possesses 14,000 and the dragonfly perhaps 30,000.

Pigment granules surround the ommatidium proximally and distally, forming a light screen that separates one ommatidium from another. The pigment granules migrate, depending upon the amount of light. In bright light the ommatidium is adapted by funneling light directly down to the rhabdome (Fig. 12), by extending the pigment screen, so that light received by one ommatidium is prevented from stimulating the rhabdome of another. Under these conditions the image produced is said to be appositional, or mosaic. The term mosaic has been misinterpreted to mean that a given ommatidium forms a separate image, even if only a part of the image. In general, however, the compound eyes function like any other eye—each photoreceptor unit represents one point in visual space. It is not obvious whether or not compound eyes have any special advantages over other eye designs, despite their universal occurrence in crustaceans and insects. However, in many arthropods the total corneal surface is greatly convex, resulting in a wide visual field. For example, the compound eye at the end of the eyestalk of a crab may cover an arc of 180° .

Nocturnal vision. The eyes of many nocturnal invertebrates are adapted for functioning in reduced light. Wolf spiders, hunters that run over the ground at night, illustrate eyes of this type. Of the eight eyes, six are indirect and are supplied with a layer of reflecting pigment, called a tapetum, which reflects light onto the photoreceptors. Like those of a cat, the eyes of these spiders appear to glow in the dark when facing a beam of light; the eyes of specimens the size of a quarter can be seen from a distance of 15 m (50 ft) or more!

Reflecting pigment is also found in the compound eyes of certain nocturnal insects and crustaceans. One type of compound eye forms a superpositional image by withdrawal of the screening pigment, so that light received by one ommatidium can strike the rhabdome of an adjacent ommatidium (Fig. 13). Under these conditions, a given rhabdome will receive light entering a number of facets so the image will appear much brighter. However, visual acuity is much lower in a superposition eye as compared to an apposition eye.

Polarization vision. Many invertebrate eyes are capable of seeing and analyzing patterns of polarized

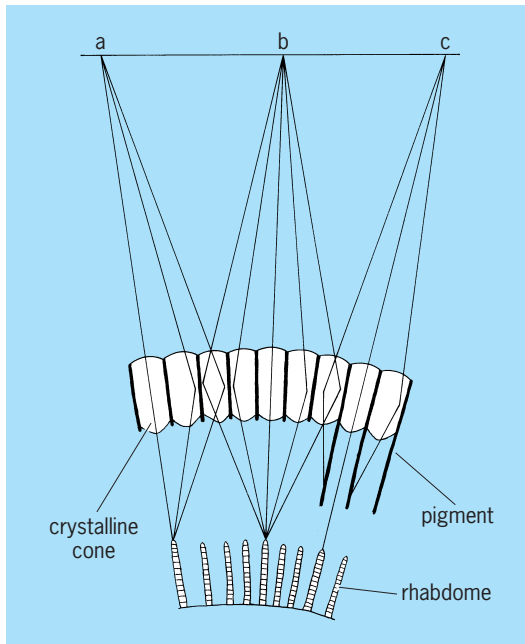


Fig. 13. Compound eye which is adapted for superposition image formation. The light from points a, b, and c is concentrated on an ommatidium. (After R. D. Barnes, *Invertebrate Zoology*, 2d ed., W. B. Saunders, 1968)

light in nature. This capacity reaches its apex in compound eyes, as well as in the simple eyes of cephalopods (octopuses, squids, and cuttlefish). Cuttlefish are known to communicate with each other with displays produced on their body surfaces that are visible only to animals that have polarization vision. Most invertebrates with polarization vision, however, use this ability to navigate with the assistance of patterns of polarization in the sky that occur naturally due to scattering of sunlight by the atmosphere. Bees and ants can find their way back to their nests or hives using only these celestial polarization cues. See EYE (VERTEBRATE).

Thomas W. Cronin; Robert D. Barnes
 Bibliography. M. A. Ali (ed.), *Photoreception and Vision in Invertebrates*, 1984; H. Autrum (ed.), *Handbook of Sensory Physiology*, vol. VII, pt. 6A-C: *Vision in Invertebrates*, 1981; P. A. Meglitsch and F. R. Schram, *Invertebrate Zoology*, 3d ed., 1991; D. G. Stavenga and R. C. Hardie (eds.), *Facets of Vision*, 1989.

Eye (vertebrate)

A sense organ that acts as a photoreceptor capable of image formation. The eye of vertebrates is constructed along a basic anatomical pattern (Fig. 1) which, in the diversification of animals, has undergone a variety of structural and functional modifications associated with different ecologies and modes of living.

Often compared with a camera (Fig. 2), the vertebrate eye is conveniently described in terms of its wall, cavities, and lens. See CAMERA.

The Wall

The wall of the eye consists of three distinct layers or tunics which, from outward to inward, are termed the fibrous, vascular, and sensory tunics.

Fibrous tunic. This continuous, outermost fibrous tunic comprises a transparent anterior portion, the cornea, and a tough posterior portion, the sclera. In the human, the cornea represents about one-sixth of the fibrous tunic, the sclera five-sixths; the transitional area where they meet is of particular functional importance and is termed the sclerocorneal junction or limbus (Fig. 1 and Fig. 3).

Cornea. The vertebrate cornea exhibits very few modifications in structure regardless of environmental influences. Its major constituent is connective tissue (both cells and fibers), regularly arranged and bordered on both anterior and posterior surfaces by an epithelium. The anterior epithelium is stratified, ectodermal in origin, and continuous with the (conjunctival) epithelium lining the eyelids. The posterior epithelium (termed the mesothelium or endothelium) is only one cell thick, perhaps mesodermal in origin, and ends at the limbus. The transparency of the cornea is attributed to the geometric organization of its connective tissue elements, its

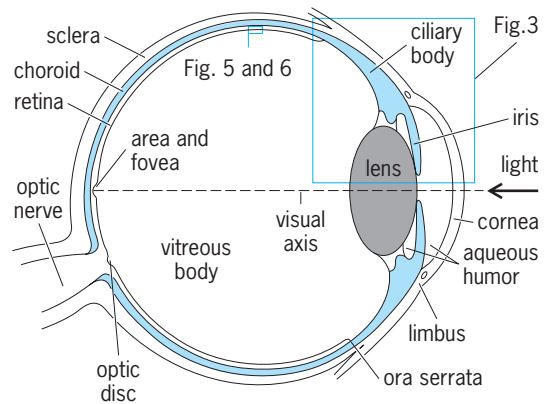
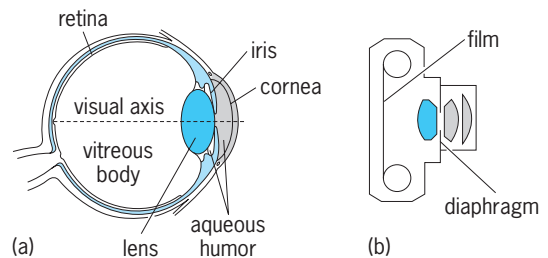


Fig. 1. Horizontal section through human eye.



Key:
 [grey box] anterior components
 [blue box] posterior components

Fig. 2. Comparison between (a) eye and (b) miniature camera. In both, the refractive part consists of anterior and posterior components, separated partially from one another by an iris diaphragm. The sensory retina corresponds to the light-sensitive film. (After E. Gardner et al., *Anatomy: A Regional Study of Human Structure*, 4th ed., Saunders, 1975)

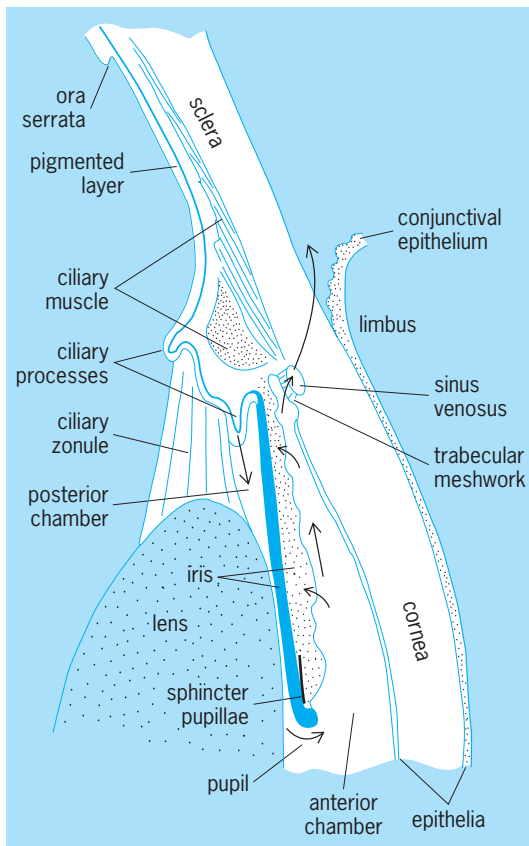


Fig. 3. Meridional section through anterior portion of human eye showing ciliary region and iridocorneal angle. The aqueous circulation (source, route, discharge) is indicated by arrows. (After E. Gardner et al., *Anatomy: A Regional Study of Human Structure*, 4th ed., Saunders, 1975)

constant state of deturgescence, and its chemical composition. It is the first ocular component traversed by the incoming light. For animals living in air, the cornea bends (refracts) this light toward the optic axis so that the visual image becomes placed on the sensory portion (retina) of the eye. The cornea provides the greatest refractive medium in these animals, whereas it appears to be optically functionless in animals living in water. Refractive errors such as myopia and astigmatism are common. See EYE DISORDERS; REFRACTION OF WAVES.

