



Charged-Particle Optics

P. W. Hawkes

CNRS, Toulouse, France

- I. Introduction
- II. Geometric Optics
- III. Wave Optics
- IV. Concluding Remarks

GLOSSARY

Aberration A perfect lens would produce an image that was a scaled representation of the object; real lenses suffer from defects known as aberrations and measured by aberration coefficients.

Cardinal elements The focusing properties of optical components such as lenses are characterized by a set of quantities known as cardinal elements; the most important are the positions of the foci and of the principal planes and the focal lengths.

Conjugate Planes are said to be conjugate if a sharp image is formed in one plane of an object situated in the other. Corresponding points in such pairs of planes are also called conjugates.

Electron lens A region of space containing a rotationally symmetric electric or magnetic field created by suitably shaped electrodes or coils and magnetic materials is known as a round (electrostatic or magnetic) lens. Other types of lenses have lower symmetry; quadrupole lenses, for example, have planes of symmetry or antisymmetry.

Electron prism A region of space containing a field in which a plane but not a straight optic axis can be defined forms a prism.

Image processing Images can be improved in various ways by manipulation in a digital computer or by optical analog techniques; they may contain latent information, which can similarly be extracted, or they may be so complex that a computer is used to reduce the labor of analyzing them. Image processing is conveniently divided into acquisition and coding; enhancement; restoration; and analysis.

Optic axis In the optical as opposed to the ballistic study of particle motion in electric and magnetic fields, the behavior of particles that remain in the neighborhood of a central trajectory is studied. This central trajectory is known as the optic axis.

Paraxial Remaining in the close vicinity of the optic axis. In the paraxial approximation, all but the lowest order terms in the general equations of motion are neglected, and the distance from the optic axis and the gradient of the trajectories are assumed to be very small.

Scanning electron microscope (SEM) Instrument in which a small probe is scanned in a raster over the surface of a specimen and provokes one or several signals, which are then used to create an image on a cathode ray tube or monitor. These signals may be X-ray intensities or secondary electron or backscattered electron currents, and there are several other possibilities.

Scanning transmission electron microscope (STEM)

As in the scanning electron microscope, a small probe explores the specimen, but the specimen is thin and the signals used to generate the images are detected downstream. The resolution is comparable with that of the transmission electron microscope.

Scattering When electrons strike a solid target or pass through a thin object, they are deflected by the local field. They are said to be scattered, elastically if the change of direction is affected with negligible loss of energy, inelastically when the energy loss is appreciable.

Transmission electron microscope (TEM) Instrument closely resembling a light microscope in its general principles. A specimen area is suitably illuminated by means of condenser lenses. An objective close to the specimen provides the first stage of magnification, and intermediate and projector lens magnify the image further. Unlike glass lenses, the lens strength can be varied at will, and the total magnification can hence be varied from a few hundred times to hundreds of thousands of times. Either the object plane or the plane in which the diffraction pattern of the object is formed can be made conjugate to the image plane.

OF THE MANY PROBES used to explore the structure of matter, charged particles are among the most versatile. At high energies they are the only tools available to the nuclear physicist; at lower energies, electrons and ions are used for high-resolution microscopy and many related tasks in the physical and life sciences. The behavior of the associated instruments can often be accurately described in the language of optics. When the wavelength associated with the particles is unimportant, geometric optics are applicable and the geometric optical properties of the principal optical components—round lenses, quadrupoles, and prisms—are therefore discussed in detail. Electron microscopes, however, are operated close to their theoretical limit of resolution, and to understand how the image is formed a knowledge of wave optics is essential. The theory is presented and applied to the two families of high-resolution instruments.

I. INTRODUCTION

Charged particles in motion are deflected by electric and magnetic fields, and their behavior is described either by the Lorentz equation, which is Newton's equation of motion modified to include any relativistic effects, or by Schrödinger's equation when spin is negligible. There

are many devices in which charged particles travel in a restricted zone in the neighborhood of a curve, or axis, which is frequently a straight line, and in the vast majority of these devices, the electric or magnetic fields exhibit some very simple symmetry. It is then possible to describe the deviations of the particle motion by the fields in the familiar language of optics. If the fields are rotationally symmetric about an axis, for example, their effects are closely analogous to those of round glass lenses on light rays. Focusing can be described by cardinal elements, and the associated defects resemble the geometric and chromatic aberrations of the lenses used in light microscopes, telescopes, and other optical instruments. If the fields are not rotationally symmetric but possess planes of symmetry or antisymmetry that intersect along the optic axis, they have an analog in toric lenses, for example the glass lenses in spectacles that correct astigmatism. The other important field configuration is the analog of the glass prism; here the axis is no longer straight but a plane curve, typically a circle, and such fields separate particles of different energy or wavelength just as glass prisms redistribute white light into a spectrum.

In these remarks, we have been regarding charged particles as classical particles, obeying Newton's laws. The mention of wavelength reminds us that their behavior is also governed by Schrödinger's equation, and the resulting description of the propagation of particle beams is needed to discuss the resolution of electron-optical instruments, notably electron microscopes, and indeed any physical effect involving charged particles in which the wavelength is not negligible.

Charged-particle optics is still a young subject. The first experiments on electron diffraction were made in the 1920s, shortly after Louis de Broglie associated the notion of wavelength with particles, and in the same decade Hans Busch showed that the effect of a rotationally symmetric magnetic field acting on a beam of electrons traveling close to the symmetry axis could be described in optical terms. The first approximate formula for the focal length was given by Busch in 1926–1927. The fundamental equations and formulas of the subject were derived during the 1930s, with Walter Glaser and Otto Scherzer contributing many original ideas, and by the end of the decade the German Siemens Company had put the first commercial electron microscope with magnetic lenses on the market. The latter was a direct descendant of the prototypes built by Max Knoll, Ernst Ruska, and Bodo von Borries from 1932 onwards. Comparable work on the development of an electrostatic instrument was being done by the AEG Company.

Subsequently, several commercial ventures were launched, and French, British, Dutch, Japanese, Swiss,

American, Czechoslovakian, and Russian electron microscopes appeared on the market as well as the German instruments. These are not the only devices that depend on charged-particle optics, however. Particle accelerators also use electric and magnetic fields to guide the particles being accelerated, but in many cases these fields are not static but dynamic; frequently the current density in the particle beam is very high. Although the traditional optical concepts need not be completely abandoned, they do not provide an adequate representation of all the properties of “heavy” beams, that is, beams in which the current density is so high that interactions between individual particles are important. The use of very high frequencies likewise requires different methods and a new vocabulary that, although known as “dynamic electron optics,” is far removed from the optics of lenses and prisms. This account is confined to the charged-particle optics of static fields or fields that vary so slowly that the static equations can be employed with negligible error (scanning devices); it is likewise restricted to beams in which the current density is so low that interactions between individual particles can be neglected, except in a few local regions (the crossover of electron guns).

New devices that exploit charged-particle optics are constantly being added to the family that began with the transmission electron microscope of Knoll and Ruska. Thus, in 1965, the Cambridge Instrument Co. launched the first commercial scanning electron microscope after many years of development under Charles Oatley in the Cambridge University Engineering Department. Here, the image is formed by generating a signal at the specimen by scanning a small electron probe over the latter in a regular pattern and using this signal to modulate the intensity of a cathode-ray tube. Shortly afterward, Albert Crewe of the Argonne National Laboratory and the University of Chicago developed the first scanning transmission electron microscope, which combines all the attractions of a scanning device with the very high resolution of a “conventional” electron microscope. More recently still, fine electron beams have been used for microlithography, for in the quest for microminiaturization of circuits, the wavelength of light set a lower limit on the dimensions attainable. Finally, there are, many devices in which the charged particles are ions of one or many species. Some of these operate on essentially the same principles as their electron counterparts; in others, such as mass spectrometers, the presence of several ion species is intrinsic. The laws that govern the motion of all charged particles are essentially the same, however, and we shall consider mainly electron optics; the equations are applicable to any charged particle, provided that the appropriate mass and charge are inserted.

II. GEOMETRIC OPTICS

A. Paraxial Equations

Although it is, strictly speaking, true that any beam of charged particles that remains in the vicinity of an arbitrary curve in space can be described in optical language, this is far too general a starting point for our present purposes. Even for light, the optics of systems in which the axis is a skew curve in space, developed for the study of the eye by Allvar Gullstrand and pursued by Constantin Carathéodory, are little known and rarely used. The same is true of the corresponding theory for particles, developed by G. A. Grinberg and Peter Sturrock. We shall instead consider the other extreme case, in which the axis is straight and any magnetic and electrostatic fields are rotationally symmetric about this axis.

1. Round Lenses

We introduce a Cartesian coordinate system in which the z axis coincides with the symmetry axis, and we provisionally denote the transverse axes X and Y . The motion of a charged particle of rest mass m_0 and charge Q in an electrostatic field \mathbf{E} and a magnetic field \mathbf{B} is then determined by the differential equation

$$(d/dt)(\gamma m_0 \mathbf{v}) = Q(\mathbf{E} + \mathbf{v} \times \mathbf{B})$$

$$\gamma = (1 - v^2/c^2)^{-1/2}, \quad (1)$$

which represents Newton's second law modified for relativistic effects (Lorentz equation); \mathbf{v} is the velocity. For electrons, we have $e = -Q \simeq 1.6 \times 10^{-19}$ C and $e/m_0 \simeq 176$ C/ μ g. Since we are concerned with static fields, the time of arrival of the particles is often of no interest, and it is then preferable to differentiate not with respect to time but with respect to the axial coordinate z . A fairly lengthy calculation yields the trajectory equations

$$\begin{aligned} \frac{d^2 X}{dz^2} &= \frac{\rho^2}{g} \left(\frac{\partial g}{\partial X} - X' \frac{\partial g}{\partial z} \right) \\ &\quad + \frac{Q\rho}{g} [Y'(B_z + X'B_X) - B_Y(1 + X'^2)] \\ \frac{d^2 Y}{dz^2} &= \frac{\rho^2}{g} \left(\frac{\partial g}{\partial Y} - Y' \frac{\partial g}{\partial z} \right) \\ &\quad + \frac{Q\rho}{g} [-X'(B_z + Y'B_Y) + B_X(1 + Y'^2)] \end{aligned} \quad (2)$$

in which $\rho^2 = 1 + X'^2 + Y'^2$ and $g = \gamma m_0 v$.

By specializing these equations to the various cases of interest, we obtain equations from which the optical properties can be derived by the “trajectory method.” It is well

known that equations such as Eq. (1) are identical with the Euler–Lagrange equations of a variational principle of the form

$$W = \int_{t_0}^{t_1} L(\mathbf{r}, \mathbf{v}, t) dt = \text{extremum} \quad (3)$$

provided that t_0 , t_1 , $\mathbf{r}(t_0)$, and $\mathbf{r}(t_1)$ are held constant. The Lagrangian L has the form

$$L = m_0 c^2 [1 - (1 - v^2/c^2)^{1/2}] + Q(\mathbf{v} \cdot \mathbf{A} - \Phi) \quad (4)$$

in which Φ and \mathbf{A} are the scalar and vector potentials corresponding to \mathbf{E} , $\mathbf{E} = -\text{grad } \Phi$ and to \mathbf{B} , $\mathbf{B} = \text{curl } \mathbf{A}$. For static systems with a straight axis, we can rewrite Eq. (3) in the form

$$S = \int_{z_0}^{z_1} M(x, y, z, x', y') dz, \quad (5)$$

where

$$M = (1 + X'^2 + Y'^2)^{1/2} g(\mathbf{r}) + Q(X'A_X + Y'A_Y + A_z). \quad (6)$$

The Euler–Lagrange equations,

$$\frac{d}{dz} \left(\frac{\partial M}{\partial X'} \right) = \frac{\partial M}{\partial X}, \quad \frac{d}{dz} \left(\frac{\partial M}{\partial Y'} \right) = \frac{\partial M}{\partial Y} \quad (7)$$

again define trajectory equations. A very powerful method of analyzing optical properties is based on a study of the function M and its integral S ; this is known as the method of characteristic functions, or eikonal method.

We now consider the special case of rotationally symmetric systems in the paraxial approximation; that is, we examine the behavior of charged particles, specifically electrons, that remain very close to the axis. For such particles, the trajectory equations collapse to a simpler form, namely,

$$\begin{aligned} X'' + \frac{\gamma \phi'}{2\hat{\phi}} X' + \frac{\gamma \phi''}{4\hat{\phi}} X + \frac{\eta B}{\hat{\phi}^{1/2}} Y' + \frac{\eta B'}{2\hat{\phi}^{1/2}} Y &= 0 \\ Y'' + \frac{\gamma \phi'}{2\hat{\phi}} Y' + \frac{\gamma \phi''}{4\hat{\phi}} Y - \frac{\eta B}{\hat{\phi}^{1/2}} X' - \frac{\eta B'}{2\hat{\phi}^{1/2}} X &= 0 \end{aligned} \quad (8)$$

in which $\phi(z)$ denotes the distribution of electrostatic potential on the optic axis, $\phi(z) = \Phi(0, 0, z)$; $\hat{\phi}(z) = \phi(z)[1 + e\phi(z)/2m_0 c^2]$. Likewise, $B(z)$ denotes the magnetic field distribution on the axis. These equations are coupled, in the sense that X and Y occur in both, but this can be remedied by introducing new coordinate axes x , y , inclined to X and Y at an angle $\theta(z)$ that varies with z ; $x = 0$, $y = 0$ will therefore define not planes but surfaces. By choosing $\theta(z)$ such that

$$d\theta/dz = \eta B/2\hat{\phi}^{1/2}; \quad \eta = (e/2m_0)^{1/2}, \quad (9)$$

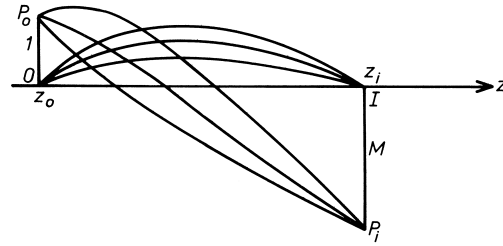


FIGURE 1 Paraxial solutions demonstrating image formation.

we find

$$\begin{aligned} x'' + \gamma \phi' x' / 2\hat{\phi} + [(\gamma \phi'' + \eta^2 B^2) / 4\hat{\phi}] / x &= 0 \\ y'' + \gamma \phi' y' / 2\hat{\phi} + [(\gamma \phi'' + \eta^2 B^2) / 4\hat{\phi}] / y &= 0. \end{aligned} \quad (10)$$

These differential equations are linear, homogeneous, and second order. The general solution of either is a linear combination of any two linearly independent solutions, and this fact is alone sufficient to show that the corresponding fields $B(z)$ and potentials $\phi(z)$ have an imaging action, as we now show. Consider the particular solution $h(z)$ of Eq. (10) that intersects the axis at $z = z_0$ and $z = z_i$ (Fig. 1). A pencil of rays that intersects the plane $z = z_0$ at some point $P_0(x_0, y_0)$ can be described by

$$x(z) = x_0 g(z) + \lambda h(z) \quad (11)$$

$$y(z) = y_0 g(z) + \mu h(z)$$

in which $g(z)$ is any solution of Eq. (10) that is linearly independent of $h(z)$ such that $g(z_0) = 1$ and λ, μ are parameters; each member of the pencil corresponds to a different pair of values of λ, μ . In the plane $z = z_i$, we find

$$x(z_i) = x_0 g(z_i); \quad y(z_i) = y_0 g(z_i) \quad (12)$$

for all λ and μ and hence for all rays passing through P_0 . This is true for every point in the plane $z = z_0$, and hence the latter will be stigmatically imaged in $z = z_i$.

Furthermore, both ratios $x(z_i)/x_0$ and $y(z_i)/y_0$ are equal to the constant $g(z_i)$, which means that any pattern of points in $z = z_0$ will be reproduced faithfully in the image plane, magnified by this factor $g(z_i)$, which is hence known as the (transverse) magnification and denoted by M .

The form of the paraxial equations has numerous other consequences. We have seen that the coordinate frame $x-y-z$ rotates relative to the fixed frame $X-Y-Z$ about the optic axis, with the result that the image will be rotated with respect to the object if magnetic fields are used. In an instrument such as an electron microscope, the image therefore rotates as the magnification is altered, since the latter is affected by altering the strength of the magnetic field and Eq. (9) shows that the angle of rotation is a function of this quantity. Even more important is the fact that the coefficient of the linear term is strictly positive in the

case of magnetic fields. This implies that the curvature of any solution $x(z)$ is opposite in sign to $x'(z)$, with the result that the field always drives the electrons toward the axis; magnetic electron lenses always have a convergent action. The same is true of the overall effect of electrostatic lenses, although the reasoning is not quite so simple.

A particular combination of any two linearly independent solutions of Eq. (10) forms the invariant known as the Wronskian. This quantity is defined by

$$\hat{\phi}^{1/2}(x_1 x'_2 - x'_1 x_2); \quad \hat{\phi}^{1/2}(y_1 y'_2 - y'_1 y_2) \quad (13)$$

Suppose that we select $x_1 = h$ and $x_2 = g$, where $h(z_o) = h(z_i) = 0$ and $g(z_o) = 1$ so that $g(z_i) = M$. Then

$$\hat{\phi}_o^{1/2} h'_o = \hat{\phi}_i^{1/2} h'_i M \quad (14)$$

The ratio h'_i/h'_o is the angular magnification M_A and so

$$MM_A = (\hat{\phi}_o/\hat{\phi}_i)^{1/2} \quad (15)$$

or $MM_A = 1$ if the lens has no overall accelerating effect and hence $\hat{\phi}_o = \hat{\phi}_i$. Identifying $\hat{\phi}^{1/2}$ with the refractive index, Eq. (15) is the particle analog of the Smith–Helmholtz formula of light optics. Analogs of all the other optical laws can be established; in particular, we find that the longitudinal magnification M_l is given by.

$$M_l = M/M_A = (\hat{\phi}_i/\hat{\phi}_o)^{1/2} M^2 \quad (16)$$

and that Abbe's sine condition and Herschel's condition take their familiar forms.

We now show that image formation by electron lenses can be characterized with the aid of cardinal elements: foci, focal lengths, and principal planes. First, however, we must explain the novel notions of real and asymptotic imaging. So far, we have simply spoken of rotationally symmetric fields without specifying their distribution in space. Electron lenses are localized regions in which the magnetic or electrostatic field is strong and outside of which the field is weak but, in theory at least, does not vanish. Some typical lens geometries are shown in Fig. 2.

If the object and image are far from the lens, in effectively field-free space, or if the object is not a physical specimen but an intermediate image of the latter, the image formation can be analyzed in terms of the asymptotes to rays entering or emerging from the lens region. If, however, the true object or image is immersed within the lens field, a different method of characterizing the lens properties must be adopted, and we shall speak of real cardinal elements. We consider the asymptotic case first.

It is convenient to introduce the solutions of Eq. (10) that satisfy the boundary conditions

$$\lim_{z \rightarrow -\infty} G(z) = 1; \quad \lim_{z \rightarrow \infty} \bar{G}(z) = 1 \quad (17)$$

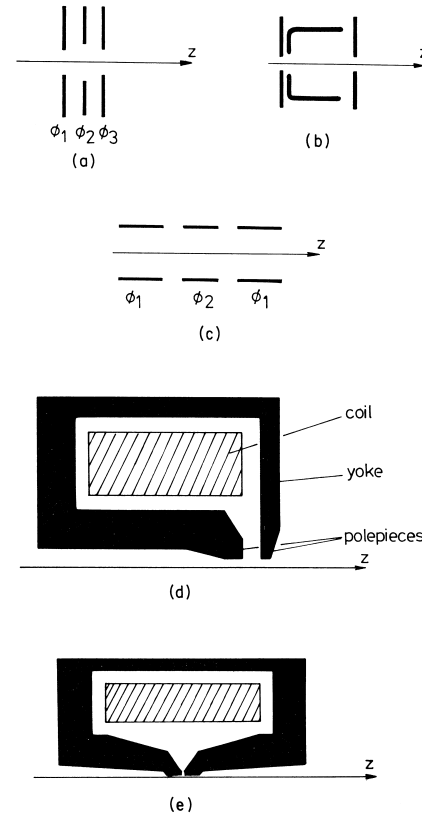


FIGURE 2 Typical electron lenses: (a–c) electrostatic lenses, of which (c) is an einzel lens; (d–e) magnetic lenses of traditional design.

These are rays that arrive at or leave the lens parallel to the axis (Fig. 3). As usual, the general solution is $x(z) = \alpha G(z) + \beta \bar{G}(z)$, where α and β are constants. We denote the emergent asymptote to $G(z)$ thus:

$$\lim_{z \rightarrow \infty} G(z) = G_i(z - z_{Fi}) \quad (18)$$

We denote the incident asymptote to $\bar{G}(z)$ thus:

$$\lim_{z \rightarrow -\infty} \bar{G}(z) = \bar{G}_o(z - z_{Fo}) \quad (19)$$

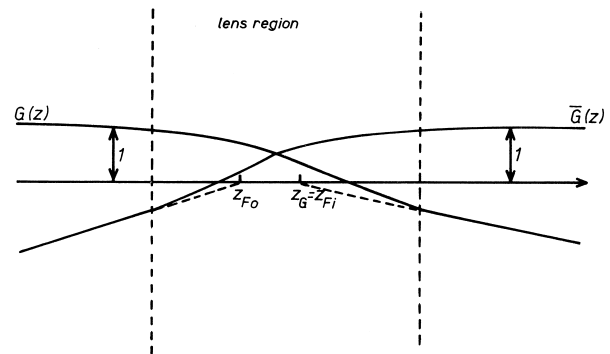


FIGURE 3 Rays $G(z)$ and $\bar{G}(z)$.

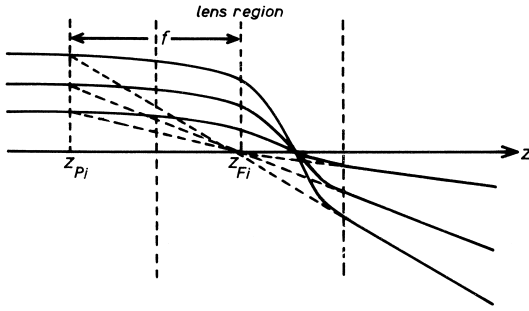


FIGURE 4 Focal and principal planes.

Clearly, all rays incident parallel to the axis have emergent asymptotes that intersect at $z = z_{Fi}$; this point is known as the asymptotic image focus. It is not difficult to show that the emergent asymptotes to any family of rays that are parallel to one another but not to the axis intersect at a point in the plane $z = z_{Fi}$. By applying a similar reasoning to $\bar{G}(z)$, we recognize that z_{Fo} is the asymptotic object focus. The incident and emergent asymptotes to $G(z)$ intersect in a plane z_{Pi} , which is known as the image principal plane (Fig. 4). The distance between z_{Fi} and z_{Pi} is the asymptotic image focal length:

$$z_{Fi} - z_{Pi} = -1/G'_i = f_i \quad (20)$$

We can likewise define z_{Po} and f_o :

$$z_{Po} - z_{Fo} = 1/\bar{G}'_o = f_o \quad (21)$$

The Wronskian tells us that $\hat{\phi}^{1/2}(G\bar{G}' - G'\bar{G})$ is constant and so

$$\hat{\phi}_o^{1/2}\bar{G}'_o = -\hat{\phi}_i^{1/2}G'_i$$

or

$$f_o/\hat{\phi}_o^{1/2} = f_i/\hat{\phi}_i^{1/2} \quad (22)$$

In magnetic lenses and electrostatic lenses, that provide no overall acceleration, $\hat{\phi}_o = \hat{\phi}_i$ and so $f_o = f_i$; we drop the subscript when no confusion can arise.

The coupling between an object space and an image space is conveniently expressed in terms of z_{Fo} , z_{Fi} , f_o , and f_i . From the general solution $x = \alpha G + \beta \bar{G}$, we see that

$$\begin{aligned} \lim_{z \rightarrow -\infty} x(z) &= \alpha + \beta(z - z_{Fo})/f_o \\ \lim_{z \rightarrow \infty} x(z) &= -\alpha(z - z_{Fi})/f_i + \beta \end{aligned} \quad (23)$$

and likewise for $y(z)$. Eliminating α and β , we find

$$\begin{bmatrix} x_2 \\ x'_2 \end{bmatrix} = \begin{bmatrix} -\frac{z_2 - z_{Fi}}{f_i} & f_o + \frac{(z_1 - z_{Fo})(z_2 - z_{Fi})}{f_i} \\ -\frac{1}{f_i} & \frac{z_o - z_{Fo}}{f_o} \end{bmatrix} \begin{bmatrix} x_1 \\ x'_1 \end{bmatrix} \quad (24)$$

where x_1 denotes $x(z)$ in some plane $z = z_1$ on the incident asymptote and x_2 denotes $x(z)$ in some plane $z = z_2$ on

the emergent asymptote; $x' = dx/dz$. The matrix that appears in this equation is widely used to study systems with many focusing elements; it is known as the (paraxial) transfer matrix and takes slightly different forms for the various elements in use, quadrupoles in particular. We denote the transfer matrix by T .

If the planes z_1 and z_2 are conjugate, the point of arrival of a ray in z_2 will vary with the position coordinates of its point of departure in z_1 but will be independent of the gradient at that point. The transfer matrix element T_{12} must therefore vanish,

$$(z_o - z_{Fo})(z_i - z_{Fi}) = -f_o f_i \quad (25)$$

in which we have replaced z_1 and z_2 by z_o and z_i to indicate that these are now conjugates (object and image). This is the familiar lens equation in Newtonian form. Writing $z_{Fi} = z_{Pi} + f_i$ and $z_{Fo} = z_{Po} - f_o$, we obtain

$$\frac{f_o}{z_{Po} - z_o} + \frac{f_i}{z_i - z_{Pi}} = 1 \quad (26)$$

the thick-lens form of the regular lens equation.

Between conjugates, the matrix T takes the form

$$T = \begin{bmatrix} M & 0 \\ -\frac{1}{f_i} & \frac{f_o}{f_i} \frac{1}{M} \end{bmatrix} \quad (27)$$

in which M denotes the asymptotic magnification, the height of the image asymptote to $G(z)$ in the image plane.

If, however, the object is a real physical specimen and not a mere intermediate image, the asymptotic cardinal elements cannot in general be used, because the object may well be situated inside the field region and only a part of the field will then contribute to the image formation. Fortunately, objective lenses, in which this situation arises, are normally operated at high magnification with the specimen close to the real object focus, the point at which the ray $\bar{G}(z)$ itself intersects the axis [whereas the asymptotic object focus is the point at which the asymptote to $\bar{G}(z)$ in object space intersects the optic axis]. The corresponding real focal length is then defined by the slope of $\bar{G}(z)$ at the object focus F_o : $f = 1/G'(F_o)$; see Fig. 5.

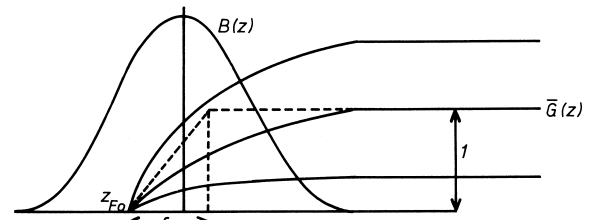


FIGURE 5 Real focus and focal length.

2. Quadrupoles

In the foregoing discussion, we have considered only rotationally symmetric fields and have needed only the axial distributions $B(z)$ and $\phi(z)$. The other symmetry of greatest practical interest is that associated with electrostatic and magnetic quadrupoles, widely used in particle accelerators. Here, the symmetry is lower, the fields possessing planes of symmetry and antisymmetry only; these planes intersect in the optic axis, and we shall assume forthwith that electrostatic and magnetic quadrupoles are disposed as shown in Fig. 6. The reason for this is simple: The paraxial equations of motion for charged particles traveling through quadrupoles separate into two uncoupled equations only if this choice is adopted. This is not merely a question of mathematical convenience; if quadrupole fields overlap and the total system does not have the symmetry indicated, the desired imaging will not be achieved.

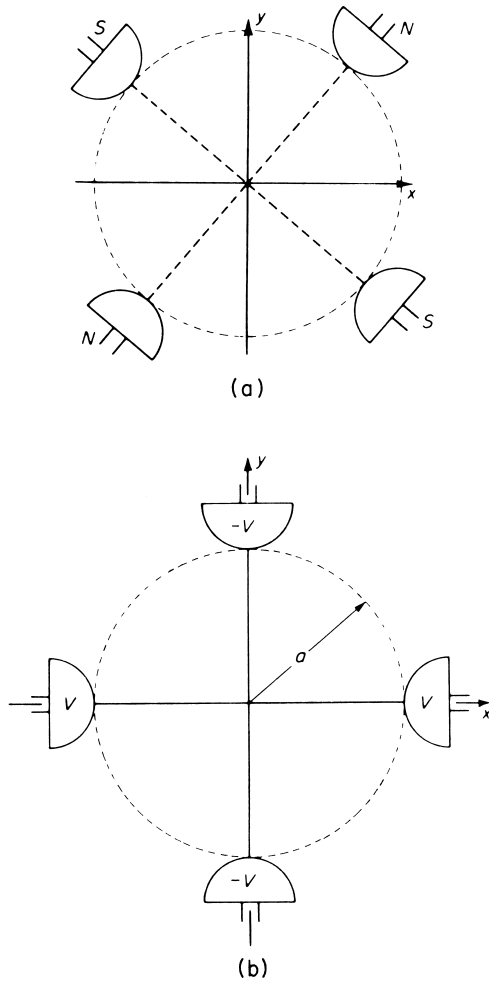


FIGURE 6 (a) Magnetic and (b) electrostatic quadrupoles.

The paraxial equations are now different in the x - z and y - z planes:

$$\begin{aligned} \frac{d}{dz}(\hat{\phi}^{1/2}x') + \frac{\gamma\phi'' - 2\gamma p_2 + 4\eta Q_2\hat{\phi}^{1/2}}{4\hat{\phi}^{1/2}}x &= 0 \\ \frac{d}{dz}(\hat{\phi}^{1/2}y') + \frac{\gamma\phi'' + 2\gamma p_2 - 4\eta Q_2\hat{\phi}^{1/2}}{4\hat{\phi}^{1/2}}y &= 0 \end{aligned} \quad (28)$$

in which we have retained the possible presence of a round electrostatic lens field $\phi(z)$. The functions $p_2(z)$ and $Q_2(z)$ that also appear characterize the quadrupole fields; their meaning is easily seen from the field expansions [for $B(z) = 0$]:

$$\begin{aligned} \Phi(x, y, z) &= \phi(z) - \frac{1}{4}(x^2 + y^2)\phi''(z) \\ &\quad + \frac{1}{64}(x^2 + y^2)^2\phi^{(4)}(z) \\ &\quad + \frac{1}{2}(x^2 - y^2)p_2(z) - \frac{1}{24}(x^4 + y^4)p_2''(z) \\ &\quad + \frac{1}{24}p_4(z)(x^4 - 6x^2y^2 + y^4) + \dots \end{aligned} \quad (29)$$

$$\begin{aligned} \Phi(r, \psi, z) &= \phi(z) - \frac{1}{4}r^2\phi'' + \frac{1}{64}r^4\phi^{(4)} \\ &\quad + \frac{1}{2}p_2r^2\cos 2\psi - \frac{1}{24}p_2''r^4\cos 2\psi \\ &\quad + \frac{1}{24}p_4r^4\cos 4\psi + \dots \end{aligned}$$

$$\begin{aligned} A_x &= -\frac{x}{12}(x^2 - 3y^2)Q_2'(z) \\ A_y &= \frac{y}{12}(y^2 - 3x^2)Q_2'(z) \\ A_z &= \frac{1}{2}(x^2 - y^2)Q_2(z) - \frac{1}{24}(x^4 - y^4)Q_2''(z) \\ &\quad + \frac{1}{24}(x^4 - 6x^2y^2 + y^4)Q_4(z) \end{aligned} \quad (30)$$

The terms $p_4(z)$ and $Q_4(z)$ characterize octopole fields, and we shall refer to them briefly in connection with the aberration correction below.

It is now necessary to define separate transfer matrices for the x - z plane and for the y - z plane. These have exactly the same form as Eqs. (24) and (27), but we have to distinguish between two sets of cardinal elements. For arbitrary planes z_1 and z_2 , we have

$$T^{(x)} = \begin{bmatrix} -\frac{z_2 - z_{Fi}^{(x)}}{f_{xi}} & \frac{(z_2 - z_{Fi}^{(x)})(z_2 - z_{Fo}^{(x)})}{f_{xo}} + f_{xi} \\ -\frac{1}{f_{xi}} & \frac{z_1 - z_{Fo}^{(x)}}{f_{xi}} \end{bmatrix}$$

$$T^{(y)} = \begin{bmatrix} -\frac{z_2 - z_{Fi}^{(y)}}{f_{yi}} & \frac{(z_2 - z_{Fi}^{(y)})(z_1 - z_{Fo}^{(y)})}{f_{yo}} + f_{yi} \\ -\frac{1}{f_{yi}} & \frac{z_1 - z_{Fo}^{(y)}}{f_{yi}} \end{bmatrix}. \quad (31)$$

Suppose now that $z = z_{x0}$ and $z = z_{xi}$ and conjugate so that $T_{12}^{(x)} = 0$; in general, $T_{12}^{(y)} \neq 0$ and so a point in the object plane $z = z_{x0}$ will be imaged as a line parallel to the y axis. Similarly, if we consider a pair of conjugates $z = z_{y0}$ and $z = z_{yi}$, we obtain a line parallel to the x axis. The imaging is hence astigmatic, and the astigmatic differences in object and image space can be related to the magnification

$$\begin{aligned} \wedge_i &:= z_{xi} - z_{yi} = \wedge_{Fi} - f_{xi}M_x + f_{yi}M_y \\ \wedge_i &:= z_{x0} - z_{y0} = \wedge_{Fo} + f_{xo}/M_x - f_{yo}/M_y, \end{aligned} \quad (32)$$

where

$$\begin{aligned} \wedge_{Fi} &:= z_{Fi}^{(x)} - z_{Fi}^{(y)} = \wedge_i(M_x = M_y = 0) \\ \wedge_{Fo} &:= z_{Fo}^{(x)} - z_{Fo}^{(y)} = \wedge_o(M_x = M_y \rightarrow \infty). \end{aligned} \quad (33)$$

Solving the equations $\wedge_i = \wedge_o = 0$ for M_x and M_y , we find that there is a pair of object planes for which the image is stigmatic though not free of distortion.

3. Prisms

There is an important class of devices in which the optic axis is not straight but a simple curve, almost invariably lying in a plane. The particles remain in the vicinity of this curve, but they experience different focusing forces in the plane and perpendicular to it. In many cases, the axis is a circular arc terminated by straight lines. We consider the situation in which charged particles travel through a magnetic sector field (Fig. 7); for simplicity, we assume that the field falls abruptly to zero at entrance and exit planes (rather than curved surfaces) and that the latter are normal to the optic axis, which is circular. We regard the plane containing the axis as horizontal. The vertical field at the axis is denoted by B_o , and off the axis, $B = B_o(r/R)^{-n}$ in the horizontal plane. It can then be shown, with the notation of Fig. 7, that paraxial trajectory equations of the form

$$x'' + k_v^2 x = 0; \quad y'' + k_H^2 y = 0 \quad (34)$$

describe the particle motion, with $k_H^2 = (1 - n)/R^2$ and $k_v^2 = n/R^2$. Since these are identical in appearance with

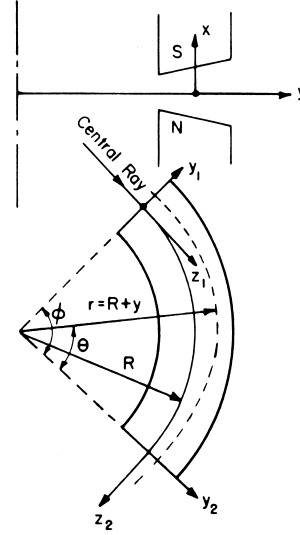


FIGURE 7 Passage through a sector magnet.

the quadrupole equations but do not have different signs, the particles will be focused in both directions but not in the same “image” plane unless $k_H = k_v$ and hence $n = \frac{1}{2}$. The cases $n = 0$, for which the magnetic field is homogeneous, and $n = \frac{1}{2}$ have been extensively studied. Since prisms are widely used to separate particles of different energy or momentum, the dispersion is an important quantity, and the transfer matrices are usually extended to include this information. In practice, more complex end faces are employed than the simple planes normal to the axis considered here, and the fringing fields cannot be completely neglected, as they are in the sharp cutoff approximation.

Electrostatic prisms can be analyzed in a similar way and will not be discussed separately.

B. Aberrations

1. Traditional Method

The paraxial approximation describes the dominant focusing in the principal electron-optical devices, but this is inevitably perturbed by higher order effects, or aberrations. There are several kinds of aberrations. By retaining higher order terms in the field or potential expansions, we obtain the family of geometric aberrations. By considering small changes in particle energy and lens strength, we obtain the chromatic aberrations. Finally, by examining the effect of small departures from the assumed symmetry of the field, we obtain the parasitic aberrations.

All these types of aberrations are conveniently studied by means of perturbation theory. Suppose that we have obtained the paraxial equations as the Euler–Lagrange equations of the paraxial form of M [Eq. (6)], which we denote

$M^{(P)}$. Apart from a trivial change of scale, we have

$$M^{(P)} = -(1/8\hat{\phi}^{1/2})(\gamma\phi'' + \eta^2 B^2)(x^2 + y^2) + \frac{1}{2}\hat{\phi}^{1/2}(x'^2 + y'^2) \quad (35)$$

Suppose now that $M^{(P)}$ is perturbed to $M^{(P)} + M^{(A)}$. The second term $M^{(A)}$ may represent additional terms, neglected in the paraxial approximation, and will then enable us to calculate the geometric aberrations; alternatively, $M^{(A)}$ may measure the change in $M^{(P)}$ when particle energy and lens strength fluctuate, in which case it tells us the chromatic aberration. Other field terms yield the parasitic aberration. We illustrate the use of perturbation theory by considering the geometric aberrations of round lenses. Here, we have

$$\begin{aligned} M^{(A)} = M^{(4)} = & -\frac{1}{4}L_1(x^2 + y^2)^2 \\ & -\frac{1}{2}L_2(x^2 + y^2)(x'^2 + y'^2) \\ & -\frac{1}{4}L_3(x'^2 + y'^2)^2 \\ & -R(xy' - x'y)^2 \\ & -P\hat{\phi}^{1/2}(x^2 + y^2)(xy' - x'y) \\ & -Q\hat{\phi}^{1/2}(x'^2 + y'^2)(xy' - x'y) \end{aligned} \quad (36)$$

with

$$\begin{aligned} L_1 = & \frac{1}{32\hat{\phi}^{1/2}} \left(\frac{\phi'^2}{\hat{\phi}} - \gamma\phi^{(4)} + \frac{2\gamma\phi''\eta^2 B^2}{\hat{\phi}} \right. \\ & \left. + \frac{\eta^4 B^4}{\hat{\phi}} - 4\eta^2 B B'' \right) \\ L_2 = & \frac{1}{8\hat{\phi}^{1/2}} (\gamma\phi'' + \eta^2 B^2) \\ L_3 = & \frac{1}{2}\hat{\phi}^{1/2}; \quad P = \frac{\eta}{16\hat{\phi}^{1/2}} \left(\frac{\gamma\phi'' B}{\hat{\phi}} + \frac{\eta^2 B^2}{\hat{\phi}} - B'' \right) \\ Q = & \frac{\eta B}{4\hat{\phi}^{1/2}}; \quad R = \frac{\eta^2 B^2}{8\hat{\phi}^{1/2}} \end{aligned} \quad (37)$$

and with $S^{(A)} = \int_{z_0}^z M^{(A)} dz$, we can show that

$$\begin{aligned} \frac{\partial S^{(A)}}{\partial x_a} &= p_x^{(A)} t(z) - x^{(A)} \hat{\phi}^{1/2} t'(z) \\ \frac{\partial S^{(A)}}{\partial y_a} &= p_y^{(A)} s(z) - x^{(A)} \hat{\phi}^{1/2} s'(z) \end{aligned} \quad (38)$$

where $s(z)$ and $t(z)$ are the solutions of Eq. (10) for which $s(z_0) = t(z_a) = 1$, $s(z_a) = t(z_0) = 0$, and $z = z_a$ denotes some aperture plane. Thus, in the image plane,

$$x^{(A)} = -(M/W) \partial S_{oi}^{(A)} / \partial x_a \quad (39)$$

where $S_{oi}^{(A)}$ denotes $\int_{z_0}^{z_i} M^{(A)} dz$, with a similar expression for $y^{(A)}$. The quantities with superscript (A) indicate the departure from the paraxial approximation, and we write

$$\begin{aligned} \Delta x_i &= x^{(A)} / M = -(1/W) \partial S_{oi}^{(A)} / \partial x_a \\ \Delta y_i &= y^{(A)} / M = -(1/W) \partial S_{oi}^{(A)} / \partial y_a \end{aligned} \quad (40)$$

The remainder of the calculation is lengthy but straightforward. Into $M^{(4)}$, the paraxial solutions are substituted and the resulting terms are grouped according to their dependence on x_o , y_o , x_a , and y_a . We find that $S^{(A)}$ can be written

$$\begin{aligned} -S^{(A)} / W = & \frac{1}{4} E r_o^4 + \frac{1}{4} C r_a^4 + \frac{1}{2} A (V^2 - v^2) \\ & + \frac{1}{2} F r_o^2 r_a^2 + D r_o^2 V + K r_a^2 V \\ & + v (d r_o^2 + k r_a^2 + a V) \end{aligned} \quad (41)$$

with

$$\begin{aligned} r_o^2 &= x_o^2 + y_o^2; & r_a^2 &= x_a^2 + y_a^2 \\ V &= x_o x_a + y_o y_a; & v &= x_o y_a - x_a y_o \end{aligned} \quad (42)$$

and

$$\begin{aligned} \Delta x_i &= x_a [C r_a^2 + 2K V + 2k v + (F - A) r_o^2] \\ &+ x_o (K r_a^2 + 2A V + a v + D r_o^2) \\ &- y_o (k r_a^2 + a V + d r_o^2) \\ \Delta y_i &= y_a [C r_a^2 + 2K V + 2k v + (F - A) r_o^2] \\ &+ x_o (k r_a^2 + a V + d r_o^2) \end{aligned} \quad (43)$$

Each coefficient A, C, \dots, d, k represents a different type of geometric aberration. Although all lenses suffer from every aberration, with the exception of the anisotropic aberrations described by k, a , and d , which are peculiar to magnetic lenses, the various aberrations are of very unequal importance when lenses are used for different purposes. In microscope objectives, for example, the incident electrons are scattered within the specimen and emerge at relatively steep angles to the optic axis (several milliradians or tens of milliradians). Here, it is the spherical (or aperture) aberration C that dominates, and since this aberration does not vanish on the optic axis, being independent of r_o , it has an extremely important effect on image quality. Of the geometric aberrations, it is this spherical aberration that determines the resolving power of the electron microscope. In the subsequent lenses of such instruments, the image is progressively enlarged until the final magnification, which may reach $100,000\times$ or $1,000,000\times$, is attained. Since angular magnification is inversely proportional to transverse magnification, the

angular spread of the beam in these projector lenses will be tiny, whereas the off-axis distance becomes large. Here, therefore, the distortions D and d are dominant.

A characteristic aberration figure is associated with each aberration. This figure is the pattern in the image plane formed by rays from some object point that cross the aperture plane around a circle. For the spherical aberration, this figure is itself a circle, irrespective of the object position, and the effect of this aberration is therefore to blur the image uniformly, each Gaussian image point being replaced by a disk of radius MCr_a^3 . The next most important aberration for objective lenses is the coma, characterized by K and k , which generates the comet-shaped streak from which it takes its name. The coefficients A and F describe Seidel astigmatism and field curvature, respectively; the astigmatism replaces stigmatic imagery by line imagery, two line foci being formed on either side of the Gaussian image plane, while the field curvature causes the image to be formed not on a plane but on a curved image surface. The distortions are more graphically understood by considering their effect on a square grid in the object plane. Such a grid is swollen or shrunk by the isotropic distortion D and warped by the anisotropic distortion d ; the latter has been evocatively styled as a pocket handkerchief distortion. Figure 8 illustrates these various aberrations.

Each aberration has a large literature, and we confine this account to the spherical aberration, an eternal pre-occupation of microscope lens designers. In practice, it is more convenient to define this in terms of angle at the specimen, and recalling that $x(z) = x_0s(z) + x_at(z)$, we see that $x'_0 = x_0s'(z_0) + x_at'(z_0)$. Hence,

$$\Delta x_i = Cx_a(x_a'^2 + y_a'^2) = \frac{C}{t_0'^3}x'_0(x_0'^2 + y_0'^2) + \dots \quad (44)$$

and we therefore write $C_s = c/t_0'^3$ so that

$$\Delta x_i = C_s x'_0(x_0'^2 + y_0'^2); \quad y_i = C_s y'_0(x_0'^2 + y_0'^2) \quad (45)$$

It is this coefficient C_s that is conventionally quoted and tabulated. A very important and disappointing property of C_s is that it is intrinsically positive: The formula for it can be cast into positive-definite form, which means that we cannot hope to design a round lens free of this aberration by skillful choice of geometry and excitation. This result is known as Scherzer's theorem. An interesting attempt to upset the theorem was made by Glaser, who tried setting the integrand that occurs in the formula for C_s , and that can be written as the sum of several squared terms, equal to zero and solving the resulting differential equation for the field (in the magnetic case). Alas, the field distribution that emerged was not suitable for image formation, thus confirming the truth of the theorem, but it has been found useful in β -ray spectroscopy. The full implications of the theorem were established by Werner Tretnner, who estab-

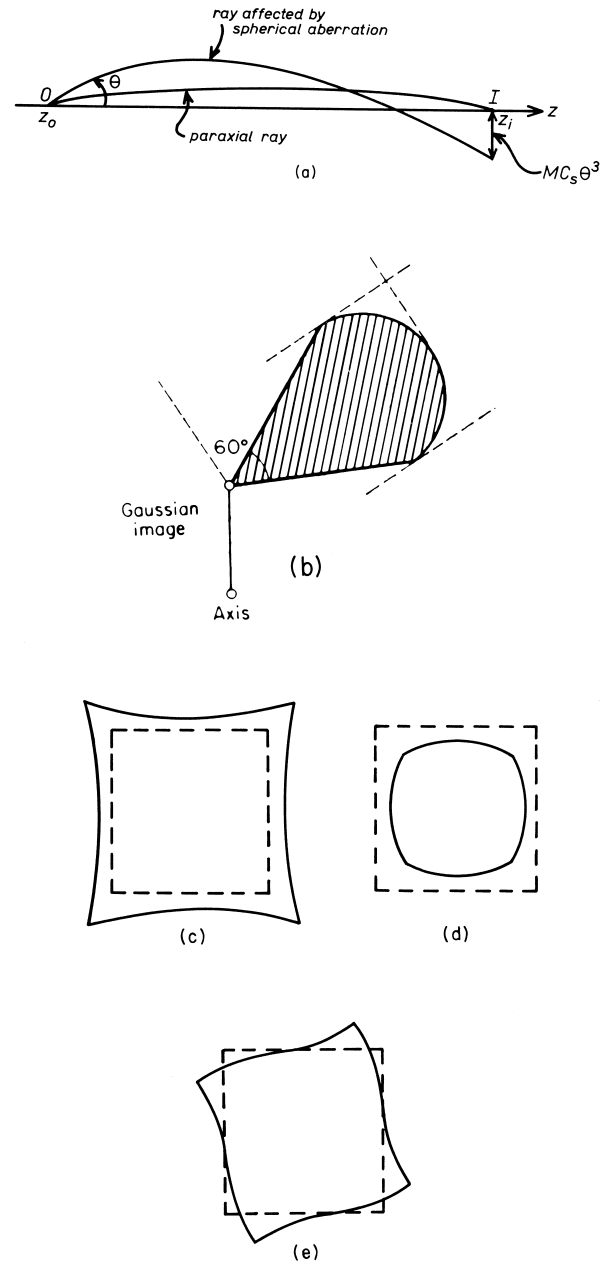


FIGURE 8 Aberration patterns: (a) spherical aberration; (b) coma; (c–e) distortions.

lished the lower limit for C_s as a function of the practical constraints imposed by electrical breakdown, magnetic saturation, and geometry.

Like the cardinal elements, the aberrations of objective lenses require a slightly different treatment from those of condenser lenses and projector lenses. The reason is easily understood: In magnetic objective lenses (and probe-forming lenses), the specimen (or target) is commonly immersed deep inside the field and only the field region downstream contributes to the image formation. The

spherical aberration is likewise generated only by this part of the field, and the expression for C_s as an integral from object plane to image plane reflects this. In other lenses, however, the object is in fact an intermediate image, formed by the next lens upstream, and the whole lens field contributes to the image formation and hence to the aberrations. It is then the coupling between incident and emergent asymptotes that is of interest, and the aberrations are characterized by asymptotic aberration coefficients. These exhibit an interesting property: They can be expressed as polynomials in reciprocal magnification m ($m = 1/M$), with the coefficients in these polynomials being determined by the lens geometry and excitation and independent of magnification (and hence of object position). This dependence can be written

$$\begin{bmatrix} C \\ K \\ A \\ F \\ D \end{bmatrix} = Q \begin{bmatrix} m^4 \\ m^3 \\ m^2 \\ m \\ 1 \end{bmatrix} \begin{bmatrix} k \\ a \\ d \end{bmatrix} = q \begin{bmatrix} m^2 \\ m \\ 1 \end{bmatrix} \quad (46)$$

in which Q and q have the patterns

$$Q = \begin{bmatrix} x & x & x & x & x \\ 0 & x & x & x & x \\ 0 & 0 & x & x & x \\ 0 & 0 & x & x & x \\ 0 & 0 & 0 & x & x \end{bmatrix}; \quad q = \begin{bmatrix} x & x & x \\ 0 & x & x \\ 0 & 0 & x \end{bmatrix}, \quad (47)$$

where an x indicates that the matrix element is a nonzero quantity determined by the lens geometry and excitation.

Turning now to chromatic aberrations, we have

$$m^{(p)} = \frac{\partial m^{(2)}}{\partial \phi} \Delta \phi + \frac{\partial m^{(2)}}{\partial B} \Delta B \quad (48)$$

and a straightforward calculation yields

$$\begin{aligned} \Delta x^{(c)} &= -(C_c x'_o + C_D x_o - C_\theta y_o) \left(\frac{\gamma \Delta \phi_o}{\hat{\phi}_o} - 2 \frac{\Delta B_o}{B_o} \right) \\ \Delta y^{(c)} &= -(C_c y_o + C_D y_o + C_\theta x_o) \left(\frac{\gamma \Delta \phi_o}{\hat{\phi}_o} - 2 \frac{\Delta B_o}{B_o} \right) \end{aligned} \quad (49)$$

for magnetic lenses or

$$x^{(c)} = -(C_c x'_o + C_D x_o) \frac{\Delta \phi_o}{\hat{\phi}_o} \quad (50)$$

with a similar expression for $y^{(c)}$ for electrostatic lenses. In objective lenses, the dominant aberration is the (axial)

chromatic aberration C_c , which causes a blur in the image that is independent of the position of the object point, like that due to C_s . The coefficient C_c also shares with C_s the property of being intrinsically positive. The coefficients C_D and C_θ affect projector lenses, but although they are pure distortions, they may well cause blurring since the term in $\Delta \phi_o$ and ΔB_o represents a spread, as in the case of the initial electron energy, or an oscillation typically at main frequency, coming from the power supplies.

Although a general theory can be established for the parasitic aberrations, this is much less useful than the theory of the geometric and chromatic aberrations; because the parasitic aberrations are those caused by accidental, unsystematic errors—imperfect roundness of the openings in a round lens, for example, or inhomogeneity of the magnetic material of the yoke of a magnetic lens, or imperfect alignment of the polepieces or electrodes. We therefore merely point out that one of the most important parasitic aberrations is an axial astigmatism due to the weak quadrupole field component associated with ellipticity of the openings. So large is this aberration, even in carefully machined lenses, that microscopes are equipped with a variable weak quadrupole, known as a stigmator, to cancel this unwelcome effect.