Sclera. This connective tissue tunic provides support for the eye and serves for the attachment (insertions) of the muscles that move it. It is cartilaginous in all elasmobranchs, anurans, reptiles (except snakes), and birds, as well as in some teleosts, urodeles, and mammals (monotremes).

Limbus. Located at the angle of the anterior chamber (Figs. 1 and 3), this small, circular transitional zone between the sclera and the cornea houses the major route for the discharge of aqueous humor from the anterior chamber: a circular venous sinus (canal of Schlemm) for transport of this fluid to the venous system, and a very porous mass of tissue (trabecular meshwork) which links this sinus with the anterior chamber. In birds and reptiles (except crocodiles and snakes), the limbus is further characterized by

the presence of a circle of flat (sometimes concave), overlapping bones, the scleral ossicles, forming a sclerotic ring. Usually 14 in number, these ossicles support a broad and strong limbus with an annular depression which is believed to aid in accommodation (changing the contour or the position of the lens for both near and far vision) in these animals.

Vascular tunic. The vascular tunic or uvea makes up the middle layer of the wall of the eye. It does not form a continuous layer around the eye but is deficient anteriorly, where the opening is termed the pupil. Beginning at the pupil, three continuous components of the uvea can easily be recognized: the iris, ciliary body, and choroid (Fig. 1).

Iris. The iris is a spongy, circular diaphragm of loose, pigmented connective tissue separating the anterior and posterior chambers and housing a hole, the pupil, in its center (Fig. 3). When heavily pigmented, the human iris appears brown; when lightly pigmented, blue. The anterior surface is lined by connective tissue cells with extensive interstices which permit the passage of aqueous humor. The posterior surface is covered by a double layer of pigmented epithelial cells, the iridial part of the retina. Situated in front of this epithelium are two muscle masses concerned with the opening (dilator pupillae) and closing (sphincter pupillae) of the pupil and, therefore, the regulation of light impinging upon the retina.

Pupil. The shape of the pupil varies considerably among vertebrates (Fig. 4). Both circular and slit-like pupils are found in all vertebrate classes. Circular forms are the most common; slitlike pupils are adaptations for nocturnality. Most slitlike pupils are oriented vertically, although horizontal slits are present in some species of all classes. In some fish, the iris possesses a single pupillary process or operculum which is capable of expanding into the pupil, relegating it to a circular slit. The greatest variety of pupillary shapes is exhibited by the amphibians.

Ciliary body. The ciliary body, triangular in section, is continuous with the root of the iris. The posterior epithelium of the iris continues along the internal surface of the ciliary body as a double layer of cells (ciliary processes) for the attachment of the ciliary zonule (suspensory ligament) of the lens. This ligament holds the lens in position and shape, and

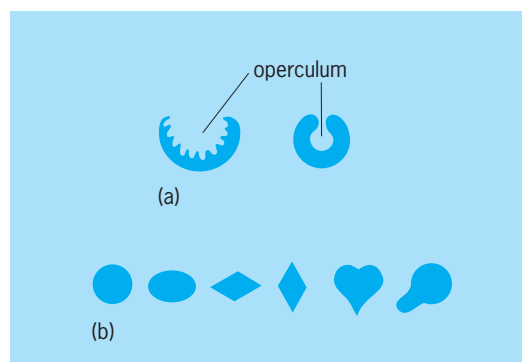


Fig. 4. Examples of various shapes of pupils in vertebrates: (a) pupillary opercula in fishes; (b) partially contracted pupil in different amphibians.

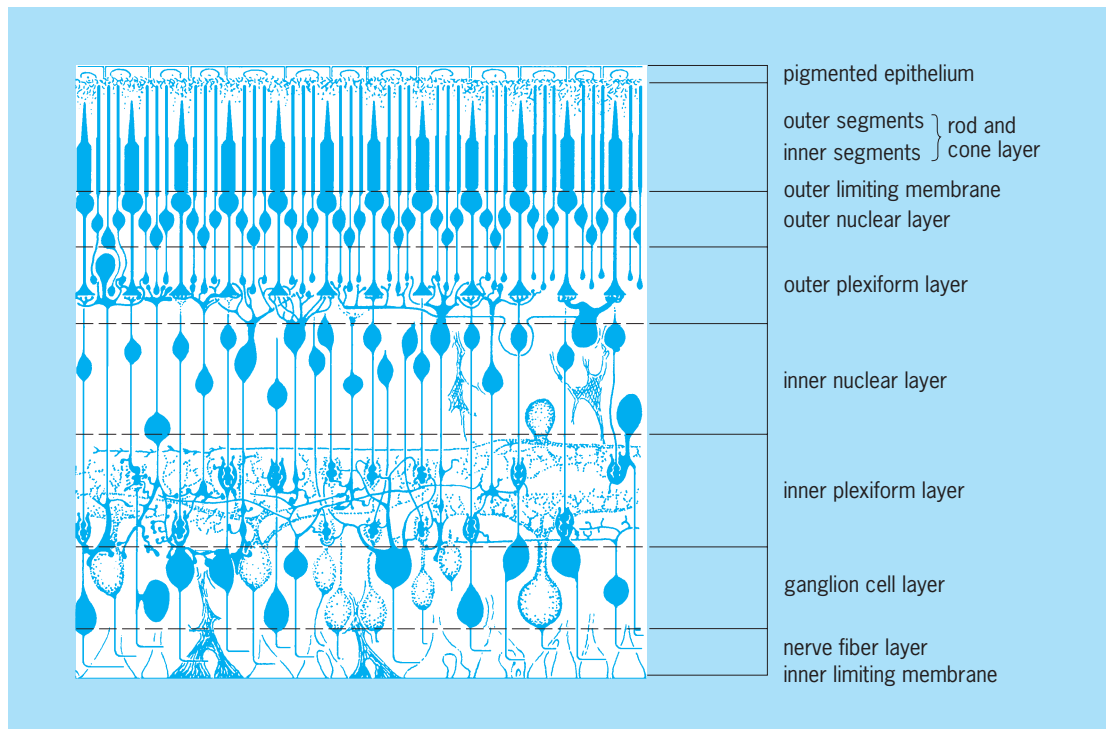


Fig. 5. Cross section of the vertebrate retina.

marks the posterior boundary of the posterior chamber. The inner layer of the ciliary epithelium contains no pigment. It produces aqueous humor which flows into the posterior chamber and thence into the anterior chamber (via the pupil) for final discharge through the sinus venosus (canal of Schlemm) via the trabecular meshwork. The continual production and removal of this fluid maintain the intraocular pressure of the eye (which is increased in glaucoma). The mass of the ciliary body, on the other hand, consists of muscular tissue (smooth in most vertebrates, skeletal in birds and reptiles) concerned with accommodation.

Choroid. The choroid is the most posterior portion of the uvea. It is directly continuous with the subepithelial portion of the ciliary body and consists primarily of blood vessels embedded within deeply pigmented connective tissue. The blood vessels are an important nutritional source for the retina (particularly the visual cells), which lies immediately internal to the choroid; the pigment serves to absorb stray light, thereby preventing internal reflections. In teleosts the choroid in the vicinity of the optic nerve develops a large capillary mass termed the choroidal gland. Although the function of this "gland" is not understood, it has been shown to be well suited for supplying the oxygen requirements of the teleost retina.