We will not give details of the aberrations of quadrupoles and prisms here. Quadrupoles have more independent aberrations than round lenses, as their lower symmetry leads us to expect, but these aberrations can be grouped into the same families: aperture aberrations, comas, field curvatures, astigmatisms, and distortions. Since the optic axis is straight, they are third-order aberrations, like those of round lenses, in the sense that the degree of the dependence on x_o , x'_o , y_o , and y'_o is three. The primary aberrations of prisms, on the other hand, are of second order, with the axis now being curved.

2. Lie Methods

An alternative way of using Hamiltonian mechanics to study the motion of charged particles has been developed, by Alex Dragt and colleagues especially, in which the properties of Lie algebra are exploited. This has come to be known as Lie optics. It has two attractions, one very important for particle optics at high energies (accelerator optics): first, interrelations between aberration coefficients are easy to establish, and second, high-order perturbations can be studied systematically with the aid of computer algebra and, in particular, of the differential algebra developed for the purpose by Martin Berz. At lower energies, the Lie methods provide a useful check of results obtained by the traditional procedures, but at higher energies they give valuable information that would be difficult to obtain in any other way.

C. Instrumental Optics: Components

1. Guns

The range of types of particle sources is very wide, from the simple triode gun with a hairpin-shaped filament relying on thermionic emission to the plasma sources furnishing high-current ion beams. We confine this account to the thermionic and field-emission guns that are used in electron-optical instruments to furnish modest electron currents: thermionic guns with tungsten or lanthanum hexaboride emitters, in which the electron emission is caused by heating the filament, and field-emission guns, in which a very high electric field is applied to a sharply pointed tip (which may also be heated). The current provided by the gun is not the only parameter of interest and is indeed often not the most crucial. For microscope applications, a knowledge of brightness B is much more important; this quantity is a measure of the quality of the beam. Its exact definition requires considerable care, but for our present purposes it is sufficient to say that it is a measure of the current density per unit solid angle in the beam. For a given current, the brightness will be high for a small area of emission and if the emission is confined to a narrow solid angle. In scanning devices, the writing speed and the brightness are interrelated, and the resulting limitation is so severe that the scanning transmission electron microscope (STEM) came into being only with the development of high-brightness field-emission guns. Apart from a change of scale with $\hat{\phi}_2/\hat{\phi}_1$ in accelerating structures, the brightness is a conserved quantity in electron-optical systems (provided that the appropriate definition of brightness is employed).

The simplest and still the most widely used electron gun is the triode gun, consisting of a heated filament or cathode, an anode held at a high positive potential relative to the cathode, and, between the two, a control electrode known as the *wehnelt*. The latter is held at a small negative potential relative to the cathode and serves to define the area of the cathode from which electrons are emitted. The electrons converge to a waist, known as the crossover, which is frequently within the gun itself (Fig. 9). If j_c is the current density at the center of this crossover and α_s is the angular spread (defined in Fig. 9), then

$$B = j_c / \pi \alpha_s^2 \quad (51)$$

It can be shown that B cannot exceed the Langmuir limit $B_{\max} = je\phi/\pi kT$, in which j is the current density at the filament, ϕ is the accelerating voltage, k is Boltzmann's constant (1.4×10^{-23} J/K), and T is the filament temperature. The various properties of the gun vary considerably with the size and position of the wehnelt and anode and the potentials applied to them; the general behavior has

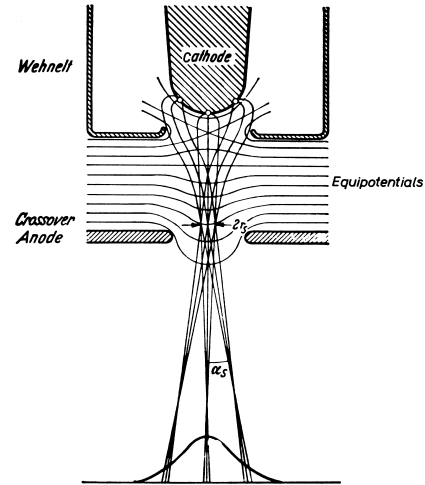


FIGURE 9 Electron gun and formation of the crossover.

been satisfactorily explained in terms of a rather simple model by Rolf Lauer.

The crossover is a region in which the current density is high, and frequently high enough for interactions between the beam electrons to be appreciable. A consequence of this is a redistribution of the energy of the particles and, in particular, an increase in the energy spread by a few electron volts. This effect, detected by Hans Boersch in 1954 and named after him, can be understood by estimating the mean interaction using statistical techniques.

Another family of thermionic guns has rare-earth boride cathodes, LaB_6 in particular. These guns were introduced in an attempt to obtain higher brightness than a traditional thermionic gun could provide, and they are indeed brighter sources; they are technologically somewhat more complex, however. They require a slightly better vacuum than tungsten triode guns, and in the first designs the LaB_6 rod was heated indirectly by placing a small heating coil around it; subsequently, however, directly heated designs were developed, which made these guns more attractive for commercial purposes.

Even LaB_6 guns are not bright enough for the needs of the high-resolution STEM, in which a probe only a few tenths of a nanometer in diameter is raster-scanned over a thin specimen and the transmitted beam is used to form an image (or images). Here, a field-emission gun is indispensable. Such guns consist of a fine tip and two (or more) electrodes, the first of which creates a very high electric field at the tip, while the second accelerates the electrons thus extracted to the desired accelerating voltage. Such guns operate satisfactorily only if the vacuum is very good indeed; the pressure in a field-emission gun must be some five or six orders of magnitude higher than that in a thermionic triode gun. The resulting brightness is

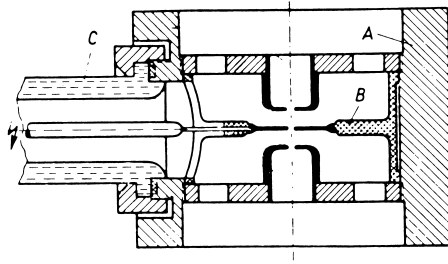


FIGURE 10 Electrostatic einzel lens design: (A) lens casing; (B and C) insulators.

appreciably higher, but the current is not always sufficient when only a modest magnification is required.

We repeat that the guns described above form only one end of the spectrum of particle sources. Others have large flat cathodes. Many are required to produce high currents and current densities, in which case we speak of space-charge flow; these are the Pierce guns and PIGs (Pierce ion guns).

2. Electrostatic Lenses

Round electrostatic lenses take the form of a series of plates in which a round opening has been pierced or circular cylinders all centered on a common axis (Fig. 10). The potentials applied may be all different or, more often, form a simple pattern. The most useful distinction in practice separates lenses that create no overall acceleration of the beam (although, of course, the particles are accelerated and decelerated within the lens field) and those that do produce an overall acceleration or deceleration. In the first case, the usual configuration is the einzel lens, in which the outer two of the three electrodes are held at anode potential (or at the potential of the last electrode of any lens upstream if this is not at anode potential) and the central electrode is held at a different potential. Such lenses were once used in electrostatic microscopes and are still routinely employed when the insensitivity of electrostatic systems to voltage fluctuations that affect all the potentials equally is exploited. Extensive sets of curves and tables describing the properties of such lenses are available.

Accelerating lenses with only a few electrodes have also been thoroughly studied; a configuration that is of interest today is the multielectrode accelerator structure. These accelerators are not intended to furnish very high particle energies, for which very different types of accelerator are employed, but rather to accelerate electrons to energies beyond the limit of the simple triode structure, which cannot be operated above ~ 150 kV. For microscope and microprobe instruments with accelerating voltages in the range of a few hundred kilovolts up to a few megavolts, therefore, an accelerating structure must be inserted

between the gun and the first condenser lens. This structure is essentially a multielectrode electrostatic lens with the desired accelerating voltage between its terminal electrodes. This point of view is particularly useful when a field-emission gun is employed because of an inconvenient aspect of the optics of such guns: The position and size of the crossover vary with the current emitted. In a thermionic gun, the current is governed essentially by the temperature of the filament and can hence be varied by changing the heating current. In field-emission guns, however, the current is determined by the field at the tip and is hence varied by changing the potential applied to the first electrode, which in turn affects the focusing field inside the gun. When such a gun is followed by an accelerator, it is not easy to achieve a satisfactory match for all emission currents and final accelerating voltages unless both gun and accelerator are treated as optical elements. Miniature lenses and guns and arrays of these are being fabricated, largely to satisfy the needs of nanolithography. A spectacular achievement is the construction of a scanning electron microscope that fits into the hand, no bigger than a pen. The optical principles are the same as for any other lens.

3. Magnetic Lenses

There are several kinds of magnetic lenses, but the vast majority have the form of an electromagnet pierced by a circular canal along which the electrons pass. Figure 11 shows such a lens schematically, and Fig. 12 illustrates a more realistic design in some detail. The magnetic flux

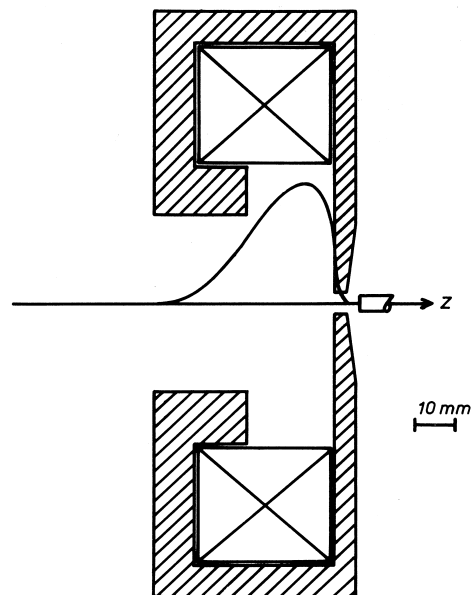


FIGURE 11 Typical field distribution in a magnetic lens.

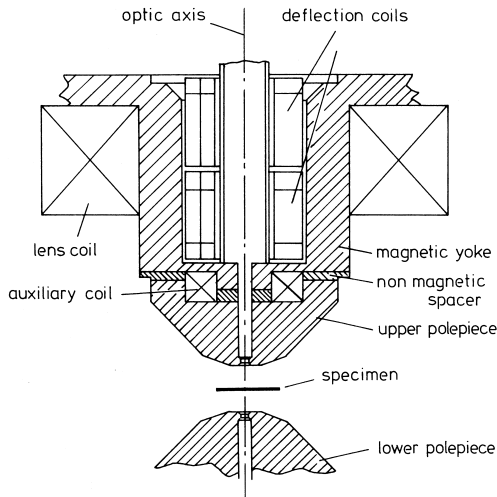


FIGURE 12 Modern magnetic objective lens design. (Courtesy of Philips, Eindhoven.)

is provided by a coil, which usually carries a low current through a large number of turns; water cooling prevents overheating. The magnetic flux is channeled through an iron yoke and escapes only at the gap, where the yoke is terminated with polepieces of high permeability. This arrangement is chosen because the lens properties will be most favorable if the axial magnetic field is in the form of a high, narrow bell shape (Fig. 11) and the use of a high-permeability alloy at the polepieces enables one to create a strong axial field without saturating the yoke. Considerable care is needed in designing the exact shape of these polepieces, but for a satisfactory choice, the properties of the lens are essentially determined by the gap S , the bore D (or the front and back bores if these are not the same), and the excitation parameter J ; the latter is defined by $J = NI/\hat{\phi}_0^{1/2}$, where NI is the product of the number of turns of the coil and the current carried by it and $\hat{\phi}_0$ is the relativistic accelerating voltage; S and D are typically of the order of millimeters and J is a few amperes per (volts)^{1/2}. The quantity NI can be related to the axial field strength with the aid of Ampère's circuital theorem (Fig. 13); we see that

$$\int_{-\infty}^{\infty} B(z) dz = \mu_0 NI \quad \text{so that} \quad NI \propto B_0$$

the maximum field in the gap, the constant of proportion-

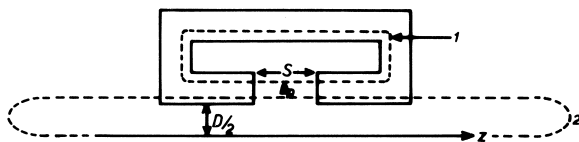


FIGURE 13 Use of Ampère's circuital theorem to relate lens excitation to axial field strength.

ality being determined by the area under the normalized flux distribution $B(z)/B_0$.

Although accurate values of the optical properties of magnetic lenses can be obtained only by numerical methods, in which the field distribution is first calculated by one of the various techniques available—finite differences, finite elements, and boundary elements in particular—their variation can be studied with the aid of field models. The most useful (though not the most accurate) of these is Glaser's bell-shaped model, which has the merits of simplicity, reasonable accuracy, and, above all, the possibility of expressing all the optical quantities such as focal length, focal distance, the spherical and chromatic aberration coefficients C_s and C_c , and indeed all the third-order aberration coefficients, in closed form, in terms of circular functions. In this model, $B(z)$ is represented by

$$B(z) = B_0 / (1 + z^2/a^2) \quad (52)$$

and writing $w^2 = 1 + k^2$, $k^2 = \eta^2 B_0^2 a^2 / 4\hat{\phi}_0$, $z = a \cot \psi$ the paraxial equation has the general solution

$$x(\psi) = (A \cos \psi + B \sin \psi) / \sin \psi \quad (53)$$

The focal length and focal distance can be written down immediately, and the integrals that give C_s and C_c can be evaluated explicitly. This model explains very satisfactorily the way in which these quantities vary with the excitation and with the geometric parameter a .

The traditional design of Fig. 12 has many minor variations in which the bore diameter is varied and the yoke shape altered, but the optical behavior is not greatly affected. The design adopted is usually a compromise between the optical performance desired and the technological needs of the user. In high-performance systems, the specimen is usually inside the field region and may be inserted either down the upper bore (top entry) or laterally through the gap (side entry). The specimen-holder mechanism requires a certain volume, especially if it is of one of the sophisticated models that permit *in situ* experiments: specimen heating, to study phase changes in alloys, for example, or specimen cooling to liquid nitrogen or liquid helium temperature, or straining; specimen rotation and tilt are routine requirements of the metallurgist. All this requires space in the gap region, which is further encumbered by a cooling device to protect the specimen from contamination, the stigmator, and the objective aperture drive. The desired optical properties must be achieved subject to the demands on space of all these devices, as far as this is possible. As Ugo Valdrè has said, the interior of an electron microscope objective should be regarded as a microlaboratory.

Magnetic lenses commonly operate at room temperature, but there is some advantage in going to very

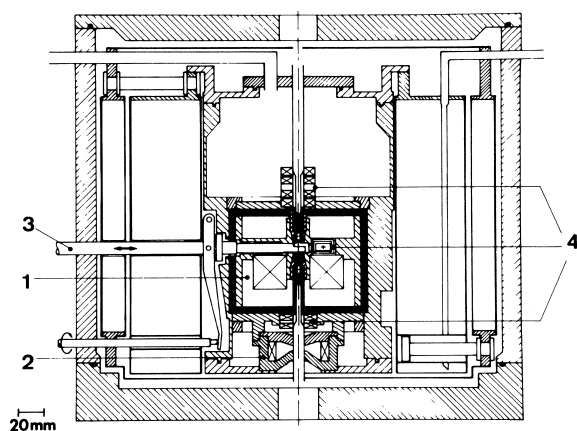


FIGURE 14 Superconducting lens system: (1) objective (shielding lens); (2) intermediate with iron circuit; (3) specimen holder; and (4) corrector device.

low temperature and running in the superconducting regime. Several designs have been explored since André Laberrigue, Humberto Fernández-Morán, and Hans Boersch introduced the first superconducting lenses, but only one has survived, the superconducting shielding lens introduced by Isolde Dietrich and colleagues at Siemens (Fig. 14). Here, the entire lens is at a very low temperature, the axial field being produced by a superconducting coil and concentrated into the narrow gap region by superconducting tubes. Owing to the Meissner–Ochsenfeld effect, the magnetic field cannot penetrate the metal of these superconducting tubes and is hence concentrated in the gap. The field is likewise prevented from escaping from the gap by a superconducting shield. Such lenses have been incorporated into a number of microscopes and are particularly useful for studying material that must be examined at extremely low temperatures; organic specimens that are irretrievably damaged by the electron beam at higher temperatures are a striking example.

Despite their very different technology, these superconducting lenses have essentially the same optical properties as their warmer counterparts. This is not true of the various magnetic lenses that are grouped under the heading of unconventional designs; these were introduced mainly by Tom Mulvey, although the earliest, the minilens, was devised by Jan Le Poole. The common feature of these lenses, which are extremely varied in appearance, is that the space occupied by the lens is very different in volume or shape from that required by a traditional lens. A substantial reduction in the volume can be achieved by increasing the current density in the coil; in the minilens (Fig. 15), the value may be $\sim 80 \text{ mm}^2$, whereas in a conventional lens, 2 A/mm^2 is a typical figure. Such lenses are employed as auxiliary lenses in zones already occupied by other elements, such as bulky traditional lenses. After the

initial success of these minilenses, a family of miniature lenses came into being, with which it would be possible to reduce the dimensions of the huge, heavy lenses used for very high voltage microscopes (in the megavolt range). Once the conventional design had been questioned, it was natural to inquire whether there was any advantage to be gained by abandoning its symmetric shape. This led to the invention of the pancake lens, flat like a phonograph record, and various single-polepiece or “snorkel” lenses (Fig. 16). These are attractive in situations where the electrons are at the end of their trajectory, and the single-polepiece design of Fig. 16 can be used with a target in front of it or a gun beyond it. Owing to their very flat shape, such lenses, with a bore, can be used to render microscope projector systems free of certain distortions, which are otherwise very difficult to eliminate.

This does not exhaust all the types of magnetic lens. For many years, permanent-magnet lenses were investigated in the hope that a simple and inexpensive microscope could be constructed with them. An addition to the family of traditional lenses is the unsymmetric triple-polepiece lens, which offers the same advantages as the single-polepiece designs in the projector system. Magnetic lens studies have also been revived by the needs of electron beam lithography.

4. Aberration Correction

The quest for high resolution has been a persistent preoccupation of microscope designers since these instruments came into being. Scherzer’s theorem (1936) was therefore a very unwelcome result, showing as it did that the principal resolution-limiting aberration could never vanish in

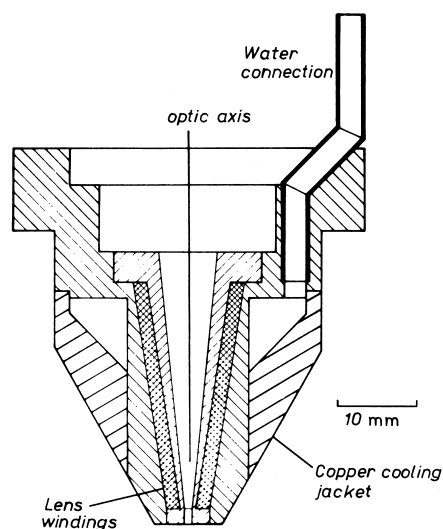


FIGURE 15 Minilens.

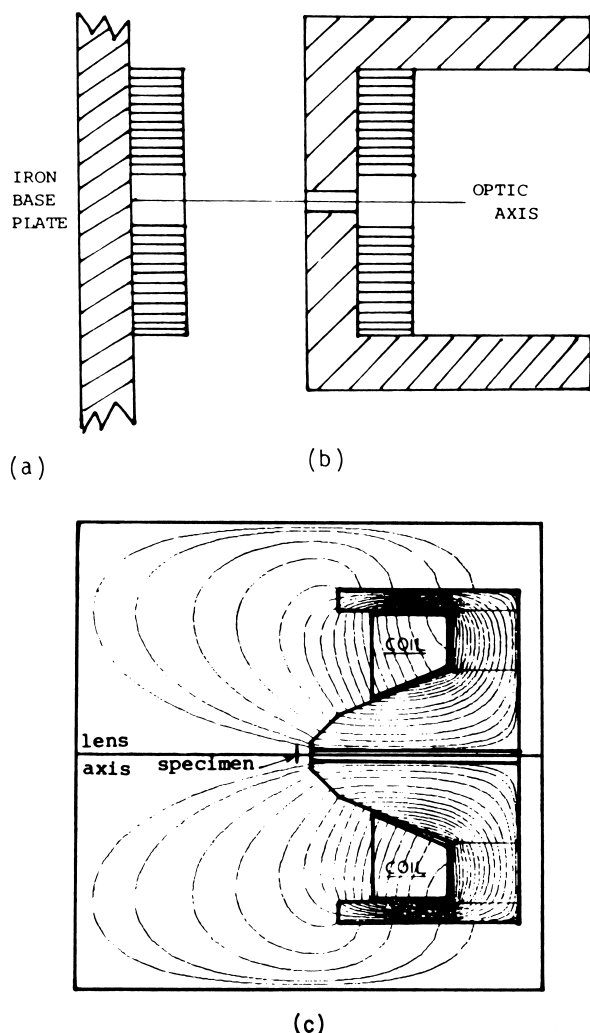


FIGURE 16 Some unconventional magnetic lenses.

round lenses. It was Scherzer again who pointed out (1947) the various ways of circumventing his earlier result by introducing aberration correctors of various kinds. The proof of the theorem required rotational symmetry, static fields, the absence of space charge, and the continuity of certain properties of the electrostatic potential. By relaxing any one of these requirements, aberration correction is in principle possible, but only two approaches have achieved any measure of success.

The most promising type of corrector was long believed to be that obtained by departing from rotational symmetry, and it was with such devices that correction was at last successfully achieved in the late 1990s. Such correctors fall into two classes. In the first, quadrupole lenses are employed. These introduce new aperture aberrations, but by adding octopole fields, the combined aberration of the round lens and the quadrupoles can be cancelled. At least four quadrupoles and three octopoles are required.

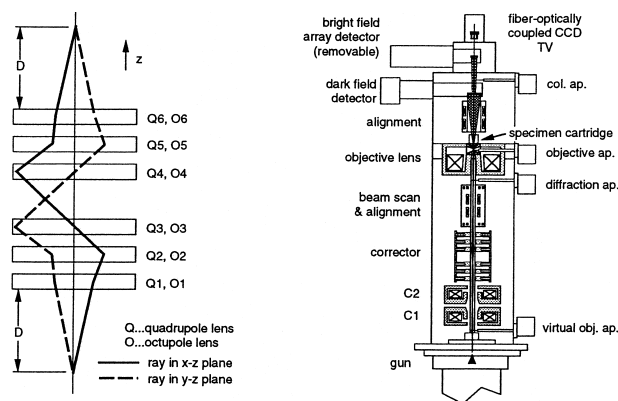


FIGURE 17 Correction of spherical aberration in a scanning transmission electron microscope. (Left) Schematic diagram of the quadrupole-octopole corrector and typical trajectories. (Right) Incorporation of the corrector in the column of a Vacuum Generators STEM. [From Krivanek, O. L., et al. (1997). Institute of Physics Conference Series 153, 35. Copyright IOP Publishing.]

A corrector based on this principle has been incorporated into a scanning transmission electron microscope by O. Krivanek at the University of Cambridge (Fig. 17). In the second class of corrector, the nonrotationally symmetric elements are sextupoles. A suitable combination of two sextupoles has a spherical aberration similar to that of a round lens but of opposite sign, and the undesirable second-order aberrations cancel out (Fig. 18). The technical difficulties of introducing such a corrector in a high-resolution transmission electron microscope have been overcome by M. Haider (Fig. 19).

Quadrupoles and octopoles had seemed the most likely type of corrector to succeed because the disturbance to the existing instrument, already capable of an unaided resolution of a few angstroms, was slight. The family of correctors that employ space charge or charged foils placed across the beam perturb the microscope rather more. Efforts continue to improve lenses by inserting one or more

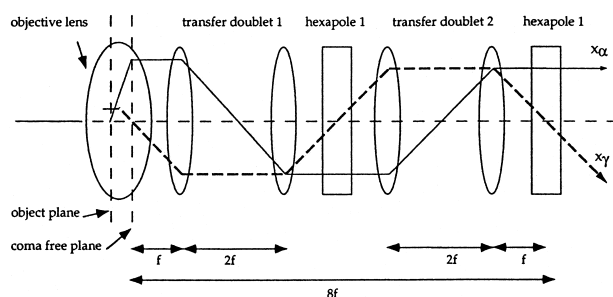


FIGURE 18 Correction of spherical aberration in a transmission electron microscope. Arrangement of round lenses and sextupoles (hexapoles) that forms a semiaplanatic objective lens. The distances are chosen to eliminate radial coma. [From Haider, M., et al. (1995). Optik 99, 167. Copyright Wissenschaftliche Verlagsgesellschaft.]

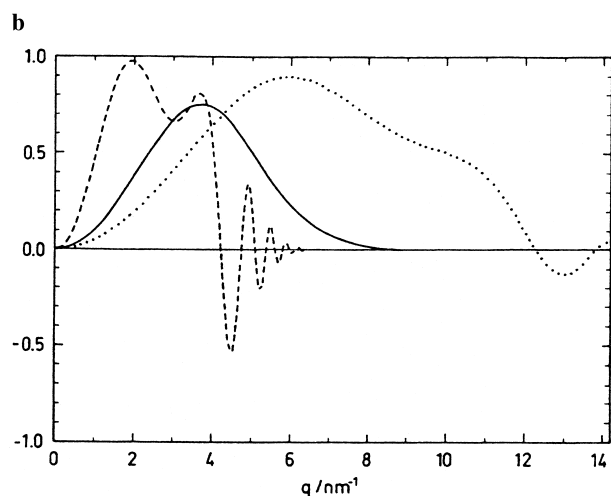
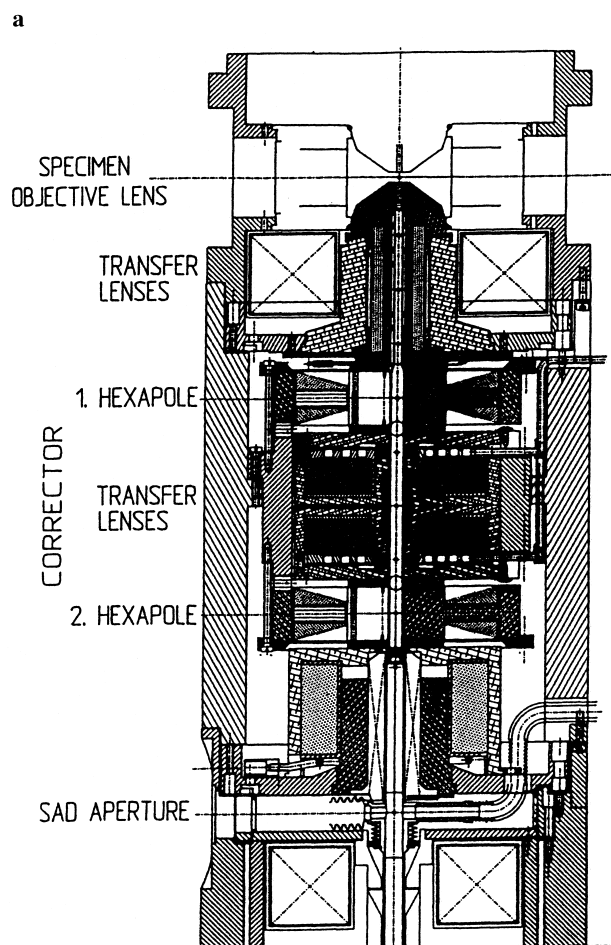


FIGURE 19 (a) The corrector of Fig. 18 incorporated in a transmission electron microscope. (b) The phase contrast transfer function of the corrected microscope. Dashed line: no correction. Full line: corrector switched on, energy width (a measure of the temporal coherence) 0.7 eV. Dotted line: energy width 0.2 eV. Chromatic aberration remains a problem, and the full benefit of the corrector is obtained only if the energy width is very narrow. [From Haider, M., et al. (1998). *J. Electron Microsc.* 47, 395. Copyright Japanese Society of Electron Microscopy.]

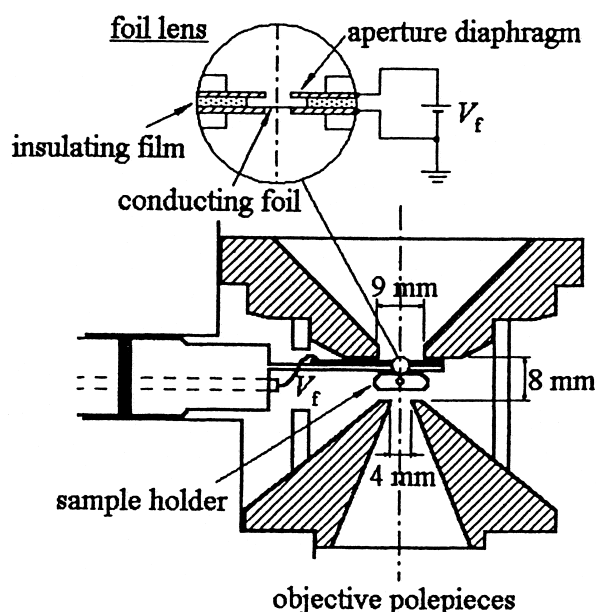


FIGURE 20 Foil lens and polepieces of an objective lens to be corrected. [From Hanai, T., et al. (1998). *J. Electron Microsc.* 47, 185. Copyright Japanese Society of Electron Microscopy.]

foils in the path of the electrons, with a certain measure of success, but doubts still persist about this method. Even if a reduction in total C_s is achieved, the foil must have a finite thickness and will inevitably scatter the electrons traversing it. How is this scattering to be separated from that due to the specimen? Figure 20 shows the design employed in an ongoing Japanese project.

An even more radical solution involves replacing the static objective lens by one or more microwave cavities. In Scherzer's original proposal, the incident electron beam was broken into short pulses and the electrons far from the axis would hence arrive at the lens slightly later than those traveling near the axis. By arranging that the axial electrons encounter the maximum field so that the peripheral electrons experience a weaker field, Scherzer argued, the effect of C_s could be eliminated since, in static lenses, the peripheral electrons are too strongly focused. Unfortunately, when we insert realistic figures into the corresponding equations, we find that the necessary frequency is in the gigahertz range, with the result that the electrons spend a substantial part of a cycle, or more than a cycle, within the microwave field. Although this means that the simple explanation is inadequate, it does not invalidate the principle, and experiment and theory both show that microwave cavity lenses can have positive or negative spherical aberration coefficients. The principal obstacles to their use are the need to produce very short pulses containing sufficient current and, above all, the fact that the beam emerging from such cavity lenses has a rather large energy

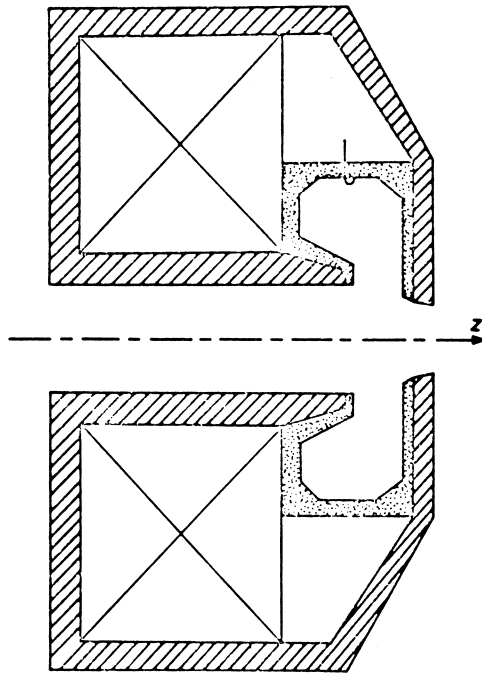


FIGURE 21 Microwave cavity lens between the polepieces of a magnetic lens. (Courtesy of L. C. Oldfield.)

spread, which makes further magnification a problem. An example is shown in Fig. 21.

Finally, we mention the possibility of *a posteriori* correction in which we accept the deleterious effect of C_s on the recorded micrograph but attempt to reduce or eliminate it by subsequent digital or analog processing of the image. A knowledge of the wave theory of electron image formation is needed to understand this idea and we therefore defer discussion of it to Section III.B.

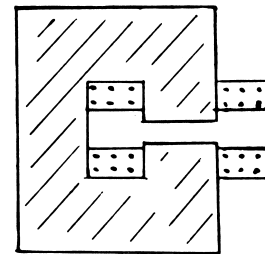
5. Prisms, Mirrors, and Energy Analyzers

Magnetic and electrostatic prisms and systems built up from these are used mainly for their dispersive properties in particle optics. We have not yet encountered electron mirrors, but we mention them here because a mirror action is associated with some prisms; if electrons encounter a potential barrier that is high enough to halt them, they will be reflected and a paraxial optical formalism can be developed to describe such mirror optics. This is less straightforward than for lenses, since the ray gradient is far from small at the turning point, which means that one of the usual paraxial assumptions that off-axis distance and ray gradient are everywhere small is no longer justified.

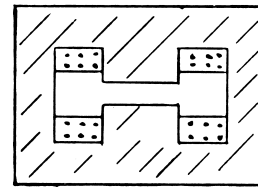
The simplest magnetic prisms, as we have seen, are sector fields created by magnets of the C-type or picture-frame arrangement (Fig. 22) with circular poles or sector poles with a sector or rectangular yoke. These or analogous

electrostatic designs can be combined in many ways, of which we can mention only a small selection. A very ingenious arrangement, which combines magnetic deflection with an electrostatic mirror, is the Castaing–Henry analyzer (Figs. 23a–23c) which has the constructional convenience that the incident and emergent optic axes are in line; its optical properties are such that an energy-filtered image or an energy spectrum from a selected area can be obtained. A natural extension of this is the magnetic Ω filter (Fig. 23d), in which the mirror is suppressed; if the particle energy is not too high, use of the electrostatic analog of this can be envisaged (Fig. 23e). It is possible to eliminate many of the aberrations of such filters by arranging the system not only symmetrically about the mid-plane ($x' - x$ in Fig. 23d), but also antisymmetrically about the planes midway between the mid-plane and the optic axis. A vast number of prism combinations have been explored by Veniamin Kel'man and colleagues in Alma-Ata in the quest for high-performance mass and electron spectrometers.

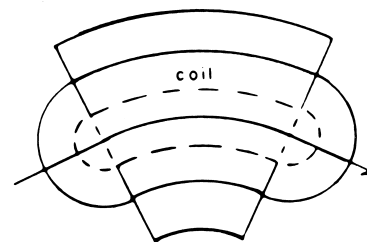
Energy analysis is a subject in itself, and we can do no more than mention various other kinds of energy or



(a)



(b)



(c)

FIGURE 22 (a) C-Type and (b) picture-frame magnets URE typically having (c) sector-shaped yoke and poles.

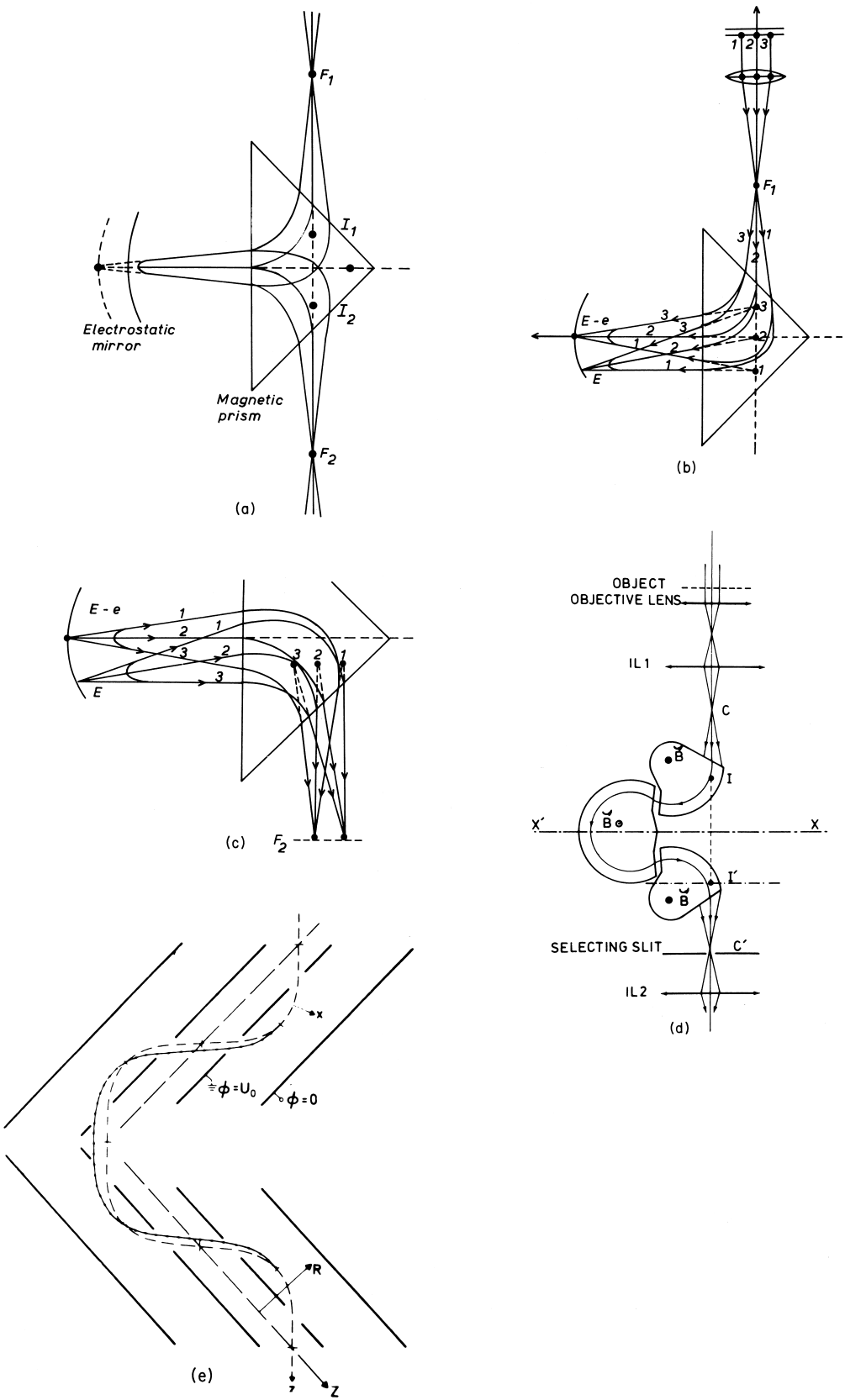


FIGURE 23 Analyzers: (a–c) Castaing–Henry analyzer; (d) Ω filter; and (e) electrostatic analog of the Ω filter.

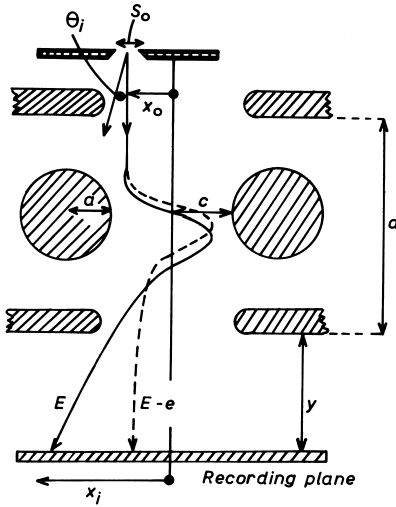


FIGURE 24 Möllenstedt analyzer.

momentum analyzers. The Wien filter consists of crossed electrostatic and magnetic fields, through which particles of a particular energy will pass undeflected, whereas all others will be deviated from their path. The early β -ray spectrometers exploited the fact that the chromatic aberration of a lens causes particles of different energies to be focused in different planes. The Möllenstedt analyzer is based on the fact that rays in an electrostatic lens far from the axis are rapidly separated if their energies are different (Fig. 24). The Ichinokawa analyzer is the magnetic analog of this and is used at higher accelerating voltages where electrostatic lenses are no longer practicable. In retarding-field analyzers, a potential barrier is placed in the path of the electrons and the current recorded as the barrier is progressively lowered.

6. Combined Deflection and Focusing Devices

In the quest for microminiaturization, electron beam lithography has acquired considerable importance. It proves to be advantageous to include focusing and deflecting fields within the same volume, and the optical properties of such combined devices have hence been thoroughly studied, particularly, their aberrations. It is important to keep the adverse effect of these aberrations small, especially because the beam must be deflected far from the original optical axis. An ingenious way of achieving this, proposed by Hajime, Ohiwa, is to arrange that the optic axis effectively shifts parallel to itself as the deflecting field is applied; for this, appropriate additional deflection, round and multipole fields must be superimposed and the result may be regarded as a “moving objective lens” (MOL) or “variable-axis lens” (VAL). Perfected immersion versions of these and of the “swinging objective

lens” (SOL) have been developed, in which the target lies within the field region.

III. WAVE OPTICS

A. Wave Propagation

The starting point here is not the Newton–Lorentz equations but Schrödinger’s equation; we shall use the nonrelativistic form, which can be trivially extended to include relativistic effects for magnetic lenses. Spin is thus neglected, which is entirely justifiable in the vast majority of practical situations. The full Schrödinger equation takes the form

$$-\frac{\hbar^2}{m_0} \nabla^2 \Psi + \frac{e\hbar}{im_0} \mathbf{A} \cdot \text{grad } \Psi + \left(-e\Phi + \frac{e^2}{2m_0} \mathbf{A}^2 \right) \Psi - i\hbar \frac{\partial \Psi}{\partial t} = 0 \quad (54)$$

and writing

$$\Psi(x, y, z, t) = \psi(x, y, z) e^{-i\omega t} \quad (55)$$

we obtain

$$-\frac{\hbar^2}{2m_0} \nabla^2 \psi + \frac{e\hbar}{im_0} \mathbf{A} \cdot \text{grad } \psi + \left(-e\Phi + \frac{e^2}{2m_0} \mathbf{A}^2 \right) \psi = E\psi \quad (56)$$

with

$$E = \hbar\omega \quad (57)$$

where $\hbar = h/2\pi$ and h is Planck’s constant. The free-space solution corresponds to

$$p = h/\lambda$$

or

$$\lambda = h/(2em_0\phi_0)^{1/2} \approx 12.5/\phi_0^{1/2} \quad (58)$$

where p is the momentum.

As in the case of geometric optics, we consider the paraxial approximation, which for the Schrödinger equation takes the form

$$-\hbar^2 \left(\frac{\partial^2 \psi}{\partial x^2} + \frac{\partial^2 \psi}{\partial y^2} \right) + \frac{1}{2} em_0 (\phi'' + \eta^2 B^2) (x^2 + y^2) \psi - i\hbar p' \psi - 2i\hbar p \frac{\partial \psi}{\partial z} = 0 \quad (59)$$

and we seek a wavelike solution:

$$\psi(x, y, z) = a(z) \exp[iS(x, y, z)/\hbar]. \quad (60)$$

After some calculation, we obtain the required equation describing the propagation of the wave function through electrostatic and magnetic fields:

$$\begin{aligned} \psi(x, y, z) = & \frac{p_0^{3/2}}{2\pi i \hbar h(z) p^{1/2}} \exp \left[\frac{i p g'(z)}{2 \hbar g(z)} (x^2 + y^2) \right] \\ & \times \iint_{-\infty}^{\infty} \psi(x_o, y_o, z_o) \exp \left\{ \frac{i p_o}{2 \hbar g(z) h(z)} \right. \\ & \times \left[(x - x_o g)^2 + (y - y_o g)^2 \right] dx_o dy_o \end{aligned} \quad (61)$$

This extremely important equation is the basis for all that follows. In it, $g(z)$ and $h(z)$ now denote the solutions of the geometric paraxial equations satisfying the boundary conditions $g(z_o) = h'(z_o) = 1$, $g'(z_o) = h(z_o) = 0$. Reorganizing the various terms, Eq. (61) can be written

$$\begin{aligned} \psi(x, y, z) = & \frac{1}{i \lambda r h(z)} \iint_{-\infty}^{\infty} \psi(x_o, y_o, z_o) \\ & \times \exp \left\{ \frac{i \pi}{\lambda h(z)} [g(z)(x_o^2 + y_o^2) \right. \\ & - 2(x_o x + y_o y) \\ & \left. + r h'(z)(x^2 + y^2)] \right\} dx_o dy_o \end{aligned} \quad (62)$$

with $\lambda = h/p_o$ and $r = p/p_o = (\phi/\phi_o)^{1/2}$.

Let us consider the plane $z = z_d$ in which $g(z)$ vanishes, $g(z_d) = 0$. For the magnetic case ($r = 1$), we find

$$\begin{aligned} \psi(x_d, y_d, z_d) = & \frac{E_d}{i \lambda h(z_o)} \iint \psi(x_o, y_o, z_o) \\ & \times \exp \left[-\frac{2i}{\lambda h(z_d)} (x_o x_d + y_o y_d) \right] dx_o dy_o \end{aligned} \quad (63)$$

with $E_d = \exp[i\pi h'(z_d)(x_d^2 + y_d^2)/\lambda h(z_d)]$, so that, scale factors apart, the wave function in this plane is the Fourier transform of the same function in the object plane.

We now consider the relation between the wave function in the object plane and in the image plane $z = z_i$ conjugate to this, in which $h(z)$ vanishes: $h(z_i) = 0$. It is convenient to calculate this in two stages, first relating $\psi(x_i, y_i, z_i)$ to the wave function in the exit pupil plane of the lens, $\psi(x_a, y_a, z_a)$ and then calculating the latter with the aid of Eq. (62). Introducing the paraxial solutions $G(z)$, $H(z)$ such that

$$G(z_a) = H'(z_a) = 1; \quad G'(z_a) = H(z_a) = 0$$

we have

$$\begin{aligned} \psi(x_i, y_i, z_i) = & \frac{1}{i \lambda H(z_i)} \iint \psi(x_a, y_a, z_a) \\ & \times \exp \left\{ \frac{i \pi}{\lambda H(z)} [G(z_i)(x_a^2 + y_a^2) \right. \\ & - 2(x_a x_i + y_a y_i) \\ & \left. + H'(z_i)(x_i^2 + y_i^2)] \right\} dx_a dy_a \end{aligned} \quad (64)$$

Using Eq. (62), we find

$$\begin{aligned} M \psi(x_i, y_i, z_i) E_i \\ = \iint \psi(x_o, y_o, z_o) K(x_i, y_i; x_o, y_o) E_o dx_o dy_o \end{aligned} \quad (65)$$

where M is the magnification, $M = g(z_i)$, and

$$\begin{aligned} E_i = & \exp \left[\frac{i \pi}{\lambda M} \frac{g_a h'_i - g'_i h_a}{h_a} (x_i^2 + y_i^2) \right] \\ E_o = & \exp \left[\frac{i \pi g_a}{\lambda h_a} (x_o^2 + y_o^2) \right] \end{aligned} \quad (66)$$

These quadratic factors are of little practical consequence; they measure the curvature of the wave surface arriving at the specimen and at the image. If the diffraction pattern plane coincides with the exit pupil, then $E_o = 1$. We write $h(z_a) = f$ since this quantity is in practice close to the focal length, so that for the case $z_d = z_a$,

$$E_i = \exp \left[-\frac{i \pi g'_i}{\lambda M} (x_i^2 + y_i^2) \right] \quad (67)$$

The most important quantity in Eq. (65) is the function $K(x_i, y_i; x_o, y_o)$, which is given by

$$\begin{aligned} K(x, y; x_o, y_o) = & \frac{1}{\lambda^2 f^2} \iint A(x_a, y_a) \\ & \times \exp \left\{ -\frac{2\pi i}{\lambda f} \left[\left(x_o - \frac{x}{M} \right) x_a \right. \right. \\ & \left. \left. + \left(y_o - \frac{y}{M} \right) y_a \right] \right\} dx_a dy_a \end{aligned} \quad (68)$$

or introducing the spatial frequency components

$$\xi = x_a/\lambda f; \quad \eta = y_a/\lambda f \quad (69)$$

we find

$$\begin{aligned} K(x, y; x_o, y_o) = & \iint A(\lambda f \xi, \lambda f \eta) \\ & \times \exp \left\{ -2\pi i \left[\left(x_o - \frac{x}{M} \right) \xi \right. \right. \\ & \left. \left. + \left(y_o - \frac{y}{M} \right) \eta \right] \right\} d\xi d\eta \end{aligned} \quad (70)$$

In the paraxial approximation, the aperture function A is simply a mathematical device defining the area of integration in the aperture plane: $A = 1$ inside the pupil and $A = 0$ outside the pupil. If we wish to include the effect of geometric aberrations, however, we can represent them as a phase shift of the electron wave function at the exit pupil. Thus, if the lens suffers from spherical aberration, we write

$$A(x_a, y_a) = a(x_a, y_a) \exp[-i\gamma(x_a, y_a)] \quad (71)$$

in which

$$\begin{aligned} \gamma &= \frac{2\pi}{\lambda} \left\{ \frac{1}{4} C_s \left(\frac{x_a^2 + y_a^2}{f^2} \right)^2 - \frac{1}{2} \Delta \frac{x_a^2 + y_a^2}{f^2} \right\} \\ &= \frac{\pi\lambda}{2} \{ C_s \lambda^2 (\xi^2 + \eta^2)^2 - 2\Delta(\xi^2 + \eta^2) \} \end{aligned} \quad (72)$$

the last term in Δ allowing for any defocus, that is, any small difference between the object plane and the plane conjugate to the image plane. All the third-order geometric aberrations can be included in the phase shift γ , but we consider only C_s and the defocus Δ . This limitation is justified by the fact that C_s is the dominant aberration of objective lenses and proves to be extremely convenient because Eq. (65) relating the image and object wave functions then has the form of a convolution, which it loses if other aberrations are retained (although coma can be accommodated rather uncomfortably). It is now the amplitude function $a(x_a, y_a)$ that represents the physical pupil, being equal to unity inside the opening and zero elsewhere.

In the light of all this, we rewrite Eq. (65) as

$$\begin{aligned} E_i \psi(x_i, y_i, z_i) &= \frac{1}{M} \iint K \left(\frac{x_i}{M} - x_o, \frac{y_i}{M} - y_o \right) E_o \\ &\quad \times \psi_o(x_o, y_o, z_o) dx_o dy_o \end{aligned} \quad (73)$$

Defining the Fourier transforms of ψ_o , ψ_i , and K as follows,

$$\begin{aligned} \tilde{\psi}_o(\xi, \eta) &= \iint E_o \psi_o \\ &\quad \times \exp[-2\pi i(\xi x_o + \eta y_o)] dx_o dy_o \\ \tilde{\psi}_i(\xi, \eta) &= \iint E_i \psi_i(Mx_i, My_i) \\ &\quad \times \exp[-2\pi i(\xi x_i + \eta y_i)] dx_i dy_i \\ &= \frac{1}{M^2} \iint E_i \psi_i(x_i, y_i) \\ &\quad \times \exp \left[-2\pi i \frac{(\xi x_i + \eta y_i)}{M} \right] dx_i dy_i \\ \tilde{K}(\xi, \eta) &= \iint K(x, y) \\ &\quad \times \exp[-2\pi i(\xi x + \eta y)] dx dy \end{aligned} \quad (74)$$

in which small departures from the conventional definitions have been introduced to assimilate inconvenient factors, Eq. (65) becomes

$$\tilde{\psi}_i(\xi, \eta) = \frac{1}{M} \tilde{K}(\xi, \eta) \tilde{\psi}_o(\xi, \eta) \quad (75)$$

This relation is central to the comprehension of electron-optical image-forming instruments, for it tells us that the formation of an image may be regarded as a filtering operation. If \tilde{K} were equal to unity, the image wave function would be identical with the object wave function, appropriately magnified; but in reality \tilde{K} is not unity and different spatial frequencies of the wave leaving the specimen, $\psi(x_o, y_o, z_o)$, are transferred to the image with different weights. Some may be suppressed, some attenuated, some may have their sign reversed, and some, fortunately, pass through the filter unaffected. The notion of spatial frequency is the spatial analog of the temporal frequency, and we associate high spatial frequencies with fine detail and low frequencies with coarse detail; the exact interpretation is in terms of the Fourier transform, as we have seen.

We shall use Eqs. (73) and (75) to study image formation in two types of optical instruments, the transmission electron microscope (TEM) and its scanning transmission counterpart, the STEM. This is the subject of the next section.

B. Instrumental Optics: Microscopes

The conventional electron microscope (TEM) consists of a source, condenser lenses to illuminate a limited area of the specimen, an objective to provide the first stage of magnification beyond the specimen, and projector lenses, which magnify the first intermediate image or, in diffraction conditions, the pattern formed in the plane denoted by $z = z_d$ in the preceding section. In the STEM, the role of the condenser lenses is to demagnify the crossover so that a very small electron probe is formed on the specimen. Scanning coils move this probe over the surface of the latter in a regular raster, and detectors downstream measure the current transmitted. There are inevitably several detectors, because transmission microscope specimens are essentially transparent to electrons, and thus there is no diminution of the total current but there is a redistribution of the directions of motion present in the beam. Electron-optical specimens deflect electrons but do not absorb them. In the language of light optics, they are phase specimens, and the electron microscope possesses means of converting invisible phase variations of amplitude variations that the eye can see.

We now examine image formation in the TEM in more detail. We first assume, and it is a very reasonable first approximation, that the specimen is illuminated by a parallel

uniform beam of electrons or, in other words, that the wave incident on the specimen is a constant. We represent the effect of the specimen on this wave by a multiplicative specimen transparency function $S(x_o, y_o)$, which is a satisfactory model for the very thin specimens employed for high-resolution work and for many other specimens. This specimen transparency is a complex function, and we write

$$S(x_o, y_o) = [1 - s(x_o, y_o)] \exp[i\varphi(x_o, y_o)] \quad (76a)$$

$$= 1 - s + i\varphi \quad (76b)$$

for small values of s and φ . The real term s requires some explanation, for our earlier remarks suggest that s must vanish if no electrons are halted by the specimen. we retain the term in s for two reasons. First, some electrons may be scattered inelastically in the specimen, in which case they must be regarded as lost in this simple monochromatic and hence monoenergetic version of the theory. Second, all but linear terms have been neglected in the approximate expression (76b) and, if necessary, the next-higher-order term ($-\frac{1}{2}\varphi^2$) can be represented by s .

The wave leaving the specimen is now proportional to S normalizing, so that the constant of proportionality is unity; after we substitute

$$\psi(x_o, y_o, z_o) = 1 - s + i\varphi \quad (77)$$

into Eq. (75). Again denoting Fourier transforms by the tilde, we have

$$\begin{aligned} \tilde{\psi}_i(\xi, \eta) &= \frac{1}{M} \tilde{K}(\xi, \eta) [\delta(\xi, \eta) - \tilde{s}(\xi, \eta) + i\tilde{\varphi}(\xi, \eta)] \\ &= \frac{1}{M} a \exp(-i\gamma)(\delta - \tilde{s} + i\tilde{\varphi}) \end{aligned} \quad (78)$$

and hence

$$\begin{aligned} \psi_i(Mx_i, My_i) &= \frac{1}{M} \iint a \exp(-i\gamma)(\delta - \tilde{s} + i\tilde{\varphi}) \\ &\quad \times \exp[2\pi i(\xi x_i + \eta y_i)] d\xi d\eta \end{aligned} \quad (79)$$

The current density at the image, which is what we see on the fluorescent screen of the microscope and record on film, is proportional to $\psi_i \psi_i^*$. We find that if both φ and s are small,

$$\begin{aligned} M^2 \psi_i \psi_i^* &\approx 1 - 2 \iint_{-\infty}^{\infty} a \tilde{s} \cos \gamma \\ &\quad \times \exp[2\pi i(\xi x + \eta y)] d\xi d\eta \\ &\quad + 2 \iint_{-\infty}^{\infty} a \tilde{\varphi} \sin \gamma \\ &\quad \times \exp[2\pi i(\xi x + \eta y)] d\xi d\eta \end{aligned} \quad (80)$$

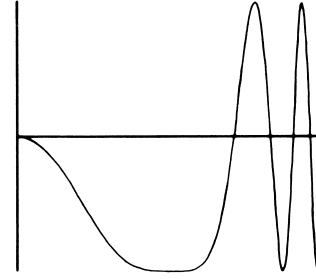


FIGURE 25 Function $\sin \gamma$ at Scherzer defocus $\Delta = (C_s \lambda)^{1/2}$.

and writing $j = M^2 \psi_i \psi_i^*$ and $C = j - 1$, we see that

$$\tilde{C} = -2a\tilde{s} \cos \gamma + 2a\tilde{\varphi} \sin \gamma \quad (81)$$

This justifies our earlier qualitative description of image formation as a filter process. Here we see that the two families of spatial frequencies characterizing the specimen, $\tilde{\varphi}$ and \tilde{s} , are distorted before they reach the image by the linear filters $\cos \gamma$ and $\sin \gamma$. The latter is by far the more important. A typical example is shown in Fig. 25. The distribution $2 \sin \gamma$ can be observed directly by examining an amorphous phase specimen, a very thin carbon film, for example. The spatial frequency spectrum of such a specimen is fairly uniform over a wide range of frequencies so that $\tilde{C} \propto \sin \gamma$. A typical spectrum is shown in Fig. 26, in which the radial intensity distribution is proportional to $\sin^2 \gamma$. Such spectra can be used to estimate the defocus Δ and the coefficient C_s very accurately.

The foregoing discussion is idealized in two respects, both serious in practice. First, the illuminating beam has been assumed to be perfectly monochromatic, whereas in reality there will be a spread of wavelengths of several parts per million; in addition, the wave incident on the specimen has been regarded as a uniform plane wave, which is equivalent to saying that it originated in an ideal point source. Real sources, of course, have a finite size, and the single plane wave should therefore be replaced by a spectrum of plane waves incident at a range of small angles to the specimen. The step from point source and monochromatic beam to finite source size and finite wavelength spread is equivalent to replacing perfectly coherent illumination by partially coherent radiation, with the wavelength spread corresponding to temporal partial coherence and the finite source size corresponding to spatial partial coherence. (We cannot discuss the legitimacy of separating these effects here, but simply state that this is almost always permissible.) Each can be represented by an envelope function, which multiplies the coherent transfer functions $\sin \gamma$ and $\cos \gamma$. This is easily seen for the temporal spatial coherence. Let us associate a probability distribution $H(f)$, $\int H(f) df = 1$, with the current density at each point in the image, the argument f being some

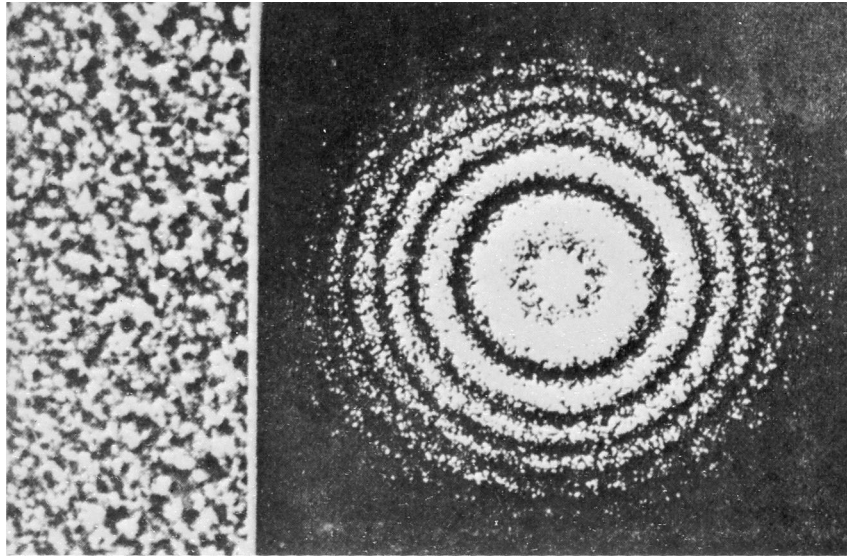


FIGURE 26 Spatial frequency spectrum (right) of an amorphous phase specimen (left).

convenient measure of the energy variation in the beam incident on the specimen. Hence, $dj/j = H(f) df$. From Eq. (80), we find

$$j = 1 - \int a \tilde{T}_s \exp[2\pi i(\xi x + \eta y)] d\xi d\eta + \int a \tilde{T}_\varphi \exp[2\pi i(\xi x + \eta y)] d\xi d\eta \quad (82)$$

where

$$T_s = 2 \int \cos \gamma(\xi, \eta, f) H(f) df \quad (83)$$

$$T_\varphi = 2 \int \sin \gamma(\xi, \eta, f) H(f) df$$

and if f is a measure of the defocus variation associated with the energy spread, we may set Δ equal to $\Delta_o + f$, giving

$$T_s = 2 \cos \gamma \int H(f) \cos[\pi \lambda f(\xi^2 + \eta^2)] df \quad (84)$$

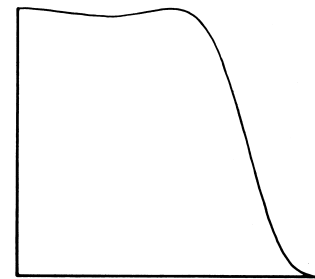
$$T_\varphi = 2 \sin \gamma \int H(f) \cos[\pi \lambda f(\xi^2 + \eta^2)] df$$

if $H(f)$ is even, and a slightly longer expression when it is not.

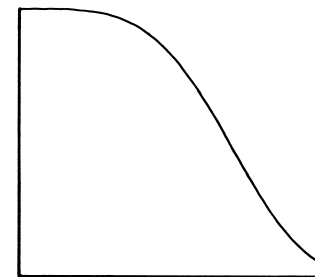
The familiar $\sin \gamma$ and $\cos \gamma$ are thus clearly seen to be modulated by an envelope function, which is essentially the Fourier transform of $H(f)$. A similar result can be obtained for the effect of spatial partial coherence, but the demonstration is longer. Some typical envelope functions are shown in Fig. 27.

An important feature of the function $\sin \gamma$ is that it gives us a convenient measure of the resolution of the

microscope. Beyond the first zero of the function, information is no longer transferred faithfully, but in the first zone the transfer is reasonably correct until the curve begins to dip toward zero for certain privileged values of the defocus, $\Delta = (C_s \lambda)^{1/2}$, $(3C_s \lambda)^{1/2}$, and $(5C_s \lambda)^{1/2}$; for the first of these values, known as the Scherzer defocus,



(a)



(b)

FIGURE 27 Envelope functions characterizing (a) spatial and (b) temporal partial coherence.

the zero occurs at the spatial frequency $(C_s\lambda^3)^{-1/4}$; the reciprocal of this multiplied by one of various factors has long been regarded as the resolution limit of the electron microscope, but transfer function theory enables us to understand the content of the image in the vicinity of the limit in much greater detail. The arrival of commercial electron microscopes equipped with spherical aberration correctors is having a profound influence on the practical exploitation of transfer theory. Hitherto, the effect of spherical aberration dictated the mode of operation of the TEM when the highest resolution was required. When the coefficient of spherical aberration has been rendered very small by correction, this defect is no longer the limiting factor and other modes of operation become of interest.

We now turn to the STEM. Here a bright source, typically a field-emission gun, is focused onto the specimen; the small probe is scanned over the surface and, well beyond the specimen, a far-field diffraction pattern of each elementary object area is formed. This pattern is sampled by a structured detector, which in the simplest case consists of a plate with a hole in the center, behind which is another plate or, more commonly, an energy analyzer. The signals from the various detectors are displayed on cathode-ray tubes, locked in synchronism with the scanning coils of the microscope. The reason for this combination of annular detector and central detector is to be found in the laws describing electron scattering. The electrons incident on a thin specimen may pass through unaffected; or they may be deflected with virtually no transfer of energy to the specimen, in which case they are said to be elastically scattered; or they may be deflected and lose energy, in which case they are inelastically scattered. The important point is that, on average, inelastically scattered electrons are deflected through smaller angles than those scattered elastically, with the result that the annular detector receives mostly elastically scattered particles, whereas the central detector collects those that have suffered inelastic collisions. The latter therefore have a range of energies, which can be separated by means of an energy analyzer, and we could, for example, form an image with the electrons corresponding to the most probable energy loss for some particular chemical element of interest. Another imaging mode exploits electrons that have been Rutherford scattered through rather large angles.

These modes of STEM image formation and others that we shall meet below can be explained in terms of a transfer function theory analogous to that derived for the TEM. This is not surprising, for many of the properties of the STEM can be understood by regarding it as an inverted TEM, the TEM gun corresponding to the small central

detector of the STEM and the large recording area of the TEM to the source in the STEM, spread out over a large zone if we project back the scanning probe. We will not pursue this analogy here, but most texts on the STEM explore it in some detail. Consider now a probe centered on a point $\mathbf{x}_0 = \boldsymbol{\xi}$ in the specimen plane of the STEM. We shall use a vector notation here, so that $\mathbf{x}_0 = (x_0, y_0)$, and similarly for other coordinates. The wave emerging from the specimen will be given by

$$\psi(\mathbf{x}_0; \boldsymbol{\xi}) = S(\mathbf{x}_0)K(\boldsymbol{\xi} - \mathbf{x}_0) \quad (85)$$

in which $S(\mathbf{x}_0)$ is again the specimen transparency and K describes the incident wave and, in particular, the effect of the pupil size, defocus, and aberrations of the probe-forming lens, the last member of the condenser system. Far below the specimen, in the detector plane (subscript d) the wave function is given by

$$\psi_d(\mathbf{x}_d, \boldsymbol{\xi}) = \int S(\mathbf{x}_0)K(\boldsymbol{\xi} - \mathbf{x}_0) \times \exp(-2\pi i \mathbf{x}_d \cdot \mathbf{x}_0 / \lambda R) d\mathbf{x}_0 \quad (86)$$

in which R is a measure of the effective distance between the specimen and the detector. The shape of the detector (and its response if this is not uniform) can most easily be expressed by introducing a detector function $D(\mathbf{x}_d)$, equal to zero outside the detector and equal to its response, usually uniform, over its surface. The detector records incident current, and the signal generated is therefore proportional to

$$\begin{aligned} j_d(\boldsymbol{\xi}) &= \int |\psi_d(\mathbf{x}_d; \boldsymbol{\xi})|^2 D(\mathbf{x}_d) d\mathbf{x}_d \\ &= \iiint S(\mathbf{x}_0)S^*(\mathbf{x}'_0)K(\boldsymbol{\xi} - \mathbf{x}_0)K^*(\boldsymbol{\xi} - \mathbf{x}'_0) \\ &\quad \times \exp[-2\pi i \mathbf{x}_d \cdot (\mathbf{x}_0 - \mathbf{x}'_0) / \lambda R] \\ &\quad \times D(\mathbf{x}_d) d\mathbf{x}_0 d\mathbf{x}'_0 d\mathbf{x}_d \end{aligned} \quad (87)$$

or introducing the Fourier transform of the detector response,

$$\begin{aligned} j_d(\boldsymbol{\xi}) &= \iint S(\mathbf{x}_0)S^*(\mathbf{x}'_0)K(\boldsymbol{\xi} - \mathbf{x}_0)K^*(\boldsymbol{\xi} - \mathbf{x}'_0) \\ &\quad \times \tilde{D}\left(\frac{\mathbf{x}_0 - \mathbf{x}'_0}{\lambda R}\right) d\mathbf{x}_0 d\mathbf{x}'_0 \end{aligned} \quad (88)$$

We shall use the formula below to analyze the signals collected by the simpler detectors, but first we derive the STEM analog of the filter Eq. (81). For this we introduce

the Fourier transforms of S and K into the expression for $\psi_d(\mathbf{x}_d, \boldsymbol{\xi})$. Setting $\mathbf{u} = \mathbf{x}_d/\lambda R$, we obtain

$$\begin{aligned}
 \psi_d(\lambda R \mathbf{u}; \boldsymbol{\xi}) &= \int S(\mathbf{x}_o) K(\boldsymbol{\xi} - \mathbf{x}_o) \exp(-2\pi i \mathbf{u} \cdot \mathbf{x}_o) d\mathbf{x}_o \\
 &= \iiint \tilde{S}(\mathbf{p}) \tilde{K}(\mathbf{q}) \\
 &\quad \times \exp[-2\pi i \mathbf{x}_o \cdot (\mathbf{u} - \mathbf{p} + \mathbf{q})] \\
 &\quad \times \exp[(-2\pi i \mathbf{q} \cdot \boldsymbol{\xi})] d\mathbf{p} d\mathbf{q} d\mathbf{x}_o \\
 &= \iint \tilde{S}(\mathbf{p}) \tilde{K}(\mathbf{q}) \delta(\mathbf{u} - \mathbf{p} + \mathbf{q}) \\
 &\quad \times \exp(2\pi i \mathbf{q} \cdot \boldsymbol{\xi}) d\mathbf{p} d\mathbf{q} \\
 &= \int \tilde{S}(\mathbf{p}) \tilde{K}(\mathbf{p} - \mathbf{u}) \exp[2\pi i \boldsymbol{\xi} \cdot (\mathbf{p} - \mathbf{u})] d\mathbf{p}
 \end{aligned} \tag{89}$$

After some calculation, we obtain an expression for $j_d(\boldsymbol{\xi}) = \int j(\mathbf{x}_d; \boldsymbol{\xi}) D(\mathbf{x}_d) d\mathbf{x}_d$ and hence for its Fourier transform

$$\tilde{j}_d(\mathbf{p}) = \int \tilde{j}(\mathbf{x}_d; \mathbf{p}) D(\mathbf{x}_d) d\mathbf{x}_d \tag{90}$$

Explicitly,

$$\begin{aligned}
 \tilde{j}_d(\mathbf{p}) &= \int |\tilde{K}(\mathbf{x}_d/\lambda R)|^2 D(\mathbf{x}_d) \delta(\mathbf{p}) \\
 &\quad - \tilde{s}(\mathbf{p}) \int q_s(\mathbf{x}_d/\lambda R; \mathbf{p}) D(\mathbf{x}_d) d\mathbf{x}_d \\
 &\quad + i \tilde{\varphi}(\mathbf{p}) \int q_\varphi(\mathbf{x}_d/\lambda R; \mathbf{p}) D(\mathbf{x}_d) d\mathbf{x}_d
 \end{aligned} \tag{91}$$

for weakly scattering objects, $s \ll 1$, $\varphi \ll 1$. The spatial frequency spectrum of the bright-field image signal is thus related to s and φ by a filter relation very similar to that obtained for the TEM.

We now return to Eqs. (87) and (88) to analyze the annular and central detector configuration. For a small axial detector, we see immediately that

$$j_d(\boldsymbol{\xi}) \propto \left| \int \tilde{S}(\mathbf{x}_o) K(\boldsymbol{\xi} - \mathbf{x}_o) d\mathbf{x}_o \right|^2 \tag{92}$$

which is very similar to the image observed in a TEM. For an annular detector, we divide $S(\mathbf{x}_o)$ into an unscattered and a scattered part, $S(\mathbf{x}_o) = 1 + \sigma_s(\mathbf{x}_o)$. The signal consists of two main contributions, one of the form $\int [\sigma_s(\mathbf{x}_o) + \sigma_s^*(\mathbf{x}_o)] |K(\boldsymbol{\xi} - \mathbf{x}_o)|^2 d\mathbf{x}_o$, and the other $\int |\sigma_s(\mathbf{x}_o)|^2 |K(\boldsymbol{\xi} - \mathbf{x}_o)|^2 d\mathbf{x}_o$. The latter term usually dominates.

We have seen that the current distribution in the detector plane at any instant is the far-field diffraction pattern

of the object element currently illuminated. The fact that we have direct access to this wealth of information about the specimen is one of the remarkable and attractive features of the STEM, rendering possible imaging modes that present insuperable difficulties in the TEM. The simple detectors so far described hardly exploit this wealth of information at all, since only two total currents are measured, one falling on the central region, the other on the annular detector. A slightly more complicated geometry permits us to extract directly information about the phase variation $\varphi(\mathbf{x}_o)$ of the specimen transparency $S(\mathbf{x}_o)$. Here the detector is divided into four quadrants, and by forming appropriate linear combinations of the four signals thus generated, the gradient of the phase variation can be displayed immediately. This technique has been used to study the magnetic fields across domain boundaries in magnetic materials.

Other detector geometries have been proposed, and it is of interest that it is not necessary to equip the microscope with a host of different detectors, provided that the instrument has been interfaced to a computer. It is one of the conveniences of all scanning systems that the signal that serves to generate the image is produced sequentially and can therefore be dispatched directly to computer memory for subsequent or on-line processing if required. By forming the far-field diffraction pattern not on a single large detector but on a honeycomb of very small detectors and reading the currents detected by each cell into framestore memory, complete information about each elementary object area can be recorded. Framestore memory can be programmed to perform simple arithmetic operations, and the framestore can thus be instructed to multiply the incoming intensity data by 1 or 0 in such a way as to mimic any desired detector geometry. The signals from connected regions of the detector—quadrants, for example—are then added, and the total signal on each part is then stored, after which the operation is repeated for the next object element under the probe. Alternatively, the image of each elementary object area can be exploited to extract information about the phase and amplitude of the electron wave emerging from the specimen.

A STEM imaging mode that is capable of furnishing very high resolution images has largely superseded the modes described above. Electrons scattered through relatively wide angles (Rutherford scattering) and collected by an annular detector with appropriate dimensions form an “incoherent” image of the specimen structure, but with phase information converted into amplitude variations in the image. Atomic columns can be made visible by this technique, which is rapidly gaining importance.

The effect of partial coherence in the STEM can be analyzed by a reasoning similar to that followed for the TEM; we will not reproduce this here.

C. Image Processing

1. Interference and Holography

The resolution of electron lenses is, as we have seen, limited by the spherical aberration of the objective lens, and many types of correctors have been devised in the hope of overcoming this limit. It was realized by Dennis Gabor in the late 1940s, however, that although image detail beyond the limit cannot be discerned by eye, the information is still there if only we could retrieve it. The method he proposed for doing this was holography, but it was many years before his idea could be successfully put into practice; this had to await the invention of the laser and the development of high-performance electron microscopes. With the detailed understanding of electron image formation, the intimate connection between electron holography, electron interference, and transfer theory has become much clearer, largely thanks to Karl-Joseph Hanszen and colleagues in Braunschweig. The electron analogs of the principal holographic modes have been thoroughly explored with the aid of the Möllenstedt biprism. In the hands of Akira Tonomura in Tokyo and Hannes Lichte in Tübingen, electron holography has become a tool of practical importance.

The simplest type of hologram is the Fraunhofer in-line hologram, which is none other than a defocused electron image. Successful reconstruction requires a very coherent source (a field-emission gun) and, if the reconstruction is performed light-optically rather than digitally, glass lenses with immense spherical aberration. Such holograms should permit high-contrast detection of small weak objects.

The next degree of complexity is the single-sideband hologram, which is a defocused micrograph obtained with

half of the diffraction pattern plane obscured. From the two complementary holograms obtained by obscuring each half in turn, separate phase and amplitude reconstruction is, in principle, possible. Unfortunately, this procedure is extremely difficult to put into practice, because charge accumulates along the edge of the plane that cuts off half the aperture and severely distorts the wave fronts in its vicinity; compensation is possible, but the usefulness of the technique is much reduced.

In view of these comments, it is not surprising that off-axis holography, in which a true reference wave interferes with the image wave in the recording plane, has completely supplanted these earlier arrangements. In the in-line methods, the reference wave is, of course, to be identified with the unscattered part of the main beam. [Figure 28](#) shows an arrangement suitable for obtaining the hologram; the reference wave and image wave are recombined by the electrostatic counterpart of a biprism. In the reconstruction step, a reference wave must again be suitably combined with the wave field generated by the hologram, and the most suitable arrangement has been found to be that of the Mach-Zehnder interferometer. Many spectacular results have been obtained in this way, largely thanks to the various interference techniques developed by the school of A. Tonomura and the Bolognese group. Here, the reconstructed image is made to interfere with a plane wave. The two may be exactly aligned and yield an interference pattern representing the magnetic field in the specimen, for example; often, however, it is preferable to arrange that they are slightly inclined with respect to one another since phase “valleys” can then be distinguished from “hills.” In another arrangement, the twin images are made to interfere, thereby amplifying the corresponding phase shifts twofold (or more, if higher order diffracted beams are employed).

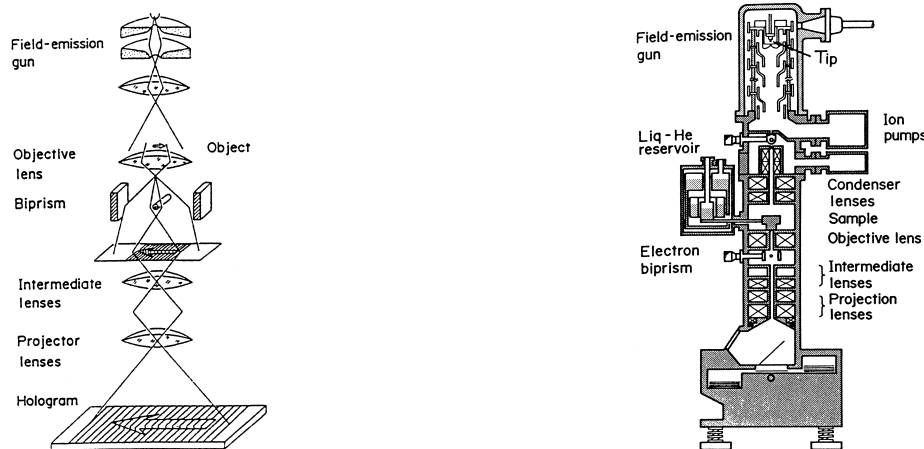


FIGURE 28 (Left) Ray diagram showing how an electron hologram is formed. (Right) Cross-section of an electron microscope equipped for holography. [From Tonomura, A. (1999). "Electron Holography," Springer-Verlag, Berlin/New York.]

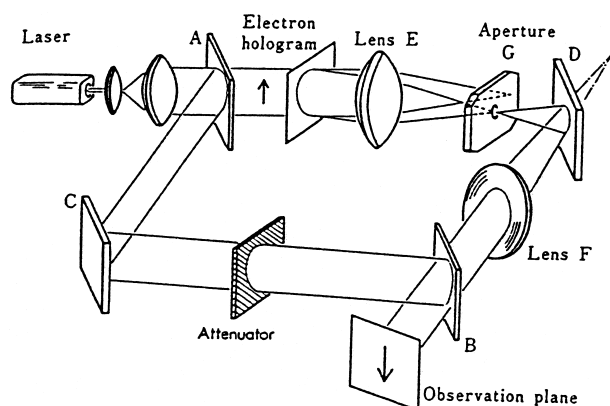


FIGURE 29 Arrangement of lenses and mirrors suitable for interference microscopy. [From Tonomura A. (1999). "Electron Holography," Springer-Verlag, Berlin/New York.]

Electron holography has a great many ramifications, which we cannot describe here, but we repeat that many of the problems that arise in the reconstruction step vanish if the hologram is available in digital form and can hence be processed in a computer. We now examine the related techniques, although not specifically in connection with holography.

2. Digital Processing

If we can sample and measure the gray levels of the electron image accurately and reliably, we can employ the computer to process the resulting matrix of image gray-level measurements in many ways. The simplest techniques, usually known as image enhancement, help to adapt the image to the visual response or to highlight features of particular interest. Many of these are routinely available on commercial scanning microscopes, and we will say no more about them here. The class of methods that allow image restoration to be achieved offer solutions of more difficult problems. Restoration filters, for example, reduce the adverse effect of the transfer functions of Eq. (81). Here, we record two or more images with different values of the defocus and hence with different forms of the transfer function and seek the weighted linear combinations of these images, or rather of their spatial frequency spectra, that yield the best estimates (in the least-squares sense) of $\tilde{\varphi}$ and \tilde{s} . By using a focal series of such images, we can both cancel, or at least substantially reduce, the effect of the transfer functions $\sin \gamma$ and $\cos \gamma$ and fill in the information missing from each individual picture around the zeros of these functions.

Another problem of considerable interest, in other fields as well as in electron optics, concerns the phase of the object wave for strongly scattering objects. We have seen that the specimens studied in transmission microscopy

are essentially transparent: The image is formed not by absorption but by scattering. The information about the specimen is therefore in some sense coded in the angular distribution of the electron trajectories emerging from the specimen. In an ideal system, this angular distribution would be preserved, apart from magnification effects, at the image and no contrast would be seen. Fortunately, however, the microscope is imperfect; contrast is generated by the loss of electrons scattered through large angles and intercepted by the diaphragm or "objective aperture" and by the crude analog of a phase plate provided by the combination of spherical aberration and defocus. It is the fact that the latter affects the angular distribution within the beam and converts it to a positional dependence with a fidelity that is measured by the transfer function $\sin \gamma$ that is important. The resulting contrast can be related simply to the specimen transparency only if the phase and amplitude variations are small, however, and this is true of only a tiny class of specimens. For many of the remainder, the problem remains. It can be expressed graphically by saying that we know from our intensity record where the electrons arrive (amplitude) but not their directions of motion at the point of arrival (phase). Several ways of obtaining this missing information have been proposed, many inspired by the earliest suggestion, the Gerchberg–Saxton algorithm. Here, the image and diffraction pattern of exactly the same area are recorded, and the fact that the corresponding wave functions are related by a Fourier transform is used to find the phase iteratively. First, the known amplitudes in the image, say, are given arbitrary phases and the Fourier transform is calculated; the amplitudes thus found are then replaced by the known diffraction pattern amplitudes and the process is repeated. After several iterations, the unknown phases should be recovered. This procedure encounters many practical difficulties and some theoretical ones as well, since the effect of noise is difficult to incorporate. This and several related algorithms have now been thoroughly studied and their reliability is well understood. In these iterative procedures, two signals generated by the object are required (image and diffraction pattern or two images at different defocus values in particular). If a STEM is used, this multiplicity of information is available in a single record if the intensity distribution associated with every object pixel is recorded and not reduced to one or a few summed values. A sequence of Fourier transforms and mask operations that generate the phase and amplitude of the electron wave has been devised by John Rodenburg.

A very different group of methods has grown up around the problem of three-dimensional reconstruction. The two-dimensional projected image that we see in the microscope often gives very little idea of the complex spatial relationships of the true structure, and techniques have therefore

been developed for reconstructing the latter. They consist essentially in combining the information provided by several different views of the specimen, supplemented if possible by prior knowledge of an intrinsic symmetry of the structure. The fact that several views are required reminds us that not all specimens can withstand the electron onslaught that such multiple exposure represents. Indeed, there are interesting specimens that cannot be directly observed at all, because they are destroyed by the electron dose that would be needed to form a discernible image. Very low dose imaging must therefore be employed, and this has led to the development of an additional class of image restoration methods. Here, the aim is first to detect the structures, invisible to the unaided eye, and superimpose low-dose images of identical structures in such a way that the signal increases more rapidly than the noise and so gradually emerges from the surrounding fog. Three-dimensional reconstruction may then be the next step. The problem here, therefore, is first to find the structures, then to align them in position and orientation with the precision needed to achieve the desired resolution. Some statistical check must be applied to be sure that all the structures found are indeed the same and not members of distinct groups that bear a resemblance to one other but are not identical. Finally, individual members of the same group are superposed. Each step demands a different treatment. The individual structures are first found by elaborate cross-correlation calculations. Cross-correlation likewise enables us to align them with high precision. Multivariate analysis is then used to classify them into groups or to prove that they do, after all, belong to the same group and, a very important point, to assign probabilities to their membership of a particular group.

IV. CONCLUDING REMARKS

Charged-particle optics has never remained stationary with the times, but the greatest upheaval has certainly been that caused by the widespread availability of large, fast computers. Before, the analysis of electron lenses relied heavily on rather simple field or potential models, and much ingenuity was devoted to finding models that were at once physically realistic and mathematically tractable. Apart from sets of measurements, guns were almost virgin territory. The analysis of in-lens deflectors would have been unthinkable but fortunately was not indispensable since even the word *microminiaturization* has not yet been coined. Today, it is possible to predict with great accuracy the behavior of almost any system; it is even possible to obtain aberration coefficients, not by evaluating the corresponding integrals, themselves obtained as a result of exceedingly long and tedious algebra, but by solving the

exact ray equations and fitting the results to the known aberration pattern. This is particularly valuable when parasitic aberrations, for which aberration integrals are not much help, are being studied. Moreover, the aberration integrals can themselves now be established not by long hours of laborious calculation, but by means of one of the computer algebra languages. A knowledge of the fundamentals of the subject, presented here, will always be necessary for students of the subject, but modern numerical methods now allow them to go as deeply as they wish into the properties of the most complex systems.

SEE ALSO THE FOLLOWING ARTICLES

ACCELERATOR PHYSICS AND ENGINEERING • HOLOGRAPHY • QUANTUM OPTICS • SCANNING ELECTRON MICROSCOPY • SCATTERING AND RECOILING SPECTROSCOPY • SIGNAL PROCESSING, DIGITAL • WAVE PHENOMENA

BIBLIOGRAPHY

- Carey, D. C. (1987). "The Optics of Charged Particle Beams," Harwood Academic, London.
- Chapman, J. N., and Craven, A. J., eds. (1984). "Quantitative Electron Microscopy," SUSSP, Edinburgh.
- Dragt, A. J., and Forest, E. (1986). *Adv. Electron. Electron. Phys.* **67**, 65–120.
- Feinerman, A. D. and Crewe, D. A. (1998). "Miniature electron optics." *Adv. Imaging Electron Phys.* **102**, 187–234.
- Frank, J. (1996). "Three-Dimensional Electron Microscopy of Macromolecular Assemblies," Academic Press, San Diego.
- Glaser, W. (1952). "Grundlagen der Elektronenoptik," Springer-Verlag, Vienna.
- Glaser, W. (1956). Elektronen- und Ionenoptik, *Handb. Phys.* **33**, 123–395.
- Grivet, P. (1972). "Electron Optics," 2nd Ed. Pergamon, Oxford.
- Hawkes, P. W. (1970). *Adv. Electron. Electron Phys.*, Suppl. 7. Academic Press, New York.
- Hawkes, P. W., ed. (1973). "Image Processing and Computer-Aided Design in Electron Optics," Academic Press, New York.
- Hawkes, P. W., ed. (1980). "Computer Processing of Electron Microscope Images," Springer-Verlag, Berlin and New York.
- Hawkes, P. W., ed. (1982). "Magnetic Electron Lenses," Springer-Verlag, Berlin and New York.
- Hawkes, P. W., and Kasper, E. (1989, 1994). "Principles of Electron Optics," Academic Press, San Diego.
- Hawkes, P. W., ed. (1994). "Selected Papers on Electron Optics," SPIE Milestones Series, Vol. 94.
- Humphries, S. (1986). "Principles of Charged Particle Acceleration." (1990). "Charged Particle Beams," Wiley-Interscience, New York and Chichester.
- Lawson, J. D. (1988). "The Physics of Charged-Particle Beams," Oxford Univ. Press, Oxford.
- Lencová, B. (1997). Electrostatic Lenses. In "Handbook of Charged Particle Optics" (J. Orloff, ed.), pp. 177–221, CRC Press, Boca Raton, FL.

- Livingood, J. J. (1969). "The Optics of Dipole Magnets," Academic Press, New York.
- Orloff, J., ed. (1997). "Handbook of Charged Particle Optics," CRC Press, Boca Raton, FL.
- Reimer, L. (1997). "Transmission Electron Microscopy," Springer-Verlag, Berlin and New York.
- Reimer, L. (1998). "Scanning Electron Microscopy," Springer-Verlag, Berlin and New York.
- Saxton, W. O. (1978). *Adv. Electron. Electron Phys., Suppl.* 10. Academic Press, New York.
- Septier, A., ed. (1967). "Focusing of Charged Particles," Vols. 1 and 2, Academic Press, New York.
- Septier, A., ed. (1980–1983). *Adv. Electron. Electron Phys., Suppl.* **13A–C**. Academic Press, New York.
- Tonomura, A. (1999). "Electron Holography," Springer-Verlag, Berlin.
- Tsuno, K. (1997). Magnetic Lenses for Electron Microscopy. In "Handbook of Charged Particle Optics" (J. Orloff, ed.), pp. 143–175, CRC Press, Boca Raton, FL.
- Wollnik, H. (1987). "Optics of Charged Particles," Academic Press, Orlando.



Diffractive Optical Components

R. Magnusson

University of Connecticut

D. Shin

University of Texas at Arlington

- I. Introduction
- II. Grating Diffraction Models
- III. Applications
- IV. Fabrication

GLOSSARY

Bragg condition A condition for constructive wave addition in a periodic lattice.

Diffraction efficiency A measure of the power density (intensity) of a diffracted wave normalized to that of the incident wave.

Diffraction grating A structure consisting of a periodic spatial variation of the dielectric constant, conductivity, or surface profile.

Evanescence wave An exponentially decaying, nonpropagating wave.

Grating equation An expression that provides the directions θ_i of the diffracted waves produced by a plane wave with wavelength λ_0 incident at θ_{in} on a periodic structure with period Λ . For a grating in air it is $\sin \theta_i = \sin \theta_{in} - i\lambda_0/\Lambda$ where i is an integer labeling the diffraction orders.

Grating fill factor Fraction of the grating period occupied by the high-refractive-index material, n_H .

Grating vector A vector, analogous to the wave vector \mathbf{k} normal to the grating fringes with magnitude $K = 2\pi/\Lambda$, where Λ is the grating period.

Monochromatic plane wave A plane wave with single frequency and wavelength.

Planar grating A spatially modulated periodic structure with the modulated region confined between parallel planes.

Plane of incidence A plane containing the wave vectors of the incident and diffracted waves.

Plane wave An electromagnetic wave with infinite transverse dimensions. It is an idealized model of a light beam where the finite transverse extent of a real beam is ignored for mathematical expediency; this works well if the beam diameter is large on the scale of the wavelength. The planes of constant phase are parallel everywhere.

Reflection grating A grating where the diffracted wave of interest emerges in the space of the cover material.

Surface-relief grating A periodic corrugated surface.

TE (TM) polarized optical wave An electromagnetic wave with its electric (magnetic) vector normal to the plane of incidence.

Transmission grating A grating where the diffracted wave of interest emerges in the space of the substrate material.

Wave vector A vector normal to the planes of constant phase associated with a plane wave. Its magnitude is $k_0 = 2\pi/\lambda_0$, where λ_0 is the wavelength.

DIFFRACTIVE OPTICAL COMPONENTS consist of fine spatial patterns arranged to control propagation of light. Lithographic patterning of dielectric surfaces, layers, or volume regions yields low-loss structures that affect the spatial distribution, spectral content, energy content, polarization state, and propagation direction of an optical wave. Applications include spectral filters, diffractive lenses, antireflection surfaces, beam splitters, beam steering elements, laser mirrors, polarization devices, beam-shaping elements, couplers, and switches. These components are widely used in lasers, fiber-optic communication systems, spectroscopy, integrated optics, imaging, and other optical systems.

I. INTRODUCTION

A. Overview

Application of fabrication methods similar to those used in the semiconductor industry enables optical components with features on the order of and smaller than the wavelength of light. Diffractive optical components possessing multiple phase levels in glass, for example, require sequential high-precision mask alignment and etching steps. Such microscopic surface-relief structures can effectively transform the phasefront of an incident optical beam to redirect or focus it. These components are compact, lightweight, and replicable in many cases. Diffractive optical elements can be used independently or in conjunction with conventional refractive elements for design flexibility and enhanced functionality. Progress

in microfabrication methods, including nanotechnology, will undoubtedly further advance this field. For example, complex three-dimensional diffractive components called photonic crystals are under rapid development. Additionally, progress in modeling and design methodologies (for example, finite-difference time domain analysis and genetic algorithm optimization methods) is occurring in parallel.

This article addresses key aspects of diffractive optics. Common analytical models are described and their main results summarized. Exact numerical methods are applied when precise characterization of the periodic component is required, whereas approximate models provide analytical results that are valuable for preliminary design and improved physical insight. Numerous examples of the applications of diffractive optical components are presented. These are optical interconnects, diffractive lenses, and subwavelength elements including antireflection surfaces, polarization devices, distributed-index components, and resonant filters. Finally, recording of gratings by laser interference is presented and an example fabrication process summarized.

B. The Generic Problem

The diffraction grating is the fundamental building block of diffractive optical components. These may be planar or surface-relief gratings made of dielectrics or metals. Figure 1 illustrates a single-layer diffraction grating and key descriptive parameters including the grating period, Λ , the grating fill factor, f , the grating refractive indices (n_H, n_L), thickness, d , and refractive indices of the cover and substrate media (n_C, n_S). The grating shown has a rectangular profile; if $f = 0.5$ the profile is said to be square. An incident plane wave at angle θ_{in} is dispersed into several diffracted waves (diffraction orders labeled by the integer i as shown) propagating both forwards and backwards.

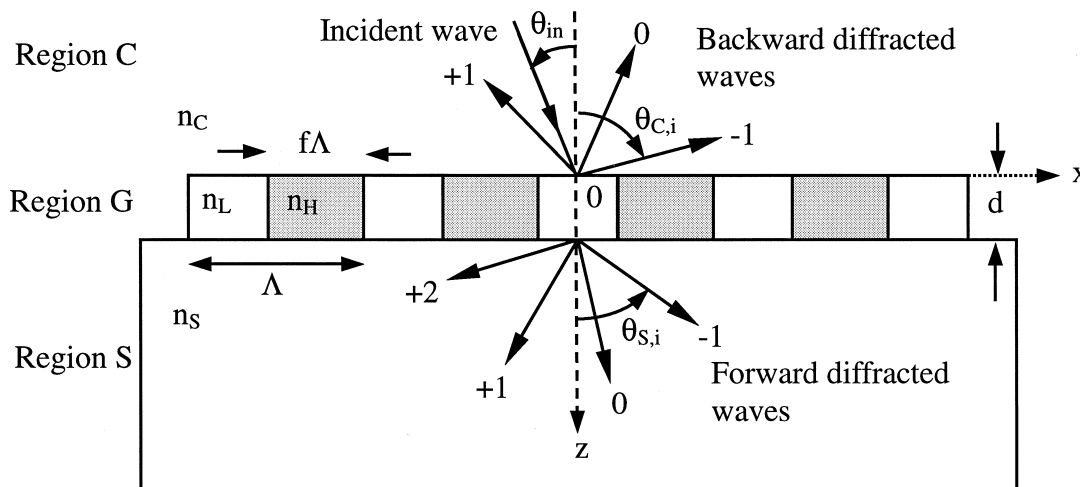


FIGURE 1 Geometry of diffraction by a rectangular grating.

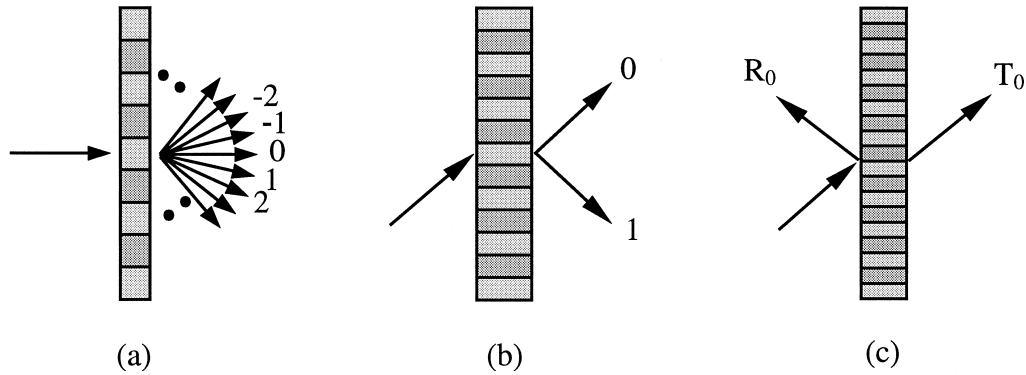


FIGURE 2 Operating regimes of transmissive diffraction gratings: (a) multiwave regime, (b) two-wave regime, and (c) zero-order regime.

The propagation angle of each diffracted wave may be obtained from the grating equation given by

$$n_P \sin \theta_{P,i} = n_C \sin \theta_{in} - i \lambda_0 / \Lambda \quad (1)$$

where λ_0 is the free-space wavelength, i is an integer, and P is either C or S depending on the region under consideration.

Exact electromagnetic analysis can be applied to find the intensity of each diffracted wave shown in Fig. 1. In such calculations, which provide purely numerical results due to the inherent complexity of the full problem, all the waves indicated as well as evanescent waves (not shown) must be included. Simplifying assumptions yield approximate results that are exceedingly useful in the understanding and design of diffractive optical components. Design procedures may additionally include numerical optimization techniques such as simulated annealing or genetic algorithms as well as iterative Fourier transform methods.

C. Operating Regimes

Figure 2 defines the operating regimes of transmissive diffractive elements. Figure 2(a) illustrates the multiwave regime where multiple output diffracted waves are gen-

erated by a grating with a large period on the scale of the wavelength. Note that reflections at refractive-index discontinuities as well as waves propagating in the $-z$ direction (generated by diffraction) are ignored in this figure. In Fig. 2(b), under Bragg conditions, only two main waves are assumed to be dominant with reflected waves neglected. Figure 2(c) defines the zero-order regime in which the grating period is smaller than the wavelength such that all higher diffracted orders are cut off (evanescent). Each of these cases can be modeled approximately resulting in analytical solutions for the diffraction efficiencies, η , associated with the various orders.

Figure 3 provides examples of reflection-type diffractive elements. Typically, a single dominant Bragg-diffracted wave prevails as shown. The waveguide reflection grating is an important element in integrated optical devices and in fiber-optical communication systems for wavelength division multiplexing and other applications. In a waveguide geometry, the incident and diffracted waves are treated as guided modes (not infinite plane waves). Figure 4 illustrates the filtering property of a bulk reflection grating (distributed Bragg mirror such as used for cavity reflection in vertical-cavity lasers); note the efficient reflection within a specific wavelength range.

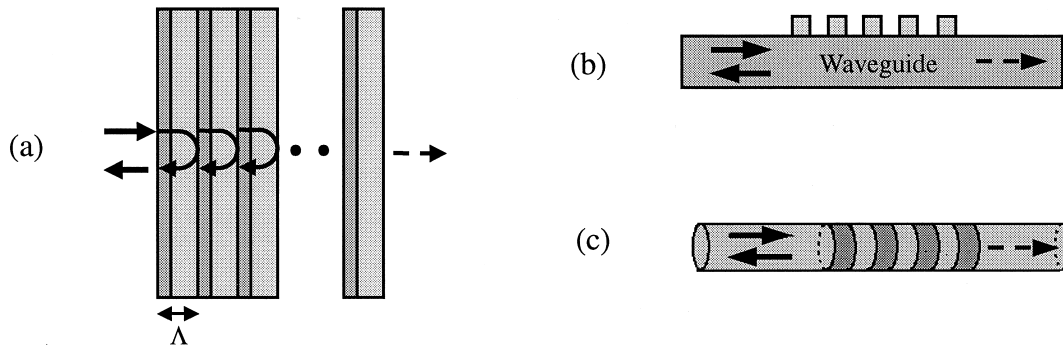


FIGURE 3 Examples of reflection gratings: (a) bulk grating, (b) waveguide grating, and (c) fiber grating.

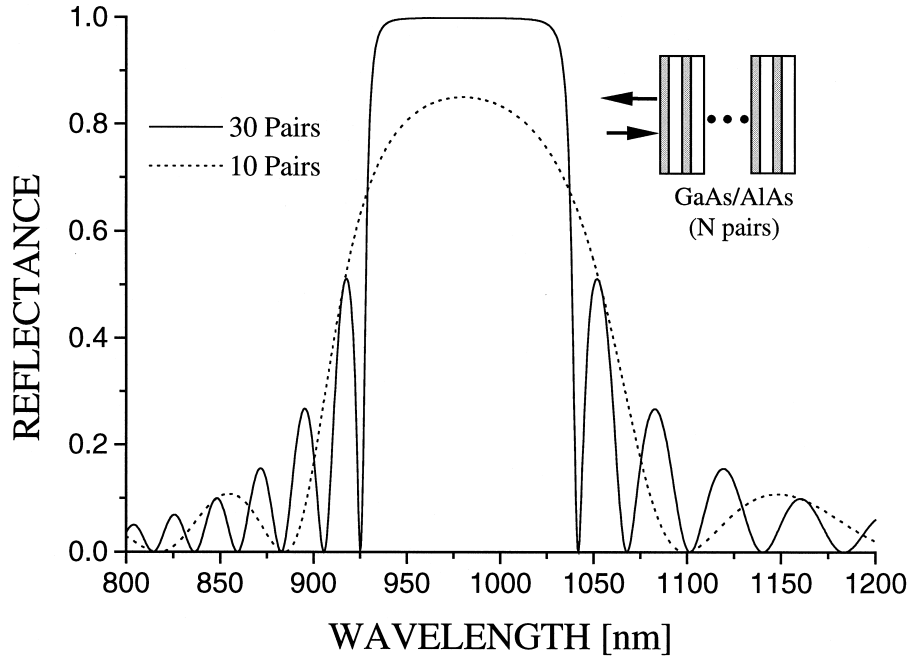


FIGURE 4 Spectral reflectance of a distributed Bragg reflector consisting of 10 and 30 pairs of alternating GaAs/AlAs layers with quarter-wavelength thicknesses at $\lambda_0 = 980$ nm. The cover region material is $\text{Al}_{0.7}\text{Ga}_{0.3}\text{As}$ and the substrate is GaAs. Refractive indices used in the simulation are $n_{\text{GaAs}} = 3.512$, $n_{\text{AlGaAs}} = 3.151$, and $n_{\text{AlAs}} = 3.007$.

D. Application Example: Monochromator

Diffraction gratings are commonly used in spectroscopy. For example, a typical monochromator is schematically illustrated in Fig. 5. For polychromatic light incident on a reflection grating as shown, each component wavelength of a specific diffracted order is spatially (angularly) sepa-

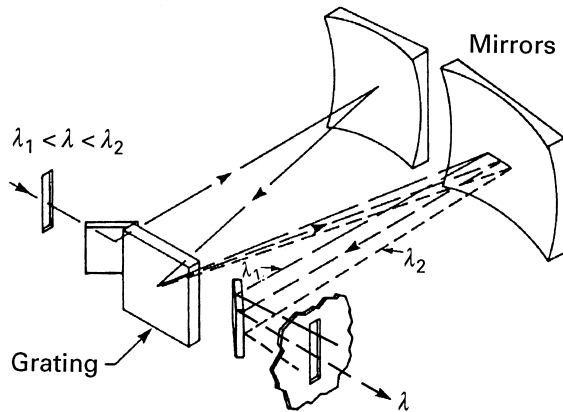


FIGURE 5 The basic elements of a monochromator. The input spectrum contains wavelengths in the range $\lambda_1 < \lambda < \lambda_2$. The grating spreads the light within the diffraction order (say $i = 1$) in which the instrument is operated. No other diffracted orders are shown in the drawing for clarity. The width of the slit defines the spectral contents of the output light. [From Oriel Corporation (1989). "Oriel 1989 Catalog," Volume II.]

rated according to the grating equation. The angular spread (dispersion) per unit spectral interval is given, differentiating the grating equation for a constant angle of incidence, by

$$\frac{d\theta}{d\lambda} = \left| \frac{i}{\Lambda \cos \theta_{C,i}} \right| \quad (2)$$

The exit slit controls the spectral content of the output light by blocking the undesirable wavelength components around the central one. A similar instrument is the spectrograph where the slit is replaced by an array of detectors. The output spectrum is picked up by the array such that a prescribed spectral band can be associated with a given detector element.

II. GRATING DIFFRACTION MODELS

A. Rigorous Coupled-Wave Analysis

Stimulated by advances in digital computers, various grating-diffraction analysis techniques, which solve the electromagnetic Maxwell's equations numerically either in differential or integral form, have been developed. Rigorous coupled-wave analysis (RCWA) is widely applied; this method is briefly summarized herein using the diffraction geometry given in Fig. 1.