Sensory tunic. The retina is the sensory tunic of the eye. It has the form of a cup closely applied to the inner portion of the choroid, and, internally, it is slightly adherent to the semisolid vitreous body. Anteriorly, the sensory layers of the retina terminate at the liplike ora serrata, where they become continuous with the nonsensory epithelium lining

the ciliary body (ciliary part of the retina). Derived from the brain, the vertebrate retina contains the light-sensitive receptors (visual cells) and a complex of well-organized impulse-carrying nerve cells (neurons), all arranged into discrete layers in an orderly and remarkably similar fashion (Fig. 5). Basically, this stratification is produced by the spatial alignment of a three-neuron linear chain oriented perpendicularly to the surface of the retina, and by the orderly synaptic connections that these neurons make with each other and with other nerve cells (Fig. 6).

The visual cells or photoreceptors represent the first neurons in this series. They are nerve cells specialized for the reception of light stimuli and the subsequent conversion of this energy into electrical impulses. Two physiological types of visual cells occur in vertebrates: rods, which operate in dim light (scotopic vision) and are involved with visual sensitivity; and cones, which function in bright light (photopic vision) and are responsible for visual acuity and color perception. Both rods and cones transmit their responses to a second group of neurons, the bipolar cells, by means of a specialized contact or junction called a synapse. The bipolar cells, in turn, carry the signals by a similar means to a third group of neurons, the ganglion cells. Single nerve processes (axons) from each of the ganglion cells then leave the retina together as the optic nerve, and convey their information to higher centers in the brain. Further complicating this circuit is the presence of two other neurons (horizontal and amacrine cells) which maintain similar positions in all vertebrate retinas and which, by their synaptic contacts, must greatly influence the responses carried along the three-neuron

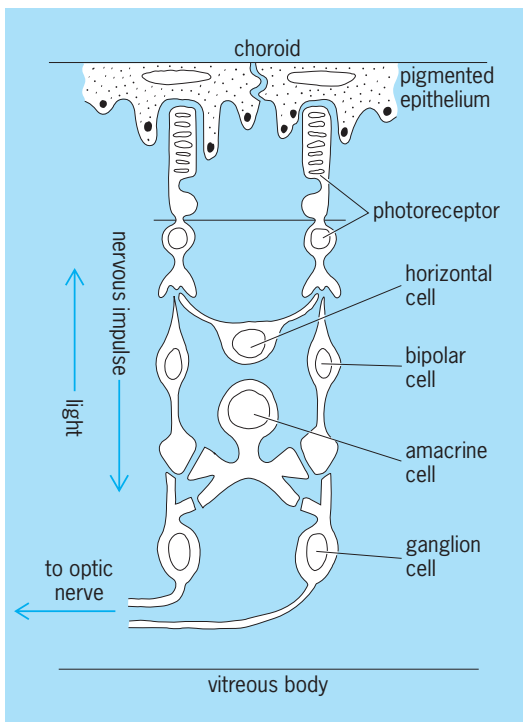


Fig. 6. Basic organization of the vertebrate retina. (After J. E. Dowling, *Synaptic organization of the frog retina: An electron microscopic analysis comparing the retinas of frogs and primates*, *Proc. Roy. Soc. London*, 170B:205–228, 1968)

chain. Horizontal cells exhibit their influence in the region of the receptor-bipolar junction, amacrines by the bipolar-ganglion cell contacts. Discrete, well-demarcated layers are formed as a result of the spatial alignment of these different cells and their synapses.

Aggregations of nuclei form three layers: outer nuclear, by nuclei from the photoreceptors; inner nuclear, by nuclei of bipolar, horizontal, and amacrine cells; and ganglion cell layer, by ganglion cell nuclei. The synaptic contacts between these nuclei form the plexiform layers. It should be emphasized that the vertebrate retina is an inverted one; that is, its receptors are situated in the deepest (most external) portion of the retina and point away from the incoming light. Light, therefore, must traverse almost the entire retina before striking the photosensitive elements. So that these relationships can be more easily understood, the retina is conveniently divided into 10 layers, beginning with its most external layer, the pigmented epithelium (Fig. 6). Finally, it should also be emphasized that the visual image is not imprinted on the retina, as on a film, but is coded and transmitted to the brain by way of the optic nerve, somewhat as in television. See TELEVISION.

Pigmented epithelium. Cuboidal in shape and heavily pigmented, this simple layer of epithelial cells forms an important barrier between the light-sensitive receptors (visual cells) and their blood supply, the choroid. As in the choroid, the pigmentation serves to absorb light and prevent its reflection. In fish, frogs, birds, and some reptiles the pigment gran-

ules effectively screen scattered light by migrating among the visual cells in the light and by receding into the cell body in the dark. In higher vertebrates, including the human, this close association between pigmented epithelium and visual cells involves the participation of the epithelium in the turnover of the rod photoreceptors, continuously phagocytizing the growing tips of the outer segments of the rods. Remnants of extruded rod tips, termed phagosomes, have been revealed in the electron microscope within the cytoplasm of the pigment epithelium in various stages of degeneration. A similar process has not been confirmed in cones. Aggregations of tightly packed disks or saccules, termed myeloid bodies, are also found in the cytoplasm of amphibians, birds, and reptiles, but their significance is not known. At the ora serrata the pigmented epithelium continues as the pigmented layer of the ciliary part of the retina.

Tapeta lucida. Many nocturnal or crepuscular vertebrates possess a reflecting layer, or tapetum lucidum, situated within the pigment epithelium (retinal tapetum) or choroid (choroidal tapetum). Densely packed crystals (guanine, riboflavin), rodlets (zinc complex), and lipid droplets or fibrils (collagenous) are responsible for the eyeshine produced by these reflectors, an adaptation for increasing light sensitivity at the expense of visual acuity. Retinal tapeta are found in bony fish, alligators, and the opossum; choroidal tapeta are characteristic of cats and ungulates.

Visual cells. The rods and cones of vertebrates generally occur as single units, but combinations of each type are frequently encountered in several vertebrate classes (Fig. 7). Double rods are scarce and are confined to some nocturnal lizards and snakes. Double cones, on the other hand, are distributed widely among vertebrates, for example, amphibians, birds, reptiles, and nonplacental mammals. Physiologically, they are regarded as one of the most obscure elements in vertebrate retinas. Whereas both double rods and double cones involve unequal, unfused members, the twin cones, which are another visual cell type unique to bony fish, exhibit identical components which are fused along the entire extent of their inner segments.

Cones appear to be adapted for photopic, or daylight, vision, based on correlations with the visual habits of the animals involved. Rods, which predominate in nocturnal vertebrates, are adapted for scotopic, or night, vision. Except for their external process, the structure of these cells does not reflect these differences. Indeed, both are organized in a very similar fashion and exhibit almost identical compartmentalization of organelles (Fig. 8) in all vertebrates. Shaped like a bipolar cell, each visual element contains a cell body with a nucleus, and two processes, one directed externally (toward the choroid) and functioning as the receptor end organ, the other directed internally (toward the vitreous body) for transmission of the visual image to the bipolar cells. Since the external processes are either rod- or cone-shaped and are easily recognized as such

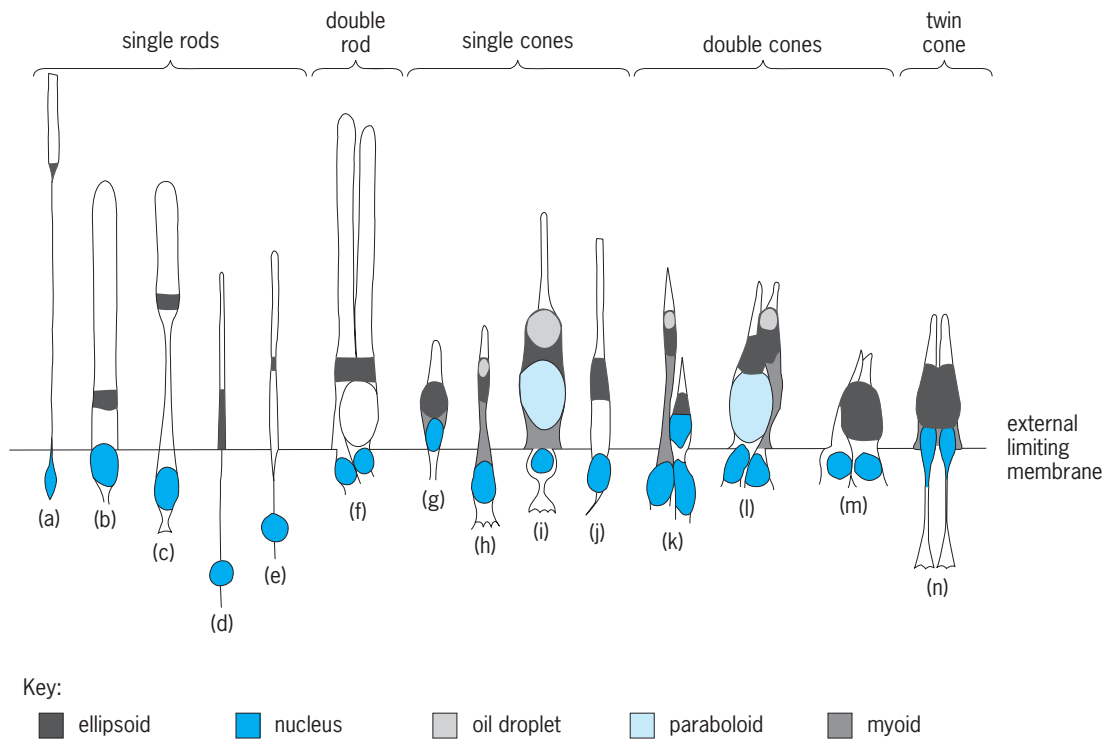


Fig. 7. Types of rods and cones found in different vertebrates: (a) goldfish, (b, c) frogs, (d) squirrel, (e) human, (f) gecko, (g) goldfish, (h) frog, (i) turtle, (j) human, (k) frog, (l) turtle, (m) snake, and (n) fish.