An incident TE-polarized plane wave may be expressed as

$$E_{y,in}(x, z) = \exp[-jk_0 n_C (x \sin \theta_{in} + z \cos \theta_{in})] \quad (3)$$

where $k_0 = 2\pi/\lambda_0$ and $j = \sqrt{-1}$. In the homogeneous regions (regions C and S) adjacent to the grating layer, the diffracted fields may be expressed as superpositions of plane waves (so-called Rayleigh expansion) and, thus, the total electric fields in each region are given by

$$E_y(x, z < 0) = E_{y,in}(x, z) + \sum_{i=-\infty}^{\infty} r_i \exp[-j(k_{x,i}x - k_{z,C,i}z)] \quad (4)$$

$$E_y(x, z > d) = \sum_{i=-\infty}^{\infty} t_i \exp[-j\{k_{x,i}x + k_{z,S,i}(z - d)\}], \quad (5)$$

where r_i and t_i are amplitudes of i -th order backward- and forward-diffracted waves, respectively. The x component of the diffracted wave vector, $k_{x,i}$, is given, from the Floquet condition, by

$$k_{x,i} = k_0(n_C \sin \theta_{in} - i\lambda_0/\Lambda) \quad (6)$$

and the z component $k_{z,P,i}$ is

$$k_{z,P,i} = \begin{cases} (k_0^2 n_P^2 - k_{x,i}^2)^{1/2}, & k_0^2 n_P^2 > k_{x,i}^2 \\ -j(k_{x,i}^2 - k_0^2 n_P^2)^{1/2}, & k_0^2 n_P^2 < k_{x,i}^2 \end{cases}, \quad \text{where } P = C \text{ or } S, \quad (7)$$

where the real values of $k_{z,P,i}$ correspond to propagating waves and the imaginary values correspond to evanescent waves.

Inside the grating region ($0 \leq z \leq d$), the total electric field may be expressed as a coupled-wave expansion

$$E_y(x, 0 \leq z \leq d) = \sum_{i=-\infty}^{\infty} S_i(z) \exp(-jk_{x,i}x), \quad (8)$$

where $S_i(z)$ is the complex amplitude of i -th diffracted wave. The total field satisfies the Helmholtz equation

$$\frac{d^2 E_y(x, z)}{dx^2} + \frac{d^2 E_y(x, z)}{dz^2} + k_0^2 \varepsilon(x) E_y(x, z) = 0, \quad (9)$$

where $\varepsilon(x)$ is the periodic relative permittivity (dielectric constant) modulation given by

$$\varepsilon(x) = \begin{cases} n_L^2, & 0 \leq x \leq f\Lambda \\ n_H^2, & \text{otherwise} \end{cases}. \quad (10)$$

Due to the periodicity, $\varepsilon(x)$ can be expressed as a Fourier series

$$\varepsilon(x) = \sum_{h=-\infty}^{\infty} \varepsilon_h \exp(jhKx), \quad (11)$$

where ε_h is the h -th Fourier harmonic coefficient and $K = 2\pi/\Lambda$ is the amplitude of grating vector. Inserting Eqs. (8) and (11) into the Helmholtz Equation (9) results in an infinite set of coupled-wave equations given by

$$\frac{d^2 S_i(z)}{dz^2} = k_{xi}^2 S_i(z) - k_0^2 \sum_{h'=-\infty}^{\infty} \varepsilon_{i-h'} S_{h'}(z) \quad (12)$$

where h' is an integer. The coupled-wave Eq. (12) can be cast into an eigenvalue problem as $d^2 S_i(z)/dz^2 = q^2 S_i(z)$. Numerical solution involves a truncated ($N \times N$) matrix equation $[\ddot{S}_i(z)] = [A][S_i(z)]$ with solutions of the form

$$S_i(z) = \sum_{m=1}^N W_{i,m} [C_m^+ \exp(-q_m z) + C_m^- \exp\{q_m(z - d)\}], \quad (13)$$

where $W_{i,m}$ and q_m are the eigenvectors and positive square roots of the eigenvalues of matrix $[A]$, respectively. The C_m^+ and C_m^- are unknown constants to be determined by the boundary conditions. From the total electric field Eqs. (4), (5), and (8) with (13), the tangential magnetic fields (x component) can be obtained by Maxwell's curl equation

$$H_x(x, z) = \frac{1}{j\omega\mu_0} \frac{dE_y(x, z)}{dz}, \quad (14)$$

where ω is the frequency and μ_0 is the permeability of free space. Finally, the unknown quantities of interest r_i and t_i as well as C_m^+ and C_m^- can be obtained by solving the boundary-condition equations obtained by matching the tangential electric and magnetic fields at boundaries $z = 0$ and $z = d$. The diffraction efficiencies of backward (C) and forward (S) diffracted waves are defined as

$$\eta_{r,i} = \text{Re} \left(\frac{k_{z,C,i}}{k_{z,C,0}} \right) |r_i|^2 \quad (15)$$

$$\eta_{t,i} = \text{Re} \left(\frac{k_{z,S,i}}{k_{z,C,0}} \right) |t_i|^2. \quad (16)$$

For lossless gratings, energy conservation yields $\sum_{i=-\infty}^{\infty} (\eta_{r,i} + \eta_{t,i}) = 1$.

The accuracy of the RCWA method depends on the number of space harmonics (N) retained in the truncation of matrices. An arbitrary grating shape with a dielectric constant modulation $\varepsilon(x, z)$ may be implemented by replacing the grating layer with a stack of digitized rectangular grating layers with $\varepsilon(x)$ only. Rigorous coupled-wave analysis is a noniterative, stable numerical technique

that is relatively straightforward in implementation. It has been successfully applied to planar and surface-relief gratings (with arbitrary grating profiles), dielectric (lossless and lossy) and metallic gratings, transmission and reflection gratings, isotropic and anisotropic gratings, multiplexed gratings, phase-shifted gratings, and photonic crystals. Extensions to include transversely finite beams (such as Gaussian laser beams) have also been implemented as found in the literature.

The calculated results presented in Figs. 4, 14, 16, and 19 were generated with RCWA.

B. Transmittance Diffraction Theory

The diffraction efficiencies of transmission gratings with low spatial frequency ($\Lambda \gg \lambda$) may be estimated using formulation based on the transmittance function. In this region, sometimes called the Raman–Nath diffraction regime, multiple diffracted output waves are typically observed and polarization effects are weak. The pertinent model is illustrated in Fig. 6. A unit-amplitude plane wave is incident upon a low-spatial-frequency grating and is diffracted into multiple forward-propagating waves. Backward propagating diffracted waves are ignored. The transmittance of the grating may be expressed as the convolution

$$T(x) = \tau(x) * \sum_{h=-\infty}^{\infty} \delta(x - h\Lambda), \quad (17)$$

where $\tau(x)$ is the transmittance of each period of the grating and h is an integer. Since $T(x)$ is a periodic function, it may be expanded into a Fourier series as

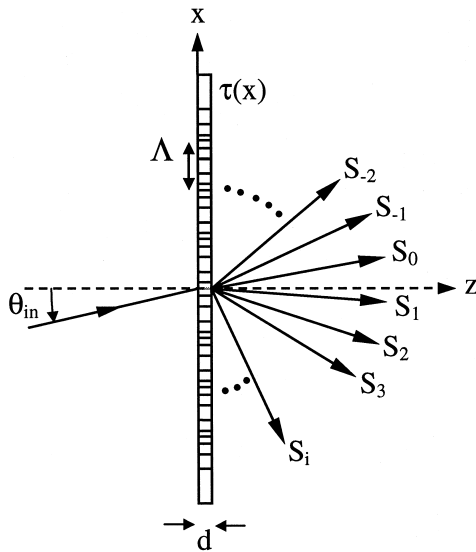


FIGURE 6 Transmittance theory model.

$$T(x) = \sum_{i=-\infty}^{\infty} S_i \exp(jiKx). \quad (18)$$

In the transmittance approach, the emerging diffracted field is well approximated by the product of the input field and the transmittance of the structure. The Fourier coefficients, S_i , are thus identified as the amplitudes of the diffracted waves, which can be calculated by

$$S_i = \frac{1}{\Lambda} \int_{-\Lambda/2}^{\Lambda/2} \tau(x) \exp(-jiKx) dx. \quad (19)$$

The diffraction efficiency of i -th diffracted order is defined as

$$\eta_i = |S_i|^2 \quad (20)$$

with the corresponding power conservation law being $\sum_{i=-\infty}^{\infty} \eta_i = 1$ for lossless structures.

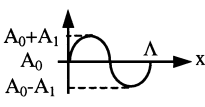
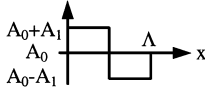
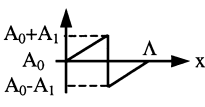
Representative results including analytical expressions of diffraction efficiencies are given in Tables Ia and Ib for common planar and surface-relief grating shapes. Absorption and amplitude transmittance gratings exhibit much lower efficiencies than phase gratings. Diffraction efficiency depends strongly on the shape of grating modulation. In particular, the phase gratings with a sawtooth (or blazed) profile (both planar and surface-relief type) can transfer the input energy completely to a particular diffracted order.

The array of dielectric cylinders shown in the inset of Fig. 7 is an example of a low-spatial-frequency grating. In the figure, diffraction-efficiency measurements are plotted for structures with two different cylinder diameters (equal to the grating period) 20 and 50 μm under a TE-polarized normally incident HeNe laser ($\lambda_0 = 632.8 \text{ nm}$) illumination. This structure generates multiple forward-propagating diffracted orders with similar diffraction efficiencies, functioning as a fan-out element in optical interconnections or as a multiple imaging device. In this experiment, 63 and 159 forward-propagating orders were observed for these gratings with $\Lambda = 20$ and 50 μm , respectively, in agreement with the number of propagating orders (N), predicted by the grating equation, $N = 2\Lambda/\lambda_0 + 1$. The calculated results in Fig. 7 found by

$$S_i = \int_0^1 \exp[-j2g(1-y^2)^{1/2}] \cos(i\pi y) dy, \quad (21)$$

where $y = x/\Lambda$ and $g = \pi \Lambda n d / \lambda_0$ show reasonably good agreement with experimental results. The agreement improves as λ_0/Λ becomes smaller. Results experimentally obtained with TM-polarized incident waves show a small but observable difference ($\sim 8\%$ higher efficiency in the region near $i = 0$ for $\Lambda = 50 \mu\text{m}$) with respect to results for TE-polarized light.

TABLE Ia Summary of Diffraction Efficiency Results from Transmittance Theory (Planar Gratings)

Grating type	Phase gratings	Absorption gratings	Transmittance gratings
$\tau(x)$	$\exp[-jk_0 \Delta n(x) d]$	$\exp[-jk_0 \{\alpha_0 + \Delta \alpha(x)\} d]$	$\tau_0 + \Delta \tau(x)$
Sinusoidal 	$J_i^2(2\gamma)$ (33.8%)	$\exp[-2\psi_0] I_i^2(\psi_1)$ (4.80%)	$\begin{cases} \tau_0^2, & i = 0 \\ (\tau_1/2)^2, & i = \pm 1 \\ 0, & i \neq \pm 1 \end{cases}$ (6.25%)
Square 	$\begin{cases} \cos^2(2\gamma), & i = 0 \\ 0, & i = \text{even} \\ (2/i\pi)^2 \sin^2(2\gamma), & i = \text{odd} \end{cases}$ (40.5%)	$\begin{cases} \exp[-2\psi_0] \cosh^2(\psi_1), & i = 0 \\ 0, & i = \text{even} \\ \exp[-2\psi_0] (2/i\pi)^2 \sinh^2(\psi_1), & i = \text{odd} \end{cases}$ (10.13%)	$\begin{cases} \tau_0^2, & i = 0 \\ 0, & i = \text{even} \\ (2\tau_1/\pi i)^2, & i = \text{odd} \end{cases}$ (10.13%)
Sawtooth 	$[2\gamma + i\pi]^{-2} \sin^2(2\gamma)$ (100%)	$\exp[-2\psi_0] \{(i\pi)^2 + (\psi_1^2)\}^{-1} \sinh^2(\psi_1)$ (1.86%)	$\begin{cases} \tau_0^2, & i = 0 \\ (\tau_1/\pi i)^2, & i \neq 0 \end{cases}$ (2.53%)
Definitions $A_0 = n_0, \alpha_0, \tau_0$ $A_1 = n_1, \alpha_1, \tau_1$	$\gamma = \pi n_1 d / \lambda_0$	$\psi_0 = \alpha_0 d, \psi_1 = \alpha_1 d$	

Note: Maximum achievable efficiency (η_{\max}) is given in the parenthesis. $\theta_{\text{in}} = 0$.

C. Two-Wave Coupled-Wave Analysis

Two-wave coupled-wave theory (Kogelnik's theory) treats diffraction from a bulk-modulated, planar structure operating in a Bragg regime. Such gratings can be realized by laser-interference recording (holographic recording) in, for example, photorefractive crystals and polymer media. The resulting spatial modulation of refractive index (n) and absorption coefficient (α) is appropriate for the Kogelnik model.

As indicated in Fig. 8, an arbitrary grating slant is allowed. This is a key feature since in practical holographic recording, the angles of incidence of the recording waves

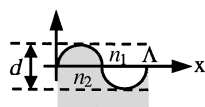
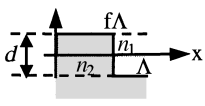
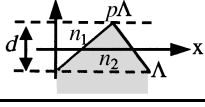
are, in general, not symmetric with respect to the surface normal. The slant angle is denoted ϕ ; $\phi = \pi/2$ defines an unslanted transmission grating, whereas $\phi = 0$ for an unslanted reflection grating.

Table II summarizes the major features and limitations of Kogelnik's theory. The assumptions made yield simple analytical solutions that are easy to apply, provide valuable physical insight, and agree well with experiment in many cases. Consequently, these results are widely cited in the diffractive optics and holography literature.

The pertinent scalar wave equation (TE polarization) is

$$\nabla^2 E_y(x, z) + k^2(x, z) E_y(x, z) = 0, \quad (22)$$

TABLE Ib Summary of Diffraction Efficiency Results from Transmittance Theory (Surface-Relief Phase Gratings)^a

Grating profile	Efficiency (η_i)	η_{\max}
Sinusoidal 	$J_i^2(g)$	33.8%
Rectangular ($0 \leq f \leq 1$) 	$\begin{cases} 1 - 4f(1-f) \sin^2(g), & i = 0 \\ [\sin(i\pi f) \sin(g)/(i\pi/2)]^2, & i \neq 0 \end{cases}$	40.5%
Triangular ($0 \leq p \leq 1$) 	$\begin{cases} p^2, & i = -\pi p/g \\ (1-p)^2, & i = \pi(1-p)/g \\ [g \sin(g + i\pi p)/(g + i\pi p)\{g - i\pi(1-p)\}]^2, & \text{otherwise} \end{cases}$	100%

^a Remarks: $\tau(x) = \exp[-jk_0 \Delta n d(x)]$, $g = \pi \Delta n d / \lambda_0$, $\Delta n = n_2 - n_1$, and $\theta_{\text{in}} = 0$.

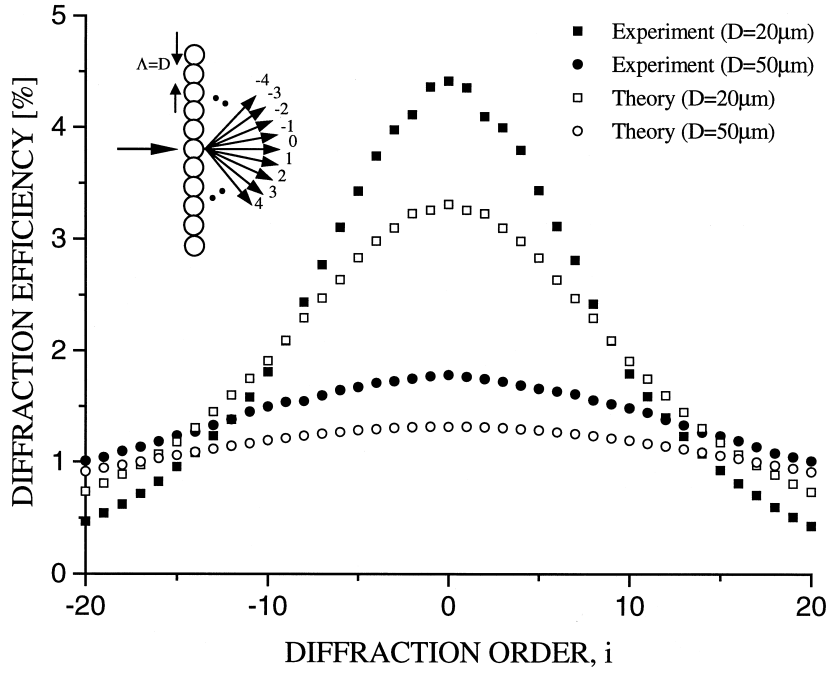


FIGURE 7 Comparison of experimental and theoretical efficiencies of diffracted orders from two ($\Lambda = 20$ and $50 \mu\text{m}$) arrays of dielectric circular cylinders. The cylinder with a refractive index of 1.71 is surrounded by a slab with an index 1.47. A spatially filtered and collimated HeNe laser ($\lambda_0 = 632.8 \text{ nm}$) is used as a source at normal incidence. [From Magnusson, R., and Shin, D. (1989). "Diffraction by periodic arrays of dielectric cylinders." *J. Opt. Soc. Am. A* **6**, 412–414.]

where $k(x, z)$ includes the spatial modulation of the optical constants for a mixed (phase and absorption modulated) grating and is given by

$$k^2(x, z) = \left(\frac{\omega}{c}\right)^2 \varepsilon(x, z) - j\omega\mu\sigma(x, z), \quad (23)$$

where $c = 1/\sqrt{\mu_0\varepsilon_0}$ is the vacuum (with permittivity ε_0 and permeability μ_0) speed of light, the dielectric constant is $\varepsilon(x, z) = n^2(x, z)$ and $\sigma(x, z)$ is the conductivity. For good dielectrics, $\sigma/\omega\varepsilon \ll 1$ and α can be approximated as (nonmagnetic media $\mu = \mu_0$) $\alpha \approx \sigma c\mu_0/2n$. Considering modulations as $n = n_0 + \Delta n(x, z)$ and $\alpha = \alpha_0 + \Delta\alpha(x, z)$ and assuming $\Delta n \ll 1$, $\Delta\alpha \ll 1$, and $\beta = 2\pi n_0/\lambda_0 \gg \alpha_0$, there results

$$k^2(x, z) = \beta^2 - 2j\alpha_0\beta + 2\beta\left(\frac{2\pi\Delta n(x, z)}{\lambda_0} - j\Delta\alpha(x, z)\right). \quad (24)$$

Taking sinusoidal modulations of the basic quantities ε and σ and converting to modulations of n and α leads to $\Delta n = n_1 \cos(\vec{K} \cdot \vec{r})$ and $\Delta\alpha = \alpha_1 \cos(\vec{K} \cdot \vec{r})$ where $\Delta n = \Delta\varepsilon/2n_0$, $\Delta\alpha = \mu_0 c \Delta\sigma/2n_0$, \vec{K} is the grating vector, and $\vec{r} = x\hat{x} + z\hat{z}$ is the position vector in Fig. 8. Thus for sinusoidal $\Delta\alpha$ and Δn

$$k^2(x, z) = \beta^2 - 2j\alpha_0\beta + 2\beta\kappa[\exp(j\vec{K} \cdot \vec{r}) + \exp(-j\vec{K} \cdot \vec{r})], \quad (25)$$

where $\kappa = \pi n_1/\lambda_0 - j\alpha_1/2$ is the coupling coefficient. The electric field inside the grating is expanded in coupled-wave form (similar to Eq. (8) in Section II.A retaining only two waves) as

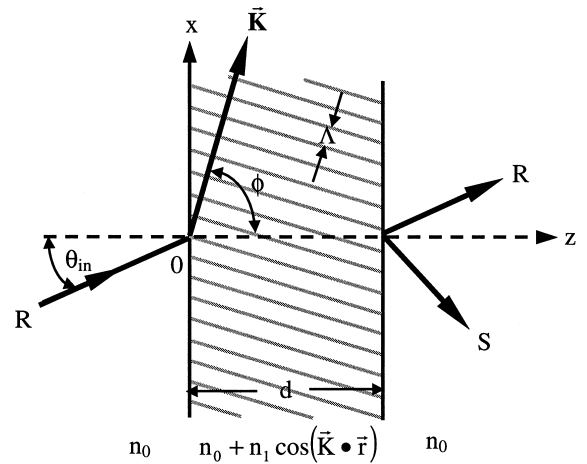


FIGURE 8 Geometry of diffraction by thick gratings (Kogelnik's model).

TABLE II Features and Limitations of Kogelnik's Theory

Features	Limitations
Sinusoidal modulation of n and α	Neglects second derivatives
Slanted gratings	No boundary reflections
Two-wave solution (plane waves)	Approximate boundary conditions (E only)
Bragg or near-Bragg solutions	
TE (H-mode, p) or TM (E-mode, p) polarizations	Neglects higher order waves
Transmission and reflection gratings	

$$E_y(x, z) = R(z) \exp(-j\vec{\rho} \bullet \vec{r}) + S(z) \exp(-j\vec{\sigma} \bullet \vec{r}) \quad (26)$$

with the wavevector of the diffracted wave expressed by the Floquet condition as $\vec{\sigma} = \vec{\rho} - \vec{K}$. The Bragg condition for the slanted grating is given by

$$\cos(\theta_{in} - \phi) = \frac{\lambda_0}{2n_0\Lambda} = \frac{K}{2\beta}. \quad (27)$$

By calculating $\nabla^2 E_y(x, z)$, including the modulations in $k^2(x, z)$, inserting both into Eq. (22), and collecting terms in $\exp(-j\vec{\rho} \bullet \vec{r})$ and $\exp(-j\vec{\sigma} \bullet \vec{r})$ leads to the coupled-wave equations

$$\begin{aligned} C_R \frac{\partial R}{\partial z} + \alpha_0 R &= -j\kappa S \\ C_S \frac{\partial S}{\partial z} + (\alpha_0 + j\vartheta) S &= -j\kappa R, \end{aligned} \quad (28)$$

where κS , κR account for the coupling of the reference and signal waves, $\alpha_0 S$, $\alpha_0 R$ account for absorption, and ϑS defines dephasing of R and S on propagation through

the grating. To solve the coupled-wave Eqs. (28), the following boundary conditions are used:

Transmission type: $R(0) = 1$, $S(0) = 0$, with $S(d)$ giving the output amplitude of interest.

Reflection type: $R(0) = 1$, $S(d) = 0$, with $S(0)$ to be determined.

Diffraction efficiency is defined, as in Section II.A, by the ratio between the diffracted power density (along the z -direction) and the input power density and is given by $\eta = (|C_S|/|C_R|)|S|^2$, where $S = S(d)$ for a transmission grating and $S = S(0)$ for a reflection grating.

Table III summarizes key results for the practical case of lossless dielectric structures or pure-phase gratings. Example calculated results are presented in Fig. 9. The diffraction efficiency of a lossless transmission grating for three slightly different grating periods is shown in Fig. 9(a); these might represent three data pages in an angularly multiplexed holographic memory crystal. The reflectance spectra of several Bragg gratings are shown in Fig. 9(b); these might correspond to reflection gratings in an optical fiber used in a wavelength-division-multiplexed system.

D. Effective Medium Theory (EMT)

For gratings with sufficiently high spatial frequency ($\Lambda \ll \lambda_0$), only zero-order forward- and backward-diffracted waves propagate with all higher orders being evanescent, as seen from the grating equation. Such

TABLE III Summary of Main Results from Kogelnik's Theory (Pure-Phase Gratings)^a

Type	Transmission gratings	Reflection gratings
Configuration		
Efficiency (η)	$\frac{\sin^2[(v^2 + \xi^2)^{1/2}]}{1 + \xi^2/v^2}$	$\left[1 + \frac{1 - \xi^2/v^2}{\sinh^2\{(v^2 - \xi^2)^{1/2}\}}\right]^{-1}$
v	$\kappa d / (C_R C_S)^{1/2}$	$j\kappa d / (C_R C_S)^{1/2}$
ξ	$\frac{\vartheta d}{2C_S}, C_S > 0$	$-\frac{\vartheta d}{2C_S}, C_S < 0$
Unslanted case (On Bragg: $\vartheta = 0$)	$\eta = \sin^2\left(\frac{\pi n_1 d}{\lambda_0 \cos \theta_{in}}\right), \phi = \pi/2$	$\eta = \tanh^2\left(\frac{\pi n_1 d}{\lambda_0 \cos \theta_{in}}\right), \phi = 0$

^a Parameters: $\kappa = \pi n_1 / \lambda_0$, $C_R = \cos \theta_{in}$, $C_S = \cos \theta_{in} - K \cos \phi / \beta$, $\beta = 2\pi n_0 / \lambda_0$, $\vartheta = K \cos(\theta_{in} - \phi) - K^2 / 2\beta$.

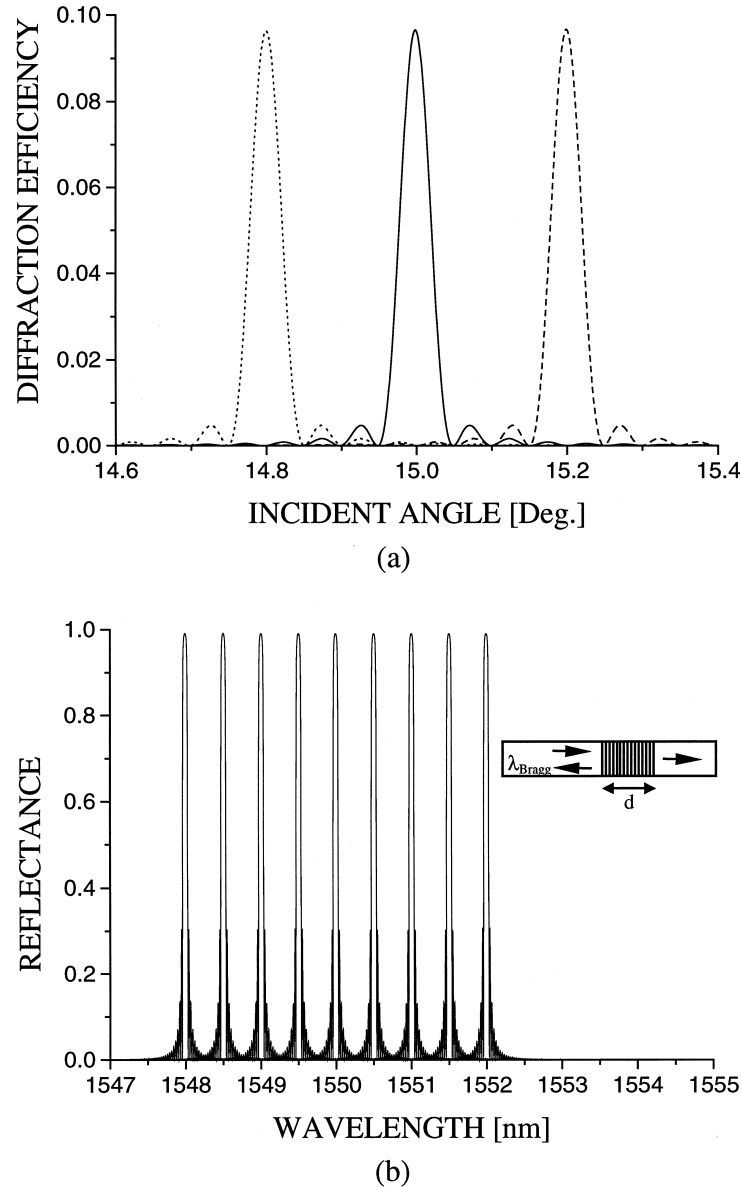


FIGURE 9 Example results from Kogelnik's theory. (a) Diffraction efficiencies of an angularly multiplexed transmission ($\phi = \pi/2$) hologram grating in LiNbO_3 as a function of incident angle. The parameters are $n_0 = 2.25$, $n_1 = 10^{-4}$, $d = 0.5$ mm, and $\lambda_0 = 514.5$ nm. Grating periods are 447.6 nm (dotted line), 441.8 nm (solid line), and 436.1 nm (dashed line). (b) Diffraction efficiencies from unslanted ($\phi = 0$) fiber Bragg gratings as a function of wavelength for several grating periods. The parameters are $n_0 = 1.45$, $n_1 = 5 \times 10^{-5}$, $d = 3$ cm, and $\Lambda = 534.5$ nm at $\lambda_0 = 550$ nm, for example.

gratings are referred to as subwavelength gratings or zero-order gratings. In this operating region, an incident electromagnetic wave cannot resolve the finely modulated layer and interacts with the structure, approximately, as if it were a homogeneous layer with its optical constants equal to the spatial average of those in the modulated layer. This is analogous to optical crystals which are periodic structures exhibiting strong diffraction effects in the x-ray region while acting as homogeneous media in the optical

spectral region. As crystals have a natural birefringence depending on the arrangement of atoms or molecules, subwavelength gratings also exhibit polarization dependence called form birefringence due to the spatial periodic modulation.

In the subwavelength region, a rectangular grating layer such as that in Fig. 1, may be approximately modeled as a homogeneous negative uniaxial layer with ordinary and extraordinary indices of refraction given by

$$n_O = [n_L^2 + f(n_H^2 - n_L^2)]^{1/2} \quad (29)$$

$$n_E = [n_L^{-2} + f(n_H^{-2} - n_L^{-2})]^{-1/2}. \quad (30)$$

Similarly, a multilevel stair-step grating can be modeled as a stack of homogeneous layers with different effective indices. A grating with a continuously-varying dielectric constant profile $\varepsilon(x, z)$, such as a sinusoidal shape for example, can be approximated as a graded-index layer along the thickness direction as

$$n_O(z) = \left[\frac{1}{\Lambda} \int_{-\Lambda/2}^{\Lambda/2} \varepsilon(x, z) dx \right]^{1/2} \quad (31)$$

$$n_E(z) = \left[\frac{1}{\Lambda} \int_{-\Lambda/2}^{\Lambda/2} \frac{1}{\varepsilon(x, z)} dx \right]^{-1/2}. \quad (32)$$

For a wave incident from the cover (region C) as shown in Fig. 1, the ordinary and extraordinary waves correspond to TE- and TM-polarized waves, respectively. Classical (anisotropic) thin-film analysis can be applied to obtain the transmission and reflection coefficients (i.e., diffraction efficiencies) of the propagating zero-order waves.

The EMT expressions given above are independent of Λ/λ_0 and are called the zero-order effective indices. They provide accurate results for weakly modulated gratings ($\Delta\varepsilon = n_H^2 - n_L^2 \ll 1$) in the limit $\Lambda/\lambda_0 \ll 1$. The applicable range in Λ/λ_0 may be extended by using the Λ/λ_0 -dependent second-order effective indices, given by

$$n_O^{(2)} = \left[n_O^2 + \frac{1}{3} \left\{ \pi \frac{\Lambda}{\lambda_0} f(1-f) \right\}^2 (n_H^2 - n_L^2)^2 \right]^{1/2} \quad (33)$$

$$n_E^{(2)} = \left[n_E^2 + \frac{1}{3} \left\{ \pi \frac{\Lambda}{\lambda_0} f(1-f) \right\}^2 \times (n_H^{-2} - n_L^{-2})^2 n_E^6 n_O^2 \right]^{1/2}. \quad (34)$$

More accurate results are found via higher-order EMT indices or by solving the Rytov eigenvalue equations numerically. Comparison of zero-order diffraction efficiencies found by EMT and RCWA is available in the literature.

Figure 10(a) shows the effective index of an etched GaAs grating, calculated using the effective medium theory, as a function of grating fill factor. This shows that an arbitrary effective index ranging from 1 (air) to 3.27 (bulk GaAs) is achievable for both polarizations by proper choice of the grating fill factor. Estimation using the second order EMT (dashed line) differs slightly from that using the zero-order EMT for this high-spatial-frequency ($\Lambda/\lambda_0 = 0.1$) grating. The maximum value of form bire-

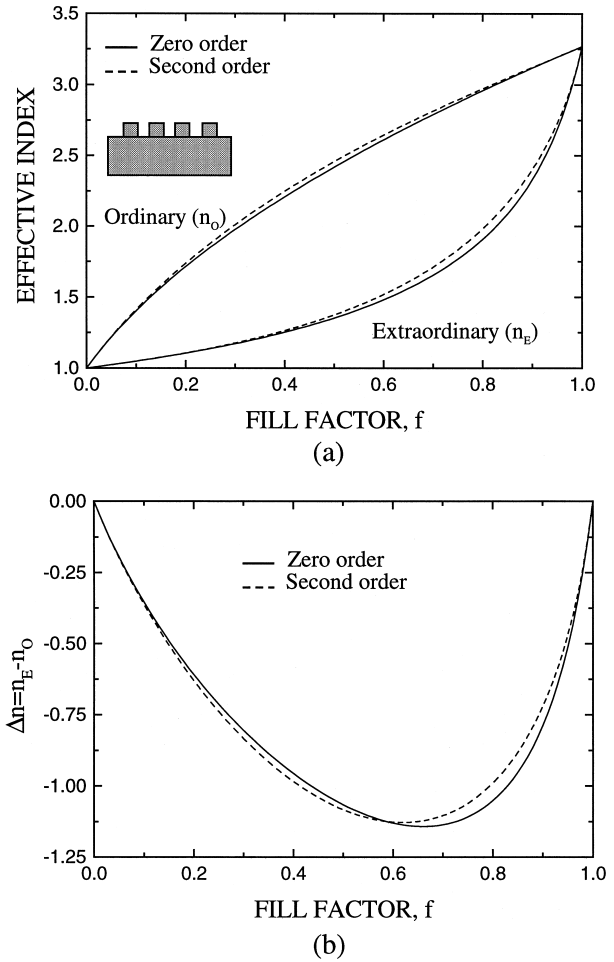


FIGURE 10 Effective medium properties of a one-dimensional, etched, GaAs subwavelength grating calculated by the zero-order (solid line) and second-order (dashed line) EMT. Parameters used are $\lambda_0 = 10.6 \mu\text{m}$, $\Lambda = 0.1\lambda_0$, and $n_{\text{GaAs}} = 3.27$. (a) Effective refractive indices for ordinary and extraordinary waves as a function of fill factor. (b) Form birefringence of the GaAs subwavelength grating as a function of fill factor.

fringence, defined by $\Delta n = n_E - n_O$, is shown (Fig. 10(b)) to be -1.128 at fill factor of 0.62; this value is much larger than that of naturally birefringent materials such as calcite with $\Delta n = -0.172$ in the visible region.

III. APPLICATIONS

A. Optical Interconnects

Due to large bandwidths and inherent parallel-processing capability, optical interconnections are of interest in communication signal routing and optical computing. Figure 11 shows example optical interconnection devices made with diffractive components including interconnections in free space (Figs. 11(a)–(d)) and in integrated and

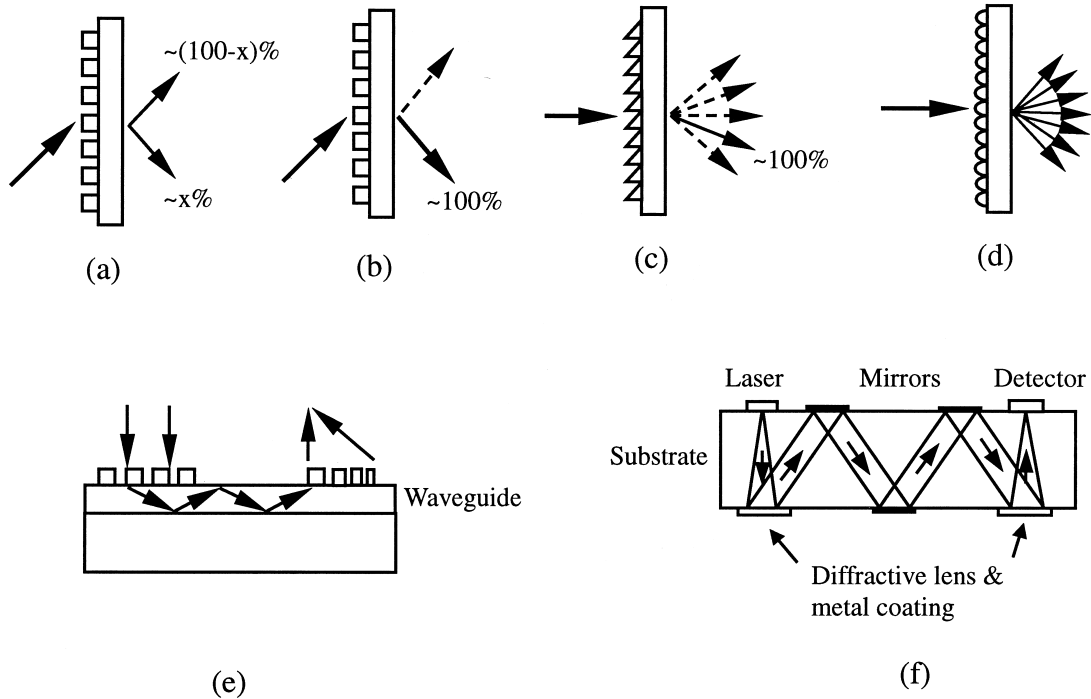


FIGURE 11 Optical interconnects: (a) beam divider (3 dB if $x = 50$), (b) beam deflector, (c) blazed grating, (d) array generator, (e) waveguide interconnect, and (f) substrate interconnect.

planar optical circuits (Figs. 11(e)–(f)). Gratings operating in a two-wave regime may be used as beam dividers (Fig. 11(a)). Under a Bragg incidence, these gratings generate two most efficient forward diffracted orders. The beam ratio (represented as x in the figure) may be controlled by the grating thickness; a beam deflector that transfers most of the energy into the first diffracted order, as shown in Fig. 11(b), is also achievable with this structure. Another type of beam deflector may be designed using a low-spatial-frequency blazed grating [Fig. 11(c)]. A low-frequency grating that generates multiple diffracted orders with similar efficiencies, as in Fig. 7, may be used as a channel divider (fanout element) [Fig. 11(d)]. In general, 1-to- N channel dividers with specified beam-division ratio can be achieved by Dammann gratings. Integrated optical input couplers are realizable using a diffractive phase matching element [Fig. 11(e)]. Similar gratings can also be used as output couplers. If a chirped grating (grating with variable local periods) is used, a focused output beam is obtainable, which is useful in optical memory readout devices. A substrate mode optical interconnect system is shown in Fig. 11(f).

B. Diffractive Lenses

Diffractive optical components are used to construct imaging and beam shaping devices. Compact, lightweight, and

fast (high numerical aperture) lenses can be obtained via radially modulated diffractive optical components. Diffractive lenses may be constructed with reference to the refractive spherical lens as explained by Fig. 12(a). The phase profile of the spherical lens is transformed based on modulo 2π operation that eliminates bulk sections of 2π phase shifts, resulting in a Fresnel lens with a set of annular rings within which the surface has continuous profile. The maximum thickness of the lens is $d_{2\pi} = \lambda_0 / [n(\lambda_0) - 1]$, which corresponds to a 2π phase shift at wavelength λ_0 . These lenses may be fabricated by diamond turning, gray level masking, or direct laser writing lithography, depending on the feature size and aspect ratio. Alternatively, the continuous phase profile may be approximated by a discrete binary or multilevel profile as shown in Fig. 12(a). Fabrication of multilevel diffractive lenses is accomplished via photolithography employing accurately placed masks and sequential lithographic and etching steps.

Figure 12(b) illustrates schematically the operation of a diffractive lens using a ray approach. The lens may be viewed as a low-spatial frequency chirped grating with local grating periods decreasing in the radial direction. The Fresnel zone radii, r_h , are defined such that rays emanating from adjacent zones add constructively at the focal point. This means that the adjacent ray paths differ by λ_0 , the design wavelength. Thus, the radius of h -th Fresnel zone is

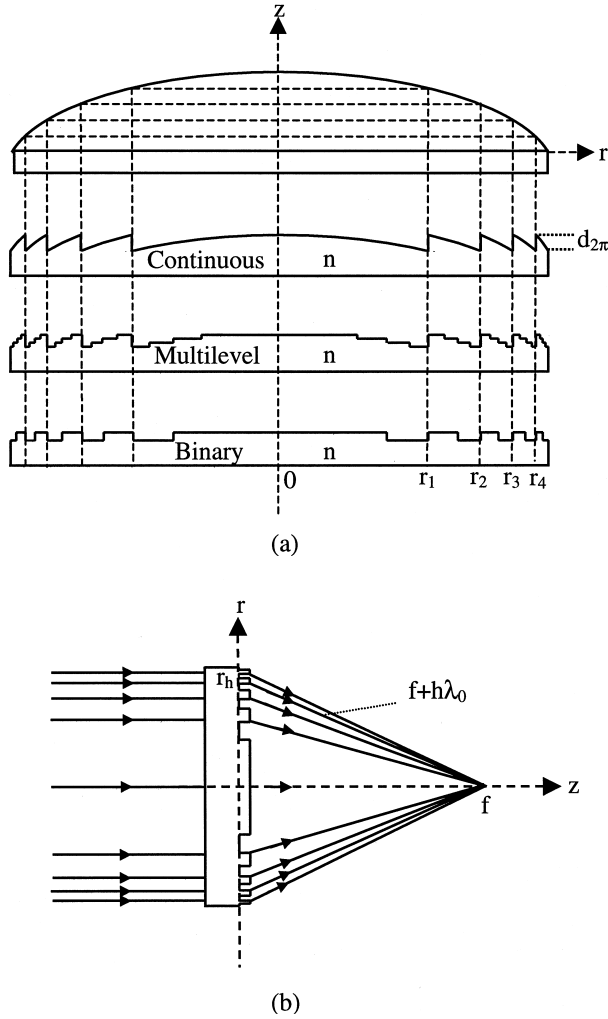


FIGURE 12 (a) A refractive spherical lens and analogous diffractive lenses. (b) Diffractive-lens focusing via ray tracing.

$$r_h = [(f + h\lambda_0)^2 - f^2]^{1/2}, \quad h = \text{integer}. \quad (35)$$

In the paraxial region ($f \gg h_{\max}\lambda_0$), it may be approximated as $r_h = [2h\lambda_0 f]^{1/2}$.

The wavelength-dependent focal length and diffraction efficiency of diffractive lenses with a continuous profile are given, through transmittance diffraction analysis, by

$$f_m(\lambda) = \frac{\lambda_0 f}{\lambda m}, \quad m = \text{integer} \quad (36)$$

and

$$\eta_m(\lambda) = \frac{\sin^2[\pi\{\alpha(\lambda) - m\}]}{[\pi\{\alpha(\lambda) - m\}]^2}, \quad (37)$$

where $\alpha(\lambda)$ is a detuning parameter defined by

$$\alpha(\lambda) = \frac{\lambda_0}{\lambda} \left[\frac{n(\lambda) - 1}{n(\lambda_0) - 1} \right]. \quad (38)$$

The expression for the focal length shows that the diffractive lens has an infinite number of focal points that correspond to the diffracted orders m , and that it is highly dispersive (i.e., f depends on λ). The diffraction efficiency can approach 100% for $\alpha(\lambda) = 1$. In contrast, the peak efficiency of a diffractive lens with a multilevel profile is given by

$$\eta = \frac{\sin^2[\pi/M]}{[\pi/M]^2} \quad (39)$$

at the operating wavelength λ_0 where M is the number of phase levels. Diffractive lenses with eight levels ($M = 8$), 4-levels ($M = 4$), and a binary phase profile ($M = 2$), for example, have maximum efficiencies of 95, 81, and 40.5%, respectively. Detailed discussion on diffractive lenses and their applications can be found in the literature.

C. Subwavelength Grating Devices

Due to advances in microfabrication technology, gratings with periods much smaller than the wavelength can be made. Various applications of subwavelength gratings are shown in Fig. 13. The operation of most subwavelength diffractive optical devices is based on the idea of an effective medium, i.e., the spatial averaging of the refractive index and the associated form birefringence described in Section II.D. Additional devices are based on combination of subwavelength gratings with thin-film and waveguide layers.

1. Antireflection (AR) Surfaces

Common AR coating techniques involve deposition of single or multiple homogeneous thin-film layers, whose refractive indices and thicknesses are chosen to reduce the Fresnel reflections at optical interfaces. Subwavelength gratings may be used as AR coatings (Fig. 13(a)) since the effective refractive index can be controlled by the grating fill factor as seen in Fig. 10(a). Single-layer AR structures can be etched into a variety of surfaces without additional thin-film deposition steps with improved adhesion and performance.

For example, an antireflection surface may be obtained by placing a rectangular subwavelength grating directly on the substrate. The commonly used quarter-wave matching technique (so called V-coating) requires that the grating layer have an effective refractive index equal to $n_{\text{eff}} = [n_c n_s]^{1/2}$ and layer thickness of $d = \lambda_0 / 4n_{\text{eff}}$ in order to be antireflective at the center wavelength of λ_0 for a normally incident wave. By connecting these with the zero-order effective index expressions (Eqs. (29) and (30)) with $n_L = n_c$ and $n_H = n_s$, an approximate grating fill factor required for antireflection is obtained as

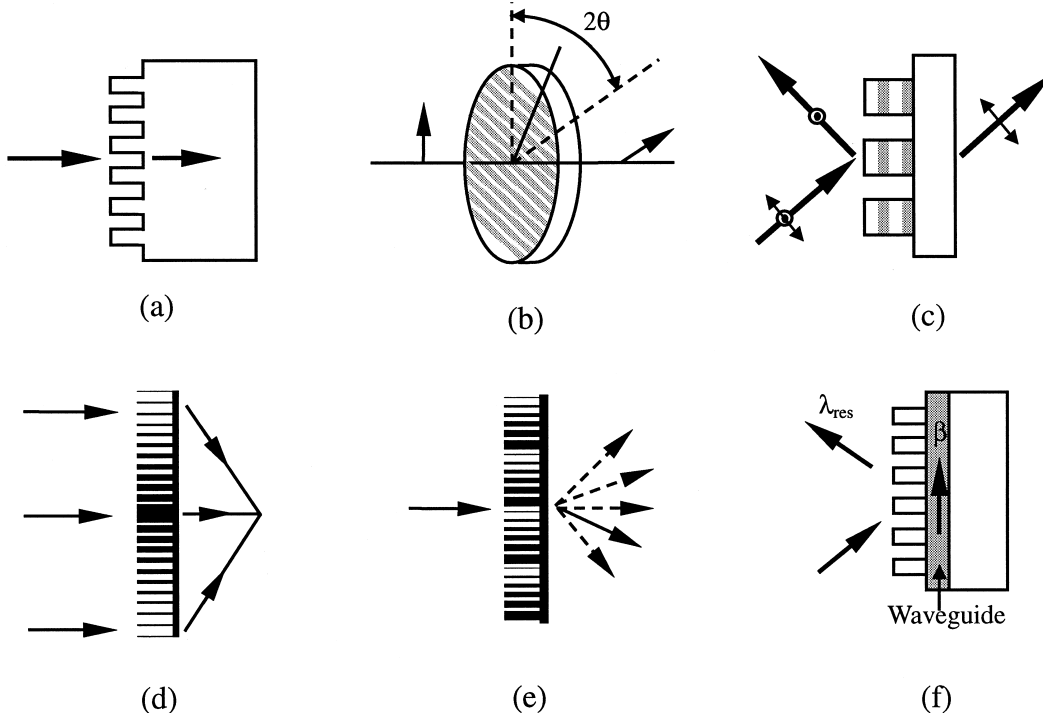


FIGURE 13 Applications of subwavelength gratings: (a) antireflection surface, (b) retardation plate, (c) polarizing beam splitter, (d) diffractive lens, (e) blazed grating, and (f) wavelength-selective reflector (filter).

$$f_{AR} = \begin{cases} n_C/(n_C + n_S) & \text{(TE)} \\ n_S/(n_C + n_S) & \text{(TM)} \end{cases} \quad (40)$$

Figure 14 shows an example spectral reflectance, calculated by RCWA, of a one-dimensional subwavelength antireflection grating formed by etching a GaAs substrate. The resulting reflectance (below $\sim 0.1\%$ at the center wavelength of $10.6 \mu\text{m}$) with the grating layer is much

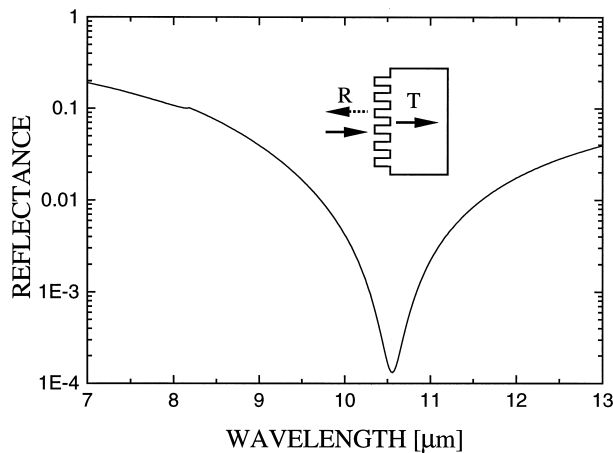


FIGURE 14 Spectral reflectance of a one-dimensional antireflection grating etched on a GaAs substrate under a normally incident wave at $10.6 \mu\text{m}$. The device parameters are $n_{\text{GaAs}} = 3.27$, $\Lambda = 2.5 \mu\text{m}$, $f = 0.19$, and $d = 1.4655 \mu\text{m}$.

lower than the bulk GaAs substrate reflection ($\sim 28\%$). A multilevel or continuously varying grating profile obtains a wider antireflection bandwidth. For example, a triangular profile simulates a layer with a graded refractive index along the z -axis and provides an improved spectral and angular AR performance. Symmetric crossed gratings (grating modulation in two orthogonal directions) may be used for polarization-independent antireflection structures as shown in Fig. 15. Applications of subwavelength AR gratings include detectors, solar cells, and optical windows.

2. Polarization Devices

Subwavelength gratings may be used as polarization components (Figs. 13(b) and (c)). The metallic wire-grid polarizer is a well-known example. The form birefringence associated with dielectric subwavelength gratings enables construction of polarization components analogous to those made by naturally birefringent materials. For example, Fig. 13(b) shows a subwavelength retardation plate. The phase retardation experienced by a normally incident, linearly polarized wave passing through a subwavelength grating is given by

$$\Delta\phi = \phi_{\text{TE}} - \phi_{\text{TM}} = -k\Delta nd, \quad (41)$$

where $\Delta n = n_{\text{TM}} - n_{\text{TE}}$ is the form birefringence and d is the thickness of the grating. Due to their large

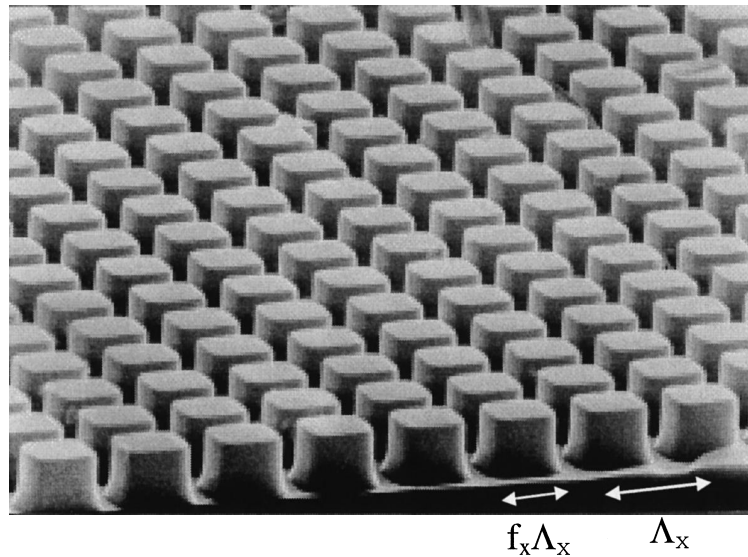


FIGURE 15 Scanning electron micrograph of a two-dimensional Si AR surface designed for operation in the infrared (8- to 12- μm) band. The grating has a 2.45- μm period and fill factors $f_x = f_y = 0.7$. The grating is quarter wavelength thick at 10.6 μm . [From Raguin, D. H., Norton, S., and Morris, G. M. (1994). "Diffractive and Miniaturized Optics" (S. H. Lee, ed.), pp. 234–261, Critical review **CR 49**, SPIE Press, Bellingham.]

birefringence (Fig. 10(b)) compared with that provided by naturally birefringent materials, compact and light weight zero-order retardation plates such as quarter-wave plates, half-wave plates, polarization rotators, and polarization compensators may be fabricated. As this artificial birefringence can be obtained using isotropic materials, low-cost devices are feasible at any operating wavelength. Note that more compact devices are possible with higher index materials at the expense of higher insertion loss due to increased reflection. Design considerations may involve optimization to maximize the birefringence and minimize insertion loss.