microscopically, the term “rod” or “cone” has been used to designate the entire cell. The external processes occupy, therefore, the rod and cone layer (Fig. 5). Each such process is divided rather evenly into a proximal inner segment and distal outer segment joined by a thin connecting stalk (Fig. 8). The outer segment represents the photosensitive receptor and

characteristically consists of a stack of membranous disks containing a light-sensitive visual pigment. A multiplicity of visual pigments occurs among vertebrates, but all contain vitamin A conjugated to a protein. The response of a vertebrate eye to light is determined by the type of visual pigment, the most well known of which are rhodopsin (rod) and iodopsin (cone). See PHOTORECEPTION; VITAMIN A.

The inner segments house most of the cellular organelles and inclusions in a very characteristic compartmental pattern. The mitochondria (respiratory organelles), for example, are always concentrated at the distal ends of the inner segments in the form of an ellipsoid, the name adopted for this entire mass, whereas the rest of the inner segment, termed the myoid, contains the Golgi apparatus and the bulk of the endoplasmic reticulum. In most vertebrates, except mammals, “myoid” is a very appropriate name for this region, because in these animals it has the ability to contract, like muscle. Unlike muscular contractions, however, these contractions are under the influence of light and in coordination with the migration of the pigment of the pigmented epithelium. In light, the pigment migrates toward the visual cells, while the rods elongate and the cones contract. These photomechanical responses shield the outer segments of the rods from the light and make those of the cones more accessible to it. The opposite reactions occur in the dark, when the rods become the dominant visual element.

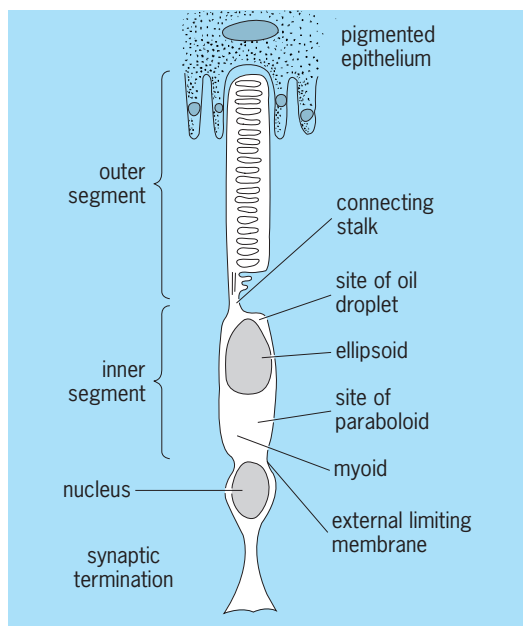


Fig. 8. Typical vertebrate visual cell showing the compartmentalization produced by the spatial arrangement of its organelles and inclusions. An oil droplet and a paraboloid are present only in some species.

In birds, reptiles (except snakes), and nonplacental mammals, the distal ends of the cone ellipsoids are further characterized by the presence of brightly colored (red, orange, yellow, or green) oil droplets

which, by their strategic position (to the incoming light) and their coloration, undoubtedly filter and modify the light destined for the outer segment. Rod myoids and some cone myoids in amphibians, birds, and reptiles likewise become compartmentalized by the presence of discrete deposits of glycogen which have a hyperboloid and paraboloid shape, respectively. The function of the rod hyperboloids and cone paraboloids, as they are called, is not known, although they have been thought to improve visual acuity or to provide an energy source for visual cell metabolism, or to perform both functions. The vertebrate retina displays a variety of other structural adaptations, some of which are believed to improve visual abilities under various environmental stresses; others are little understood and appreciated. Among the adaptations that have been shown to be characteristic features of some vertebrate groups are specialized retinal areas for acute vision, with or without a foveal depression and supplementary nutritive devices for the retina, for example, the pecten oculi and the conus papillaris.

Areas and foveae. An important adaptation for improving visual detail in vertebrates is the formation of circumscribed thickenings (circular, oval, or linear) of the retina resulting from localized increases in the number of visual cells and the other retinal neurons associated synaptically with them. Such thickenings, termed areas of acute vision, appear in some members of all vertebrate classes and reach their greatest development in birds, in which one to three distinct areas may be found in the same retina. Only a single area occurs in humans; it is colored yellow and is called the macula.

These areas most commonly develop a depression or pit of various depth and size called a fovea. Formed by a lateral displacement of the cellular components of the retina, foveae most often lie in the central or lateral region, or both, of the fundus (fovea centralis or fovea lateralis), and in birds are positioned so as to lie directly in the visual axis for either monocular (fovea centralis) or binocular (fovea lateralis) vision (Fig. 9). In the human only a single fovea is present; it is situated in the center of the fundus and contains only cones.

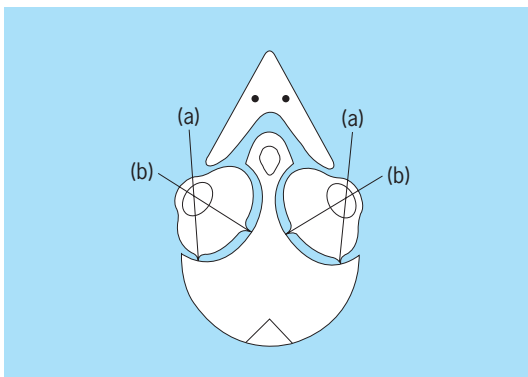


Fig. 9. Visual axes of bifoveate retinas: (a) for the lateral fovea and binocular vision; and (b) for the central fovea and monocular vision.

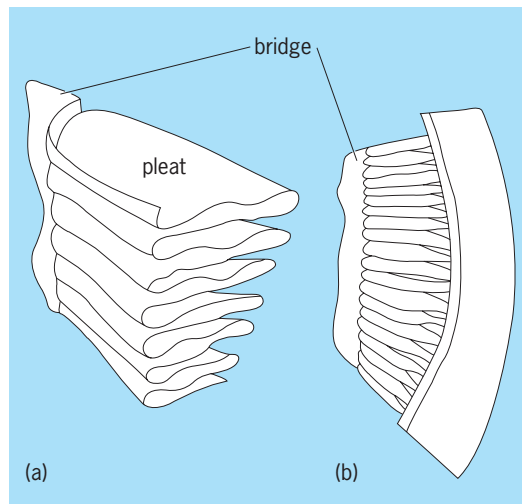


Fig. 10. Pleated pecten of the domestic chicken: (a) pleats and "bridge" in three-dimensional view; (b) longitudinal view. (After G. K. Smelser, ed., *The Structure of the Eye*, Academic Press, 1961)

Pecten and conus papillaris. In addition to the choroid, the nutritional supply to the vertebrate retina draws on several different sources, the most bizarre of which are the heavily pigmented, vascular bodies arising from the optic nerve head and projecting into the vitreous body of birds (pecten oculi) and lizards (conus papillaris). In lieu of retinal or vitreal vessels, these homologous organs are believed to provide nourishment for the inner layers of the vertebrate retina. The avian pecten has been the most extensively studied and has been found to assume a conical, vaned, or pleated appearance. The pleated type (Fig. 10) is by far the most common and varies considerably in size, shape, and number of pleats. The largest and most pleated pectens are found in active diurnal birds with high visual acuity, for example, hawks and eagles. Nocturnal birds with poor vision tend to have smaller and simpler ones. Because of its unique intravitreal location, the pecten has also been implicated in many other functions, particularly those concerned with vision. For example, the pecten is believed to cast shadows upon the retina that aid in the detection of movement and perhaps also, by acting as a sextant, enable the bird to estimate the position of the Sun and thereby aid in navigation. See VISION.