Subwavelength gratings may be used to form multi-function polarization elements. The inset of Fig. 16 shows a diffractive polarizing mirror. A lateral subwavelength corrugation yields a polarization-dependent effective refractive index, while a stack of thin-film layers with alternating high/low refractive index simulates a distributed Bragg reflector. By taking each layer-thickness to be quarter wavelength for one polarization (TE polarization in this example), the structure reflects this polarization while it transmits the orthogonal polarization. The device in Fig. 16 can be used as a polarizing mirror at 1.55 μm . Similar structures can be designed to operate as polarizing beam splitter as shown in Fig. 13(c).

3. Distributed-Index Devices

The controllability of the effective refractive index (by the fill factor f) of a subwavelength grating enables applica-

tions such as phase plates. A refractive surface, such as shown in Fig. 17(a), may be modeled as a single, planar, laterally graded index layer with a fixed thickness if the x -dependent phase accumulation provided by the refractive surface profile (with fixed refractive indices n_1 and n_2) is translated into an equivalent refractive index

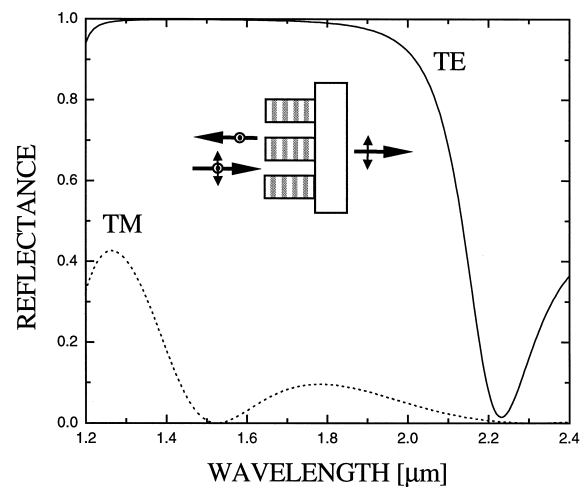


FIGURE 16 Diffractive polarization-dependent mirror. The grating consists of four pairs of quarter-wavelength thick Si-SiO₂ stack deposited on a fused silica substrate and etched with a period of 600 nm and a fill factor of 0.5. Refractive indices are $n_{\text{Si}} = 3.48$ and $n_{\text{SiO}_2} = 1.44$. [From Tyan, R.-C., Salveker, A. A., Chou, H.-P., Cheng, C.-C., Scherer, A., Sun, P.-C., Xu, F., and Fainman, Y. (1997). "Design, fabrication, and characterization of form-birefringent multilayer polarizing beam splitter." *J. Opt. Soc. Am. A* **14**, 1627–1636.]

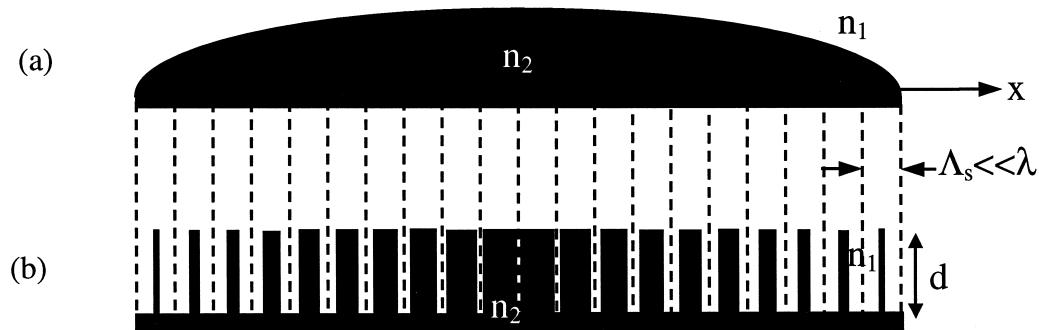


FIGURE 17 Construction of subwavelength distributed-index devices. (a) Refractive surface (center zone of diffractive lens). (b) Analogous structure realized by binary fill-factor modulated subwavelength gratings.

distribution. A subwavelength grating can be used to approximate this graded-index layer using the effective medium property, as shown in Fig. 17(b). That is, the structure is first replaced by a subwavelength binary grating layer with a sampling period of Λ_s . The grating thickness is set to $d = \lambda_0 / (n_2 - n_1)$ to obtain a phase difference up to 2π by varying the fill factor in a local period. Then, in each local grating period, the fill factor is chosen, with the help of EMT, such that the overall effective index distribution approximates the graded-index profile required.

In general, an arbitrary phase function may be achieved by tailoring the local refractive index with two-dimensional fill-factor-modulated subwavelength gratings. Because the resulting structure is a single-layer binary-modulated surface, planar technology can be used for fabrication. Applications include diffractive lenses (Fig. 13(d)), blazed gratings (Fig. 13(e)), phase compensators, and other beam-shaping devices.

Figure 18 shows a scanning electron micrograph (SEM) of one period (along x) in a two-dimensional blazed

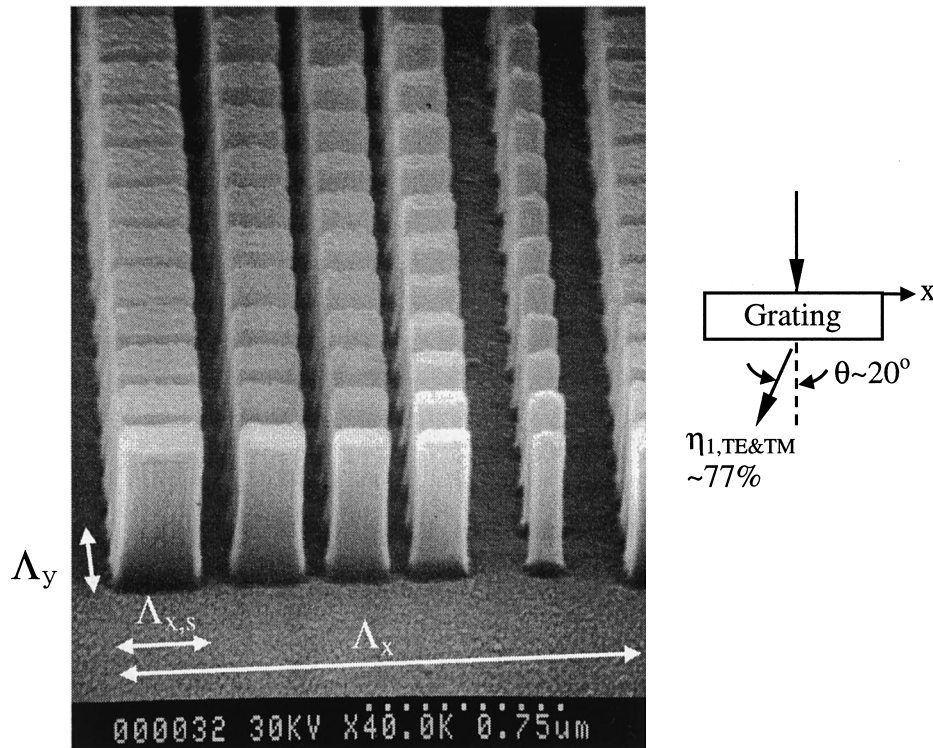


FIGURE 18 Scanning electron micrograph of a blazed grating operating at $\lambda_0 = 632.8$ nm realized by two-dimensional fill-factor-modulated subwavelength grating. The global grating periods are $\Lambda_x = 1900$ nm and $\Lambda_y = 380$ nm. The sampling period ($\Lambda_{x,s}$) is 380 nm. [From Lalanne, P., Astilean, S., Chavel, P., Cambril, E., and Launois, H. (1998). "Blazed binary subwavelength gratings with efficiencies larger than those of conventional echelette gratings." *Opt. Lett.* **23**, 1081–1083.]

grating made of TiO_2 . The grating has global periods of $\Lambda_x = 1900$ nm and $\Lambda_y = 380$ nm such that an incident HeNe laser beam ($\lambda_0 = 632.8$ nm) is diffracted into five forward diffracted orders. Each x -period (Λ_x) is subdivided into five fill-factor modulated subwavelength gratings with a sampling period ($\Lambda_{x,s}$) equal to 380 nm to simulate the blazed grating profile. The two-dimensional area fill factors used in each local subwavelength grating are 0.31, 0.47, 0.53, 0.65, and 0.77. The grating was fabricated by e-beam lithography and reactive-ion etching. For a normally incident HeNe laser beam, this device redirect most of its energy into a first-order wave propagating at an angle $\sim 20^\circ$ nearly independent of the polarization state. Measured diffraction efficiencies were 77 and 78% for TE and TM polarizations, respectively.

4. Guided-Mode Resonance (GMR) Devices

Thin-film structures containing waveguide layers and diffractive elements exhibit the guided-mode resonance effect. When an incident wave is phase-matched, by the periodic element, to a leaky waveguide-mode, it is reradiated in the specular-reflection direction as it propagates along the waveguide and constructively interferes with the directly reflected wave. This resonance coupling is manifested as rapid spectral or angular variations of

the diffraction efficiencies of the propagating waves. When zero-order gratings are used, in particular, a high-efficiency resonance reflection can be obtained, which can be a basis for high-efficiency filtering. (Fig. 13(f))

In Fig. 19, the calculated (using RCWA) spectral reflectance from a double-layer waveguide-grating structure, SEM picture of which is shown in the inset, is given as a dotted line. The device includes HfO_2 waveguide layer and a holographically recorded photoresist grating. A high-efficiency resonant reflection with a linewidth of ~ 2.2 nm is shown to occur near 860 nm for a normally incident TE-polarized wave. This resonance is induced by phase matching to the TE_0 waveguide mode via the first evanescent diffracted order.

The resonance wavelength may be estimated by the phase matching condition

$$2\pi n_C \sin \theta_{in} / \lambda_0 - i2\pi / \Lambda = \beta(\lambda_0), \quad (42)$$

where $\beta(\lambda_0)$ is the mode-propagation constant that can be obtained by solving the eigenvalue equation of the waveguide-grating structure. The resonance spectral linewidth is typically narrow (on the order of nanometers or less) and can be controlled by the modulation amplitude, fill factor, grating thickness, and the refractive-index contrast of the device layers. The resonance response (location and linewidth) is polarization dependent due to

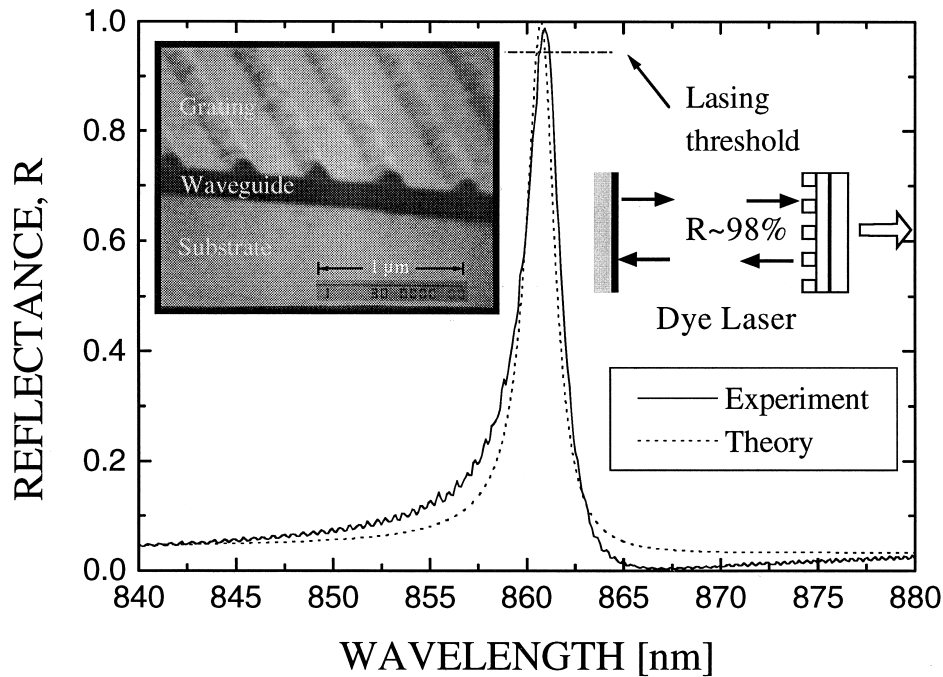


FIGURE 19 Theoretical (dotted line) and experimental (solid line) spectral response from a double-layer waveguide grating for a normally-incident TE-polarized wave. The device consists of holographically-recorded photoresist grating ($n = 1.63$, $\Lambda = 487$ nm, $f = 0.3$, $d = 160$ nm) on top of e-beam evaporated HfO_2 waveguide layer ($n = 1.98$, $d = 270$ nm) on a fused silica substrate ($n = 1.48$). [From Liu, Z. S., Tibuleac, S., Shin, D., Young, P. P., and Magnusson, R. (1998). "High-efficiency guided-mode resonance filter." *Opt. Lett.* **23**, 1556–1558.]

inherent difference in modal characteristics of excited TE- and TM-waveguide modes.

The experimentally obtained high efficiency of $\sim 98\%$ in Fig. 19 (solid line) confirms the feasibility of guided-mode resonance devices. Potential applications of these devices include laser resonator frequency-selective polarizing mirrors, laser cavity tuning elements, electro-optic modulators and switches, tunable filters, mirrors for vertical-cavity lasers, wavelength-division multiplexing, and chemical and biosensors. For example, the fabricated double-layer filter in Fig. 19 was used to realize a GMR laser mirror. The flat output mirror of the dye laser with broadband output (800–920 nm) was replaced with the GMR filter and the birefringent tuning element removed. Lasing occurred at a wavelength of ~ 860 nm. The laser power was ~ 100 mW when pumped with an Ar⁺ laser emitting a power of ~ 5 W at a 514-nm wavelength. The linewidth of the output laser beam was measured as ~ 0.3 nm. This linewidth was set by the GMR filter linewidth at the threshold reflectance for laser oscillation to occur; in this case at $\sim 95\%$ reflectance value in Fig. 19.

IV. FABRICATION

A. Photolithography

For diffractive components with features that are larger than the wavelength, ordinary photolithography may be

applied for fabrication. Typically, several steps are needed including deposition of a photoresist layer on a thin film or substrate, UV light exposure of the photoresist through a prefabricated mask defining the diffractive element, development, and processing. The photoresist, a light-sensitive organic liquid, is spin-coated on the substrate resulting in a film with thickness ($\sim 1 \mu\text{m}$) that depends on the spin rate. If the UV-exposed regions remain (vanish) after development, the resist is said to be negative (positive). Photoresist selection takes account of resolution, sensitivity, adhesion, and other factors. Typically, the resist is UV-exposed for a few seconds followed by development with appropriate chemical solutions. A postbake step at $100\text{--}200^\circ\text{C}$ to increase etching resistance may follow.

B. Laser Interference Recording

Two-wave laser interference (Fig. 20(a)) may be used to record periodic elements with subwavelength features ($\Lambda \sim 100$ nm). The UV laser (for example helium cadmium, argon-ion, or excimer laser) beam is focused down with an objective lens and passed through a matching pin-hole to spatially filter the beam. The central part of the emerging spherical wave, which has nearly planar phase fronts, illuminates the sample directly. A part of the wave is reflected towards the sample as shown. These two waves interfere to produce a standing, periodic intensity pattern with period

$$\Lambda = \lambda_0 / (\sin \theta_1 + \sin \theta_2) \quad (43)$$

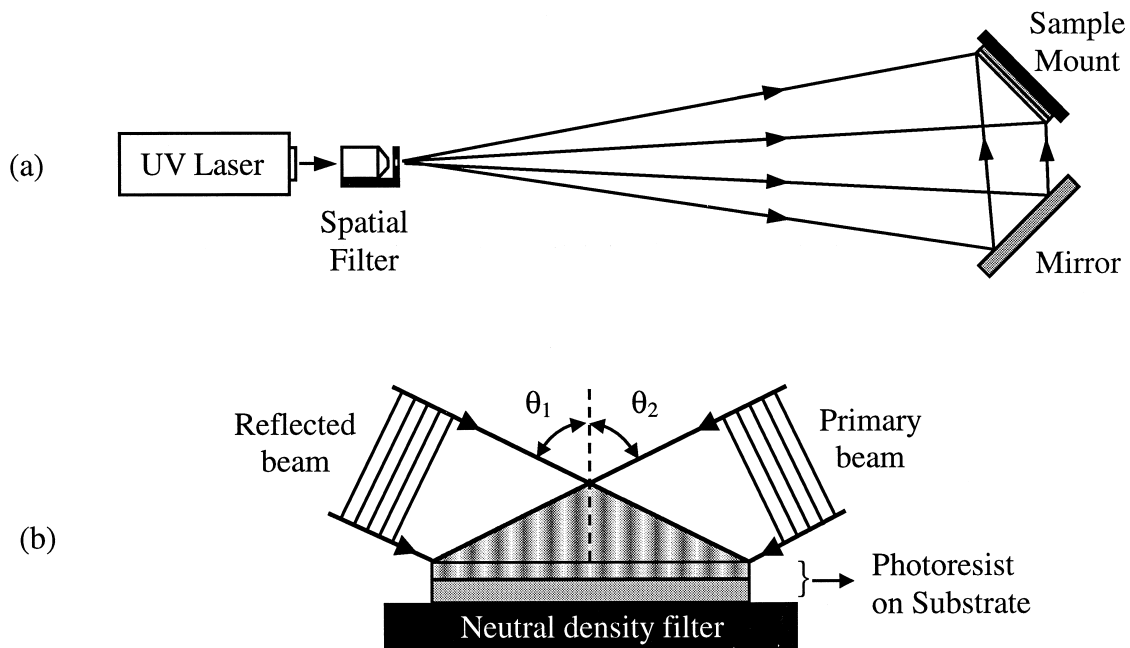


FIGURE 20 (a) Laser interference grating recording system. (b) Details of interference light pattern.

indicated by Fig. 20(b). For symmetrically incident ($\theta_1 = \theta_2 = \theta$) recording beams, the period is $\Lambda = \lambda_0/2 \sin \theta$. The interference pattern has maximum contrast if the two recording beams are polarized in parallel (electric-field vector normal to the plane of incidence). This pattern must be kept stationary with respect to the photoresist during exposure. Heavy vibration-isolating, air-suspended optics tables are used for this; additional measures should be taken to screen off air currents and acoustic noise. The single beam system in Fig. 20 is often used in practice; two-beam systems with separate collimated beams are also in use.

From the expression for the period, it is seen that the minimum period obtainable with this approach is $\Lambda = \lambda_0/2$ for counterpropagating recording beams. It can be further reduced by increasing the refractive index of the input space by placing a glass block or a prism on the sample in which case the period is reduced by the refractive index of the prism. Laser interference recording, thus, allows fabrication of high-quality, large-area gratings with small periods. Noise patterns may arise from interface reflections during exposure. These can be minimized by use of index-matching liquids and neutral density filters (that absorb the light passing through preventing its return to the sample region) as indicated in Fig. 20(b).

C. Other Fabrication Techniques

Many additional fabrication methods are described in the references. Electron-beam writing is a high-resolution

direct-write approach that is widely applied. Since the equipment is expensive and relatively slow, e-beam writing is particularly appropriate for making masks that can serve as masters for subsequent copying. Direct laser writing with a focused UV laser beam is another promising method. The generation of the spatial pattern defining the diffractive element is often followed by a variety of processing steps. These may include reactive-ion etching, wet chemical etching, ion-beam milling, etc. Finally, thin-film deposition may be needed with methods such as sputtering, thermal evaporation, e-beam evaporation, chemical vapor deposition, and molecular beam epitaxy being representative examples.

D. An Example Process

Figure 21 shows a process used for fabricating a buried waveguide grating. A thin film of silicon nitride (thickness ~ 200 nm) is e-beam evaporated on a fused silica substrate followed by thermal evaporation of a thin aluminum layer (thickness ~ 50 nm). A grating with period $\Lambda = 300$ nm is recorded using argon-ion ($\lambda_0 = 364$ nm) laser interference in photoresist. The developed photoresist grating is dry etched (RIE) to remove any residual resist in the grating troughs through which the Al layer is wet etched. The silicon nitride film is then dry etched through the metal mask and the Al grating removed. On sputtering deposition of a layer of silicon dioxide, a silicon nitride/silicon dioxide waveguide grating with a square-wave profile results; it is a waveguide grating as its average index of refraction exceeds that of the surrounding media.

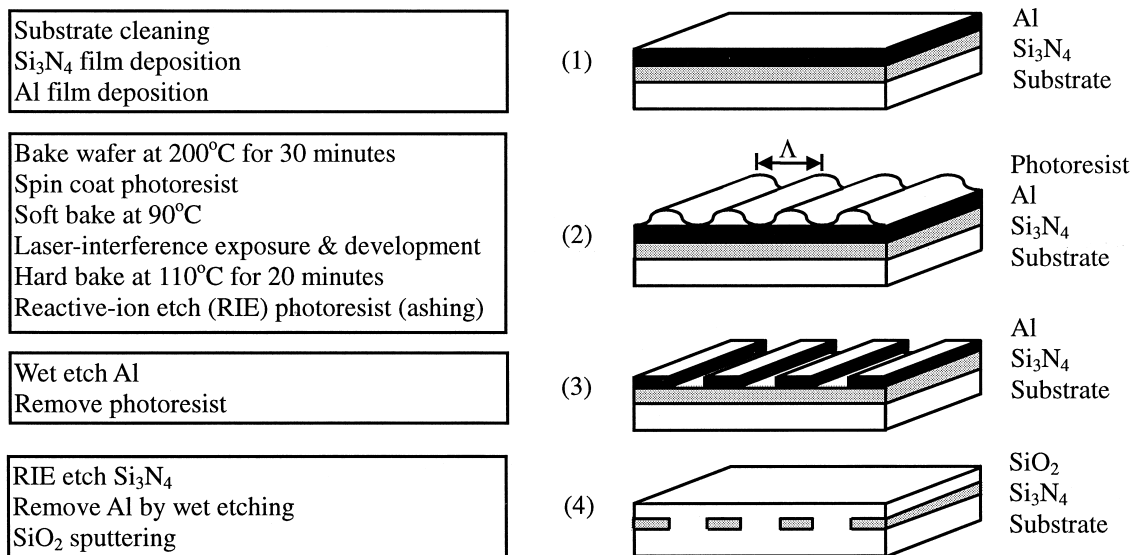


FIGURE 21 An example procedure to fabricate an embedded-grating device.

SEE ALSO THE FOLLOWING ARTICLES

HOLOGRAPHY • IMAGING OPTICS • LASERS • LASERS,
OPTICAL FIBER • MICROOPTICS • OPTICAL DIFFRACTION
• OPTICAL WAVEGUIDES (PLANAR) AND WAVEGUIDE
DEVICES

BIBLIOGRAPHY

- Born, M., and Wolf, E. (1980). "Principles of Optics" (5th edition), Pergamon, Oxford.
- Gaylord, T. K., and Moharam, M. G. (1985). Analysis and applications of optical diffraction by gratings. *Proc. IEEE* **73**, 894–937.
- Goodman, J. W. (1996). "Introduction to Fourier Optics" (2nd edition), McGraw Hill, New York.
- Herzig, H. P. (ed.) (1997). "Micro Optics: Elements, Systems and Applications," Taylor and Francis, London.
- Joannopoulos, J. D., Meade, R. D., and Winn, J. N. (1995). "Photonic Crystals: Molding the Flow of Light," Princeton University Press, Princeton, NJ.
- Kogelnik, H. (1969). Coupled wave theory for thick hologram gratings. *Bell Sys. Tech. J.* **48**, 2909–2947.
- Lee, S. H. (ed.). (1994). "Diffractive and Miniaturized Optics," Critical review **CR 49**, SPIE Press, Bellingham.
- Mait, J. N., and Prather, D. W. (ed.) (2001). "Selected Papers on Sub-wavelength Diffractive Optics," Milestone series **MS 166**, SPIE Optical Engineering Press, Bellingham.
- Martellucci, S., and Chester, A. N. (ed.) (1997). "Diffractive Optics and Optical Microsystems," Plenum, New York.
- Moharam, M. G., Grann, E. B., Pommet, D. A., and Gaylord, T. K. (1995). Formulation of stable and efficient implementation of the rigorous coupled-wave analysis of binary gratings. *J. Opt. Soc. Am. A* **12**, 1068–1076.
- Nishihara, H., Haruna, M., and Suhara, T. (1989). "Optical Integrated Circuits," McGraw-Hill, New York.
- Petit, R. (ed.) (1980). "Electromagnetic Theory of Gratings," Springer Verlag, Berlin.
- Sinzinger, S., and Jahns, J. (1999). "Micro-optics," Wiley-VCH, Weinheim.
- Solymer, L., and Cooke, D. J. (1981). "Volume Holography and Volume Gratings," Academic Press, London.
- Turunen, J., and Wyrowski, F. (ed.) (1997). "Diffractive Optics for Industrial and Commercial Applications," Akademie Verlag, Berlin.



Holography

Clark C. Guest

University of California, San Diego

- I. Introduction
- II. Basic Principles
- III. Classification of Holograms
- IV. Recording Materials
- V. Applications
- VI. History of Holography

GLOSSARY

Diffraction Property exhibited by waves, including optical waves. When part of a wavefront is blocked, the remaining portion spreads to fill in the space behind the obstacle.

Diffraction orders When a wave passes through regularly spaced openings, such as a recorded holographic interference pattern, the diffracted waves combine to form several beams of different angles. These beams are called diffraction orders.

Fringes Regular pattern of bright and dark lines produced by the interference of optical waves.

Hologram Physical record of an interference pattern. It contains phase and amplitude information about the wavefronts that produced it.

Holograph Although this is an English word meaning signature, it is often improperly used as a synonym for hologram.

Holography Process of recording holograms and reproducing wavefronts from them.

Index of refraction Property of transparent materials related to their polarizability at optical frequencies. The

speed of light in a vacuum divided by the speed of light in a material gives the index of refraction for the material.

Interference When two waves are brought together they may be in phase, in which case their amplitudes add, or they may be out of phase, in which case their amplitudes cancel. The reinforcement and cancellation of wavefronts due to their relative phase is called interference. Waves that do not have the same phase structure will add in some regions of space and cancel in others. The resulting regions of high and low amplitude form an interference pattern.

Parallax Difference, due to perspective, in a scene viewed from different locations.

Planewave Wave configuration in which surfaces of constant phase form parallel flat planes. All the light in a planewave is travelling the same direction, perpendicular to the surface of the planes.

Reconstructed beam Light diffracted by a hologram that reproduces a recorded wavefront.

Reconstructing beam Beam that is incident on a hologram to provide light for the reconstructed beam.

Reconstruction Either the process of reading out a

recorded hologram, or the wavefront produced by reading out a hologram.

Refractive index See “index of refraction.”

Spherical wave Wave configuration in which surfaces of constant phase form concentric spherical shells or segments of shells. Light in an expanding spherical wave is propagating radially outward from a point, and light in an converging spherical wave is propagating radially inward toward a point.

Surface relief pattern Ridges and valleys on the surface of a material.

Wavefront Surface of constant phase in a propagating wave.

HOLOGRAPHY is the technology of recording wavefront information and producing reconstructed wavefronts from those recordings. The record of the wavefront information is called a hologram. Any propagating wave phenomenon such as microwaves or acoustic waves in a candidate for application of the principles of holography, but most interest in this field has centered on waves in the visible portion of the electromagnetic spectrum. Therefore, this article will concentrate on optical holography.

I. INTRODUCTION

Although holography has many applications, it is best known for its ability to produce three-dimensional images. A hologram captures the perspective of a scene in a way that no simple photograph can. For instance, when viewing a hologram it is possible by moving one's head to look around objects in the foreground and see what is behind them. Yet holograms can be recorded on the same photographic film used for photographs.

Two questions naturally occur: What makes holograms different from photographs, and how can a flat piece of film store a three-dimensional scene? Answering these questions must begin with a review of the properties of light. As an electromagnetic wave, light possesses several characteristics: amplitude, phase, polarization, color, and direction of travel. When the conditions for recording a hologram are met, there is a very close relationship between phase and direction of travel. The key to the answers to both our questions is that photographs record only amplitude information (actually, they record intensity, which is proportional to the square of the amplitude) and holograms record both amplitude and phase information. How ordinary photographic film can be used to record both amplitude and phase information is described in Section II of this article.

Another interesting property of holograms is that they look nothing like the scene they have recorded. Usually, a

hologram appears to be a fairly uniform gray blur, with perhaps a few visible ring and line patterns randomly placed on it. In fact, all the visible patterns on a hologram are useless information, or noise. The useful information in a hologram is recorded in patterns that are too small to see with the unaided eye; features in these patterns are about the size of a wavelength of light, one two-thousandth of a millimeter.

One useful way to think of a hologram is as a special kind of window. Light is reflected off the objects behind the window. Some of the light passes through the window, and with that light we see the objects. At the moment we record the hologram, the window “remembers” the amplitude and direction of all the light that is passing through it. When the hologram is used to play back (reconstruct) the three-dimensional scene, it uses this recorded information to reproduce the original pattern of light amplitude and direction that was passing through it. The light reaching our eye from the holographic window is the same as when we were viewing the objects themselves. We can move our heads around and view different parts of the scene just as if we were viewing the objects through the window. If part of the hologram is covered up, or cut off, the entire scene can still be viewed, but through a restricted part of the window.

There are actually many different types of holograms. Although photographic film is the most widely used recording material, several other recording materials are available. The properties of a hologram are governed by the thickness of the recording material and the configuration of the recording beams. The various classifications of holograms will be discussed in Section III. Holograms can be produced in materials that record the light intensity through alterations in their optical absorption, their index of refraction, or both. Materials commonly used for recording holograms are discussed in Section IV.

Holograms have many uses besides the display of three-dimensional images. Applications include industrial testing, precise measurements, optical data storage, and pattern recognition. A presentation of the applications of holography is given in Section V.

II. BASIC PRINCIPLES

Photographs record light-intensity information. When a photograph is made, precautions must be taken to ensure that the intensities in the scene are suitable; they must be neither too dim nor too bright. Holograms record intensity and phase information. In addition to the limits placed on the intensity, the phase of light used to record holograms must meet certain conditions as well. These phase conditions require that the light is coherent. There are two types of coherence, temporal and spatial; the light used

for recording holograms must have both types of coherence. Temporal coherence is related to the colors in the light. Temporally coherent light contains only one color: it is monochromatic. Each color of light has a phase associated with it; multicolored light cannot be used to record a hologram because there is no one specific phase to record. Spatial coherence is related to the direction of light. At any given point in space, spatially coherent light is always travelling in one direction, and that direction does not change with time. Light that constantly changes its direction also constantly changes its relative phase in space and therefore is unsuitable for recording holograms.

Temporal and spatial coherence are graded quantities. We cannot say that light is definitely coherent or definitely incoherent; we can only say that a given source of light is temporally and spatially coherent by a certain amount. Ordinary light, from a light bulb, for example, is temporally and spatially very incoherent. It contains many different colors, and at any point in space it is changing directions so rapidly that our eyes cannot keep up; we just see that on average it appears to come from many different directions. The temporal coherence of ordinary light can be improved by passing it through a filter that lets only a narrow band of colors pass. Spatial coherence can be improved by using light coming from a very small source, such as a pinhole in an opaque mask. Then we know that light at any point in space has to be coming from the direction of the pinhole. Ordinary light that has been properly filtered for color and passed through a small aperture can be used to record holograms. However, light from a laser is naturally very temporally and spatially coherent. For this reason, practically all holograms are recorded using laser light.

The simplest possible hologram results from bringing together two linearly polarized optical planewaves. Imagine that a planewave is incident at an angle θ_1 on a flat surface. At a particular moment in time we can plot the electric field of the planewave at positions on that surface. This is done in Fig. 1a. If a second planewave is incident on the same surface at a different angle θ_2 , its electric field can also be plotted. This, along with the combined field from both planewaves is plotted in Fig. 1b. Both planewaves incident on the surface are, of course, travelling forward at the speed of light. In Fig. 1, parts c–e show the electric field at the observation surface for each planewave and for the combined field when the waves have each travelled forward by one-fourth, one-half, and three-quarters of a wavelength, respectively. The interesting thing to notice is that the locations on the observation plane where the total electric field is zero remain fixed as the waves travel. Locations midway between zero electric field locations experience an oscillating electric field. To an observer, locations with constant zero electric field appear dark, and locations with oscillating elec-

tric field appear bright. These alternating bright and dark lines, called fringes, form the interference pattern produced by the planewaves. Likewise, locations with zero electric field will leave photographic film unexposed, and locations with oscillating electric field will expose film. Thus, the interference pattern can be recorded.

The interference fringes resulting from two planewaves will appear as straight lines. Wavefronts do not have to be planewaves to produce an interference pattern, not do the wavefronts have to match each other in shape. Interference between arbitrary wavefronts can appear as concentric circles or ellipses or as wavy lines.

The distance L from one dark interference fringe to the next (or from one bright fringe to the next) depends on the wavelength λ of the light and the angle θ between the directions of propagation for the wavefronts:

$$L = \lambda / [2 \sin(\theta/2)]. \quad (1)$$

For reasonable angles, this fringe spacing is about the size of a wavelength of light, around one two-thousandth of a millimeter. This explains why holograms appear to be a rather uniform gray blur: the useful information is recorded in interference fringes that are too small for the eye to see.

Variations in the amplitudes of the recording beams are also recorded in the hologram and contribute to accurate reproduction of the reconstructed wavefront. During hologram recording, locations of zero electric field will occur only if the two beams are equal in amplitude. If one beam is stronger than the other, complete cancelation of their electric fields is impossible, and the depth of modulation, or contrast, of the recorded fringes is decreased. In practical terms, this means that all of the hologram is exposed to some extent.

There are two steps to the use of ordinary photographs, taking (and developing) the photograph and viewing the photograph. Similarly, there are two steps to using a hologram: recording (and developing) the hologram and reconstructing the holographic image. As we have just seen, recording a hologram amounts to recording the interference pattern produced by two coherent beams of light. Reconstructing the holographic image is usually accomplished by shining one of those two beams through the developed hologram. Through a wave phenomenon known as diffraction, the recorded interference fringes redirect some of the light in the reconstructing beam to form a replica of the second recording beam. This replica, or reconstructed, beam travels away from the hologram with the same variation in phase and amplitude that the original beam had. Thus, for the eye, or for common image recording instruments such as a photographic camera or a video camera, the reconstructed wavefront is indistinguishable from the original wavefront and therefore possesses all the

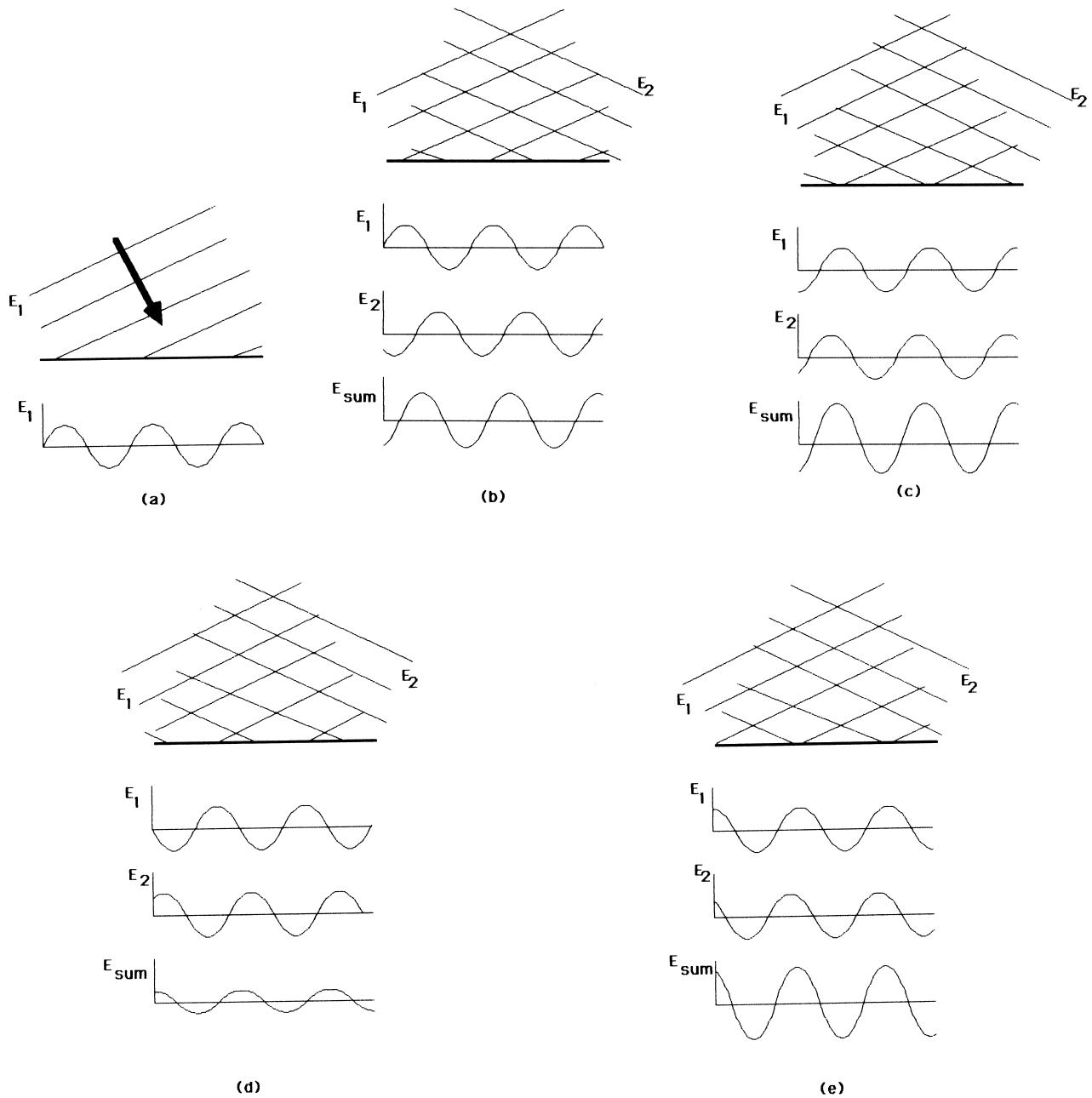


FIGURE 1 (a) The electric field of a planewave incident on a surface. The electric field due to two planewaves incident on a surface at different angles as the waves are (b) initially, (c) after one-quarter wavelength of propagation, (d) after one-half wavelength of propagation, and (e) after three-quarters wavelength of propagation.

visual properties of the original wavefront, including the three-dimensional aspects of the scene.

The interference fringe spacing that is recorded depends on the angle between the recorded beams. During reconstruction, this fringe spacing is translated, through diffraction, back into the proper angle between the reconstructing and the reconstructed beams. During reconstruction, locations on the hologram with high fringe contrast divert

more optical energy into the reconstructed beam than locations with low fringe contrast, thereby reproducing the amplitude distribution of the original beam.

For most holograms, even if the contrast of the recorded fringes is large, not all of the light in the reconstructing beam can be diverted into the reconstructed image. The ratio of the amount of optical power in the reconstructed image to the total optical power in the reconstructing beam

is the diffraction efficiency of the hologram. For common holograms recorded on photographic film, the maximum diffraction efficiency is limited to 6.25%. The rest of the light power, 93.75%, passes through the hologram unaffected or ends up in beams that are called higher diffraction orders. These higher diffraction orders leave the hologram at different angles and are generally not useful. As will be noted in Section III, certain types of holograms, notably thick phase holograms, are able to eliminate the undiffracted light and the higher diffraction orders; they can produce a diffraction efficiency close to 100%.

Certain assumptions that are made in the explanation given above should now be discussed. First, it is assumed that both beams used for recording are the same optical wavelength, that is, the same color. This is necessary to produce a stationary interference pattern to record. It is also assumed that the wavelength of the reconstructing beam is the same as that of the recording beams. This is often the case, but is not necessary. Using a reconstructing beam of a different wavelength changes the size of the reconstructed image: a shorter wavelength produces a smaller image, and a longer wavelength produces a larger image. Oddly, the dimensions of the image parallel to the plane of the hologram scale linearly with wavelength, but the depth dimension of the image scales proportional to the square of the wavelength, so three-dimensional images reconstructed with a different wavelength will appear distorted. Also, if the change in wavelength is large, distortions will occur in the other dimensions of the image, and the diffraction efficiency will decrease.

Another assumption is that the two recording beams have the same polarization. If beams with different polarizations or unpolarized beams are used, the contrast of the fringes, and therefore the diffraction efficiency of the hologram, is decreased. Beams that are linearly polarized in perpendicular directions cannot interfere and so cannot record a hologram. The polarization of the reconstructing beam usually matters very little in the quality of the reconstruction. An exception to this is when the hologram is recorded as an induced birefringence pattern in a material. Then the polarization of the reconstructing beam should be aligned with the maximum variation in the recording material index of refraction.

It is also possible that the reconstructing beam is not the same as one of the recording beams, either in its phase variation, its amplitude variation, or both. If the structure of the reconstructing beam differs significantly from both of the recording beams, the reconstructed image is usually garbled. One particular case where the image is not garbled is if the reconstructing beam has the same relative phase structure as one of the recording beams, but approaches the hologram at a different angle. In this case, provided the angle is not too large and that the recording

behaves as a thin hologram (see Section III for an explanation of thin and thick holograms), the reconstructed image is produced.

Photographic film is assumed to be the recording material used in the example above. Many other materials can be used to record holograms, and these are the subject of Section IV. Photographic film gives a particular type of recording, classified as a thin absorption hologram. The meaning and properties of different classifications of holograms are dealt with in Section III.

III. CLASSIFICATION OF HOLOGRAMS

Many different types of holograms are possible. Holograms are classified according to the material property in which the interference pattern is recorded, the diffraction characteristics of the hologram, the orientation of the recording beams with respect to the hologram, and the optical system configuration used for recording and reconstructing the hologram.

Holograms are recorded by exposing an optically sensitive material to light in the interference pattern produced by optical beams. For example, exposing photographic film to light triggers a chemical reaction that, after development, produces a variation in the optical absorption of the film. Portions of the film exposed to high optical intensity become absorbing, and unexposed portions of the film remain transparent. Other materials also exhibit this characteristic of changing their optical absorption in response to exposure to light. Holograms that result from interference patterns recorded as variations in material absorption are known as amplitude holograms.

There are also materials whose index of refraction changes in response to exposure to light. These materials are usually quite transparent, but the index of refraction of the material increases or decreases slightly where it is exposed to light. Holograms that result from interference patterns recorded as index-of-refraction variations are known as phase holograms. During reconstruction, light encountering regions with a higher index of refraction travels more slowly than light passing through lower index regions. Thus, the phase of the light is modified in relation to the recorded interference pattern.

It is not correct to assume that amplitude holograms can reconstruct wavefronts with only amplitude variations and phase holograms can reconstruct wavefronts with only phase variations. In Section II it was explained that wavefront direction (i.e., phase) is recorded by interference fringe spacing, and wavefront amplitude is recorded by interference fringe contrast. In reality, both amplitude and phase types of holograms are capable of recording wavefront amplitude and phase information.

Holograms are also classified as being “thin” or “thick”. These terms are related to the diffraction characteristics of the hologram. A thin hologram is expected to produce multiple diffraction orders. That is, although only two beams may have been used for recording, a single reconstructing beam will give rise to several reconstructed beams, called diffraction orders. Another property associated with thin holograms is that if the angle at which the reconstructing beam approaches the hologram is changed, the hologram continues to diffract light, with little change in diffraction efficiency. The diffraction orders will rotate in angle as the reconstructing beam is rotated. Thick holograms, on the other hand, produce only a single diffracted beam; a portion of the reconstructing beam may continue through the hologram in its original direction as well. Also, noticeable diffraction efficiency for thick holograms occurs only if the reconstructing beam is incident on the hologram from one of a discrete set of directions, called the Bragg angles. If the beam is not at a Bragg angle, it passes through the hologram and no diffracted beam is produced. The property of thick holograms that diffraction efficiency falls off if the reconstructing beam is not at a Bragg angle is called angular selectivity. Many thick holograms can be recorded in the same material and reconstructed separately by arranging for their Bragg angles to be different.

The terms thin and thick were originally applied to holograms based solely on the thickness of the recording material. The situation is, in fact, more complicated. Whether a particular hologram displays the characteristics associated with being thick or thin depends not only on the thickness of the recording material, but also on the relative sizes of the optical wavelength and the interference fringe spacing and on the strength of the change produced in the absorption or refractive index of the material.

The next category of hologram classification has to do with the arrangement of the recording beams (and therefore the reconstructing and reconstructed beams) with respect to the recording material. When two planewave beams produce interference fringes, the fringes form a set of planes in space. The planes lie parallel to the bisector of the angle between the beams, as shown in Fig. 2a. If the recording material is arranged so that both recording beams approach it from the same side, fringes are generally perpendicular to the material surfaces, as shown in Fig. 2b, and a transmission-type hologram is formed. During readout of the transmission hologram, the reconstructing and the reconstructed beams lie on opposite sides of the hologram, as in Fig. 2c. Alternatively, the recording material can be arranged so that the recording beams approach it from opposite sides. In this case, the fringes lie parallel to the surfaces of the material, as shown in Fig. 2d, and a reflection-type hologram is formed. For a reflection hologram, the reconstructing and the reconstructed

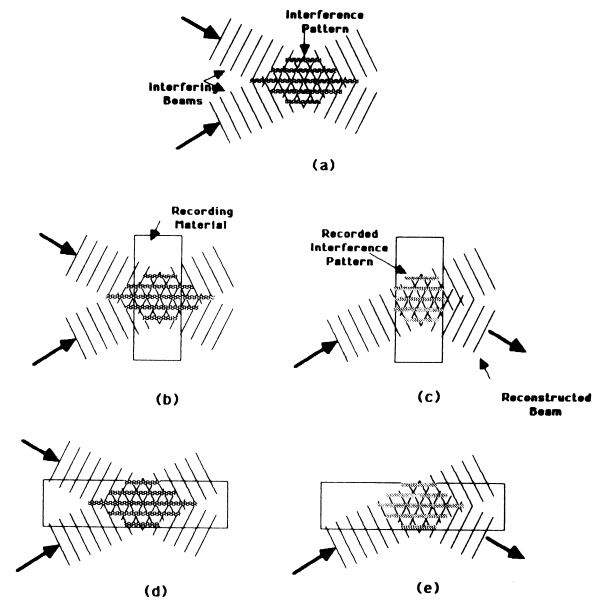


FIGURE 2 (a) The interference produced by two planewave beams. (b) The orientation of the recording material for a transmission hologram. (c) The configuration used for reconstruction with a transmission hologram. (d) The orientation of the recording material for a reflection hologram. (e) The configuration used for reconstruction with a reflection hologram.

beams lie on the same side of the hologram, portrayed in Fig. 2e.

The three hologram classification criteria discussed so far—phase or absorption, thick or thin, and transmission or reflection—play a role in determining the maximum possible diffraction efficiency of the hologram. Table I summarizes the diffraction efficiencies for holograms possessing various combinations of these characteristics. Because the fringes for reflection holograms lie parallel to the surface, there must be an appreciable material thickness to record the fringes; therefore, thin reflection holograms are not possible and are absent from Table I. (Often, holograms use reflected light but have fringes perpendicular to the material surfaces and are properly classified

TABLE I Maximum Diffraction Efficiencies of Hologram Classes

Thickness	Modulation	Configuration	Maximum efficiency (%)
Thin	Absorption	Transmission	6.25
Thin	Phase	Transmission	33.90
Thick	Absorption	Transmission	3.70
Thick	Absorption	Reflection	7.20
Thick	Phase	Transmission	100.00
Thick	Phase	Reflection	100.00

as transmission holograms.) Notice that phase holograms are generally more efficient than absorption holograms, and thick phase holograms are able to achieve perfect efficiency; all of the reconstructing beam power is coupled into the reconstructed beam. Keep in mind that the figures in the table represent the absolute highest diffraction efficiencies that can be achieved. In practice, hologram diffraction efficiencies are often substantially lower.

The final classification criterion to be discussed has more to do with the configuration of the optical system used for recording than with the hologram itself. For recording a hologram, coherent light is reflected from or transmitted through an object and propagates to the position of the recording material. A second beam of coherent light produces interference with light from the object, and the interference pattern is recorded in the material. If the object is very close to the recording material or is imaged onto the recording material, then an image-plane hologram is formed. If the separation between the object and the recording material is a few times larger than the size of the object or the material, a Fresnel hologram is formed. If the object is very far from the recording material, a Fraunhofer hologram is recorded. Another possible configuration is to place a lens between the object and the recording material such that the distance between the lens and the material is equal to the focal length of the lens. This arrangement produces a Fourier transform hologram, so called because during reconstruction the Fourier transform of the reconstructed wave must be taken (with a lens) to reproduce the recorded image. These previous configurations use a planewave as the second recording beam. If an expanding spherical wave is used instead, with the source of the spherical wave near the location of the object, a quasi-Fourier transform hologram, or lensless Fourier transform hologram, is formed. This type of hologram shares many of the properties of the true Fourier transform hologram.

Other classes of holograms also exist, such as polarization holograms, rainbow holograms, synthetic holograms, and computer-generated holograms. However, these are specialized topics that are best dealt with separately from the central concepts of holography.

IV. RECORDING MATERIALS

There are many materials that can be used for recording holographic interference patterns. Some materials record the pattern as a change in their optical absorption; this yields absorption holograms. Other materials record the patterns as changes in their index of refraction or as a relief pattern on their surface; these materials produce phase holograms. The thickness of the recording layer is also im-

portant to the characteristics of the hologram, as discussed in Section III.

Practical concerns related to holographic recording materials include the recording resolution, the material sensitivity as a function of optical wavelength, and the processing steps required to develop the hologram. The spacing of interference fringes can be adjusted by changing the angle between the beams: a small angle gives large fringes, and a large angle gives fringes as small as one-half the wavelength of the light used. An ideal recording material would have a resolution of at least 5000 fringes per millimeter.

Some materials are sensitive to all visible wavelengths, and others are sensitive to only a portion of the spectrum. The wavelength sensitivity of the recording material must be matched to the light source used. Sensitive recording materials are generally desirable, since high sensitivity reduces the recording exposure time and the amount of optical power required of the source. Long exposures are undesirable because of the increased chance that a random disturbance will disrupt the coherence of the recording beams.

Many recording materials require some chemical processing after exposure to develop the holographic pattern. Some materials develop when heat is applied, and a few materials require no processing at all: the hologram is immediately available. Of course, the need for developing complicates the system and introduces a delay between recording the hologram and being able to use it. Important characteristics of the most common hologram recording materials are summarized in [Table II](#).

Silver halide photographic emulsions are the most common recording material used for holograms. They are a mature and commercially available technology. The emulsion may be on a flexible acetate film or, for greater precision, a flat glass plate. Photographic emulsions have a very high sensitivity and respond to a broad spectral range. The ordinary development procedure for photographic emulsions causes them to become absorptive at the locations that have been exposed to light. Thus, an absorption hologram is produced. Alternate developing procedures employing bleaches leave the emulsion transparent but modulate the index of refraction or the surface relief of the emulsion. These processes lead to phase holograms. Many photographic emulsions are thick enough to produce holograms with some characteristics of a thick grating.

Dichromated gelatin is one of the most popular materials for recording thick phase holograms. Exposure to light causes the gelatin molecules to crosslink. The gelatin is then washed, followed by dehydration in alcohol. Dehydration causes the gelatin to shrink, causing cracks and tears to occur in the regions of the gelatin that are not crosslinked. The cracks and tears produce a phase change

TABLE II Holographic Recording Materials

Material	Modulation	Sensitivity (J/cm ²)	Resolution (line pairs/mm)	Thickness (μ m)
Photographic emulsion	Absorption or phase	$\sim 5 \times 10^{-5}$	~ 5000	< 17
Dichromated gelatin	Phase	$\sim 7 \times 10^{-2}$	> 3000	12
Photoresist	Phase	$\sim 1 \times 10^{-2}$	~ 1000	> 1
Photopolymer	Phase	$\sim 1 \times 10^{-2}$	3000	3–150
Photoplastic	Phase	$\sim 5 \times 10^{-5}$	> 4100	1–3
Photochromic	Absorption	~ 2	> 2000	100–1000
Photorefractive	Phase	~ 3	> 1000	5000

in light passing through those regions. Phase holograms recorded in dichromated gelatin are capable of achieving diffraction efficiencies of 90% or better with very little optical noise. The primary limitations of dichromated gelatin are its very low sensitivity and the undesirable effects of shrinkage during development.

Photoresists are commonly used in lithography for fabrication of integrated circuits but can be utilized for holography too. Negative and positive photoresists are available. During developing, negative photoresists dissolve away in locations that have not been exposed to light, and positive photoresists dissolve where there has been exposure. In either case, a surface relief recording of the interference pattern is produced. This surface relief pattern can be used as a phase hologram either by passing light through it or by coating its surface with metal and reflecting light off it. The photoresist can also be electroplated with nickel, which is then used as the master for embossing plastic softened by heat. The embossing process can be done rapidly and inexpensively and therefore is useful for mass producing holograms for use in magazines and other large-quantity applications.

Photopolymers behave in a fashion similar to photoresists, but instead of dissolving away during development, exposure of a photopolymer to light induces a chemical reaction in the material that changes its index of refraction or modulates its surface relief. Some photopolymers require no development processing, and others must be heated or exposed to ultraviolet light.

Photoplastics are noted for their ability to record and erase different holographic patterns through many cycles. The photoplastics are actually multilayer structures. A glass substrate plate is coated with a conductive metal film. On top of this is deposited a photoconductor. The final layer is a thermoplastic material. For recording, a uniform static electric charge is applied to the surface of the thermoplastic. The voltage drop due to the charge is divided between the photoconductor and the thermoplastic. The structure is then exposed to the holographic interference pattern. The voltage across the illuminated portions of the

photoconductor is discharged. Charge is then applied a second time to the surface of the device. This time excess charge accumulates in the regions of lowered voltage. The device is now heated until the thermoplastic softens. The electrostatic attraction between the charge distribution on the surface of the thermoplastic and the conductive metal film deforms the plastic surface into a surface relief phase hologram. Cooling the plastic then fixes it in this pattern. The hologram may be erased by heating the plastic to a higher temperature so that it becomes conductive and discharges its surface.

Photochromics are materials that change their color when exposed to light. For example, the material may change from transparent to absorbing for a certain wavelength. This effect can be used to record absorption holograms. Furthermore, the recording process can be reversed by heating or exposure to a different wavelength. This allows patterns to be recorded and erased. These materials, however, have very low sensitivity.

Photorefractive materials alter their refractive index in response to light. These materials can be used to record thick phase holograms with very high diffraction efficiency. This recording process can also be reversed, either by uniform exposure to light or by heating. These materials, too, have rather low sensitivity, but research is continuing to produce improvements.

V. APPLICATIONS

Holography is best known for its ability to reproduce three-dimensional images, but it has many other applications as well. Holographic nondestructive testing is the largest commercial application of holography. Holography can also be used for storage of digital data and images, precise interferometric measurements, pattern recognition, image processing, and holographic optical elements. These applications are treated in detail in this section.

An image reconstructed from a hologram possesses all the three-dimensional characteristics of the original scene.

The hologram can be considered a window through which the scene is viewed. As described in Section II, a hologram records information about the intensity and direction of the light that forms it. These two quantities (along with color) are all that the eye uses to perceive a scene. The light in the image reconstructed by the hologram has the same intensity and direction properties as the light from the original scene, so the eye sees an image that is nearly indistinguishable from the original. There are two important aspects in which the holographic reconstruction of a scene differs from the original: color and speckle. Most holograms are recorded using a single wavelength of light. The image is reconstructed with this same wavelength, so all objects in the scene have the same color. Also, because of the coherence properties of laser light, holographic images do not appear smooth: they are grainy, consisting of many closely spaced random spots of light called speckles. Attempts to produce full-color holograms and eliminate speckle will be described in this section.

The simplest method of recording an image hologram is shown in Fig. 3a. Light reflecting off the object forms one recording beam, and light incident on the film directly from the laser is the other beam needed to produce the interference pattern. The exposed film is developed and then placed back in its original position in the system. The object that has been recorded is removed from the system, and the laser is turned on. Light falling on the developed hologram reconstructs a virtual image of the object that can be viewed by looking through the hologram toward the original position of the object, as shown in Fig. 3b.

Many variations on the arrangement described above are possible. Often, the portion of the beam directed to-

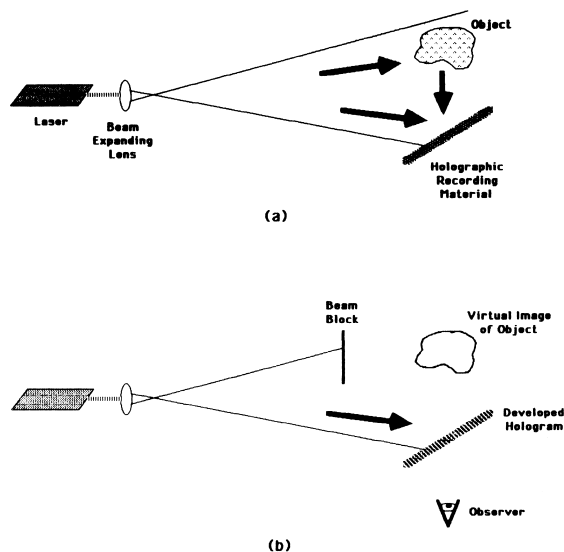


FIGURE 3 A simple optical system for (a) recording and (b) viewing a transmission hologram.

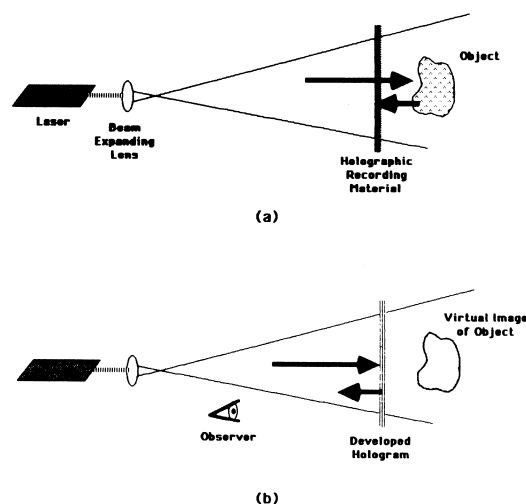


FIGURE 4 A simple optical system for (a) recording and (b) viewing a reflection hologram.

ward the object is split into several beams that are arranged to fall on the object from different directions so that it is uniformly illuminated. Attention to lighting is especially important when the scene to be recorded consists of several separate objects that are different distances from the film. Another recording arrangement is to have the laser beam pass through the film plane before reflecting from the object, as shown in Fig. 4a. A reflection hologram is formed with this arrangement and can be viewed as shown in Fig. 4b. An important advantage of the reflection hologram is that a laser is not needed to view it: a point source of white light will work quite well. Use of white light to reconstruct a hologram is not only more convenient, but it eliminates speckle too. This is possible because the reflection hologram also acts as a color filter, efficiently reflecting light of only the proper wavelength. It is important, however, that the source be very small (as compared to its distance from the hologram), otherwise the reconstructed image will be blurred.

A hologram that reconstructs an image containing all the colors of the original scene would, obviously, be a desirable accomplishment. The principal obstacle in achieving this goal is that holographic interference patterns result only when very monochromatic light is used. Recording a hologram with a single color and reconstructing the image with several colors does not work. Reconstructing a hologram with a wavelength other than the one used to record it changes the size and direction of the reconstructed image, so the variously colored images produced with white-light illumination are not aligned. It is possible to record three different holograms of the same object on one piece of film, using a red, a blue, and a green laser. The three holograms can be reconstructed simultaneously with the appropriate colors, and a reasonable representation of

the colors of the original object is produced. The problem, however, is that each reconstructing beam illuminates not only the holograms it is intended for, but the other two holograms as well. This leads to numerous false-color images mingled with the desired image. The false images can be eliminated with a thick hologram recording material by using the angular selectivity of thick holograms.

Another approach to color holography is the stereogram. Sets of three separate holograms (one each for red, green, and blue) are recorded for light coming from the scene at various angles. A projection system using a special screen is used for viewing. Light from different angular perspectives is directed at each eye of the viewer, thus providing the parallax information needed to yield a three-dimensional image. This system requires that the viewer is positioned rather exactly with respect to the screen.

Interferometry is a means of making precise measurements by observing the interference of optical wavefronts. Since holograms record phase information about the light from an object, they are useful in making before and after comparisons of the deformation of objects in response to some change in their environment. A typical arrangement for holographic interferometry is shown in Fig. 5. The first step is to record and develop a hologram of the object. The hologram is replaced in the system, and the image of the object is reconstructed from it. The object itself is subjected to some change and illuminated as it was to record the hologram. When viewing the object through the hologram, light from the object and from the holographic reconstruction of the object will interfere. If the surface of the object has deformed, bright and dark fringes

will appear on the surface of the object. Each fringe corresponds to a displacement in the surface of the object by one-quarter of the optical wavelength. Thus, this is a very sensitive measurement technique.

The same principles apply to transparent objects that undergo an index-of-refraction change. Changes in index of refraction alter the optical path length for light passing through the object. This leads to interference fringes when the object is viewed through a holographic recording of itself. Such index-of-refraction changes can occur in air, for example, in response to heating or aerodynamic flows.

A slight modification to the technique described above can provide a permanent record of the interference pattern. A hologram of the original object is recorded but not developed. The object is subjected to some change, and a second recording is made on the same film. The film is now developed and contains superimposed images of both the original and the changed object. During reconstruction, wavefronts from the two images will interfere and show the desired fringe pattern.

Another variation of holographic interferometry is useful for visualizing periodic vibrations of objects. Typically, the object is subjected to some excitation that causes its surface to vibrate. A hologram of the vibrating object is made with an exposure time that is longer than the period of the vibration. A time integration of the lightwave phases from locations on the object is recorded on the film. Portions of the object that do not move, vibrational nodes, contribute the same lightwave phase throughout the exposure. This produces constructive interference, and these locations appear bright in the reconstruction. Bright and dark fringes occur elsewhere on the object, with the number of fringes between a particular location and a vibrational node indicating the amplitude of the vibration.

Deformations on the surface of an object in response to applied pressure or heating can often be used to make determinations concerning changes within the volume of the object under the surface. For this reason, holographic interferometry is useful for holographic nondestructive testing; characteristics of the interior of an object can be determined without cutting the object apart. If interference fringes are concentrated on one area of a stressed object, that portion of the object is probably bearing a greater portion of the stress than other locations; or a pattern of fringes may indicate a void in the interior of the object or a location where layers in a laminated structure are not adhering.

Holography plays an important role in the field of optical data processing. A common optical system using a holographic filter is shown in Fig. 6. One application of this system is pattern recognition through image correlation. This is useful for detecting the presence of a reference object in a larger scene. The first step is to record

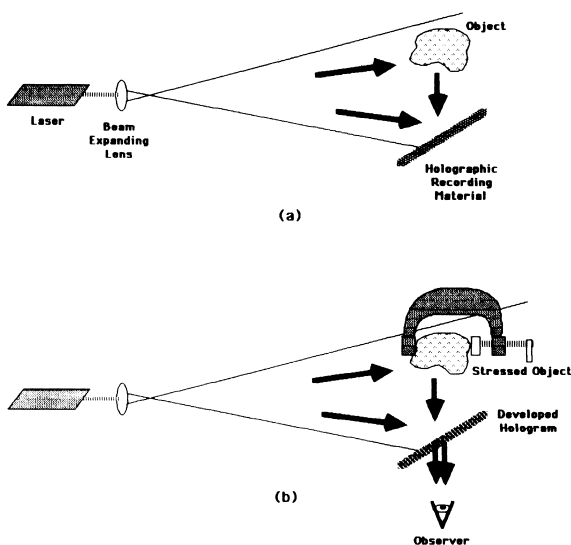


FIGURE 5 (a) For real-time holographic interferometry, a hologram of the object is recorded. (b) Then the stressed object is viewed through the developed hologram.

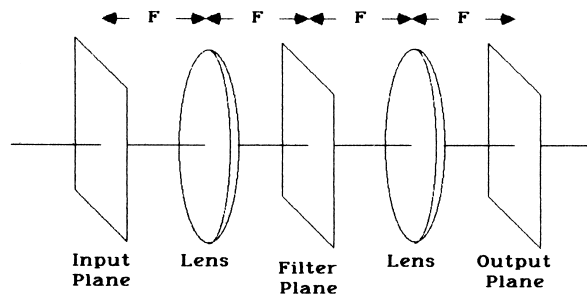


FIGURE 6 A simple, yet versatile, optical system for using holographic filters for image processing.

a Fourier hologram of the object that is to be detected. A transparency of the object is placed in the input plane and is holographically recorded onto photographic film in the filter plane; a reference beam is also incident on the film. After exposure, the film is removed from the system, developed, and then replaced in the filter plane. The scene believed to contain occurrences of the reference object is placed in the input plane and illuminated. A portion of the light distribution in the output plane will represent the correlation of the input image with the reference image. This means that wherever the reference image occurs in the input scene, a bright spot of light will be present at the corresponding position on the output plane. As a specific example, the reference object might be a printed word. A hologram of the word would be placed in the filter plane. Input to the system could be pages of printed text. For each occurrence of the selected word within the text, a bright spot would be present at the output plane. In the simple form described, the system can detect occurrences of the reference object that are very close to the same size and rotational orientation as the object used to record the hologram. However, research is being conducted to produce recognition results without regard to scaling or rotation of the object and has already met with considerable success.

Other image-processing operations can be performed using the same optical system. The hologram is recorded to represent a frequency-domain filter function rather than a reference object. For example, blurring of an image produced by an unfocused optical system can be modeled as a filtering of the spatial frequencies in the image. Much of the detail can be restored to a blurred image by placing a hologram representing the inverse of the blurring function in the filter plane and a transparency of the blurred image in the input plane. Holographic filters can also be used to produce other useful image-processing operations, such as edge enhancement.

Holography is also attractive for data storage. Because holograms record information in a distributed fashion, with no one point on the hologram corresponding to a particular part of the image, data recorded in holographic

form are resistant to errors caused by defects in the recording material. Returning to the analogy of the hologram as a window through which to view an image, if part of the hologram is obscured or destroyed, it is still possible to recover all the data simply by looking through the remaining usable portion of the window. Obviously, if part of an ordinary photograph containing data is destroyed, the data on that portion of the photograph is irrevocably lost. The data stored in a holographic system can be pages of text, images, or digitally encoded computer data. If thick recording materials are used, many pages of data can be stored in one piece of material by utilizing the angular selectivity of thick holograms.

Holograms can be recorded to provide the same functions as refractive optical elements such as lenses and prisms. For example, if a hologram is recorded as the interference of a planewave and a converging spherical wave, when the developed hologram is illuminated with a planewave it will produce a converging spherical wave, just as a positive lens would. Holographic optical elements (HOEs), as they are called, have two principal disadvantages with respect to the more common refractive elements. HOEs work as designed only for one wavelength of light, and because they usually have a diffraction efficiency less than 100%, not all of the light is redirected as desired. However, HOEs also have advantages over refractive elements in certain applications. First, the optical function of a HOE is not linked to its physical shape. A HOE may be placed on a curved or angled surface. For example, HOEs are used on the surface of the visor on pilots' helmets to serve as a projection lens for instrumentation displays. Also, HOEs can be created that provide the function of refractive elements that would be very difficult to fabricate, such as an off-axis segment of a lens. Several HOEs can be recorded on the same piece of material to give the function of multiple refractive elements located at the same spatial position. Another popular configuration is to record a HOE on the surface of an existing lens. The lens takes care of most of the refraction required by an optical system, and the HOE provides small additional corrections to reduce aberrations. A second advantage of HOEs is their small physical volume and weight. A HOE recorded on film can provide the same function as a thick, heavy piece of glass. In particular, lenses with a large aperture and short focal length can be quite massive when fabricated in glass, while those characteristics are readily produced using a HOE. Finally, there are applications in which the limited range of working wavelengths of HOEs is desirable. Reflective holograms are often used in these applications. Information presented in the intended color can be reflectively imaged toward a viewer, while the scene behind the hologram, containing mainly other colors, is also visible through the hologram without distortion.

The applications just cited are only a sampling of the uses of holography. It is a field that is the subject of ongoing research and that enjoys the continuing discovery of new forms and applications.

VI. HISTORY OF HOLOGRAPHY

The fundamental principle of holography, recording the phase of a wavefront as an intensity pattern by using interference, was devised by Dennis Gabor in 1948 to solve a problem with aberrations in electron microscopy. He also coined the word hologram from Greek roots meaning whole record. The results of Gabor's work of translating electron diffraction patterns into optical diffraction patterns and then removing aberrations from the resulting image were less than satisfactory. Other researchers following his lead were not significantly more successful, so this first application of holography soon faded from the scene.

Other researchers experimented with optical holography for its own sake during the early and mid-1950s, but results were generally disappointing. At that time, holography suffered from two important disadvantages. First, no truly interesting application had been found for it. The rather poor-quality images it produced were laboratory curiosities. Second, coherent light to record and view holograms was not readily available. The laser was not available until the early 1960s. Coherent light had to be produced through spectral and spatial filtering of incoherent light.

In 1955 an interesting, and eventually very successful, application for holography was uncovered. Leith and Upatnieks, working at the University of Michigan's Radar Laboratory, discovered that holography could be used to reconstruct images obtained from radar signals. Their analysis of this technique led them to experiment with holography in 1960, and by 1962 they had introduced an important improvement to the field. Gabor's holograms

had used a single beam to produce holograms, recording the interference of light coming only from various locations on the object. This led to reconstructed images that were seriously degraded by the presence of higher diffraction orders and undiffracted light. Leith and Upatnieks introduced a second beam for recording. The presence of this second beam separated the reconstructed image from other undersired light, thus significantly improving the image quality.

Also in the early 1960s, Denisjuk was producing thick reflection holograms in photographic emulsions. These holograms have the advantage that they can be viewed with a point source of white (incoherent) light, since the reflection hologram acts also as a color filter.

Advances in practical applications of holography began to occur in the mid-1960s. In 1963, Vander Lugt introduced the form of holographic matched filter that is still used today for pattern recognition. Powell and Stetson discovered holographic interferometry in 1965. Other groups working independently introduced other useful forms of this important technique.

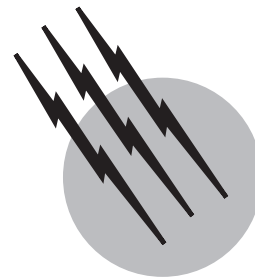
The first form of thin transmission hologram that can be viewed in white light, the rainbow hologram, was developed by Benton in 1969. Another form of hologram that can be viewed in white light, the multiplex hologram, was produced by Cross in 1977.

SEE ALSO THE FOLLOWING ARTICLES

COLOR SCIENCE • LASERS • OPTICAL DIFFRACTION • PHOTOGRAPHIC PROCESSES AND MATERIALS

BIBLIOGRAPHY

- Caulfield, H. J., ed. (1979). "Handbook of Optical Holography," Academic Press, New York.
- Goodman, J. W. (1968). "Introduction to Fourier Optics," McGraw-Hill, New York.



Light-Emitting Diodes

Robert Weissman

Agilent Technologies, Inc.

Gina Christenson

Patrick N. Grillo

Yu-Chen Shen

Dawnelle Wynne

LumiLeds Lighting LLC

- I. Material and Device Physics
- II. Fabrication Techniques and Processes
- III. Infrared LEDs
- IV. Visible LEDs

GLOSSARY

Bandgap The energy difference between the top of the highest filled band (valence band) and the bottom of the lowest empty band (conduction band).

Efficiency The fraction of input energy that is converted into useful output energy. Internal quantum efficiency is the number of photons generated inside the light-emitting diode (LED) per number of input electrons. External quantum efficiency is the number of photons that escape from the device per number of input electrons. Extraction efficiency is the ratio of the external quantum efficiency to the internal quantum efficiency, or the fraction of photons generated in the semiconductor that escape from the device. Wall-plug efficiency is the fraction of electrical power into the LED that is converted to optical power out of the device. Luminous efficiency is the power sensed by the human eye per electrical power input into the device.

Epitaxial growth Growth of a crystal on a host crystal such that they have similar structures.

Heterostructure Two or more semiconductors with different physical properties (such as their bandgaps) grown on top of each other.

Lattice constant A number specifying the size of the basic unit in a crystal.

N-type semiconductor (p-type semiconductor) A semiconductor containing impurity atoms in which the electron (hole) is the primary charge carrier.

Photon A quantum of electromagnetic radiation with energy $h\nu$, where h is Planck's constant and ν is the frequency of the electromagnetic radiation.

Radiative transition An electronic transition between energy bands that produces photons and possibly phonons (quantized lattice vibrations). A nonradiative transition is an electronic transition that produces only phonons.

LIGHT-EMITTING DIODES (LEDs) are semiconductor devices that emit light when a bias voltage is applied and current flows through the device. Injected electrons in

the semiconductor conduction band recombine with holes in the valence band, and photons are emitted. LEDs can be made with emission wavelengths ranging from the mid-infrared to the ultraviolet (encompassing the entire visible spectrum from red to blue). Compared to other types of light sources, LEDs use less power, are more compact, are mechanically rugged, and have longer operating lifetimes. Today, LEDs are used in a vast array of products in the home, office, factory, and field. Applications include illumination, information display, computer interconnects, sensors, and process control. A wide variety of package configurations is available to satisfy the diverse product specifications.

This chapter reviews the current state-of-the-art of LED technology used for commercially available products. The chapter begins with a discussion of the material and device physics of LEDs. A discussion of light generation in and extraction from LEDs follows. The engineering of materials and device structures is explained so the reader can better understand the various approaches in use today. Next the processes used to fabricate, test, and package LEDs are discussed. The major techniques for growing light-emitting semiconductor materials are outlined, and the basic processing steps required to fabricate LED chips are explained. The important topic of LED reliability is also addressed.

Many types of semiconductor materials and devices have been developed to make infrared and visible LEDs. The most important types of materials and devices, ranging from simple p–n junctions to high-performance, double-heterostructure LEDs, are discussed. Major LED-based products, such as indicator lamps, displays, illuminators, fiberoptic transmitters, emitter-detector products, and optocouplers are described.

I. MATERIAL AND DEVICE PHYSICS

A. Semiconductor Physics

The materials physics of LEDs can be understood through the band theory of solids, which is based on the time-independent Schrodinger equation. The potential energy generated by the protons and electrons in the crystal lattice can be input into the Schrodinger equation to determine the electronic structure. The electronic structure in a uniform crystal lattice consists of allowed and disallowed bands of energy. Electrons reside in the allowed bands of energy, and electronic transitions occur between the allowed bands. In an LED, photons are emitted when electrons in an energy band transition to a band with lower energy. The energy (and thus the color of the light) is determined by the energy gap, E_g , between the two bands. Simple mod-

els of band theory, such as the Kronig–Penney model, the Ziman model, and the Feynman model, give good qualitative pictures of the dependence of the bands on lattice spacing and the width and height of potential barriers. More sophisticated theories such as density functional theory and tight-binding methods predict more accurate electronic structures. These theories can include the effects of defects, surfaces, and impurities.

In order for LEDs to emit light, it is important that the electrons are able to move. At the temperature $T = 0$ K, when the electrons occupy the lowest states, some materials are conducting, and some are insulating. Conductive materials have energy bands that are not completely filled, while insulating materials have energy bands that are completely filled. At nonzero temperatures, the probability that a given energy state, E , is occupied is given by the Fermi function

$$f(E) = [1 + e^{\beta(E-E_f)}]^{-1}$$

where $\beta = 1/(k_B T)$ (where k_B is Boltzmann's constant), and E_f is the material-dependent Fermi level. At $T = 0$ K, the Fermi function indicates that energy levels less than or equal to E_f are occupied. At $T > 0$ K, electrons are thermally excited to a higher energy level. If the energy gap between the filled and unfilled bands is small enough for electrons to make transitions through thermal excitation, an insulating material at $T = 0$ K can become a conducting material at finite temperatures. Such materials are called semiconductors. In a semiconductor at $T = 0$ K, the highest energy band that is filled is called the valence band; the empty band above it is called the conduction band.

In a semiconductor, conduction occurs through the movement of electrons excited to the conduction band and through the movement of holes in the valence band. Holes are positively charged particles that represent the absence of electrons. The dynamical equations of an electron in a crystal in an applied electric field can be rewritten to resemble the equations of motion of a free electron in an electric field. To compensate for the effects of the crystal, the mass of the electron is replaced with an effective mass, m^* . The effective mass can be negative, representing a hole. A hole with a negative effective mass behaves in the same way as a positive charge and moves in the opposite direction as an electron in an electric field. Both electrons and holes contribute to the current flow.

Doping can be used to change the conduction properties of a semiconductor. By adding elements from neighboring columns on the periodic table, the number of conducting electrons and/or holes can be altered. If group IV silicon (with four valence electrons) is doped with a group V impurity (with five valence electrons) such as antimony (Sb), arsenic (As), or phosphorus (P), each group V atom will bond covalently to four silicon atoms. Since only four

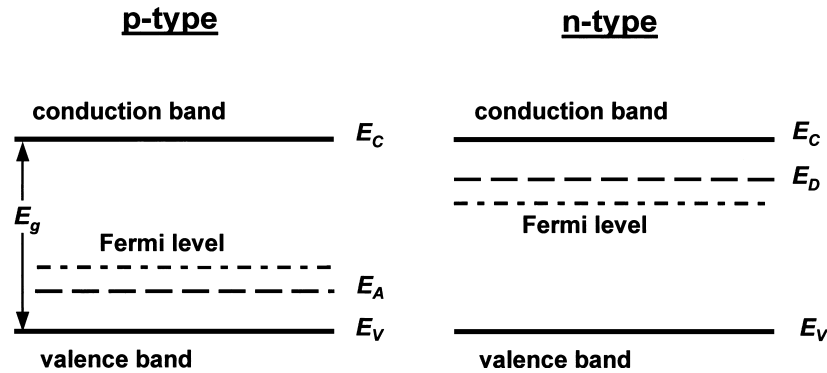


FIGURE 1 Band diagrams (energy vs. position) for a p-type and an n-type semiconductor. In the p-type semiconductor, the acceptor level, E_A , is slightly above the valence band, E_V . In the n-type semiconductor, the donor level, E_D , is slightly below the conduction band, E_C . The bandgap, E_g is the difference between E_C and E_V . The position of the Fermi level changes with temperature and doping.

electrons are necessary to complete the outer shell of electrons, the fifth electron on the group V atom will be a spare electron that is loosely bound to the atom. The amount of energy needed to ionize this spare electron is far less than for the four other dopant electrons which are involved in the bonding process, and at $T = 0$ K the electron will reside in an electronic level that is closer to the conduction band than to the valence band. This electronic level is called a donor level with energy E_D (see Fig. 1). Typically, $E_C - E_D$ is on the order of 0.01 eV (as compared to $1.0 \text{ eV} < E_g < 3.5 \text{ eV}$), and at higher temperatures the electron may have enough energy to jump into the conduction band. This electron has then been donated by an impurity (donor) atom and contributes to conduction. A semiconductor doped with donor atoms is called an n-type semiconductor.

On the other hand, if we had doped the silicon with a group III atom, a bonding electron will be missing since the group III atom can only contribute three valence electrons to the outer shell. Because there are only seven electrons in the outer shell (eight are needed to complete the four bonds), an acceptor electronic level with energy E_A is created near the valence band (see Fig. 1). An electron with extra energy can jump into the acceptor level to complete the outer shell of electrons (valence shell), creating a hole in the valence band. This hole has then been contributed to the valence band by an impurity (acceptor) atom. A semiconductor doped with impurity atoms that have insufficient electrons to complete the valence shell is called a p-type semiconductor.

Most inorganic semiconductors that are used for LEDs crystallize into a zincblende or a wurtzite crystal structure. A zincblende crystal can be visualized as a cube in which two atoms are located near each corner and two atoms are located near the center of each face of the cube. In the silicon crystal described above, both of the atoms at these

locations are silicon atoms. On the other hand, in a III–V semiconductor, one of these two atoms is a group III element such as gallium, aluminum, or indium, and the other atom is a group V element such as arsenic or phosphorous. A wurtzite crystal is similar to a zincblende crystal, except that the atoms are located in slightly different places.

In III–V compound semiconductors, common donor atoms that reside in group V sublattices are S, Se, and Te (group VI impurities), while common acceptor atoms residing on the group III sublattice include Cd, Zn, Mg, and Be. The donor states in III–V semiconductors are typically on the order of 5 meV below the conduction band edge, while the acceptor states are somewhat deeper (approximately 30 meV above the valence band edge). Group IV impurities (C, Si, and Ge) can act as either donors or acceptors depending on whether they reside on group III or group V sites, respectively.

Crystal defects, surfaces, and some impurities create deep levels in the bandgap. These levels are separated from the conduction and valence bands by more than approximately $5 k_B T$. Deep levels act as neither good acceptor nor good donor levels because they are not usually ionized at room temperature. Deep levels provide an alternative mechanism for the recombination of holes and electrons and thus affect radiative efficiency. Transitions where phonons are not involved are called direct transitions; indirect transitions involve phonons and are less radiatively efficient.

Figure 2 shows the energy band diagram of a direct bandgap semiconductor where E_Γ is lower in energy than E_X , and spontaneous emission of a photon is the most likely path for recombination. In Fig. 2, the Γ valley of the conduction band is directly above the holes where the momentum change in the transition is nearly zero. Electrons in the X valley of the conduction band have a different momentum and thus cannot directly recombine with

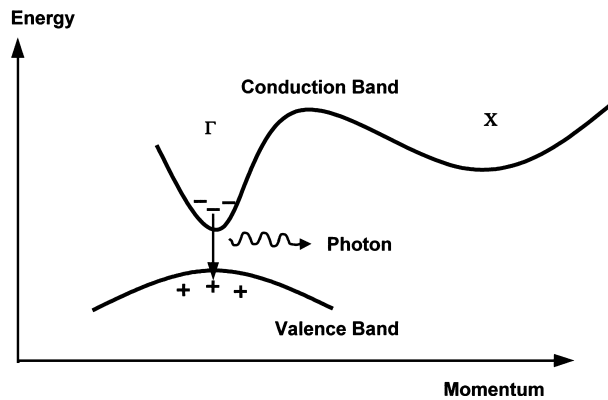


FIGURE 2 Energy vs. momentum band diagram for a direct bandgap semiconductor. ($E_X - E_\Gamma \gg k_B T$). The electrons (–) in the Γ valley of the conduction band recombine with the holes (+) in the valence band, and a photon is emitted to conserve energy.

valence band holes. In an indirect bandgap semiconductor, E_X is lower in energy than E_Γ (see Fig. 3), and recombination occurs through deep level states and involves phonons. If $E_X \sim E_\Gamma$, the indirect and direct transitions are both likely. Deep level states allow for alternate nonradiative paths. Thus, radiative efficiency increases with fewer deep level states. Electrons in the Γ valley may recombine nonradiatively with holes through deep electronic levels induced by defects. This process is called the Hall–Shockley–Read (HSR) recombination. Since indirect transitions are less radiatively efficient than direct transitions, most III–V LED semiconductors are made from direct gap materials.

The wavelength (color) of light, λ , emitted from the device is determined by the bandgap ($E_g = hc/\lambda$), where c is the speed of light and h is Planck's constant. The

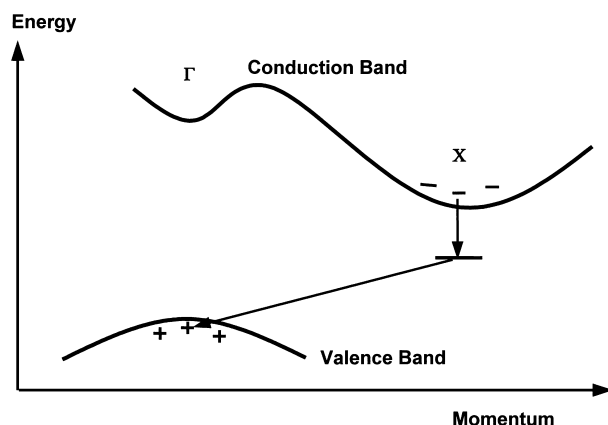


FIGURE 3 Energy vs. momentum band diagram for an indirect bandgap semiconductor ($E_\Gamma - E_X \gg k_B T$). Phonons are involved in the recombination of the holes from the valence band with the electrons in the X valley of the conduction band.

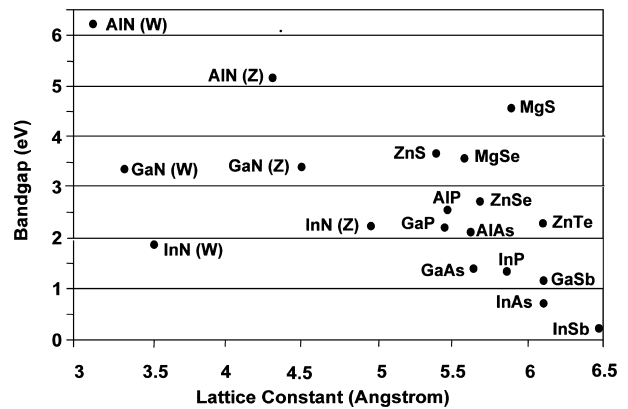


FIGURE 4 Bandgap vs. lattice constant for various III–V materials used in the manufacturing of red, yellow, green, and blue LEDs. W, wurtzite crystal structure; Z, zincblende crystal structure.

bandgap in turn is controlled by the lattice constant, a , of the crystal. A diagram of bandgap vs. lattice constant is given in Fig. 4 for III–V semiconductors commonly used for LEDs. For LEDs consisting of different III–V alloys grown on each other, the lattice constants must be similar to avoid creating defects. If the ratio of the difference of the lattice constants to the substrate lattice constant, $\Delta a/a$, is $\leq 0.1\%$, the crystal structures are said to be lattice-matched. If $\Delta a/a$ is too large, the crystal structure will be too strained, and defects will appear to relax the strain. These defects can create deep levels, which can decrease LED quantum efficiency.

Light-emitting diodes can be made from organic as well as inorganic materials. One type of organic LED is the polymer LED. A polymer is a chain of covalently bonded molecular units or monomers. Opposite ends of a monomer have an incomplete electron orbital that readily bonds to another monomer unit. In the condensed phase, polymers consist of extended chain structures. The polymers used for LED fabrication are *conjugated* polymers, which means that the bonds connecting the monomers are an alternating sequence of single and double (or triple) bonds. This alternating sequence of bonds is responsible for both the metallic and semiconducting character of this type of solid. Overlap of the P_z (and P_y) orbitals of the double (or triple) bonds results in the formation of a delocalized π -electron system along the polymer backbone. Bonding and anti bonding π -orbitals form the equivalent of valence and conduction bands, respectively, allowing conjugated polymers to support the positive and negative charge carriers with high mobilities along the polymer chain. However, the poor overlap of electron orbitals on neighboring long-chain molecules limit the carrier mobility of the polymer solid to low values, typically much less than for inorganic crystalline semiconductors.

Luminescent or photoabsorptive functional units (known as side chains) can be attached to the polymer backbone. The energy gap of polymers can be easily tuned by attaching electron-accepting side chains (which shift the luminescence to higher energies) or electron-donating side chains (which shift the luminescence to lower energies).

Another type of LED can be made from molecular organic materials. The molecules in the layers of this material are much smaller than the long-chain polymers. Thin films of these materials are chosen such that the mobility of one carrier (hole or electron) is much higher than that of the other carrier. Due to their typically low electron affinities, the majority of organic materials are hole-transporting solids. To tune the emission wavelength or to increase the luminous efficiency, the light-emitting layer is doped with organic dye molecules, allowing the transfer of energy from the host to the guest molecule.

Both the thin polymeric and molecular organic LED layers are typically amorphous to increase the radiative efficiency. Strongly coupled crystalline organic systems have many phonon modes associated with the crystal lattice that may interfere with the radiative recombination of the excitons (electron-hole pairs). However, even in amorphous polymeric or molecular organic films, there is strong exciton–phonon coupling energy (energy lost to phonon excitation), and the resulting spectral width of the emitted light is large, on the order of 100 nm. This presents a challenge to obtaining saturated colors from organic LEDs. In comparison, the spectral width of $\text{Ga}_x\text{In}_{1-x}\text{N}$ -based inorganic LEDs is on the order of 25 nm, and the spectral width of $(\text{Al}_x\text{Ga}_{1-x})_{0.5}\text{In}_{0.5}\text{P}$ -based LEDs is on the order of 10 nm. Due to the low intermolecular coupling and disorder of these amorphous organic materials, their carrier mobilities are lower than in most crystalline solids (10^{-2} to $10^{-8} \text{ cm}^2/\text{V s}$).

B. Device Physics

All III–V LEDs are p–n junction devices. A p–n junction device is fabricated by placing a p-type semiconductor next to an n-type semiconductor. Figure 1 shows a p-type semiconductor on the left and an n-type semiconductor on the right. Note that the Fermi level shifts with impurity concentration and with temperature. When the p- and n-type semiconductors are in contact, the valence and conduction bands adjust so that the Fermi level is constant in energy as a function of position (see Fig. 5). In the junction region (the transition region between the p and n materials), the conduction and valence bands bend monotonically, as shown in Fig. 5. If a forward bias, V_f , is applied as shown in Fig. 6, the barrier to conduction is lowered by eV_f , where e is the electronic charge ($1.60 \times 10^{-19} \text{ C}$). The current flowing from the p-type to the n-type

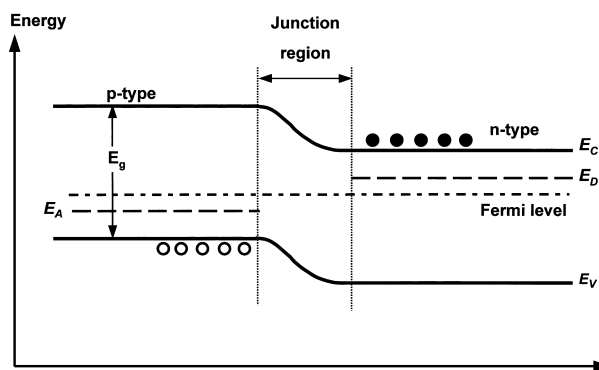


FIGURE 5 The thermal equilibrium band structure of a p-type semiconductor in contact with an n-type semiconductor. The conduction and valence bands bend in the junction region such that the Fermi level remains constant. The filled circles represent electrons and the unfilled circles represent holes.

material remains fixed at I_0 (there is no barrier); however, the current flowing from the n-type to the p-type material changes according to Boltzmann statistics. The total current is $I = I_0[\exp(eV_f/k_B T) - 1]$ and increases with increasing forward voltage. With a forward bias of V_f , the number of electrons moving into the p-type material is increased by $\exp(eV_f/k_B T)$, and typically an electron can be injected about $1 \mu\text{m}$ to 1 mm into the p-material where it recombines with holes. In the same way, holes are injected into the n-type material and recombine with the electrons to produce photons and phonons. If a reverse bias is applied, the barrier increases by eV_r , where V_r is the applied reverse voltage. The total current decreases with increasing V_r to the value $-I_0$.

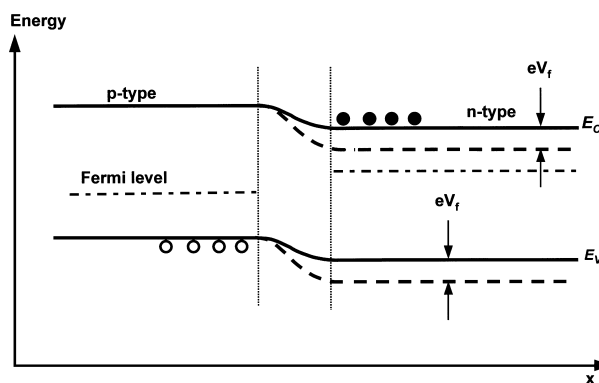


FIGURE 6 The band structure of a p–n junction when a bias voltage, V_f , is applied. The barrier to electrons moving from the n-type semiconductor to the p-type semiconductor (and the barrier to holes moving in the opposite direction) decreases by eV_f , allowing more electron injection into the p-type semiconductor and hole injection into the n-type material. The filled circles represent electrons and the unfilled circles represent holes.

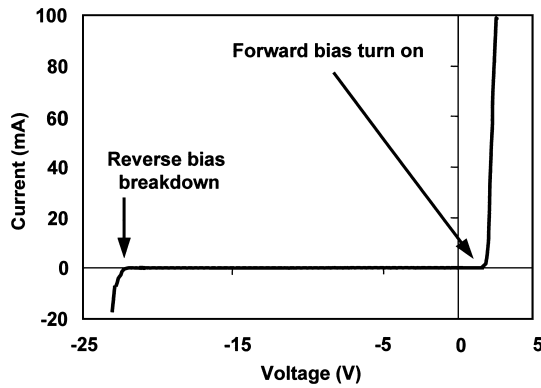


FIGURE 7 Current–voltage characteristic for a typical LED. In this example, the forward turn-on voltage is approximately 2 Volts, and the reverse bias breakdown is -22.5 Volts.

As the relationship $I = I_0[\exp(eV_f/k_B T) - 1]$ indicates, LEDs exhibit the typical rectifying current–voltage (I–V) relationship that is characteristic of a p–n junction diode (see Fig. 7). That is, LEDs have a sharp turn-on under forward bias, which varies from approximately 1.0 volts to approximately 3.5 volts, depending on the LED material. For bias values greater than this forward turn-on voltage, the LED easily conducts current in the forward direction, and typically exhibits a series resistance that is approximately a few ohms to several tens of ohms. In contrast, very little current flows under reverse bias until the LED reaches reverse bias breakdown. For the I–V characteristic in Fig. 7, the forward turn-on voltage is approximately 2 volts, and the reverse bias breakdown voltage is -22.5 volts. These values are typical for $(\text{Al}_x\text{Ga}_{1-x})_{0.5}\text{In}_{0.5}\text{P}$ LEDs, although they may be higher or lower for other materials, depending on parameters such as bandgap and doping concentrations.

In its simplest form, an LED consists of a p–n junction in a uniform composition material. Such an LED is referred to as a p–n homojunction LED, the band structure of which is illustrated in Fig. 8a. Since the homojunction LED is simple and inexpensive to manufacture, it dominates low-flux lighting applications such as seven-segment displays (see Section IV), and indoor message panels such as “EXIT” signs. More recent LED device designs, such as the double heterostructure (DH) LED or the multi-well active layer (MWAL) LED device design shown in Figs. 8b and c, respectively, have several advantages over the homojunction LED.

To better understand the advantages provided by the DH or MWAL device designs, we must first introduce several parameters that describe the ability of an LED to produce and emit light. The first of these parameters is the internal quantum efficiency, η_{int} . The internal quantum efficiency is simply a measure of the ability of the LED to

generate light or to turn electron-hole pairs into photons inside the LED. Although the LED efficiency is mathematically a unitless quantity, η_{int} can intuitively be thought of in units of photon/electron inside the LED. The internal quantum efficiency of an LED is determined by the radiative, τ_r , and nonradiative, τ_n , lifetimes of electrons and holes within the semiconductor, where the nonradiative lifetime is inversely related to the concentration of defects in the crystal. If the radiative and nonradiative lifetimes are known, one can compute the internal quantum efficiency of the LED through the equation:

$$\eta_{\text{int}} = \frac{\tau_n}{(\tau_r + \tau_n)}$$

Although the internal quantum efficiency is the most fundamental measure of the ability of an LED to generate light, the internal quantum efficiency is not commonly used to specify LED performance since the carrier lifetimes are very difficult to measure. In practice, the problem with trying to quantify the internal quantum efficiency is that not every photon generated inside the LED is emitted by the LED. Although it is difficult to count the number of photons that are generated inside the LED, it is relatively straightforward to count the number of photons that escape from the LED. The external quantum efficiency, η_{ext} , is defined as the ratio of the number of photons that escape from the LED to the number of electrons that are injected into the LED. η_{ext} is then related to η_{int} through the extraction efficiency, $C_{\text{ex}} = (\eta_{\text{ext}}/\eta_{\text{int}})$, where the extraction efficiency is the fraction of generated photons that escape from the LED.

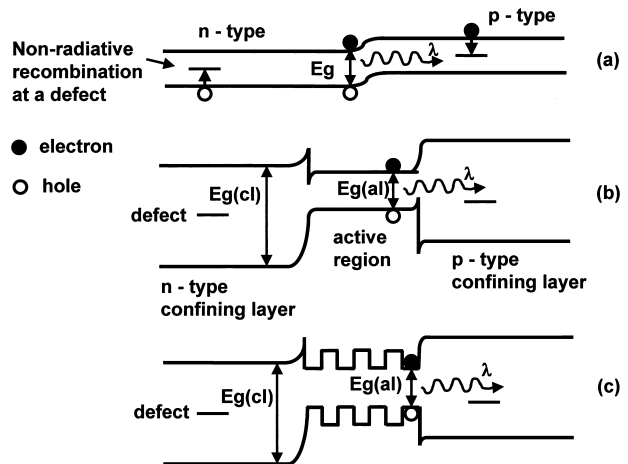


FIGURE 8 Energy band diagrams for: (a) a simple homojunction LED, (b) a double heterostructure LED, and (c) a multi-well active layer LED. Note that carriers can easily recombine at a defect in the homojunction LED, but the band structure of the double heterostructure and multi-well active layer LEDs prevents carriers from diffusing to defects outside of the active regions of these device designs. For parts (b) and (c), $E_g(\text{al})$ is the bandgap of the active region, and $E_g(\text{cl})$ is the bandgap of the confining layers.

Although the mathematical relationship between η_{ext} and η_{int} is trivial, the physical relationship between these parameters is very complicated, as will be described shortly. Before discussing this relationship in more detail, there are two other measures of LED performance that will be briefly introduced. The first concept is the wall plug efficiency, η_{wp} , which describes the ability of the LED to convert electrical input energy into the desired form of output energy (visible, infrared, or ultraviolet light). η_{wp} is related to η_{ext} through the forward voltage drop, V_f , across the diode, and the photon energy of the LED, $h\nu$, such that $\eta_{\text{wp}} = (h\nu/eV_f)\eta_{\text{ext}}$. Note that $h\nu = E_g$, and $V_f = V_j + I_f R_s$, where V_j is the voltage drop across the p-n junction of the LED, I_f is the forward current, and R_s is the series resistance of the LED. The proportionality factor, e , is the electronic charge previously introduced, and is necessary to convert the forward voltage from volts to electron volts. Typically, the energy $eV_f < E_g$, so $\eta_{\text{wp}} < \eta_{\text{ext}}$, although this is not always the case. From the above relationship, it can be seen that any increase in series resistance will decrease the ability of the LED to convert electrical energy into light energy, and η_{wp} will decrease. η_{wp} is measured in units of optical watt/electrical watt.

In addition to its influence on η_{wp} , the effects of series resistance must be properly accounted for when designing an LED circuit. The importance of this effect is partially illustrated in Fig. 9 for two $(\text{Al}_x\text{Ga}_{1-x})_{0.5}\text{In}_{0.5}\text{P}$ LEDs where the series resistance of LED 2 is seen to be greater than that of LED 1, as noted by the smaller slope of the I-V curve for LED 2. As a result, for a fixed voltage, the LED current will vary from LED to LED, due to differences in series resistance. Since the LED brightness is determined by the LED current, it is important to use a current source to drive LEDs, as opposed to the voltage sources that are used for most other electronic devices such as transistors. To further minimize the effect of variations in LED series

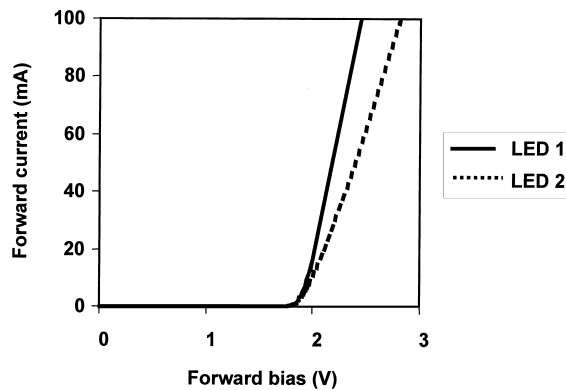


FIGURE 9 Forward-bias current-voltage curves for two LEDs with different series resistances.

resistance on circuit performance, an external resistor is often placed in series with the LEDs. If the value of this external series resistor is significantly greater than the series resistance of the LEDs themselves, the LED circuit will be further stabilized.

The final type of efficiency that is used to describe LED performance is the luminous efficiency, η_L . The luminous efficiency is a measure of the ability of the LED to generate light that can be seen by the human eye. η_L is related to η_{wp} as $\eta_L = 683 V(\lambda) \eta_{\text{wp}}$ where $V(\lambda)$ is the normalized C.I.E. curve (normalized eye response to light of wavelength λ), measured in units of lumens/watt.

We now return to the relationship between internal quantum efficiency and external quantum efficiency. As the preceding discussion implies, LED performance can be improved by increasing the internal quantum efficiency and/or the extraction efficiency. To illustrate this point, we consider a simple p-n homojunction LED, where the LED performance is limited by several mechanisms. First of all, the internal quantum efficiency of these devices is limited since the injected electrons and holes can diffuse through the semiconductor to any defects in the semiconductor and can recombine at these defects through nonradiative recombination, which does not generate light (see Fig. 8). By using a DH or MWAL structure, the injected electrons and holes are confined to the same spatial region, which is of limited volume. By confining the electrons and holes to the same spatial region, they tend to recombine more quickly than they would in a homojunction LED. Furthermore, in DH and MWAL LEDs, carrier confinement prevents carriers from reaching the defects that are outside of the active region. Both of these mechanisms tend to improve the internal quantum efficiency of DH and MWAL LEDs relative to homojunction LEDs.

In addition to these limitations on internal quantum efficiency, the extraction efficiency of homojunction LEDs is hindered by internal self-absorption in the LED. Since the composition of a homojunction LED is uniform, the bandgap of the LED is also uniform, and the photons emitted by the LED can be re-absorbed before they escape from the LED. Again, the DH and MWAL LED device designs have a major advantage over the homojunction LED with respect to self-absorption in that light generation in these DH and MWAL LEDs occurs in semiconductor materials that have a lower bandgap than the surrounding semiconductor materials. Once light escapes from the active region, it can therefore usually travel through the adjacent materials without being absorbed. So, in addition to the internal quantum efficiency benefits of the DH or MWAL LED, these device designs also result in increased light extraction and higher external quantum efficiency than the homojunction LED.