The Cavities

Three cavities or chambers are present within the vertebrate eye: anterior, posterior, and vitreous (Fig. 1). The anterior and posterior chambers are continuous with one another at the pupil and are filled with a fluid, the aqueous humor, which is continuously secreted by the ciliary body and removed by the sinus venosus (canal of Schlemm). The eye is normally maintained in a distended state by the (intraocular) pressure created by this fluid. The vitreous cavity, on the other hand, is filled with a semisolid material, the vitreous body, which is fixed in amount and relatively permanent. Its consistency is not uniform

in all vertebrates, however. In mammals and owls, for example, the vitreous body is completely gelatinous, whereas in many birds (pigeon, chicken, turkey) and fish (tuna), it is partly fluid and partly gelatinous.

The Lens

The lens is a transparent body, supported by thin suspensory fibers (ciliary zonule) and by the vitreous body behind and by the iris in front (Figs. 1 and 3). It is completely cellular, the anterior cells forming a thin epithelium, and the posterior cells, much elongated, forming the so-called lens fibers. The entire lens is surrounded by an elastic capsule which serves for the attachment of the ciliary zonule.

In all vertebrates the lens functions in accommodation, either by moving backward and forward or by changing its shape. Some vertebrate lenses are colored yellow, for example, in the lamprey, some snakes and diurnal geckos, most squirrels, and tree shrews and humans. In ground squirrels and prairie dogs they are almost orange in color. Colored lenses, like colored oil droplets, are considered as intraocular filters; they represent adaptations for increasing visual acuity by enhancing color contrasts. In humans lenticular color increases with age. An opacity of the lens is a cataract.

David B. Meyer; Ronan O'Rahilly

Electrophysiology of Rods and Cones

Visual information perceived by the vertebrate eye is fed to the brain in the form of coded electrical impulses that are initiated by the light-sensitive, visual-pigment-containing outer segments of the rods and cones. Light striking the outer segments is absorbed by these pigments, which form an integral component of their membranous discs resulting—in the case of rhodopsin, for example—in the isomerization of the 11-*cis*-retinal chromophore to all *trans*-retinal. The outcome of this photolytic process is a change in electrical activity at the plasma membrane enclosing the outer segments, and a sudden and drastic decrease in its permeability (particularly to Na^+). The net result is a hyperpolarization response, or increased negativity of membrane potential, with the rods exhibiting slightly slower but somewhat longer-lasting responses than cones. Hyperpolarization generates a membrane current that spreads to the inner segment and finally to the synaptic terminal, where it regulates the release of neurotransmitter and thus controls the flow of information from the visual cells to other retinal cells (bipolars, horizontals, other photoreceptors).

Whereas the outer-segment plasma membrane is depolarized in the dark and is permeable to Na^+ and other ions (K^+ , Ca^{2+} , Cl^-), the light-adapted hyperpolarized membrane exhibits closed ionic channels (with resultant decrease in intracellular Na^+) in response to a decrease in the level of cyclic guanosine 3',5'-monophosphate (cyclic GMP) owing to its rapid hydrolysis. Cyclic GMP is directly responsible for regulating the permeability of the plasma membrane by opening these channels (in the light). Its

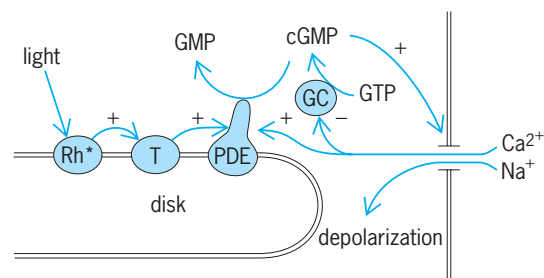


Fig. 11. Summary of the interactions occurring in the rod outer segments during light activation. Within the disc membrane, photolytic rhodopsin (Rh^*) activates transducin (T), which in turn activates phosphodiesterase (PDE), leading to the breakdown of cyclic GMP (cGMP) to inactive GMP. In the dark, cyclic GMP maintains open ionic channels in the outer-segment plasma membrane, permitting Na^+ to enter the cell. In the light, decreased cyclic GMP results in a closure of these channels, a decrease in intracellular Na^+ levels, and a hyperpolarization of the cell. A decrease in Ca^{2+} concentration acts primarily as a negative feedback, leading to an increase in cyclic GMP levels. (After J. E. Dowling, *The Retina: An Approachable Part of the Brain*, Harvard University Press, 1987)

concentration is controlled by a peripheral membrane enzyme, phosphodiesterase, which in turn is activated by transducin, an intracellular messenger protein generated by a photolytic intermediate of rhodopsin (Fig. 11). Since one molecule of photoactivated rhodopsin can react with many molecules of transducin, an amplification of the visual cells' response is produced, the final amplitude being enhanced by degradation of cyclic GMP by phosphodiesterase and subsequent inhibitory feedback (cones) from the horizontal cells.

Since photoreceptors are depolarized in the dark, their axon terminals continually release a transmitter (for example, acidic amino acids such as l-glutamate) that hyperpolarizes (inhibits) the bipolar cell, and since this cell is hyperpolarized in the dark, it is prevented from releasing its excitatory transmitter at the ganglion cell synapse so that the synapse is not excited. In the light, on the other hand, hyperpolarization of the visual cells causes a decrease in the amount of inhibitory transmitter released at the bipolar synapse, leading to a depolarization of the latter, which in turn increases the amount of excitatory transmitter released at the bipolar-ganglion synapses and affecting the ganglion cells.

A change in the light energy taking place across the retina also initiates a transient action potential, the electroretinogram, which is recorded as a difference in potential between the cornea and the back of the eye and exhibits three principal deflections: A, B, and C waves. The C wave has been attributed to the pigment epithelium but may be triggered by the photoreceptors, which have also been implicated as the site of origin of the A wave. David B. Meyer

Accommodation

Ocular adjustments for the sharp focusing of objects viewed at different distances is termed accommodation. Each vertebrate class has its own mechanism of accomplishing accommodation, and most of

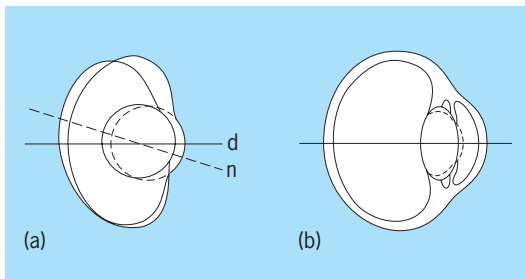


Fig. 12. Methods of accommodation in vertebrates. (a) In the teleost (bony fish) eye, on the approach of an object the lens returns to its normal relaxed (forward) position for near vision; it accommodates for distant vision. (b) In the mammalian eye, on accommodation for near vision, the lens becomes more curved. The visual axes for near (n) and distant (d) vision are indicated.

these mechanisms involve specific muscular actions. Two basic methods of accommodation are used by vertebrates: movement of the lens, either forward (elasmobranchs, amphibians, snakes) or backward (teleosts, lampreys), and changes in the shape of the lens (reptiles except snakes, birds, and mammals) [Fig. 12].

Cyclostomes, fish, and amphibians have developed specialized muscles for lenticular movement—the cornealis, retractor lentis, protractor lentis, respectively—whereas higher vertebrates, including humans, utilize the musculature of the ciliary body to change the shape of the lens. In birds and reptiles (except snakes) lens shape change is accomplished by a squeezing of the lens at the equator. In mammals, contraction of the ciliary muscles relaxes the tension that normally exists in the ciliary zonule, thereby permitting the lens to relax and change shape (become more spherical for nearer objects).

Adaptations

The preference of an animal for a certain type of illumination has produced a variety of ocular adaptations which have served to improve its chances for survival. On the basis of their visual habits, for example, vertebrates are classified into three major groups: diurnal, nocturnal, and arrhythmic.

Diurnal vertebrates. Members of this group are active primarily during daylight, and thus their eyes are suited to sharp vision at the expense of light sensitivity. Lizards, birds, and primates, which feed on small objects (seeds and insects), require such vision, as do animals that depend upon great speed and dexterity to capture prey or elude enemies. Such vertebrates possess large eyes with a large retinal surface housing chiefly cones, as well as rather flattened lenses and a long “throw” of the image from the lens to the retina (Fig. 13a and b).