While self-absorption in the active region has an impact on extraction efficiency as described above, absorption in other regions of the LED can also play a major role in LED efficiency. One of the most dominant factors with respect to LED efficiency and internal absorption relates to the choice of substrate. For example, the LED may be grown on a substrate where the bandgap energy of the substrate is less than the photon energy of the light emitted by the LED. Such an LED is referred to as an absorbing substrate (AS) LED. Examples of such absorbing substrate LEDs include $\text{GaAs}_{1-x}\text{P}_x$, $\text{Al}_x\text{Ga}_{1-x}\text{As}$, or $(\text{Al}_x\text{Ga}_{1-x})_{0.5}\text{In}_{0.5}\text{P}$ LEDs grown on a GaAs substrate. Any light that is emitted toward the substrate in such a device will be absorbed by the substrate and lost. Since light is emitted isotropically from the active region of an LED, at least 50% of the light generated in an absorbing substrate LED will not be emitted by the LED but will be absorbed in the substrate. The choice of an absorbing substrate thus places an upper limit of 50% on the external quantum efficiency of the device, and, in reality, the limitations imposed by an absorbing substrate far exceed this 50% loss. In contrast to these absorbing substrate devices, the LED may be fabricated on a transparent substrate where the bandgap of the substrate is greater than the bandgap of the light-emitting active region. Although free carrier absorption and other absorption mechanisms may still occur within such transparent substrate (TS) LEDs, significant light output gains can be achieved by using a transparent substrate as opposed to an absorbing substrate. Examples of transparent substrate devices include $\text{Al}_x\text{Ga}_{1-x}\text{As}$ LEDs grown on $\text{Al}_y\text{Ga}_{1-y}\text{As}$ pseudo-substrates where $y > x$, $\text{Ga}_x\text{In}_{1-x}\text{N}$ LEDs grown on sapphire or SiC substrates, and $(\text{Al}_x\text{Ga}_{1-x})_{0.5}\text{In}_{0.5}\text{P}$ LEDs that have been wafer-bonded to GaP substrates.

Although the use of a transparent substrate can significantly increase LED light output, other factors must be considered which may significantly limit the LEDs external quantum efficiency. For example, the index of refraction of most LED materials is approximately 3 to 4, which is significantly greater than the index of refraction of air ($n_{\text{air}} = 1$). According to Snell's law, the critical angle for total internal reflection inside an LED is therefore approximately 15° , and any light ray that hits the surface of the LED at an angle greater than approximately 15° will be reflected inside the chip. For an absorbing substrate LED, this light will be absorbed by the substrate and will not be emitted by the LED. In a conventional absorbing substrate LED, it can be shown that only 2 to 3% of the light generated in the LED chip will be emitted by the chip, with the remainder of the light being absorbed in the LED. If this were the end of the story, LEDs would be forever limited to low flux applications. Fortunately, considerable improvements in LED performance can be achieved through a number of methods including encapsulation of the bare chip in an epoxy lens ($>2\times$ improvement), the use of a thick window layer above the active region ($3\times$ improvement), the use of a transparent substrate or distributed Bragg reflector ($2\times$ improvement), chip shaping or surface texturing techniques (1.5 to $2\times$ improvement), antireflection coatings, etc. By combining all of these approaches, the LED external quantum efficiency can approach or exceed 50%. Such improvements in efficiency have enabled LEDs to move from the low-flux realm to high-flux applications such as traffic signals, automotive lighting, and full-color outdoor displays. These applications will be discussed more fully in Section IV.

ulation of the bare chip in an epoxy lens ($>2\times$ improvement), the use of a thick window layer above the active region ($3\times$ improvement), the use of a transparent substrate or distributed Bragg reflector ($2\times$ improvement), chip shaping or surface texturing techniques (1.5 to $2\times$ improvement), antireflection coatings, etc. By combining all of these approaches, the LED external quantum efficiency can approach or exceed 50%. Such improvements in efficiency have enabled LEDs to move from the low-flux realm to high-flux applications such as traffic signals, automotive lighting, and full-color outdoor displays. These applications will be discussed more fully in Section IV.

II. FABRICATION TECHNIQUES AND PROCESSES

The fabrication of LEDs requires semiconductor material to be grown, processed in wafer form, diced into individual chips, and packaged. The growth, processing, and packaging techniques vary according to the material system and application. The basic techniques are described below.

A. Material Growth

There are several ways to form the semiconductor material from which an LED is made. Most binary III–V semiconductors can be formed from a melt where a large boule is grown and sliced into wafers. To form the n- and p-type materials, diffusion or ion implantation can be used. These techniques cannot be used for the multilayer structures or alloy materials that most high-performance LEDs require. In order to produce more complex structures, a sequence of crystalline semiconductor layers is deposited on a substrate. The crystal structure of the layers and the substrate should be nearly identical to prevent defects that limit the performance of the LED, as discussed in the first section of this chapter. Early fabrication methods included vapor-phase epitaxy (VPE) and liquid-phase epitaxy (LPE). As semiconductor structures have become increasingly complex, newer growth methods have been developed, including organo-metallic vapor-phase epitaxy (OMVPE) and molecular beam epitaxy (MBE). OMVPE is also known as metal-organic vapor-phase epitaxy (MOVPE), organo-metallic chemical vapor deposition (OMCVD), and metal-organic chemical vapor deposition (MOCVD).

Vapor-phase epitaxy consists of combining group III halides with group V halides or hydrides at high temperatures. For example, GaCl can react with AsH_3 to produce GaAs with HCl as a by-product. Since it is difficult to grow complex structures and many of the newer compound materials using VPE, VPE is used to grow only p–n homojunction devices.

Liquid-phase epitaxy is a simple near-equilibrium process that can produce very pure, thick layers and as a result is still commonly used for producing LEDs. The growth system consists of a graphite holder with wells for various melts. Each melt contains group III elements, group V elements, and dopants in the proper proportions for a defined layer in the semiconductor structure. A graphite slider bar holds a substrate on which the semiconductor is to be grown. The slider moves from one well to another, putting the substrate sequentially in contact with the melts in the wells. At each well, the temperature is ramped down to supersaturate the melt and cause material to be deposited on the substrate. LPE cannot be used to grow nitride or phosphide alloys and is not well suited for the production of quantum well and superlattice structures, so other techniques have been developed for these applications.

Organo-metallic VPE growth on a substrate involves the use of gases that dissociate at high temperatures and deposit constituent atoms on the surface of the substrate. It is the most complex growth technique, but it is also the most versatile, as nearly all III–V material growth can be accomplished with OMVPE. The gases are typically a combination of hydrides and organo-metallic compounds, such as triethylindium and triethylgallium. Dopants can be introduced from various solid or gaseous sources. The process is commonly performed at atmospheric pressure, although low-pressure OMVPE is used when very abrupt interfaces are desired. One drawback to OMVPE is the use of highly toxic gases that make it more dangerous than other growth methods.

Molecular beam epitaxy growth is performed under ultra-high vacuum (UHV) and results in very tight thickness control, making it ideal for superlattices or quantum well structures. Beams of atoms are evaporated from an elemental source in an effusion cell. A shutter in front of the cell controls the material that is allowed to condense on a heated substrate. Since the group V elements are much more volatile than the group III atoms, the growth rate is determined by the group III flux. One advantage of MBE is the ability to use UHV surface analysis techniques to monitor the growth *in situ*.

Several growth methods combine elements of OMVPE and MBE to allow UHV surface analysis methods to be used with a wider variety of materials. Because of difficulties in growing phosphides with solid-source MBE, gas-source MBE was developed in which group V molecular precursors such as arsine and phosphine are used in place of the elemental constituents. Group III molecular precursors are used in metal-organic MBE (MOMBE). Chemical beam epitaxy (CBE) uses the same group III and group V sources as OMVPE but at ultra-high vacuum as in MBE, where UHV surface analysis techniques are available.

The method of fabrication of molecular organic and polymeric thin films for use in organic LEDs may be quite different from that discussed above for inorganic LEDs. Molecular organic LEDs are fabricated by evaporating a sequence of layers in a high-vacuum environment. Polymeric LED layers are deposited by spin- or dip-coating a solution of soluble polymer onto a substrate. Most organic materials have high mechanical flexibility that allows for compatibility with a large number of substrates.

B. Chip Fabrication

The fabrication of LED chips from a semiconductor material is fairly simple when compared to other semiconductor devices. Generally, the line dimensions are very large compared to silicon integrated circuits so high-resolution photomasking is not required, and there are fewer photomasking levels and processing steps. Most LEDs are produced using conventional wafer fabrication techniques. These involve multiple sequences of evaporating or sputtering metals or dielectrics on the surface of the wafer, patterning them with photolithography, and etching them to form simple structures.

In order to contact the p- and n-type materials for operation of the LED, a conducting layer must be deposited on both material types to form ohmic contacts. Unless the doping in the semiconductors is very high, metals typically form Schottky barriers when evaporated or sputtered onto semiconductor material. In order to make the contacts ohmic, annealing at high temperatures is often combined with doping in the semiconductor or metal to remove the interfacial barrier. Some systems, especially organic LEDs, use a transparent conducting oxide such as indium tin oxide as the ohmic contact in order to maximize the amount of light that is emitted from the structure. Photolithography followed by etching is used to pattern one or both of the contacts. For structures in which the substrate is conducting, one contact is on the top surface of the chip and the other contact is on the substrate side of the chip. A few semiconductors are typically grown on insulating substrates, such as $\text{Ga}_x\text{In}_{1-x}\text{N}$ grown on sapphire. In these cases, either the substrate is removed and a back contact can be used, or a hole is etched through the top semiconductor layer and a contact is put on the layer beneath, resulting in both n-type and p-type contacts on the top side of the structure. Although this complicates the fabrication process, both by adding processing steps and increasing the difficulty of packaging, it enables flip chip technology, as described in the next section.

After the wafer fabrication is complete, the wafer needs to be diced to separate the individual chips. The method again depends upon the type of semiconductor. The preferred methods of dicing include sawing with a

narrow, diamond-impregnated blade or diamond scribing and breaking. There are typically 15,000 chips per 2-inch-diameter wafer, which is extremely high compared to silicon integrated circuits that have much larger chip sizes. In addition, many of the substrates and semiconductors are more difficult to dice than silicon. As a result, the dicing process is more complicated and stringent for LEDs, although the rest of the fabrication process is simpler.

C. Packaging

The basic LED packaging process involves attaching the chip to a leadframe, wire bonding the contact pads on the chip to leads on the package, and encapsulating the chip in a transparent encapsulant for protection (see Fig. 10). To attach the chip to the package, silver-based conductive epoxy is typically used. If the chip has a conducting substrate, the back contact then automatically contacts one of the two leads of the package. A wire bonder thermosonically attaches a wire between the chip's top contact and the other package lead. If the substrate is insulating, a wire bond is needed to attach each of the chip's two top contacts to the leadframe. For some applications where the radiation pattern is important, the transparent encapsulation has a lens included to focus the light.

As LEDs are used in more sophisticated applications, the package types are becoming more specialized as in the case of flat packages that conform to surfaces for contour lighting. Flip chip devices are also becoming common. All contacts for the flip chip structure are on the top side of the device. The chip is flipped upside down and put on solder bumps that connect the contacts to the package. This

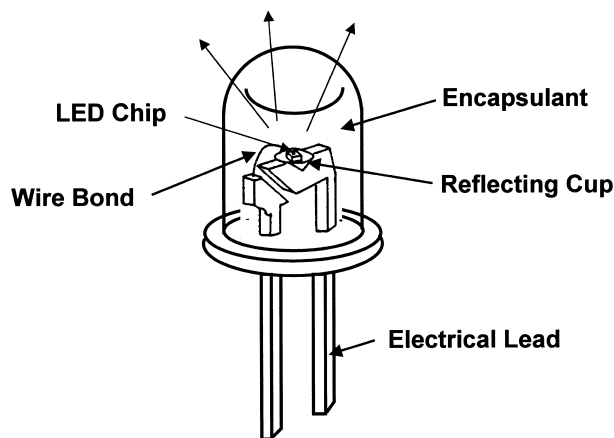


FIGURE 10 Schematic illustration of a discrete LED lamp. The LED chip is attached to a reflecting cup on an electrical lead. A wire bond connects the LED to a second electrical lead. The chip is then encased in an encapsulant. [From Haitz, R. H., Craford, M. G., and Weissman, R. H. (1995). *In* "Handbook of Optics," 2nd ed., Vol. 1, McGraw-Hill, New York. With permission.]

decreases the need for wire bonds and allows for higher power operation and better reliability performance.

Light-emitting diodes typically undergo performance testing to separate the devices according to performance and color. This testing includes measurements of the color, light output, and current-voltage characteristics of the LED. Emitters for high-speed applications, such as fiberoptics or optocouplers, typically are tested for optical rise and fall times.

D. Reliability Issues

The recombination process that produces light in an LED is not a destructive process; therefore, the reliability of LEDs is typically much better than that of incandescent bulbs. LEDs typically last many years as compared to between 1000 and 10,000 hours (approximately one year) for the incandescent bulb. This is a valuable advantage for applications in which the integrity of the light source is important or replacing a bulb is difficult, such as traffic signals.

There are, however, several sources of degradation that are dependent on environmental, packaging, or radiation conditions. For example, low-temperature operation may result in excessive stress due to the difference in thermal coefficients of expansion between the package and the semiconductor. This stress may lead to defects in the semiconductor that will reduce the life of the LED. As a result, new packaging materials that match the low-temperature properties or thermal expansion coefficient of the semiconductor materials have been developed.

Semiconductors such as $\text{Al}_x\text{Ga}_{1-x}\text{As}$ with high aluminum concentrations may degrade rapidly in high-temperature, high-humidity environments. The aluminum undergoes a hydrolization process that produces an absorbing oxide that affects the performance of the LED. To prevent this degradation, a passivation layer or transparent native oxide may be added on top of the LED to prevent exposure to the oxidizing environment.

In order to ensure process stability and long life, some products such as high-power infrared LEDs implement a constant high-current and high-temperature burn-in. Using this technique, infantile failures can be accelerated and defective units removed from the population. Samples of LEDs are often put through reliability testing with a variety of time, temperature, humidity, bias, and other operating conditions to assess the long-term performance of the devices.

Organic LEDs suffer from temperature and atmospheric stability problems. The technology is evolving quickly, and both the reliability and performance are improving as more research effort is put into material fabrication and packaging. Molecular organic materials are characterized

by weak van der Waals bonding. They are soft, susceptible to physical damage, and recrystallize at low temperatures (recrystallization has been observed at room temperature). Polymeric organic electronic materials have greater mechanical stability, due to their strong intrachain covalent bonding.

In general, polymeric and molecular organic materials are known to readily react with water and oxygen. Exposure to atmosphere and other gases can significantly alter their electronic properties. Studies show that the stability of these materials can be improved with the use of sealed packaging. Another problem with polymeric and molecular organic LEDs is that the operating conditions together with the high reactivity of the cathode (needed to ensure efficient injection of electrons into the organic thin film) cause the generation and subsequent growth of “dark spots” within both LED types. Clarification of the nature and origin of dark spot defects is a major priority for future research.

These are a few of the important environmental and material-related reliability issues that have been documented for LEDs. There are also issues with radiation-dependent degradation as well as package degradation. Application and package requirements must be examined for each product at the design stage to avoid such problems. When this is done properly, the lifetimes of LEDs far outlast many of the competing technologies.

III. INFRARED LEDs

A. Materials

The $\text{Al}_x\text{Ga}_{1-x}\text{As}$ material system is used extensively to fabricate infrared LEDs. $\text{Al}_x\text{Ga}_{1-x}\text{As}$ has a direct bandgap for $0 \leq x \leq 0.38$, making the alloy system suitable for LEDs with emission wavelengths between 0.78 and 1.00 μm . As discussed in Section I, double heterostructure (DH) device structures are commonly used to make high-efficiency and high-speed LEDs. DH structures are easily fabricated using the $\text{Al}_x\text{Ga}_{1-x}\text{As}$ material system since layers with different composition can be grown lattice-matched to a GaAs substrate over the entire alloy range. (Recall Fig. 4, where GaAs and AlAs are shown to have the same lattice constant.) Also, $\text{Al}_x\text{Ga}_{1-x}\text{As}$ can be easily grown in thin layers with very abrupt interface transitions between adjacent layers.

Figure 11 shows the cross section of a typical $\text{Al}_x\text{Ga}_{1-x}\text{As}$ LED. Light is generated in the $\text{Al}_{0.08}\text{Ga}_{0.92}\text{As}$ layer where injected holes and electrons recombine. The carriers are confined in one dimension by the $\text{Al}_{0.25}\text{Ga}_{0.75}\text{As}$ barrier layers placed above and below the active layer. The efficiency of such LEDs depends on the $\text{Al}_{0.08}\text{Ga}_{0.92}\text{As}$

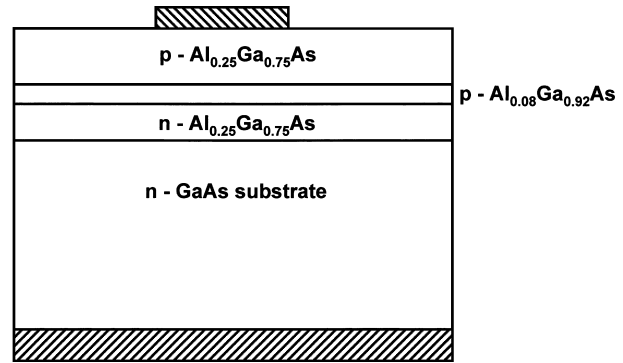


FIGURE 11 Cross section of a typical double heterostructure $\text{Al}_x\text{Ga}_{1-x}\text{As}$ LED. The composition of the light-emitting layer is chosen to emit light at 820 nm in the infrared portion of the spectrum.

layer thickness and doping concentration and on the contact geometry. Devices of this type can have external efficiencies of 5 to 20% and rise/fall times of 20 to 50 ns. The efficiency can be increased by removing the absorbing GaAs substrate.

Many infrared LED devices require well-defined, small-diameter emitting regions in order to couple the emitted light into a detector or fiber. One typical design for producing a small emitting area is given in Fig. 12. Both top and cross sectional views are shown. The reduced emitting region is achieved by three-dimensional carrier confinement. The DH layers confine injected carriers in the perpendicular direction, and the patterned n-GaAs layer is used to control current flow in the lateral direction. Because of the presence of the p-n junction, current is restricted to flow only through the diamond-shaped central region of the device.

A second widely used material system is the quaternary alloy $\text{Ga}_x\text{In}_{1-x}\text{As}_y\text{P}_{1-y}$. DH structures incorporating $\text{Ga}_x\text{In}_{1-x}\text{As}_y\text{P}_{1-y}$ alloys are used to make LEDs emitting at wavelengths from 1.0 to 1.6 μm . Such LEDs are widely deployed as optical sources for fiberoptic links operating at speeds up to 200 Mb/s and for link lengths up to 2 km. By appropriate selection of the x and y values, layers of $\text{Ga}_x\text{In}_{1-x}\text{As}_y\text{P}_{1-y}$ can be grown lattice-matched on InP substrates. Generally, lattice matching is necessary in this system to avoid defects, which can negatively impact LED performance and limit the operating lifetime. Lattice matching to InP requires that $x \approx 0.46y$.

The structure of a $\text{Ga}_x\text{In}_{1-x}\text{As}_y\text{P}_{1-y}$ DH LED used for optical fiber communications is shown in Fig. 13. The x and y parameters in the $\text{Ga}_x\text{In}_{1-x}\text{As}_y\text{P}_{1-y}$ light-emitting layer are chosen to produce a peak emission wavelength of 1.3 μm , where silica glass used for light guiding has minimum wavelength dispersion. InP barrier layers placed on either side of the emitting layer are used to confine

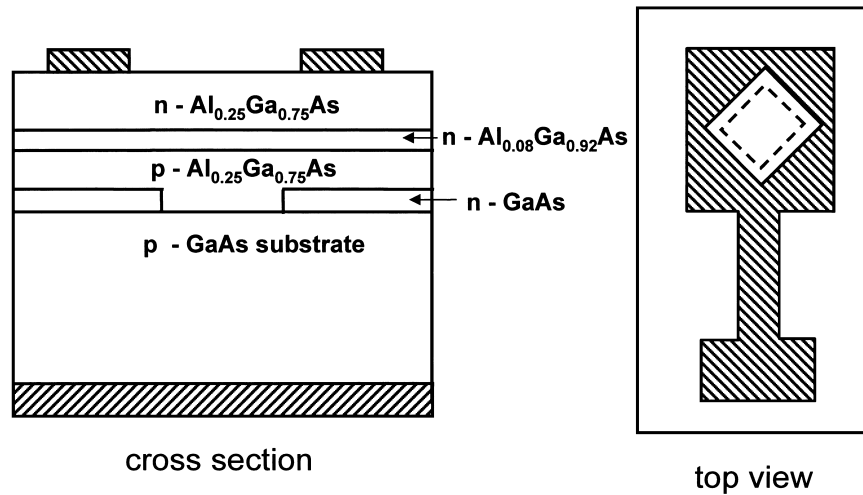


FIGURE 12 High-performance infrared LED with a small emitting area defined by the diamond-shaped pattern in the n-type GaAs layer. The DH structure confines carriers in the perpendicular direction, and the buried p–n junction controls current flow in the normal direction.

electrons and holes, resulting in efficient light generation. The bottom $\text{In}_x\text{Ga}_{1-x}\text{As}$ layer is used to lower contact resistance on the p-side of the device.

The LED in Fig. 13 shows several features commonly found in infrared LEDs. The InP substrate has been shaped to form an etched dome. The dome serves to increase light extraction and to collimate the light beam for effective coupling of emitted light into a glass fiber. A limited-area backside contact is used to control current spreading in the lateral direction. This results in a small-diameter light-emitting region, typically 10 to 30 μm in diameter. Alternatively a shallow etched mesa on the backside of the chip can be used to restrict current flow in the lateral direction.

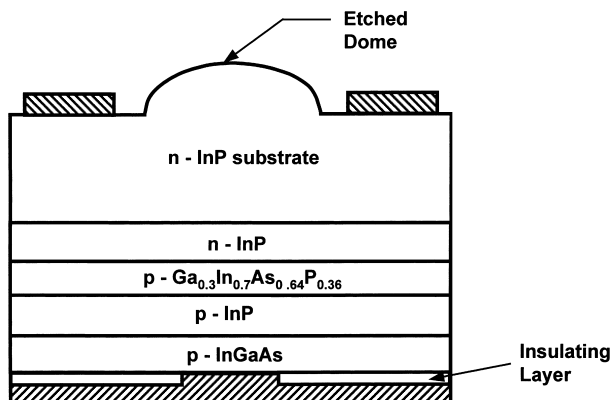


FIGURE 13 Cross section of a $\text{Ga}_x\text{In}_{1-x}\text{As}_y\text{P}_{1-y}$ DH LED used for optical fiber emitters. The InP etched lens is used to collimate light into the core of an optical fiber. The patterned insulating layer is used to control the emitting area size.

B. Applications

Today, infrared LEDs are used in a wide variety of applications, where they are typically paired with photodetectors. Infrared flux emitted by an LED source is transmitted a distance ranging from a few millimeters to many meters, and at the receiving end a photodiode or phototransistor detects the flux and converts it to an electrical signal. Silicon devices are most commonly used to detect the light emitted by $\text{Al}_x\text{Ga}_{1-x}\text{As}$ DH LEDs. Photodetectors made from $\text{In}_x\text{Ga}_{1-x}\text{As}$ are used with $\text{Ga}_x\text{In}_{1-x}\text{As}_y\text{P}_{1-y}$ LEDs. Such emitter-detector pairs are used for the remote control of televisions and other consumer electronic products. Other applications for infrared LED-detector devices include intrusion alarms, blood gas analyzers, bar code readers, limit switches, and shaft encoders.

Infrared LEDs are also used in optocouplers, where an infrared LED and a phototransistor are mounted face to face, with the region between filled with a clear insulating material. Typically the emitter-detector pair is spaced a few millimeters apart and is packaged in a dual-in-line semiconductor type package. Optocouplers are used at the interface between the line voltage side of an electronic system and the low-voltage circuit functions of the system. Pulsing the LED on and off causes an output pulse to be produced. Optocouplers reject common-mode noise.

Infrared LEDs are a primary light source used for short-distance and low-speed fiberoptic links. In this application, a glass or plastic fiber waveguide is used for the transmission medium. LED sources operate in several wavelength windows where fiber absorption is low and wavelength dispersion is minimized. These are the 0.82 to 0.86 μm and 1.3 to 1.5 μm regions. Today, short-distance fiberoptic

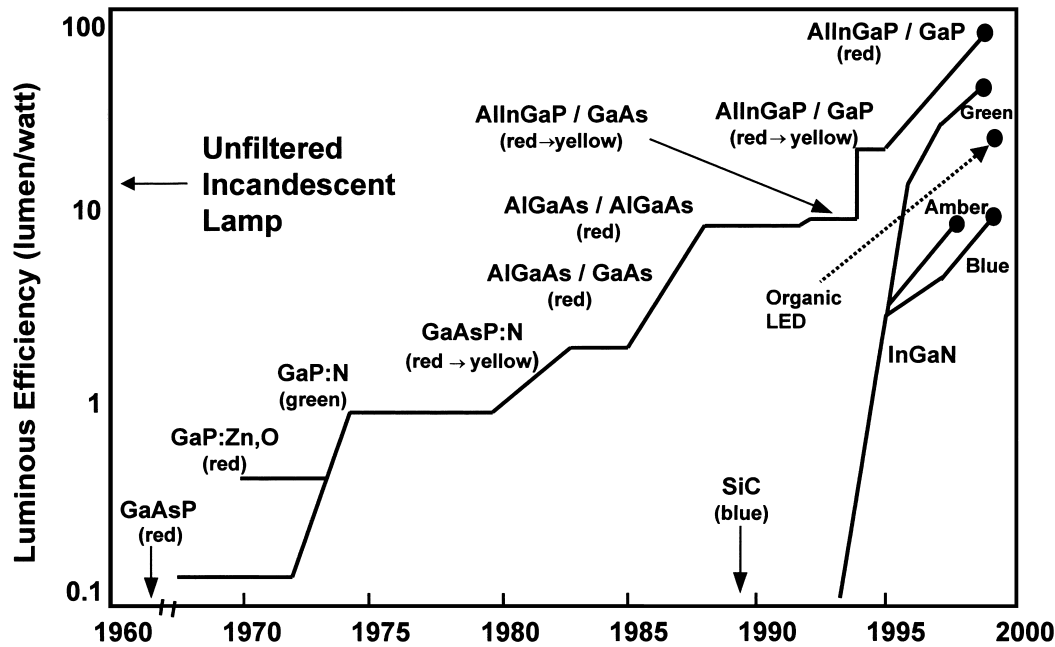


FIGURE 14 Evolution in performance of visible LEDs with time.

links are used in a variety of applications including data transmission, industrial process control, consumer electronics, and medical equipment.

IV. VISIBLE LEDs

A. Materials

1. Inorganic LED Materials

The evolution in performance of visible LED materials with time is shown in Fig. 14 along with the performance of unfiltered incandescent lamps for comparison. Over the last 35 years LEDs have improved in performance to the point where the best LEDs have a higher luminous efficiency than tungsten incandescent lamps.

A widely used alloy for visible LEDs is the ternary $\text{GaAs}_{1-x}\text{P}_x$ system, including one of its binary components GaP. The room-temperature bandgaps of GaAs and GaP are 1.4 and 2.3 eV, respectively. Mixing GaAs and GaP, and thereby adjusting the ratio of arsenic to phosphorus, the bandgap of the resulting ternary compound can be adjusted to any value between 1.4 and 2.3 eV. The $\text{GaAs}_{1-x}\text{P}_x$ alloy is grown on either a GaP or GaAs substrate depending on whether the alloy is closer in composition to GaP or GaAs. $\text{GaAs}_{1-x}\text{P}_x$ LEDs are typically grown by VPE, so devices made with these alloys are limited to simple, low-efficiency homojunction devices.

The resulting band structure with varying As to P ratio is illustrated in Fig. 15. Note that with increasing phosphorus

content the direct Γ valley moves upward in energy space faster than the indirect X valley. At a composition of about 45% GaP and 55% GaAs, the direct and indirect valleys are about equal in energy, and electrons in the conduction band can scatter from the direct valley into the indirect valley (see also Section I of this article). While the direct valley electrons still undergo rapid radiative recombination, the indirect valley electrons have a long radiative

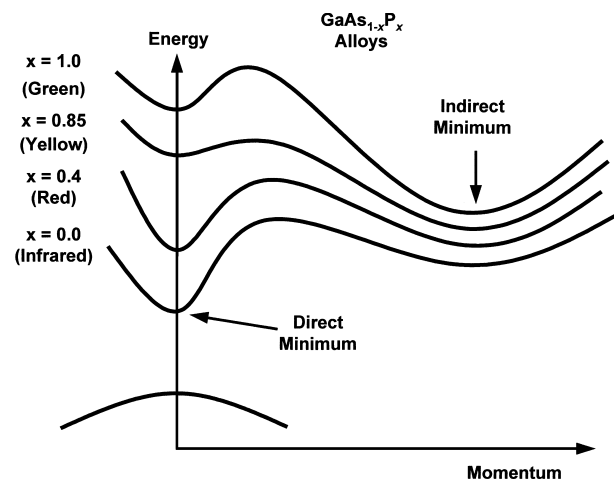


FIGURE 15 Energy band diagram for various alloys of the $\text{GaAs}_{1-x}\text{P}_x$ material system showing the direct and indirect conduction band minima for various alloy compositions. [From Haitz, R. H., Craford, M. G., and Weissman, R. H. (1995). In "Handbook of Optics," 2nd ed., Vol. 1, McGraw-Hill, New York. With permission].

lifetime. To undergo radiative recombination, they must be scattered back to the direct valley through interaction with a phonon with the correct energy and momentum. Therefore, near this crossover between direct and indirect valleys (called the direct–indirect crossover), the radiative efficiency drops off drastically, and for compositions with greater than 45% phosphorus the probability for direct radiative recombination is very small.

The introduction of so-called isoelectronic impurities into indirect bandgap semiconductors can result in increased radiative recombination. The impurity is called “isoelectronic” because it has the same number of electrons in its outer shell as the atom it replaces. In other words, both atoms are in the same group of the Periodic Table. Substituting a nitrogen atom for a phosphorus atom produces an isoelectronic trap in the $\text{GaAs}_{1-x}\text{P}_x$ alloys (see Fig. 16). The stronger electronegativity of nitrogen relative to phosphorus results in the capture of an electron from the conduction band. The negatively charged defect can attract a free hole from the valence band to form a loosely bound electron-hole pair or “exciton.” This electron-hole pair has a high probability to recombine radiatively. The energy of the emitted light is less than the bandgap energy, E_g , since the nitrogen-trap energy level is within the forbidden energy gap.

Nitrogen doping is widely used for $\text{GaAs}_{1-x}\text{P}_x$ alloys with $x = 0.65$ to $x = 1.0$. The resulting light sources cover the wavelength range from 635 nm (red) to 572 nm (yellow-green). Another isoelectronic defect in GaP is formed by ZnO pairs (zinc on a gallium site and oxygen on a phosphorus site). The ZnO electron trap is farther away in energy from the conduction band, resulting in

longer wavelength emission in the red region of the spectrum (see Fig. 16).

Green-emitting LEDs can also be constructed in GaP without nitrogen doping. Because GaP is an indirect bandgap material, these devices depend upon phonons to conserve momentum (see Section I). GaP devices without nitrogen doping have an advantage in that the emission is shorter in wavelength (565 nm) and therefore appears to be a purer green.

Another alloy used for visible LEDs is $\text{Al}_x\text{Ga}_{1-x}\text{As}$. As stated in Section III on infrared LEDs, the entire alloy range from $x = 0$ to $x = 1$ can be grown nearly lattice-matched on GaAs substrates. This means that defects in the epitaxial layers can be minimized, and alloy compositions can be grown without the use of transition layers (as are required with the $\text{GaAs}_{1-x}\text{P}_x$ system). This allows the growth of very abrupt transitions in composition and bandgap (heterojunctions), the advantages of which were previously described in Section I.

The $\text{Al}_x\text{Ga}_{1-x}\text{As}$ alloy with composition $x = 0.35$ is used for commercial red LED fabrication. Double heterostructure LEDs are formed using a p-type $\text{Al}_{0.35}\text{Ga}_{0.65}\text{As}$ active layer sandwiched between p- and n-type $\text{Al}_{0.75}\text{Ga}_{0.25}\text{As}$ confining layers. Improved external quantum efficiency can be achieved by growing a thick p-type window layer on the top side of the structure and removing the absorbing GaAs substrate. Figure 14 shows the improved performance of the $\text{Al}_x\text{Ga}_{1-x}\text{As}$ LEDs over $\text{GaAs}_{1-x}\text{P}_x$ and GaP LEDs.

A third alloy used for visible LEDs is $(\text{Al}_x\text{Ga}_{1-x})_{0.5}\text{In}_{0.5}\text{P}$, which can be grown lattice-matched to GaAs substrates. The ratio of Al to Ga can be changed without

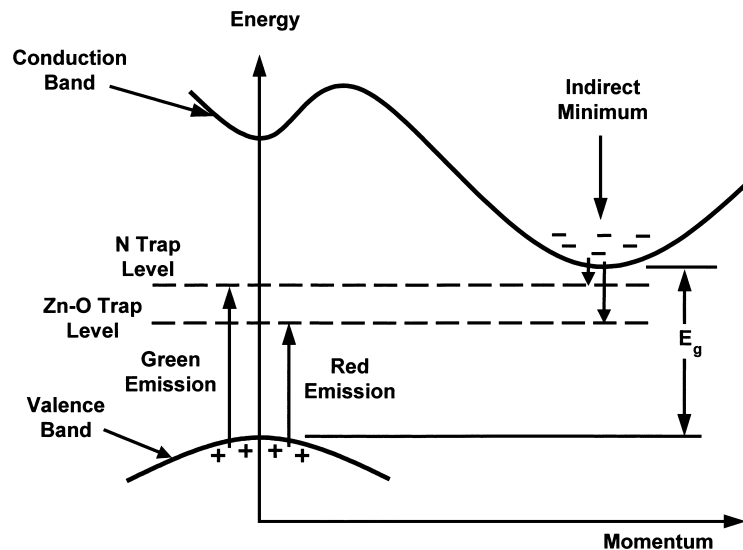


FIGURE 16 Formation of excitons (electron-hole pairs) by the addition of isoelectronic dopants N and ZnO to an indirect semiconductor. The excitons have a high probability to recombine radiatively. [From Haitz, R. H., Craford, M. G., and Weissman, R. H. (1995). *In* “Handbook of Optics,” 2nd ed., Vol. 1, McGraw-Hill, New York. With permission.]

affecting the lattice match, just as in the $\text{Al}_x\text{Ga}_{1-x}\text{As}$ system, since AlP and GaP have nearly the same lattice constant. This enables the growth of heterostructures that have the same efficiency advantages as described above for $\text{Al}_x\text{Ga}_{1-x}\text{As}$. When lattice-matched to GaAs, $(\text{Al}_x\text{Ga}_{1-x})_{0.5}\text{In}_{0.5}\text{P}$ LEDs cover a wide range of the visible spectrum from approximately 653 nm (red) for $x = 0$ to 555 nm (green) for $x \cong 0.53$. Device efficiency decreases considerably in the green, because this composition is close to the direct-indirect bandgap crossover.

For $(\text{Al}_x\text{Ga}_{1-x})_{0.5}\text{In}_{0.5}\text{P}$ LEDs, multi-well active layer structures as well as single and double heterostructures have been grown. A top window layer of $\text{Al}_x\text{Ga}_{1-x}\text{As}$ or GaP is used for current spreading and light extraction. The layer can be grown relatively thick (up to 50 μm) to maximize light extraction from the edge of the chip. $\text{Al}_x\text{Ga}_{1-x}\text{As}$ has the advantage that it is lattice matched, but the disadvantage that it absorbs some of the light in the case of yellow and green emitters. GaP is transparent to the emitted light, but it is not lattice matched. Fortunately, this does not appear to introduce defects in the active region. Very-high-performance devices are commercially available using the $(\text{Al}_x\text{Ga}_{1-x})_{0.5}\text{In}_{0.5}\text{P}$ materials system for red, orange, and yellow LEDs (see Fig. 14).

Despite the many successes in the red to yellow color range, as recently as 1990 the best green LEDs were quite dim, and blue LEDs were still not available. These colors are essential to products such as full-color displays and white-light products, and the absence of bright blue and green LEDs limited the LED market to low-light or red and yellow applications. Extensive research in the development of blue and green LEDs has been ongoing for decades in order to solve this problem. Unfortunately, many of the materials with a bandgap large enough to produce these short wavelengths, such as SiC- or ZnSe-based semiconductors, are either indirect-bandgap semiconductors (low efficiency) or have other fabrication or reliability issues that currently prevent their use in commercial products.

In the early 1990s, a significant breakthrough for blue and green LED technology came in the development of the $\text{Ga}_x\text{In}_{1-x}\text{N}$ material system, which had largely been ignored due to the lack of an appropriate substrate and difficulties in achieving p-type doping. One reason for the difficulty in achieving p-type doping of $\text{Ga}_x\text{In}_{1-x}\text{N}$ alloys is the high energies of the acceptor states ($E_A - E_V \geq 160$ meV), which are generally not ionized at room temperature. After this problem was solved, substrate material remained the primary concern. $\text{Ga}_x\text{In}_{1-x}\text{N}$ epitaxial growth is performed on sapphire (Al_2O_3) or silicon carbide (SiC) substrates that are well suited for the extreme thermal and corrosive environment experienced during nitride epitaxial growth. Cost and availability currently limit the use of silicon carbide substrates. For both sapphire and silicon carbide substrates, large differences in lattice constant

(up to 20%) and coefficient of thermal expansion between the substrate and semiconductor layers create problems regarding the growth of high-quality films. The growth of a thin nucleation layer (typically a nitride alloy) on the substrate before the $\text{Ga}_x\text{In}_{1-x}\text{N}$ layer growth improves the quality of the $\text{Ga}_x\text{In}_{1-x}\text{N}$ layers considerably. The defect density present in the $\text{Ga}_x\text{In}_{1-x}\text{N}$ layers is still high (as much as six orders of magnitude higher than in other III-V semiconductors), but it does not affect the radiative efficiency as much as for other visible LED materials (see Fig. 14 for a comparison of the efficiency of visible LED materials). This lack of dependence of radiative efficiency on the defect density in nitride materials is poorly understood and continues to be a topic of intense research.

The binary compounds that make up the $\text{Ga}_x\text{In}_{1-x}\text{N}$ system are direct bandgap materials that exhibit a range of bandgap energies from 1.9 eV (InN) to 3.4 eV (GaN). If Al is added to the system, the bandgap can be extended to 6.3 eV, and the available emission colors range continuously from red, through the visible, and into the ultraviolet with a single material system. Because of the difficulties associated with incorporating increasing percentages of indium into the $\text{Ga}_x\text{In}_{1-x}\text{N}$ alloy, in practice $\text{Ga}_x\text{In}_{1-x}\text{N}$ is used commercially only for high-efficiency green- and blue-emitting LEDs.

2. Organic LED Materials

Organic LED technology has seen a wave of interest since the first demonstration of efficient electroluminescence at Kodak in 1987. This was followed by the synthesis of the first polymer LED at Cambridge University in 1992. The device physics of organic LEDs is similar to that of inorganic LEDs. Holes and electrons are injected from opposite contacts into the organic layer sequence and transported to the light-emitting layer. Carrier recombination leads to the formation of excitons which, after a characteristic lifetime of approximately 1 to 10 ns, either decay radiatively to generate light or nonradiatively to generate heat. The best reported luminescence efficiencies of organic LEDs exceed that of incandescent light bulbs, and lifetimes on the order of 10,000 hours have been reported. A significant advantage of the organic materials is their low index of refraction, between 1.5 and 1.7. Thus, light extraction into air is much greater than for inorganic LEDs that have an index of refraction of approximately 3.5 (see Section I).

Single-layer and heterojunction organic LEDs have been fabricated. Single-layer organic LEDs have a low efficiency, due in part to the low probability for exciton formation in the thin film. Heterojunction organic LEDs are much more efficient, since the carrier confinement provided by the heterointerface increases the probability of exciton formation.

Organic LEDs have advantages over inorganic LEDs in that organic LEDs can easily be fabricated into thin, relatively large displays for use in automobile radios, navigation aids, or heads-up displays. Light emission from organic LEDs is achieved from fluorescent dyes placed within a conductive matrix. By using different dyes, organic LEDs can be fabricated in each of the three primary colors—red, green, and blue—and full-color displays using organic LEDs can be fabricated either by using pixels with distinct colors or by using an organic LED that emits white light, and using discrete red, green, and blue color filters to achieve the desired color.

While organic LEDs offer advantages over inorganic LEDs in terms of design flexibility, organic LEDs have certain disadvantages compared to inorganic LEDs, particularly in the areas of flux density and efficiency. One particular difficulty associated with organic LEDs is in the area of ohmic contact formation. Typical contact layers for organic LEDs include transparent metal-oxide layers such as indium tin oxide, and such layers can contribute significantly to the forward voltage drop across the device, resulting in turn-on voltages in the range of 3 to 5 volts for organic LEDs.

B. Applications

A substantial fraction of visible LEDs produced today are the low-efficiency homojunction type. The efficiency of these diodes is adequate for many signaling applications, and their low price makes them attractive. High-efficiency visible LEDs are serious contenders for applications where incandescent light bulbs are now used. The inherent advantages of LEDs over incandescent bulbs include high reliability, low power consumption, styling flexibility (they are small), and a tolerance for harsh environmental conditions.

1. Discrete Devices

The simplest LED product is an indicator lamp, as previously shown in Fig. 10. There are many variations of this basic lamp design to satisfy various product applications. Typical variations include size, shape, and radiation pattern. A number of lamps can be combined into a single package to illuminate a rectangular area.

Displays can be either numeric or alphanumeric. Numeric displays are usually made up of a nearly rectangular arrangement of seven elongated segments in a figure-eight pattern. Selectively switching these segments generates all ten digits from zero to nine. Often a decimal point, comma, or other symbols are added. There are two ways LEDs are used to display alphanumeric information: either by using more than seven elongated segments (i.e., 14, 16, or 24)

or by using an array of LED chips in a 5×7 dot matrix format.

A common type of LED display is called a stretched segment display. Figures 17 and 18 illustrate this type of display. Each character segment is formed by encapsulating an LED chip in an epoxy-filled diffusing cavity. Reflections off the white plastic cavity walls and additional scattering within the cavity result in a uniformly lit segment. The top area of the cavity is typically 20 to 30 times larger than the LED surface area. On/off contrast is enhanced by coloring the epoxy and the outer surfaces of the display package.

Instrument and computer applications of these discrete devices include indicator lamps, numeric and alphanumeric displays, and LED backlighting of liquid crystal displays. Higher performance LEDs are often used for the backlighting application. Consumer electronics applications (audio equipment, home appliances, clocks, radios, toys, etc.) include discrete indicator lamps and numeric and alphanumeric displays.

2. System Uses

System applications include display, signage, and traffic signal applications. Displays and signs can be “electronic billboards” as large as several meters or as small as several inches high and wide. Outdoor applications require the use of high-performance LED technologies, such as $\text{Al}_x\text{Ga}_{1-x}\text{As}$, $(\text{Al}_x\text{Ga}_{1-x})_{0.5}\text{In}_{0.5}\text{P}$ or $\text{Ga}_x\text{In}_{1-x}\text{N}$, for better visibility in sunlight. Indoor applications generally use the lower performance technologies, including $\text{GaAs}_{1-x}\text{P}_x$ doped with nitrogen (red-emitting) and GaP doped with nitrogen (green-emitting). Traffic signal applications have begun to incorporate red-emitting $(\text{Al}_x\text{Ga}_{1-x})_{0.5}\text{In}_{0.5}\text{P}$ and $\text{Al}_x\text{Ga}_{1-x}\text{As}$ LEDs for stop lights and are moving toward incorporating amber $(\text{Al}_x\text{Ga}_{1-x})_{0.5}\text{In}_{0.5}\text{P}$ and blue-green $\text{Ga}_x\text{In}_{1-x}\text{N}$ LEDs to produce a completely LED-based signal head.

For vehicular lighting applications, LEDs are used for their ruggedness, reliability, small size (providing styling flexibility), fast turn-on time, and low power use. $\text{Al}_x\text{Ga}_{1-x}\text{As}$ and $(\text{Al}_x\text{Ga}_{1-x})_{0.5}\text{In}_{0.5}\text{P}$ LEDs are incorporated in automobile trunks and “spoilers” to serve as brake lights and are being used on trucks and automobiles as side markers, running lights, and turn signals. LEDs are also used as indicators and display devices on the interior of vehicles. Dashboard displays generally use the nitrogen-doped $\text{GaAs}_{1-x}\text{P}_x$ and nitrogen-doped GaP technologies, although high-performance technologies are sometimes used to backlight liquid crystal displays.

High-efficiency LEDs are being developed for solid state lighting applications. A major goal is the development of white LED light sources that will expand the

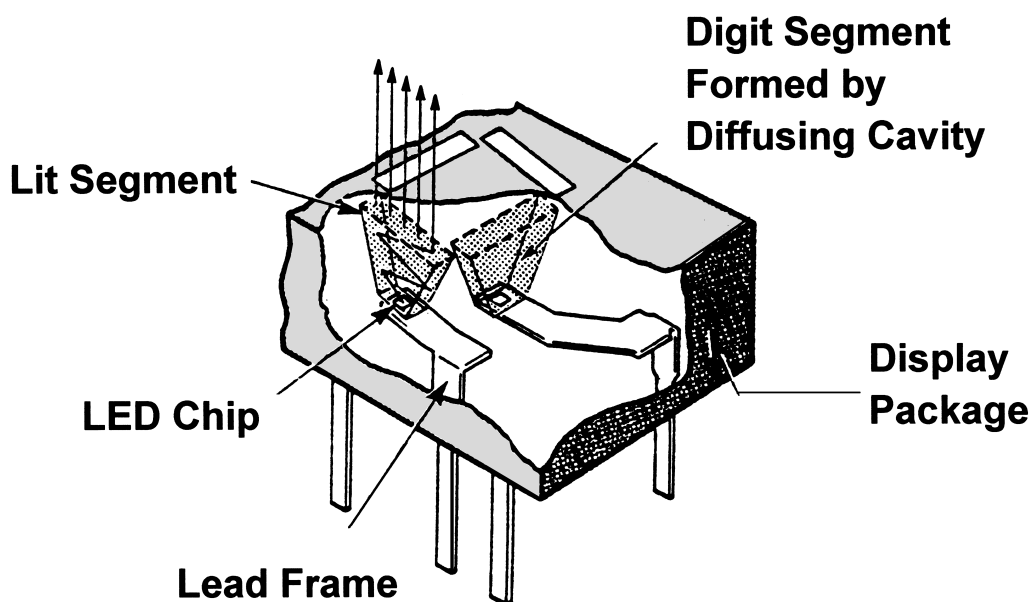


FIGURE 17 Cutaway of a seven-segment numeric LED display showing how light from a small LED chip is stretched to a large character segment using a diffusing cavity. Characters 0 to 9 are created by turning on appropriate combinations of segments. [From Haitz, R. H., Craford, M. G., and Weissman, R. H. (1995). *In* "Handbook of Optics," 2nd ed., Vol. 1, McGraw-Hill, New York. With permission.]

market for LEDs into areas that have not previously been possible. One way to create white light is to mix the light from three LEDs emitting the primary colors (red, green, and blue). Using this method, all colors in the spectrum, including white, may be emitted by adjusting the proportion of light emitted from each LED. This is a valuable property for decorative lighting or outdoor video screens where color control is essential. A second common method for

producing white light is the use of a blue diode to excite a phosphor inside the lamp package that converts the blue light into white light. A third method employs multiple layers of semiconductors that emit different colors, stacked on top of each other. The bottom layer is excited when a bias voltage is applied and emits light. Part of this light passes through the other layer(s) and part excites the layer(s) on top to produce one or more different colors of light. The mixing of the emitted light from all of the layers results in white light. There is a huge demand for white LEDs due to styling possibilities and energy cost savings over traditional lighting methods. A few of the many applications include illumination, decorative lighting, and backlit displays. Because of the huge potential market, several lighting companies and LED manufacturers have formed alliances to accelerate the adoption of LED technology. The LED sector of the lighting market is expected to grow very quickly in the next few years as LEDs become brighter and are used in a wider variety of products.

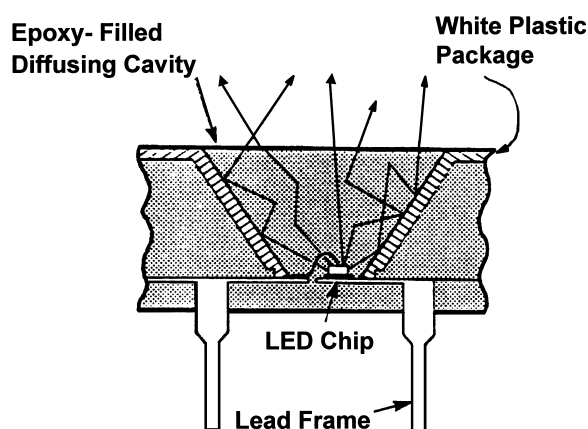


FIGURE 18 Cross section through one segment of the seven-segment numeric display shown in Fig. 17. An LED chip is placed at the bottom of a highly reflective, white plastic cavity which is filled with clear plastic containing a diffusant. Light is scattered within the cavity to produce uniform emission at the segment surface. [From Haitz, R. H., Craford, M. G., and Weissman, R. H. (1995). *In* "Handbook of Optics," 2nd ed., Vol. 1, McGraw-Hill, New York. With permission.]

SEE ALSO THE FOLLOWING ARTICLES

CRYSTAL GROWTH • EPITAXIAL TECHNOLOGY FOR INTEGRATED CIRCUIT MANUFACTURING • LASERS, OPTICAL FIBER • LASERS, SEMICONDUCTOR • MOLECULAR BEAM EPITAXY, SEMICONDUCTORS • PHOTONIC BANDGAP MATERIALS • POLYMERS, PHOTORESPONSIVE • SEMICONDUCTOR ALLOYS

BIBLIOGRAPHY

- Bergh, A., and Dean, P. (1976). "Light Emitting Diodes," Clarendon Press, Oxford.
- Craford, M. G., and Sterenka, F. M. (1994). "Light emitting diodes," *In* "Encyclopedia of Applied Physics," Vol. 8, pp. 485–514. VCH Publishers, New York.
- Gillessen, K., and Schairer, W. (1987). "Light Emitting Diodes—An Introduction," Prentice Hall, International, Englewood Cliffs, NJ.
- Haitz, R. H., Craford, M. G., and Weissman, R. H. (1995). "Light emitting diodes." *In* "Handbook of Optics," 2nd ed., Vol. 1, pp. 12.1–12.39, McGraw-Hill, New York.
- Mueller, G., vol. ed. (2000). "Semiconductors and Semimetals," Vols. 64 and 65: "Electroluminescence I and II," Academic Press, San Diego, CA.
- Nakamura, S., and Chichibu, S. F., eds. (2000). "Introduction to Nitride Semiconductor Blue Lasers and Light Emitting Diodes," Taylor & Francis, London.
- Solymar, L., and Walsh, D. (1990). "Lectures on the Electrical Properties of Materials," 4th ed. Oxford University Press, New York.
- Streetman, B. G. (1990). "Solid State Electronic Devices," 3rd ed. Prentice Hall International, Englewood Cliffs, NJ.
- Stringfellow, G. B., and Craford, M. G., vol. eds. (1997). "Semiconductors and Semimetals," Vol. 48: "High Brightness Light Emitting Diodes," Academic Press, San Diego, CA.
- Sze, S. M. (1981). "Physics of Semiconductor Devices," 2nd ed. Wiley, New York.



Imaging Optics

Matt Young

Colorado School of Mines

- I. Lenses
- II. Human Eye
- III. Camera
- IV. Projection System
- V. Magnifying Glass
- VI. Microscope
- VII. Telescope

GLOSSARY

Aberration The departure of a ray or a wavefront from the path predicted by the paraxial theory.

Airy disk The diffraction pattern of a circular aperture, as seen in the image plane of a well-corrected lens. The radius of the Airy disk is a measure of the resolution limit of a well-corrected lens.

Aperture stop An opening, usually circular, that limits the total radiant power entering a lens system. The iris is the aperture stop of the eye.

Conjugate distances The object and image distances. The object and image points are called *conjugate points*.

Diffraction-limited Capable of achieving the resolution limit imposed by diffraction.