Nocturnal vertebrates. These vertebrates are active chiefly at night. Their eyes, therefore, are modified for greater sensitivity to light at the expense of sharp vision. Vision in these animals, for example, bats, is not as important as sound and touch. Nocturnal eyes are characterized by a rod-rich retina, a larger and

broader lens with a larger pupil, and a noticeable decrease in the distance, or throw, from the lens to the retina (Fig. 13c). Some nocturnal eyes (for example, in owls, Fig. 13d) are tubular in shape, which results in a small retinal image but an increased throw from the lens to the retina.

Arrhythmic vertebrates. Arrhythmic, or 24-h, vertebrates are equally active by day and by night. Among them are the bony fish, frogs, crocodiles, and larger land mammals (elephants, bears, lions, wolves). Because the eyes of these vertebrates must function equally well under both extreme conditions of illumination, they have combined those ocular adaptations which are not mutually exclusive, for example, duplex retinas (rods and cones), visual cell and pigment migration, and mobile pupils. Their lenses are generally intermediate in size (between nocturnal and diurnal), and the throw to the retina is also intermediate in length (Fig. 13e).

Development

The development of the human eye is summarized in Fig. 14. At 3 weeks, a groove (the optic sulcus) appears on the internal surface of the brain wall (neural plate) (Fig. 14a). The sulcus rapidly becomes deeper (Fig. 14b), and the communication (rostral neuropore) between the cavity of the brain and the outside of the embryo becomes narrower and soon closes. The optic sulcus has now been transformed into a cavity, which, together with its wall, is known as the optic vesicle (Fig. 14c and d). The right and left optic vesicles communicate widely with that part of the forebrain termed the diencephalon (Fig. 14e). The optic vesicle (Fig. 15) induces the overlying surface ectoderm to form a lens disk (Fig. 14f), and the subjacent wall of the vesicle constitutes a retinal disk. During the next few days the lens and retinal disks become invaginated to form a lens pit and the optic cup, respectively (Fig. 14g). Subsequently, the cavity of the optic cup becomes obliterated, and the two retinal layers maintain contact with each other. The outer layer of the cup becomes the pigmented epithelium of the retina, and the inner, invaginated layer develops into the sensory retina (Fig. 14b). By this time (5 weeks) the lens has developed into the lens vesicle and has lost contact with

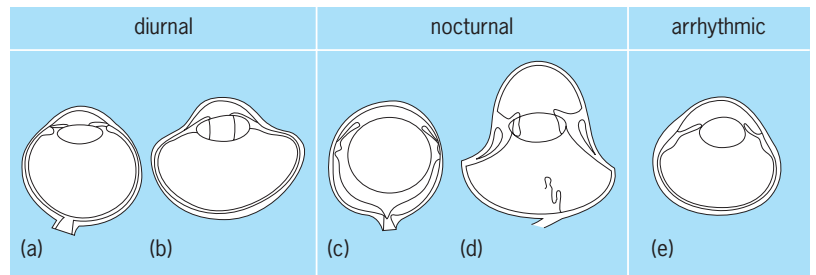


Fig. 13. Intraocular proportions of various species in relation to habit; (a) human, (b) pigeon, (c) mouse, (d) owl, and (e) dog. A tubular eye, such as that shown in e, is also found in birds of prey.

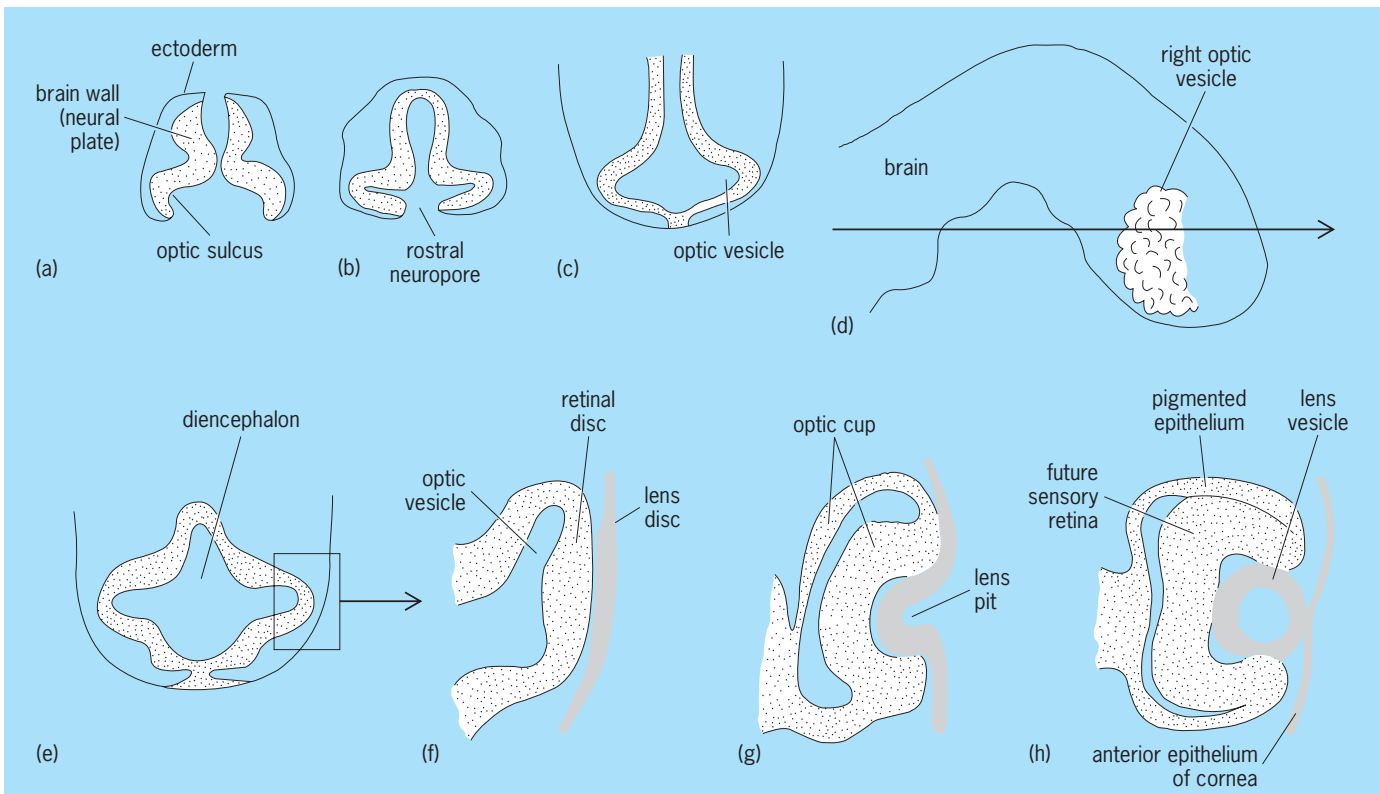


Fig. 14. Development of the eye in the human embryo at (a, b, c) 3¹/₂ weeks; (d, e, f) 4 weeks; (g) 4¹/₂ weeks; and (h) 5 weeks. (After R. O’Rahilly, *Contr. Embryol. Carnegie Inst.*, 38:1–42, 1966)

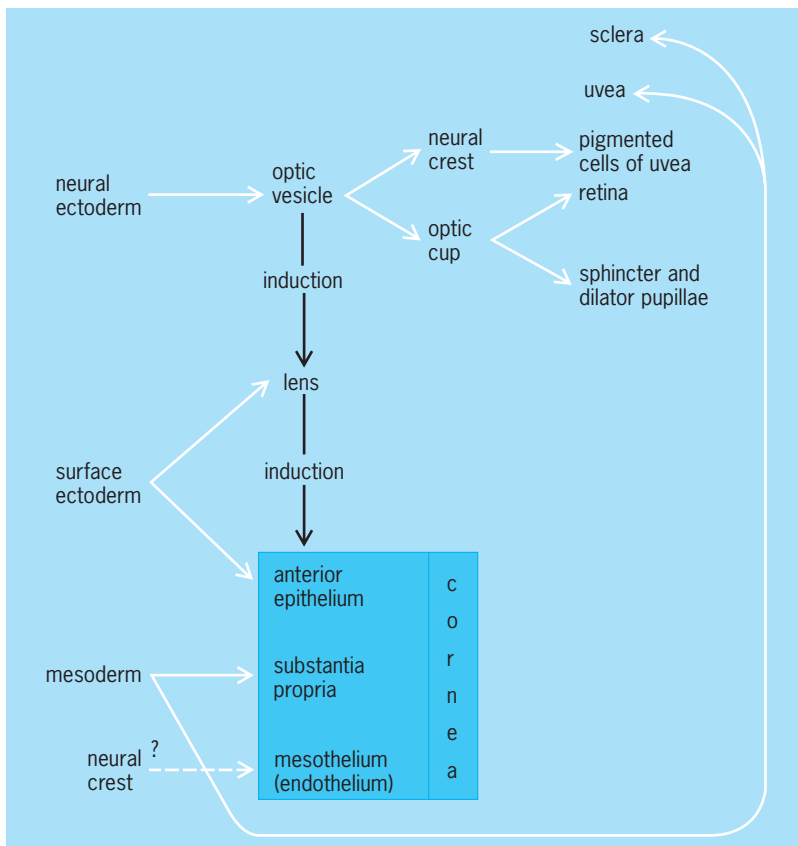


Fig. 15. Simplified scheme showing the chief embryological sources of the ocular tissues.

the surface ectoderm, which becomes the anterior epithelium of the cornea (Fig. 14b). Mesodermal invasion contributes to the rest of the cornea, the uvea, and the sclera. The posterior cells of the lens vesicle grow forward and obliterate the lens cavity. A simplified scheme of the derivation of the ocular tissues is provided in Fig. 15. See NERVOUS SYSTEM (VERTEBRATE); NEURULATION.

David B. Meyer; Ronan O’Rahilly

Bibliography. M. A. Ali and M. A. Klyne, *Vision in Vertebrates*, 1985; H. Davson, *Physiology of the Eye*, 5th ed., 1991; J. E. Dowling, *The Retina: An Approachable Part of the Brain*, 1987; J. Heggie, *The Use of the Eyes in Movement*, 1985; M. H. Wake (ed.), *Hyman’s Comparative Vertebrate Anatomy*, 3d ed., 1992.

Eye disorders

Disorders of form and function that affect the eye. Disorders of form affect the way the eye looks or feels; disorders of function affect vision. Normal vision requires proper performance from the entire visual system, from the precorneal tear film to the occipital cortex. See EYE (VERTEBRATE); VISION.

External eye. The eyelids are susceptible to a variety of disorders, many of which are common to the facial skin in general. The eyelids contain many glands (Meibomian, Zeis, Moll, Wolfring, Krause) that are susceptible to acute or chronic infection. Acute infection produces a hordeolum, or sty, which is a

localized nodule up to several millimeters in diameter. Moderate discomfort and lid swelling may develop rapidly. The lids are easily distensible, and a soft, puffy swelling may involve the entire upper or lower lid. Chronic infection results in a granulomatous nodule, or chalazion, which involutes spontaneously; however, this regression may take several months to a year or more. In many cases, surgical removal is indicated.

Conjunctival and corneal disorders. The conjunctiva is a mucous membrane that lines the inner surface of the eyelids and covers the outer surface of the anterior third of the globe except for the corneal epithelium, with which it is contiguous. Because the conjunctiva is connected to the inner surface of the eyelids and the anterior surface of the globe, it forms a barrier between the outside world and the retrobulbar space. Conjunctivitis is an inflammation of the conjunctiva, regardless of cause, which may include viral, bacterial, or chlamydial infection; mechanical or chemical trauma; allergy; and so on. The symptoms of conjunctivitis depend on the cause and generally include conjunctival redness (hyperemia), swelling (chemosis), mild to moderate discomfort, and tearing. If the conjunctiva is the only structure involved, vision is usually affected minimally or not at all.

Infectious conjunctivitis is often termed pink eye because of the bright red appearance of the conjunctiva. Bacterial infections are usually responsible for pink eye in children. Infectious conjunctivitis in adults is commonly caused by viral agents. Chlamydial infections are responsible for inclusion conjunctivitis and ophthalmia neonatorum in the United States and trachoma in many arid regions of the third world. Repeated chronic chlamydial infections characteristic of trachoma cause progressive conjunctival scarring and shortening, resulting in an inturning of the eyelids, and trauma to the cornea by eyelashes (trichiasis), exposure to the elements, a lack of lubrication, and abrasion by the scarred inner surface of the eyelids. In many villages in northern Africa, for example, trachoma is the leading cause of blindness in adults. Like endemic infections in any part of the world, trachoma is treated by behavior modification and specific therapy.

The cornea and conjunctiva are adjacent. The cornea is the window through which light passes to the retina; harm to its form or function is likely to result in visual impairment ranging from mild image distortion to blindness. Trachoma is a case in point. The initial infection involves the conjunctiva, but ultimately the disease results in blindness due to a scarred, vascularized, opaque cornea. *See VISUAL IMPAIRMENT.*

If the cornea and conjunctiva are both involved, the condition is referred to as keratoconjunctivitis. Adenovirus infection causes keratoconjunctivitis. If lubrication by the lacrimal glands or accessory lacrimal glands is interrupted by infection, trauma, or autoimmune disease, one of the dry-eye syndromes may ensue. The most common of these syndromes frequently affects the aging population, causing a

nonspecific loss in the quantity and quality of the tears. The symptoms include a mild foreign-body sensation, light sensitivity (photophobia), and gritty feeling in the eyes that progresses throughout the day. Paradoxically, increased tearing (epiphora) can be a symptom of dry eyes. In such cases, tear function is poor because the normal mucoid and oily components of the precorneal tear film are insufficient. Localized irritation results and leads to reflex increase in the secretion of the lacrimal or watery component of the tear film. Keratoconjunctivitis sicca is an autoimmune disease in which the lacrimal and salivary glands are diffusely infiltrated with an inflammatory reaction. Therapy may be difficult in severe cases.

Glaucoma. Glaucoma is a serious eye disorder with many subclassifications. Combined, they affect 1–2% of the population of the United States and Europe, and as much as 10–12% of the population in certain isolated regions of the world. Glaucoma is characterized by intraocular pressure sufficiently high to cause characteristic damage to the optic nerve. Acute glaucoma is associated with pain, redness of the eye, and sudden, often reversible, visual loss. Chronic glaucoma generally causes no symptoms other than progressive peripheral visual loss. Chronic glaucoma is often called the silent thief of vision because the person may not notice a problem until most of the vision has been destroyed. If untreated, the glaucomatous eye may lose all visual functions. *See GLAUCOMA.*

Cataract. The leading cause of blindness in the world is cataract, an opacity within the crystalline lens. Cataracts can be caused by trauma, chronic drug use, or endocrine imbalance, but the overwhelming majority have no specific cause and are associated with aging. Except for rare cataracts caused by medication or associated with systemic disorders, there is no non-surgical treatment. Effective therapy almost always involves surgical removal of the lens, during which a prosthetic intraocular lens is usually implanted. The typical symptom of cataract formation is a gradual painless loss of vision in one or both eyes. Because the process is largely limited to a physical clouding of a single refractive structure, surgery frequently restores excellent vision. *See CATARACT.*

Retinal disorders. The retina is the sense organ of the eye. Its rods and cones, or photoreceptors, and their supporting cellular elements receive light energy and convert it to electrical impulses which travel along the optic nerve to the brain. Degeneration of the retina produces a painless distortion or loss of vision. The function may be disrupted in many ways. The retina itself may be affected by by infectious processes. Cytomegalovirus retinitis, for instance, is a frequent cause of visual loss in individuals with acquired immune deficiency syndrome (AIDS). Primary retinal degeneration may occur in many forms. Retinitis pigmentosa tends to manifest itself early in life and progresses from peripheral to central visual loss; prognosis depends in part on the hereditary pattern of the disease. Involuntary macular degeneration is fairly common in the elderly. It affects the central vision first and rarely spreads

to the peripheral vision. Between these two extremes is a very large number of progressive retinal dystrophies and degenerations, involving central or peripheral vision or both. Although the ultimate visual loss is variable, it is often severe.

Retinal detachment occurs when the sensory layers of the retina are separated from their supporting foundations. It is classified as rhegmatogenous (caused by a retinal hole) or nonrhegmatogenous. Rhegmatogenous retinal detachment occurs spontaneously or following trauma. Nonrhegmatogenous retinal detachment occurs as a final stage of such pathologic conditions as retinopathy of prematurity or diabetic retinopathy. The symptoms of retinal detachment are a painless and sudden segmental or total visual loss in one eye. Treatment is aimed at reestablishing the connection between the sensory-neural retina and its supporting structures. In some complicated cases, or when there has been pathologic involvement of the vitreous body, the vitreous may need to be removed surgically. This is accomplished with vitrectomy instruments that amputate and aspirate small amounts of the vitreous gel.

Diabetic retinopathy is a common cause of severe retinal disease. Research has developed many methods of interrupting the progression of diabetic retinopathy. The most effective include laser therapy for localized retinal edema, hemorrhage, and neovascularization, and vitrectomy for late-stage diabetic retinopathy with vitreous hemorrhage and retinal traction.

Michael Drake

Bibliography. M. L. Callahan (ed.), *Current Therapy in Emergency Medicine*, 1987; F. W. Newell, *Ophthalmology: Principles and Concepts*, 8th ed., 1996; D. Vaughan and T. Asbury, *General Ophthalmology*, 15th ed., 1998.

Eyeglasses and contact lenses

Lenses worn to improve vision. The primary function of the eye is to focus light energy on the retina. In the optics of a normal eye, parallel rays of light are refracted into sharply focused points of light on the retina. Refraction is necessary because all natural light rays are divergent. See EYE (VERTEBRATE); REFRACTION OF WAVES.

The degree of divergence (or minus power) varies with distance: the closer the observer to the light source, the greater the divergence. The unit of measure of divergence is the diopter; the divergence of light rays in diopters is the reciprocal of the distance of the observer from the light source. For example, if the observer is 0.1 m (4 in.) from the light source, the light rays reaching that observer's eye have a divergence of 10 diopters; 6 m (20 ft) from the source the light rays have a divergence of one-sixth of one diopter. For practical purposes, the human eye cannot distinguish rays with one-sixth of one diopter from parallel rays, and therefore 6 m (20 ft) is referred to as optical infinity. The eye is approximately $1/460$ m (1.67 cm or 0.65 in.) and therefore requires about 60 diopters of converging (plus) power to bring light

rays parallel at the cornea into focus on the retina. See GEOMETRICAL OPTICS; VISION.

Myopia. If the eye is longer than 1/60 m (0.65 in.), or if its converging media (cornea and lens) impart more than 60 diopters to parallel light rays, they come into focus in the vitreous cavity in front of the retina, and the image on the retina is made up of blurred circles rather than focused points of light. In either case the eye possesses too much converging or plus power to form a good image from parallel light rays.

This may be remedied in two ways. First, the observer may move closer to the object. At the correct distance from a myopic eye, the divergence due to proximity to the object (the minus power of the rays reaching the eye) neutralizes the excess plus power of the eye's refracting system, thereby giving a sharp image on the retina. This explains the term nearsighted. A myopic person clearly sees things close to the eye.

A second method of correcting nearsightedness is to introduce a minus lens into the path of light striking the eye. This can be done with concave lenses, thus adding divergence to the light entering the eye. This divergence combines with the excess convergence of the eye to give clear images.

Hyperopia. The converse is true for farsightedness. The farsighted or hyperopic eye is too short or too weak in refracting power. The convergence applied to incoming light is insufficient to focus it at the retina, and blurred circles form rather than focused points of light. To correct this, plus lenses that add convergence to incoming light are used.

Presbyopia. Accommodation is the ability of the eye to change its refractive power for sharp viewing across different distances. When objects are held close for reading or detailed examination, the refracting (plus) power of the normal eye should increase to focus the more divergent light on the retina. Ability to accommodate decreases with age, and by the mid-forties or early fifties it is often difficult to read print at a normal distance. This condition is known as presbyopia and is corrected with reading glasses, or bifocals in which the distance prescription is in the upper segment and plus power in a lower segment.

Contact lenses. In addition to eyeglasses worn in frames, contact lenses can be used to assist the eye in forming a sharp image. The contact lens rests directly on the precorneal tear film; its front surface becomes the initial refracting surface of the eye. The space between the lens and the cornea is filled with the precorneal tear film. Contact lenses can be either divergent or convergent.

Contact lenses are classified as hard, soft, or gas-permeable hard according to the materials that they are made of, their rigidity or flexibility, and their permeability to oxygen and other gases. Traditional hard contact lenses are the most rigid. They float on the precorneal tear film, through which the cornea must eliminate most of its waste. Hard contact lenses tend to be smaller in diameter than the cornea and to move freely over it if fitting properly. These lenses

provide excellent vision correction where suitable, but because they require an active precorneal tear film, they tend to be tolerated better at a younger age.

Soft lenses are usually larger in diameter than hard contact lenses and are very thin and flexible. They are made of different materials than the hard lenses and are permeable to oxygen and other gases. Soft contact lenses tend to take the shape of the cornea, and if it has an irregular curvature, may not correct the vision as effectively as hard contact lenses. Soft contact lenses tend to be more comfortable and better tolerated, particularly by older people whose precorneal tear film would not support a hard contact lens.

A gas-permeable contact lens is an improvement of the hard contact lens. It is more rigid than a soft lens but more permeable to gases than the older hard contact lenses. Such lenses may give satisfactory correction where the corneas are irregular or for other reasons are not corrected satisfactorily with soft contact lenses. They tend to be better tolerated than ordinary hard contact lenses.

Contact lenses are not entirely harmless, and have been associated with a specific immune response known as giant papillary conjunctivitis and with dangerous infections capable of destroying an eye. In some users, extended-wear contact lenses may be left in place for several weeks to months at a time. Disposable contact lenses are also available, which can be left in the eye for 1-2 weeks. For many individuals, contact lenses are an effective alternative to spectacles. *See* LENS (OPTICS).

Michael Drake

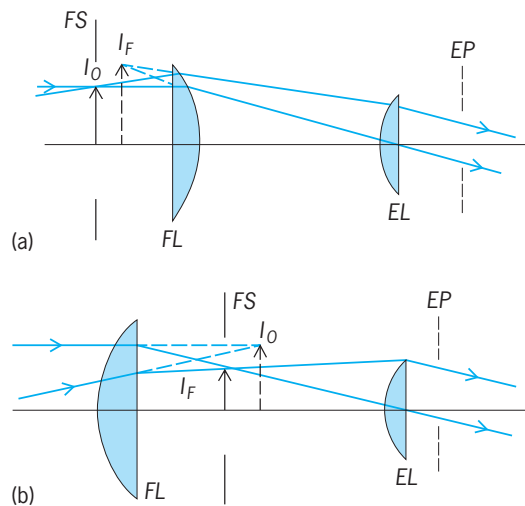
Bibliography. D. D. Michaels, *Visual Optics and Refraction: A Clinical Approach*, 3d ed., 1985; M. L. Rubin, *Optics for Clinicians*, 2d ed., 1974, revised 1993.

Eyepiece

A lens or optical system which offers to the eye the image originating from another system (the objective), at a suitable viewing distance. The image can be virtual. *See* OPTICAL IMAGE.

In modern instruments, most eyepieces (also called oculars) are not independently corrected for all errors. They are designed to balance out certain residual aberrations of the objective or (as in the microscope) of a group of objectives, for instance, chromatic difference of magnification. *See* ABERRATION (OPTICS).

The Ramsden eyepiece consists of two planoconvex lenses, the field lens and the eye lens, with their plane sides out. Both of these lenses have the same power and focal length; their separation is equal to their common focal length. The field lens is near the



Eyepieces: (a) Ramsden and (b) Huygens. FL = field lens; EL = eye lens; FS = field stop; EP = exit pupil or eye point; I_O = image (real or virtual) formed by the preceding system; I_F = image (virtual or real) formed by the preceding system and the field lens.

image formed by the objective. The Ramsden acts as a field stop. The Kellner eyepiece is a Ramsden eyepiece with an achromatic eye lens. *See* GEOMETRICAL OPTICS; LENS (OPTICS).

The Huygens eyepiece also consists of two planoconvex lenses, but the plane sides of both lenses face the eye. The focal length of the field lens is in general three times that of the eye lens, and the separation is twice the focal length of the eye lens. The Huygens eyepiece is usually corrected for astigmatism and distortion. A type of Huygens eyepiece in which the eye lens is achromatized to compensate for the color errors of the objective is called a compensating eyepiece. The field lens may also be achromatized.

Huygens and Ramsden eyepieces (see *illus.*) are in general used for low magnifying powers. For higher powers, eyepieces assume more complex forms to correct astigmatism, field curvature, and distortion for larger field angles.

The orthoscopic eyepiece consists of a single lens, made up of three cemented elements, to which a planoconvex lens is added.

Complex types of oculars enable corrections to be made over a field angle of 70° or more on the image side. They are used especially in military instruments.

Max Herzberger

Bibliography. M. Laikin, *Lens Design*, 1990; D. Malacara and Z. Malacara, *Handbook of Lens Design*, 1994; Optical Society of America, *Handbook of Optics*, 2 vols., 2d ed., 1994; W. J. Smith, *Modern Optical Engineering*, 3d ed., 2000; W. J. Smith and Genesee Optics Software, *Modern Lens Design: A Resource Manual*, 1992.