F-number The ratio of the focal length of a lens to the diameter of the exit pupil of that lens. The F-number of a well-corrected lens is the reciprocal of twice its numerical aperture.

Magnifying power The ratio of the angular subtense of an image to that of the object viewed with the naked

eye, usually from the distance of 25 cm. Used when magnification is not a meaningful measure, as with a magnifying glass, a microscope, or a telescope.

Numerical aperture The sine of the half-angle between a marginal ray and the axis, as measured at a focal point of a lens. The numerical aperture of a well-corrected lens is the reciprocal of twice its F-number.

Paraxial approximation The approximation that all rays are so nearly parallel to the axis that the small-angle approximation may be used. In geometric optics, imaging is perfect in the paraxial approximation.

Point spread function The image of an isolated point.

Principal plane Either of two planes that fully describe a lens for paraxial raytracing. The intersection of a principal plane with the axis of a lens is called the *principal point*.

Principal ray A ray that passes through the principal points of a lens.

Pupil 1: An image of the aperture stop of an optical system. 2: The black circular opening in the center of the iris (of the eye), which permits light to pass to the retina.

Resolution limit The center-to-center distance between the images of two points, usually of equal intensity, that can barely be resolved. The reciprocal of the resolution limit is called the *resolving power*.

Thin lens A lens that consists of one or more elements and whose overall thickness is very much smaller than one or both of its focal lengths or conjugate distances; a lens that may be considered infinitesimally thin. Distinguished from a *thick lens*.

OPTICS has been critical to every scientific and technological revolution beginning with Copernicus and ending, so far, with microelectronics and telecommunications. Optical instruments continue to be used to magnify both the images of distant galaxies and also subwavelength features in biology and microelectronics. Optical imaging systems are used both to enhance visual observation and also to acquire, process, or transmit images by photography, video photography, or digital electronics.

This chapter describes the principles that underlie optical systems for imaging the very large to the very small, but excluding holographic systems. Most of the instruments use lenses or mirrors, though one uses a probe. Most project images, though two scan point by point. Some of the instruments are used visually, though some use computers or video sensors. Most of the instruments can be described with ray optics, but their ultimate limitations can be described only by wave optics. Hence, the bulk of this chapter is devoted to geometrical optics, with bits from physical optics drawn in when necessary.

More specifically, the chapter describes lenses and lens optics, lens aberrations, the human eye, cameras and projectors, microscopes, and telescopes.

I. LENSES

A. Thin Lenses

A *lens* consists of one or more pieces of glass or other refracting material with (usually) spherical surfaces. Figure 1 shows in cross-section a lens that consists of a single piece of glass that has rotational symmetry about its axis. We assume for the moment that the lens is infinitesimally thin; such a lens is called a *thin lens*.

The lens in Figure 1 projects an image of an object at O to the point O' . The plane that passes through O and is perpendicular to the axis of the lens is called the *object plane*, and the corresponding plane that passes through O' is called the *image plane*. The object is located at a distance l (the *object distance*) from the lens, and the image, a distance l' (the *image distance*). Because O is located

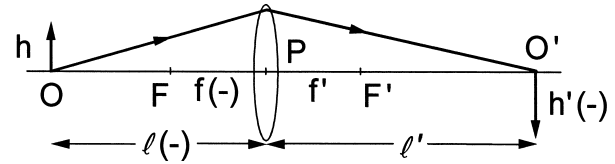


FIGURE 1 A thin lens, showing the conjugate points O and O' , the primary and secondary focal points F and F' , the primary and secondary focal lengths f and f' , the object and image distances l and l' , the object and image heights h and h' , and the principal point P .

to the left of the center P of the lens, l is negative, as is indicated on the figure by the parenthetical $(-)$. The distances l and l' are related by the *lens formula*,

$$\frac{1}{l'} - \frac{1}{l} = \frac{1}{f'}, \quad (1)$$

where f' is the *secondary focal length* of the lens and is equal to the value of l' when $l = -\infty$. The *primary focal length* f is similarly the value of l when $l' = +\infty$. Physically, f' is the image distance that corresponds to an infinitely distant object, whereas f is the object distance that yields an infinitely distant image. When a small source, often called a *point source*, is projected to ∞ , the lens projects a bundle of parallel rays, and the beam is said to be *collimated*.

The sign convention used here differs from that found in many elementary texts but has advantage in that it uses Cartesian coordinates for object and image distances and for radii, whereas elementary texts use a sign convention that works well only for simple lenses.

When the object or the image is infinitely distant, the object and image points O and O' become the *primary focal point* F and the *secondary focal point* F' . The planes that pass through those points are called the *focal planes* and are special cases of the image planes. If $f' > 0$, the lens is called a *positive lens*; otherwise, it is a *negative lens*. The primary focal point lies to the right of a negative lens, and the secondary focal point, to the left.

The primary and secondary focal lengths are equal to each other in magnitude: $f' = -f$, provided that the lens is surrounded by materials with the same index of refraction. If it is not, then $f'/f = -n'/n$, where n' is the index of refraction on the right side of the lens and n is the index of refraction on the left side of the lens. The secondary focal length of a thin lens in air is given implicitly by the *lens-maker's formula*,

$$\frac{1}{f'} = (n - 1) \left(\frac{1}{R_1} - \frac{1}{R_2} \right), \quad (2)$$

where R_1 is the radius of curvature of the first surface and R_2 , of the second, and where radii of curvature are measured from the surface, so the radius of curvature is

positive if the center of the surface is to the right of a surface.

The object and image distances are called *conjugate distances*, and the object and image points O and O' are called *conjugate points*. If an object lies to the left of F , then its image lies to the right of the lens and is called a *real image* (presuming here that the lens is positive). If an object lies between F and the lens, however, its image lies behind the lens and is called a *virtual image*. If that virtual image is used as the object of a second lens, it is called a *virtual object*. By contrast, the image formed by a negative lens is always virtual, unless the object is a virtual object located to the right of the lens. Figure 2 compares and contrasts real and virtual images and objects. If the object located at O extends a height h above the axis of the lens, then the height h' of the image is given by

$$\frac{h'}{h} \equiv m = \frac{l'}{l}, \quad (3)$$

where m is called the *magnification* or, sometimes, the *transverse magnification*; m is negative for a positive lens projecting a real image.

It is sometimes convenient to write the object and image distances l' and l in terms of the focal lengths f' and f , and the magnification m :

$$l' = f'(1 - m), \quad l = f[1 - (1/m)] \quad (4)$$

If $m = -1$ (a real, inverted image), then $l' = 2f'$, and $l = -2f$. This case is called *unit magnification*, since the magnitude of m is equal to 1.

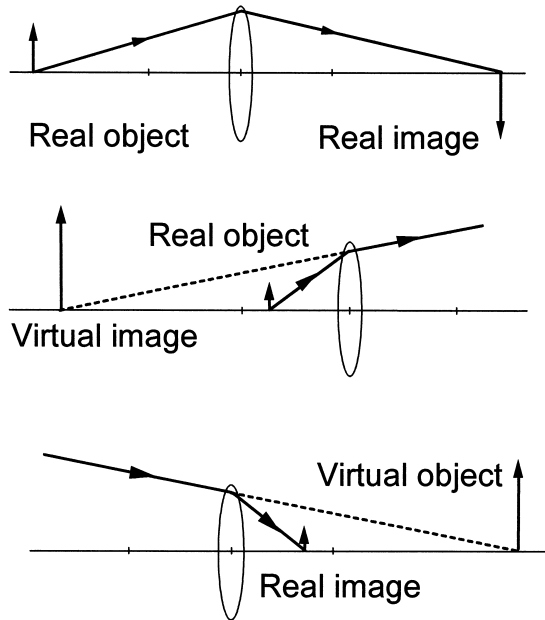


FIGURE 2 Three thin, positive lenses, showing real and virtual objects and images.

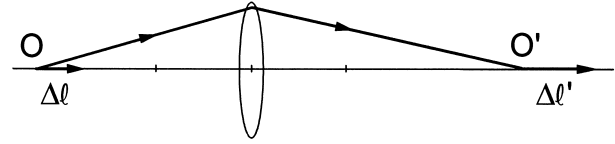


FIGURE 3 A thin lens, showing the object and image lengths Δl and $\Delta l'$, from which the longitudinal magnification is calculated.

If the object extends a distance Δl (Figure 3) along the axis, then the image extends a distance $\Delta l'$ along the axis, and

$$\frac{\Delta l'}{\Delta l} \equiv \mu = m^2, \quad (5)$$

where μ is the *longitudinal magnification* and is always positive. Longitudinal magnification is defined only when $\Delta l \ll l$ and $\Delta l' \ll l'$. It is not a useful concept when the object or the image is near a focal point, because then the image or object is likely to be very long.

A single spherical interface between dissimilar materials focuses light according to the relation

$$\frac{n'}{l'} - \frac{n}{l} = \frac{n' - n}{f'}, \quad (6)$$

where

$$\frac{n'}{f'} = \frac{n' - n}{R}. \quad (7)$$

The quantity n'/f' is called the *power* of the surface.

A spherical mirror in air may be visualized as a spherical refracting surface for which $n' = -1$; hence, the equation that relates the conjugate distances for a mirror is

$$\frac{1}{l'} + \frac{1}{l} = \frac{1}{f'}, \quad (8)$$

where $f' = R/2$.

B. Newton's Form of the Lens Equation

Newton derived an alternate form of the lens equation,

$$xx' = -f'^2, \quad (9)$$

where x and x' are measured from the respective focal points rather than from the principal points and follow the same sign convention as the conjugate distances l and l' (Figure 4). Newton's form is most often useful when one of the conjugate distances is close to the focal length, in which case the other conjugate distance is large and Newton's x parameter is very nearly equal to the related conjugate distance. Newton's formulation also yields alternate expressions for the magnification:

$$m = -\frac{x'}{f'} = -\frac{f}{x}. \quad (10)$$

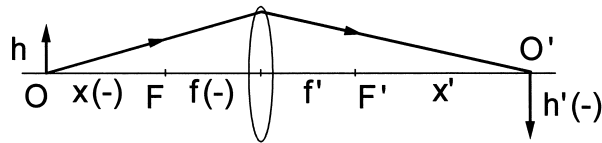


FIGURE 4 A thin lens, showing the parameters x and x' used in Newton's formulation of the lens equation.

C. Thick Lenses

A lens that may not be considered infinitesimally thin is called a *thick lens*. It is not always obvious when a lens must be considered thick; indeed, that may depend on how precisely we want to measure its focal length or its conjugate distances. In general, however, a thick lens does not satisfy the condition that its thickness t is much less than l and l' .

A thick lens may consist of a single piece of glass that is relatively thick compared to the focal length of the lens, for example, or it may consist of two or more thin lenses that are separated by a finite distance, or of two or more thick lens elements. In either case, we define the *principal planes* of the lens by the geometrical construction shown in Figure 5. Specifically, the *secondary principal plane* is constructed by tracing through the lens a ray that is parallel to the axis of the lens and then extending that ray or the emergent ray, as necessary, until the incident and emergent rays cross. The plane that passes through the intersection of those rays and is perpendicular to the axis is the secondary principal plane. Its intersection with the axis is the *secondary principal point* P' . The *primary principal plane* and the *primary principal point* P are constructed analogously. The principal planes reduce a complicated lens to a thin lens in this sense: One conjugate distance is measured from the appropriate principal point, and the thin lens equation is applied to calculate the other conjugate distance, which is then measured from the other principal point (see also Ray Tracing, below).

Conjugate distances are now measured from the respective principal points, not from the surfaces of the lens and

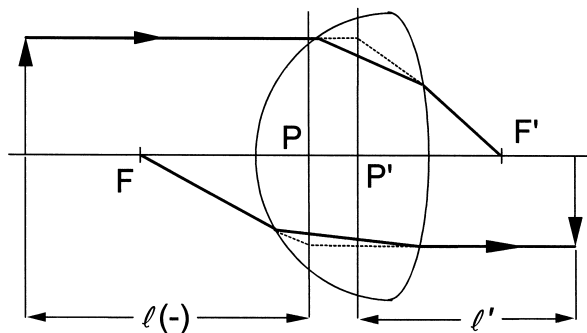


FIGURE 5 The construction of the principal planes P and P' of a thick lens.

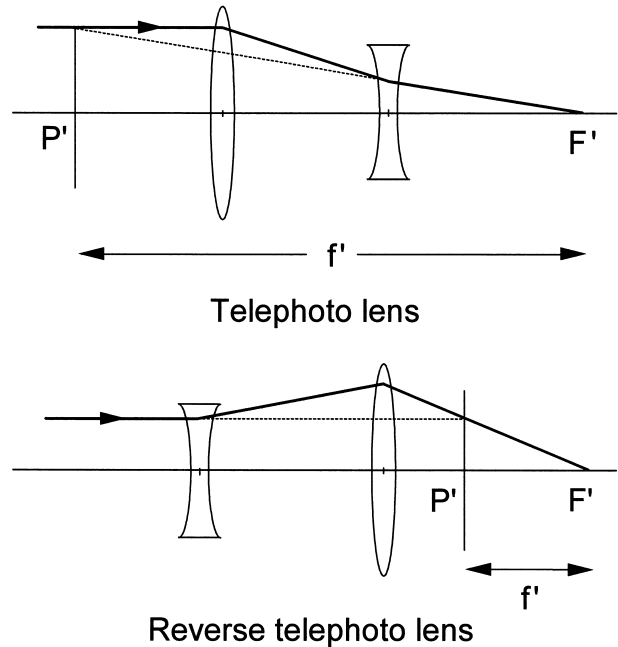


FIGURE 6 A telephoto lens, with a long effective focal length, and a reverse telephoto lens, with a long working distance.

not from the geometrical center of the lens. Specifically, the distance between the last surface of the lens and the secondary focal point is not, in general, equal to the focal length of the lens. That distance is called the *working distance* or, perhaps confusingly, the *back focal length*.

Figure 6 shows two special thick lenses: (a) the telephoto lens and (b) the reverse telephoto lens, and the construction of their secondary principal planes. The telephoto lens has a comparatively long focal length but a short physical length. It is used when compactness, weight, and freedom from vibration are important. The reverse telephoto lens has a relatively short focal length and a long working distance, and is used when it is necessary to insert another element, such as a mirror, between the lens and the image plane. Many objectives, especially short focal length objectives, for 35-mm cameras are reverse telephoto lenses.

D. Lens Aberrations

All the preceding equations are *paraxial*; that is, they assume that all angles, such as the angle of incidence of any ray on any refracting surface, are small. If that condition holds, then the sine or tangent of any angle may be replaced by that angle itself, and the geometrical image of a point is a point. In reality, however, the paraxial condition is never exactly satisfied, and the geometrical image of a point is not a point but rather a relatively compact distribution of points, sometimes wholly displaced from the paraxial location. For example, Figure 7 traces rays across a planar

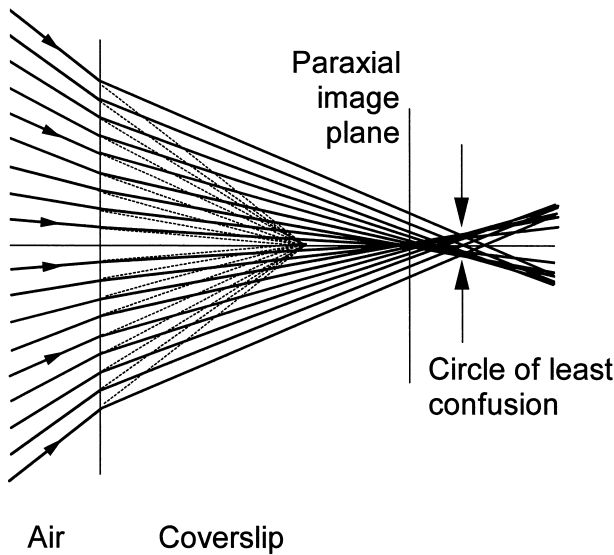


FIGURE 7 Spherical aberration caused by a planar interface.

interface such as a microscope coverslip. The incident rays converge exactly to a point, but rays that are incident on the interface at finite angles do not cross the axis precisely at the paraxial image point. Instead, the image of the point is blurred by *spherical aberration*. A ray that strikes the lens at its periphery is called a *marginal ray*. The *transverse spherical aberration* TA of the marginal ray (shown in Figure 8) is often taken as the measure of the third-order spherical aberration of a lens.

All real lenses made from spherical surfaces suffer to some degree from spherical aberration. Additionally, if the object point is distant from the axis of the lens, or *off-axis*, the image may suffer from other aberrations known as *astigmatism*, *coma*, *distortion*, and *field curvature*. These are the *third-order aberrations*, or *Seidel aberrations*, and are calculated by using the approximation that the sine of an angle θ is given by $\sin \theta \cong \theta - \theta^3/3!$. *Fifth-order aberrations* are calculated by adding the term $\theta^5/5!$ to the series. Finally, because the index of refraction of the lens is a function of wavelength, the focal length of any lens

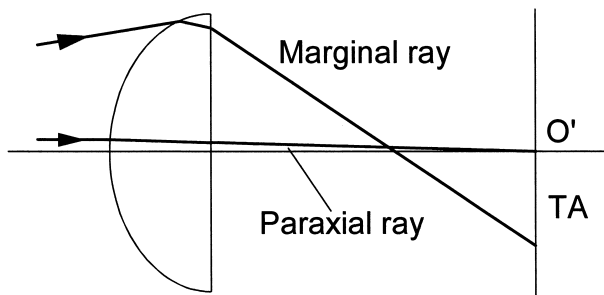


FIGURE 8 The transverse spherical aberration TA as measured in the paraxial image plane.

varies slightly with wavelength; the resulting aberration is called *chromatic aberration*.

Spherical aberration appears both on the axis and off the axis. It does not depend on the distance off-axis. In the absence of other aberrations, an image of a point is not compact but rather is surrounded by a diffuse halo. The diameter of the halo is a measure of the spherical aberration; the spherical aberration of a thin lens increases as the cube of the diameter of the lens. The other third-order aberrations begin to appear off-axis and become more severe as the distance from the axis increases.

Coma appears as a cometlike image, where the head of the comet is the paraxial image point and the tail is the aberration. The tail points away from the axis of the lens and is three times longer than it is wide. The length of the comatic image increases as the square of the lens diameter and proportionally to the distance of the image from the axis of the lens.

Astigmatism occurs because an off-axis bundle of rays strikes the lens asymmetrically. This gives rise to a pair of line images, one behind the plane of best focus and one before it. The axial distance between these two images, or *astigmatic difference*, is often taken as a measure of the astigmatism of a lens. The astigmatic difference is approximately proportional to the focal length of a thin lens; it increases as the square of the distance off-axis and proportionally to the diameter of the aperture stop.

Even in the absence of other aberrations, the image surface of a lens is not truly a plane but rather a curved surface. The aberration is known as *field curvature*, and the focal surface is called the *Petzval surface*. The Petzval surface curves toward a thin, positive lens. Its radius of curvature is the reciprocal of the *Petzval sum* $\Sigma(1/n_i f'_i)$, where the sum is taken over all the lens elements. The Petzval sum can be made 0 under certain circumstances; that is, a lens can be made *flat-fielded*.

If the magnification m is a function of distance from the axis, then an image will not be rendered rectilinearly. The resulting aberration is called *distortion*. If a square, for example, is imaged in the shape of a pincushion (its sides bowed inward), we speak of *pincushion distortion*. Pincushion distortion results when m is greater off-axis than on-axis and is also called *positive distortion*. If, on the other hand, a square is imaged in the shape of a barrel (its sides bowed outward), we speak of *barrel distortion*. Barrel distortion results when m is less off-axis than on-axis and is also called *negative distortion*.

Chromatic aberration manifests itself in two ways. Since the lens-maker's equation (2) depends on the index of refraction n , and n in turn depends on wavelength, the location of an image point depends on wavelength. The distance between a red image point and a blue image point is called *longitudinal chromatic aberration* and is one measure of chromatic aberration. Alternatively,

lateral chromatic aberration may be defined as the radius of the image in the plane of best focus. Longitudinal chromatic aberration does not depend on the diameter of the lens, whereas lateral chromatic aberration is proportional to the ratio of longitudinal chromatic aberration and the F-number of the lens and therefore varies in proportion to the lens diameter.

Aberrations may be reduced by *bending* a lens, that is, by adjusting the radii of curvature of lens elements so that, for example, angles of incidence are minimized. Astigmatism, however, is only weakly influenced by bending the elements. Similarly, one aberration can sometimes be balanced against another. In Figure 7, for example, spherical aberration can be partially compensated by defocusing and locating the image plane at the waist instead of the paraxial image plane, that is, by balancing spherical aberration and defocusing. The aberration at the waist is one-fourth that in the paraxial image plane. Finally, adjusting the position of the aperture stop can also reduce coma, distortion, and astigmatism.

Unfortunately, it is not possible to compensate aberrations over a wide range of conjugate distances and angles, so it is important to use a given lens only for its intended purpose: for example, a photographic objective, which is designed for a distant object and usually is not flat-fielded, should not be used in place of a copying or enlarging lens, which is designed for nearby objects and requires both a flat object plane and a flat image plane. Similarly, many microscope objectives are designed for use with a coverslip of a certain thickness and with a fixed magnification. Using a microscope objective that has a specified magnification of $40\times$ or more with no coverslip or with the wrong magnification will often result in unacceptable aberration. Such aberration blurs an image or diminishes contrast, or both.

Spherical aberration is the only third-order aberration that is not 0 for imaging monochromatic light on the axis of the lens. The spherical aberration of a thin lens is sometimes of interest, for example, for collimating a beam or for coupling the radiation from one optical fiber into another, especially if the light is monochromatic. The transverse spherical aberration TA (Figure 8) of a planoconvex lens may be calculated from the formula

$$TA = l' B / [64n(n-1)(f'/D)^3], \quad (11)$$

where

$$B = \frac{n+2}{n-1} q^2 + 4(n-1)pq + (3n+2)(n-1)p^2 + \frac{n^3}{n-1}, \quad (12)$$

$$p = \frac{l' + l}{l' - l}, \quad q = \frac{r_2 - r_1}{r_2^2 + r_1^2}, \quad (13)$$

and D is the diameter of the lens. The quantity p is called the *position factor* of the lens, and q is called the *shape factor*. The sign of TA is irrelevant, but for calculating p and q it is important to remember that l is negative if the object is real and r_2 is negative if the lens is double-convex. The ratio f'/D in Eq. (11) has a special name: the *F-number* of the lens. The F-number determines the maximum angle of the cone of rays that converges to the image point and is useful in diffraction theory and photography. The transverse spherical aberration of a thin lens of given F-number is proportional to the image distance l' ; thus, for example, if you double both the focal length and the image distance of a lens while keeping its F-number constant, you also double the transverse spherical aberration.

For imaging near unit magnification, a symmetrical, double-convex lens shows the least spherical aberration of any thin, spherical lens. For imaging with one conjugate near ∞ , the thin, spherical lens that shows the least spherical aberration is nearly planoconvex and is oriented with the more steeply curved side toward the long conjugate; in practice, a planoconvex lens oriented in the same way serves nearly as well.

E. Diffraction

Diffraction is a consequence of the wave nature of light. It is not an aberration, but it nevertheless prevents an image of a point from being itself a true point. Rather, the image of a point is a diffraction pattern, known as the *point spread function*. If the geometrical image of the point is free of aberrations, the point spread function is an *Airy disk*. The point spread function in the presence of significant aberration is always broader and less contrasty than the Airy disk shows the intensity of the Airy disk as a function of the distance from the geometrical image point. The intensity of the Airy disk is 0 at the radius $1.22\lambda l'/D$, where D is the diameter of the lens; this value is often called the radius of the Airy disk. The term “relatively free of aberrations” means, roughly, that the geometrical image of the point has to be of the same order as the Airy disk radius or smaller. A lens or system that is relatively free of aberrations is called *diffraction-limited*.

According to the *Rayleigh criterion*, a diffraction-limited system can distinguish two equally intense points provided that their geometrical images are separated by at least one Airy disk radius. That distance is known as the *Rayleigh limit* and defines the *resolution limit* RL of a diffraction-limited system. For a diffraction-limited lens

$$RL = 1.22\lambda l'/D, \quad (14)$$

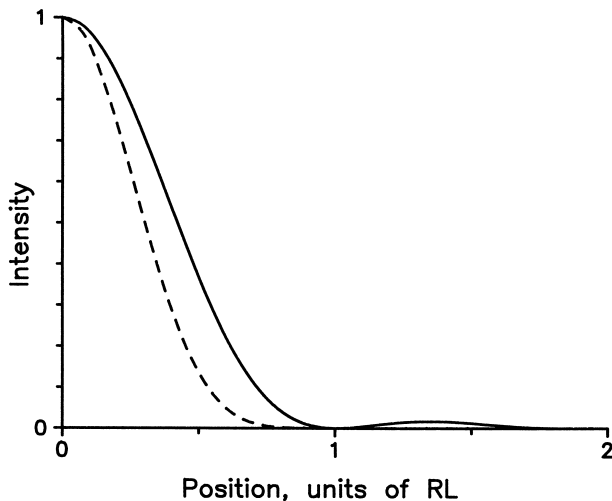


FIGURE 9 Diffraction-limited images of an isolated point. Solid curve: conventional imaging. Dashed curve: scanning confocal microscope. [From Young, M. (2000). *Optics and Lasers, Including Fibers and Optical Waveguides*, 5th ed. Springer, Berlin, Heidelberg, New York. Copyright © 2000 Springer-Verlag.]

where $l'/D = (f'/D)(1 - m)$ is the *effective F-number* of the lens. Sometimes, especially in photography, resolution is described as the number of points that can be resolved in a unit length of an image. Resolution described in this way is called *resolving power*. Usually expressed in lines per millimeter, resolving power is the reciprocal of the resolution limit.

Many lenses are not diffraction-limited; their resolution limit is instead defined by their aberrations and may be many times larger than the Rayleigh limit. Such lenses are called *aberration-limited*. Most thin lenses and most camera lenses, for example, are aberration-limited. Microscope objectives, by contrast, are often diffraction-limited or nearly so. In general, a lens that has a short focal length is more likely to be diffraction-limited than a lens that has a longer focal length and the same effective F-number. This is so because the Rayleigh limit is a function of the F-number only, whereas the aberrations of a lens scale with its focal length. A thin lens may thus be diffraction-limited if its focal length is very short or if its diameter is very small.

According to the wave theory of aberrations the converging wave that propagates toward an image point should in principle have exactly spherical wavefronts. The maximum distance between the wavefront and a true sphere is called the *wavefront aberration*. If the wavefront aberration exceeds one-quarter wavelength, then the lens is aberration-limited. The point spread function is broader than the Airy disk, and the intensity of the peak is correspondingly reduced. The ratio of the peak intensity of the point spread function to that of the Airy disk is called the

Strehl ratio and may be used as a measure of the image quality. A Strehl ratio of 0.8 approximately corresponds to a wavefront aberration of one-quarter wavelength, and a lens with a Strehl ratio greater than 0.8 may be taken to be diffraction-limited.

A thin lens used on or very near the axis may similarly be considered diffraction-limited if the normalized transverse aberration $Y \equiv TA/RL$ is less than 1.6. Figure 10 shows Y as a function of the effective F-number of thin lenses that have a focal length of 1 cm; Y is plotted as a function of effective F-number for both planoconvex lenses with the object at ∞ and symmetrical, double-convex lenses used at unit magnification. The value of Y for any other focal length may be found by multiplying the appropriate value on the graph by the focal length of the lens in centimeters. Given an effective F-number, you can always find a value of f' that will yield diffraction-limited imagery on the axis of the lens, provided that f' is short enough.

F. Raytracing

To construct the image in an optical system that consists of one or more thin lenses, we use the paraxial approximation and trace two or more of the three special rays shown in Figure 11. Ray 1 enters the system parallel to the axis, and the lens directs it through the secondary focal point F' . Ray 3 passes through the primary focal point F , and the lens directs it parallel to the axis of the lens. Ray 2, the ray that passes through the primary principal plane, is called the *principal ray*. A ray that passes through the edge of the lens is called a *marginal ray*, and a ray that passes through the center of the aperture stop (see below) is called the *chief ray*. We will have no occasion here to distinguish between the chief ray and the principal ray. For a discussion of pupils, see the telescope, below.

The principal ray is undeviated by the lens, as long as the lens is infinitesimally thin, because a very thin lens at its center is no more than a slab of glass with parallel surfaces. The three rays that emerge from the lens cross at the image of the tip of the arrowhead. Because other features of the arrow are similarly mapped into the image plane, we may construct the image of the arrow as a whole by constructing a line segment that is perpendicular to the axis of the lens and passes through the tip of the arrowhead. If a lens cannot be considered a thin lens, then we replace the lens by its principal planes. Ray 2, the principal ray, is directed toward the primary principal plane and emerges from the secondary principal plane unchanged in direction.

II. HUMAN EYE

The optical system of the human eye is sketched in Figure 12. It consists of a transparent *cornea*, which

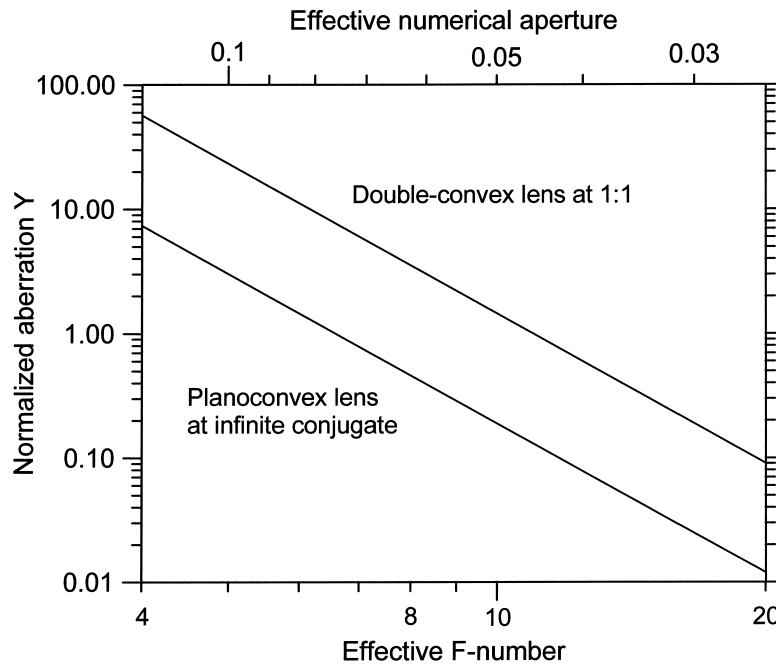


FIGURE 10 The normalized aberrations of two thin lenses with 1-cm focal lengths as functions of effective F-number. A lens is diffraction-limited if $Y < 1.6$. [From Young, M. (2000). Spherical Aberration of a Thin Lens: A Normalized Parameter Relating It to the Diffraction Limit and Fiber Coupling Efficiency. "Engineering and Laboratory Notes," Supplement to *Optics and Photonics News*, vol. 11, no. 5, pp. 1–4. See also Corrections, vol. 11, no. 8, p. 4. Contribution of the National Institute of Standards and Technology, not subject to copyright.]

performs most of the focusing, with a lens, or *crystalline lens*, close behind it. For our purposes, we may consider the cornea and the lens to be in contact and equivalent to a thin lens a few millimeters behind the cornea. The eye is filled with gelatinous substances that have an index of refraction of 1.34, very nearly that of water, or 1.33.

The power n'/f' of a lens in visual optics is usually expressed in *diopters* (D), or reciprocal meters. That is, the power of a lens is in diopters if f' is in meters. The power of the optical system of the human eye is approximately 60 D when the eye is focused at infinity. The corresponding focal length is 22 mm.

The human eye focuses on nearby points by *accommodation*, that is, by muscular contractions that change

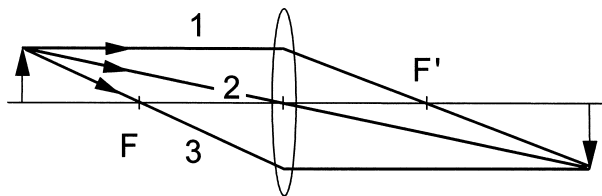


FIGURE 11 Three rays used for constructing an image. (1) A ray incident from $-\infty$. (2) A ray directed at the primary principal point of the lens. (3) A ray directed at the primary focal point of the lens.

the shape of the crystalline lens and increase its power, rather than by moving the lens back and forth along its axis. It is conventional to assume that the lens has a range of accommodation of 4 D, which gives it a *shortest distance of distinct vision*, or *near point*, $d_v = 1/4 \text{ m} = 25 \text{ cm}$, though very young eyes may have several times this range and older eyes substantially less.

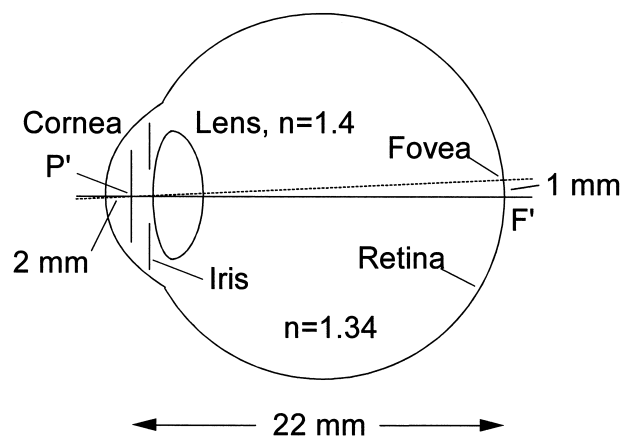


FIGURE 12 A schematic drawing of the human eye. P' represents the principal points, which are nearly coincident, and F' is the secondary focal point.

The rear surface of the eye supports a thin membrane of light receptors called the *retina*. The retina has two types of receptors: *cones*, which operate in bright light and can discriminate colors, and *rods*, which operate in relative darkness and are color-blind. Though their spectral characteristics differ slightly, both receptors are sensitive primarily to light whose wavelength lies in the *visible spectrum* between 400 and 700 nm. Figure 13 shows the luminous efficiency of the eye for cone, or *photopic*, vision and for rod, or *scotopic*, vision.

The eye *adapts* to different light intensities by switching from cones to rods as the light intensity decreases and also by adjusting the diameter of a variable aperture known as the *iris*. The iris controls the diameter of the opening, or *pupil*, through which light is allowed into the eye; it thus serves as the *limiting aperture*, or *aperture stop*, of the eye. The iris involuntarily varies the diameter of the pupil from 2 mm or less in bright light to 8 mm or more in relative darkness. By using *neural inhibition* to turn off the cones and switching to rod vision, however, the eye can adapt to a range of intensities of perhaps 10^6 to 1.

We will be concerned here with photopic vision only. The cones are concentrated in the center of the visual field, in a region called the *macula lutea*. They are most heavily concentrated in the center of the macula, in a region called the *fovea centralis*. The diameter of the macula is approximately 2 mm, and that of the fovea is 300 μm .

When we fix our gaze on an object, we are focusing its image onto the fovea. The cones in the fovea are approximately 5 μm in diameter and closely packed. Their diameter is about equal to the diffraction limit for a 2-mm pupil, and the eye is diffraction-limited when the pupil diameter is 2 mm. When the pupil diameter increases,

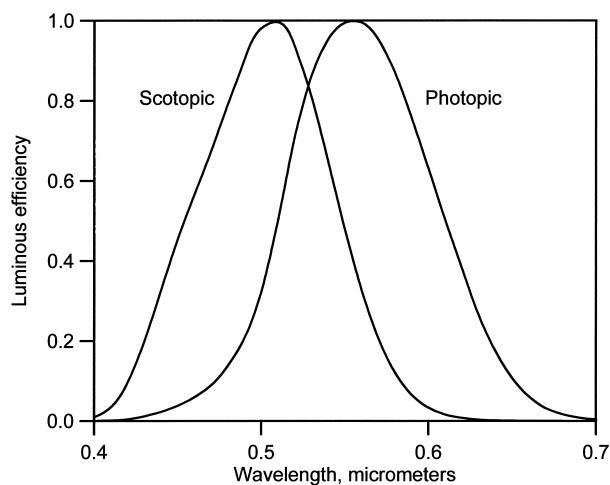


FIGURE 13 The luminous efficiency of the human eye as a function of wavelength. Photopic: bright light, typically outdoor illumination. Scotopic: near darkness.

however, the eye becomes aberration-limited, and its resolution limit increases. In bright outdoor light, the eye can resolve an angle of approximately 0.3 mrad, which translates to a resolution limit of 0.1 mm at a distance of 25 cm. In interior light, when the pupil is 5 mm in diameter, the resolution limit of the eye is a factor of 2 to 3 larger, that is, poorer, than in bright outdoor light.

The human eye may suffer from focusing error. For example, if the power of the optical system is too high for the length of the eyeball or if the eyeball is too long, the patient cannot focus accurately at great distances and is called a *myope*. The condition is called *myopia*, or nearsightedness, and is corrected with negative spectacle lenses or contact lenses. Similarly, if the power of the optical system is too low, the patient cannot focus accurately at short distances and is called a *hyperope*. The condition is called *hyperopia* (sometimes *hypermetropia*), or farsightedness, and is corrected with positive spectacle lenses or contact lenses. If the cornea is not a sphere, but rather a section of a spheroid, then vertical features may focus in different planes from horizontal features. The condition is called *astigmatism* but differs in its origin from the astigmatism of the lens designer. Finally, with age the crystalline lens loses its ability to accommodate, which causes an inability to focus at short distances, even when the eye is otherwise well corrected. This condition, which is sometimes confused with hyperopia, is called *presbyopia* and, like hyperopia, is corrected with positive spectacle lenses or contact lenses. Eventually, the crystalline lens may become wholly opaque and is said to have a *cataract*. Cataracts may be corrected with a surgical procedure that replaces the defective lens with a plastic implant.

Optical systems such as magnifying glasses, microscopes, and telescopes are analyzed under the assumption that the pupil of the eye is 5 mm in diameter and that the near point of the eye is 25 cm. Inasmuch as these are arbitrary assumptions, the resulting magnifying powers are only nominal.

III. CAMERA

The basic camera, like the eye, consists of a lens and a photosensitive surface. Also like the eye, the camera controls the intensity of the light falling onto the photosensitive surface by adjustments of the diameter of the aperture stop, usually an *iris diaphragm* in a professional camera. Unlike the eye, the camera captures only instantaneous images, and the camera varies the duration of the exposure to light, or the *exposure time*, depending on the intensity of the light and the sensitivity of the photosensitive surface. Also unlike the eye, the camera usually has a flat photosensitive

surface and roughly uniform resolution across the entire field of view. Finally, the camera has a *field stop* that limits the image plane to a fixed rectangle. Nowadays, there are two major types of cameras: *film cameras* and *digital cameras*.

A. Film Cameras

Most popular professional cameras today use the 35-mm *format*, which produces an image that is in reality 24×36 mm. The image is recorded on photographic film, which consists of a light-sensitive *emulsion* coated onto a flexible backing. The emulsion is actually a suspension of light-sensitive silver halide particles in a thin layer of gelatin. After the film is exposed, it has to be processed chemically to render the image visible. The size and distribution of the particles, or *grains*, determines the resolving power of the film. Much film for scientific use is black and white. Color film is a stack of black-and-white emulsions that are specially sensitized to different colors and separated by colored filter layers.

Common films have resolving powers in the range of 50 to 100 lines/mm. Resolution is measured by photographing a target that consists of alternating black-and-white lines of equal thickness, where one line means a black line plus the adjacent white line.

The exposure of the film is determined by the F-number of the lens and the exposure time, sometimes called the *shutter speed*. The F-number of a lens is the ratio of its focal length to its diameter, and most cameras have adjustable F-numbers. Lenses on professional cameras are calibrated in specific values, called *F-stops*, of the F-number. The response of photographic film to light is logarithmic, so each F-stop differs by a constant ratio, specifically a factor of $\sqrt{2}$, from the previous F-stop; typical F-stops are 2, 2.8, 4, 5.6, 8, 11, 16, and 22. Changing by one F-stop changes the exposure of the film by a factor of 2. Exposure times are therefore likewise calibrated in factors of 2: 1/500, 1/250, 1/125, 1/60, 1/30, 1/15 s. Many cameras today are fully or partly automated, so the exposure times and F-stops may be mostly invisible to the user.

Camera lenses, or *photographic objectives*, are rarely diffraction-limited. A high-quality 35-mm camera lens may display at most a resolving power of 60 or 70 lines/mm at an F-number of 8 (written F/8). At higher and lower F-numbers, the resolving power is usually slightly lower, and in the corner of the film, often much lower. Resolving power may suffer further if the exposure time is long and the camera is not held rigidly. Common films are thus adequate to resolve all the detail in the image.

The resolving power of the film or of the lens (whichever is lower) determines the *depth of focus* of the camera, or

the distance $\pm\delta'$ by which an image point may depart from the image plane and still be considered in acceptable focus. δ' may be related to the lens and film parameters by

$$\delta' = (l'/D)/RP \quad (15)$$

Similarly, if we hold the film fixed and consider a range of object distances, we find that the range that is in acceptable focus, or the *depth of field* $\pm\delta$, is given by $\pm\delta'/m^2$, where m^2 is the longitudinal magnification. When the camera is focused at ∞ (and longitudinal magnification is not a useful concept), the nearest object distance that is in sharp focus, or the *hyperfocal distance*, is given by

$$H = f'^2 RP / FN, \quad (16)$$

where we write the equation to show the dependence on F-number *FN* explicitly. The preceding two equations are only approximately correct for diffraction-limited systems, in which *RP* is the reciprocal of the Rayleigh limit.

B. Digital Cameras

These replace the film with a digital receptor called a *CCD array*. A CCD (for charge-coupled device) array is a rectangular array of photosensitive elements or *pixels*. A typical array might contain 1300×1000 pixels, which require approximately 1.3 MB of digital memory. A 24×36 -mm color slide, for comparison, might produce about 60 lines/mm, which translates to 120 pixels/mm, or 4000×3000 pixels. The slide has roughly three times the resolution of the CCD array, but to duplicate the slide digitally would require about an order of magnitude (3^2) more memory.

C. Video Cameras

Typical video cameras have CCD arrays with 780×480 pixels and capture a new picture every 1/30 s. In the United States each complete picture, or *frame*, is divided into two *fields* that consist of alternate horizontal lines (1, 3, 5, . . . and 2, 4, 6, . . .) and are *interlaced*. Thus, the rate at which the picture flickers when it is played back is 60 Hz, which is generally higher than the rate at which the eye can perceive flickering intensities. It is difficult to vary the exposure time of a video camera if the recorded pictures are to be played back at the usual rate of 60 Hz, so the video camera adjusts the exposure by varying the gain of the electronics as well as the diameter of the iris diaphragm.

IV. PROJECTION SYSTEM

Figure 14 shows a *projection system* such as a slide projector. An object, such as a 35-mm slide, is back-lighted by

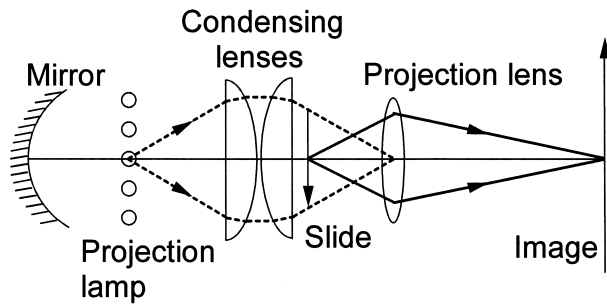


FIGURE 14 Typical projection system, such as a slide projector.

a *projection lamp*, which is usually a tungsten lamp with several tightly coiled filaments arranged to form a square. The lamp or the projector may also include a concave mirror to capture backward-going rays as well. A well-corrected *projection lens* projects an image of the slide onto a distant screen. Because much of the light from the lamp passes through the slide and does not (or would not) enter the aperture of the projection lens, the projector includes a *condensing lens*, which projects an image of the filament into the aperture of the projection lens. The condensing lens is not especially well corrected inasmuch as it is used to concentrate light into the projection lens rather than to project a high-quality image. Nevertheless, most condensing lenses are *aspheric* to reduce spherical aberration and ensure uniform illumination in the image plane.

In principle, the image of the filament precisely fills the aperture of the projection lens. If the image is larger than the aperture of the projection lens, then the lens is said to be *overfilled*, and light from the extremities of the filament is wasted. Conversely, if the projection lens is *underfilled*, then its full aperture is not used, and resolution may suffer. The principle of *filling the aperture* is important in other optical systems as well (see below).

Other projection systems may image the filament directly onto the film plane; it is important then to use a projection lamp that has a *ribbon filament* for uniform illumination of the film plane. Such a layout is useful for motion-picture projectors, where the fast-moving film is unlikely to be burned by the heat of the lamp.

V. MAGNIFYING GLASS

We assume here that the near point of the eye is $d_v = 25$ cm. A small object of height h located at the near point subtends an angle h/d_v , as measured from the primary principal point of the eye. If we want to see the object with more detail than with the naked eye, that is, magnify the object, we may place a positive lens, or *magnifying glass*, just in front of the eye and use it to examine this object. In prin-

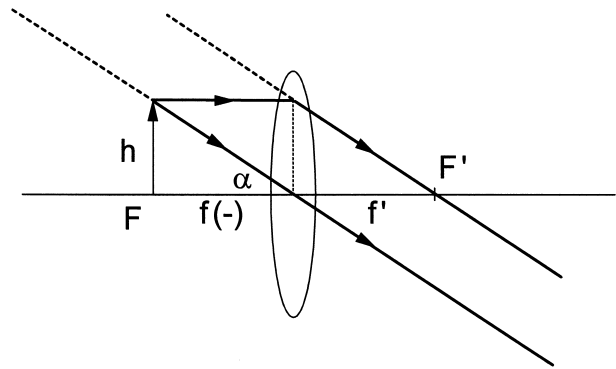


FIGURE 15 An object of height h located at the focal point of a magnifying glass, where it subtends an angle α .

ciple, the object is located at the focal point of the lens, so a virtual image is projected to ∞ (Figure 15). This image subtends an angle $\alpha = h/f'$ at the primary principal point of the lens. The ratio of h/f' to h/d_v is the *magnifying power*

$$MP = d_v/f' \quad (17)$$

of the magnifying glass. As long as $f' < d_v$, the magnifying power will exceed 1. A magnifying glass may be used to achieve a magnifying power up to approximately 10, provided that the lens is well corrected.

Many observers adjust the lens so that the image lies at the near point of the eye. The magnifying power is then increased slightly, but it also depends on the position of the eye with respect to the magnifying glass. Thus, magnifying power should be regarded as a nominal specification, not an exact number.

Magnifying power is not to be confused with magnification: It is an arbitrary number that depends on the value chosen for d_v and on the exact location of the image. For example, if the image is located at d_v rather than ∞ , and if the eye is in contact with the magnifying glass, then $MP = 1 + d_v/f'$. Similarly, if the magnifying glass is used by a myope, then d_v may be substantially less than 25 cm.

VI. MICROSCOPE

A. Conventional Microscope

Sometimes called a *compound microscope* to distinguish it from the magnifying glass, or simple microscope, a microscope consists of an *objective lens* (*objective*, for short) that projects a real, magnified image and an *eyepiece* with which that image is examined (Figure 16). The eyepiece is usually no more than a specialized magnifying glass. If the magnification of the objective is m_o and the magnifying

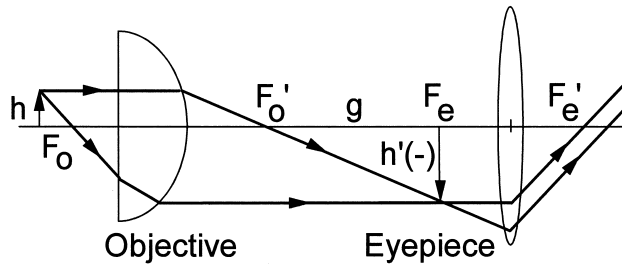


FIGURE 16 A microscope, showing the tube length g , or the distance between F'_o and F_e .

power of the eyepiece is MP_e , then the magnifying power MP of the microscope is the product

$$MP = m_o \times MP_e. \quad (18)$$

In the past, many microscopes were designed with an *optical tube length* g equal to 160 mm; g is the distance between the secondary focal point of the objective and the image point. Some manufacturers, however, use a *mechanical tube length* of 140 or 160 mm, for example, so g is only approximately equal to 160 mm. *Metallurgical microscopes*, which are designed to have a long working distance and to use no coverslip, often have optical tube lengths of 210 mm. Finally, many modern microscopes set the tube length equal to ∞ (that is, project the image to ∞) and use an *auxiliary lens* to project a real image to a finite distance. If the auxiliary lens has the correct focal length, typically 180 mm, then the image is displayed with the nominal magnification of the objective. The auxiliary lens, because it works at a high F-number, need not be highly corrected for aberrations.

If a microscope objective with a magnification of perhaps 40 or more is to be used outside the microscope for which it was designed, as for expanding a laser beam, it is important to use the proper tube length and to use a coverslip if that is called for. The coverslip may be located anywhere on the short-conjugate side of the objective and need not be in contact with the specimen. The thickness and index of refraction of the coverslip are important; most coverslips are 170- μm thick and have an index of 1.522. The coverslip produces aberrations when it is located in a diverging or converging beam (see Figure 7), and the objective is designed to correct these aberrations. Similarly, older objectives were not completely corrected for chromatic aberration, and the eyepieces were used to effect an additional correction. More modern objectives are better corrected and require different eyepieces. It is therefore important to use an objective with the proper eyepiece when you want optimum performance.

Most microscope objectives are diffraction-limited or very nearly so. They are characterized by their magnification and their *numerical aperture*. Numerical aper-

ture is the sine of the half-angle subtended by the aperture stop at the location of the object. It is related to F-number by the relation $NA = 1/(2 \times FN)$, where FN is the F-number of the lens. Though it looks like a paraxial approximation, this relation is exact for a well-corrected lens. The resolution limit of a microscope objective is therefore

$$RL = 0.61\lambda/NA. \quad (19)$$

When the eye can clearly resolve the image projected by the objective lens, there is no need to increase the magnifying power of the eyepiece. Worse, if the magnifying power is increased beyond this limit, then a single Airy disk in the image plane covers more than one cone on the retina. Each cone therefore receives reduced power, and the perceived image becomes dim. The magnifying power that achieves a diffraction-limited resolution on the retina is called the *useful magnifying power*; magnifying power in excess of that value is called *empty magnifying power*. The useful magnifying power of a microscope is given by

$$MP_u \cong 300 \times NA, \quad (20)$$

depending on the assumption concerning the angular resolution limit of the eye. In general, the magnifying power should not exceed the useful magnifying power by more than a factor of 2. The useful magnifying power will be lower if a detector with a smaller angular resolution limit than the human eye is used.

B. Microscope Illuminators

Most microscopes are used with white-light sources. An image of that source may be projected directly into the object plane of the objective; this is *critical illumination*, so called because of an early conceptual error. Critical illumination requires a very uniform light source and may be impractical, for example, if it heats the specimen. Most microscopes today use *Köhler illumination*, as illustrated in Figure 17. Köhler illumination differs from

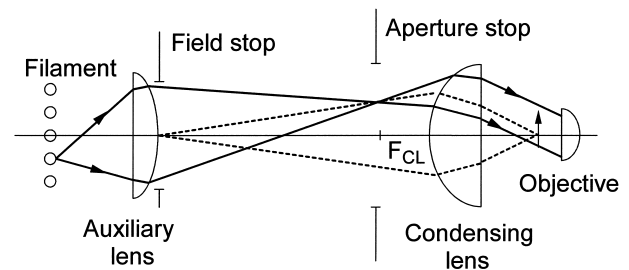


FIGURE 17 A microscope illuminator using Köhler illumination. The aperture stop is located at the primary focal point of the condensing lens.

the condensing system in a projector in that an image of the source is not projected into the aperture stop of the objective. Rather, the image of the source is projected to ∞ , so each point on the lamp gives rise to a collimated beam in the object plane. This device allows control of the numerical aperture of the condenser independently of the diameter of the illuminated area. The numerical aperture of the condenser influences the image quality of the microscope.

Specifically, if the numerical aperture of the condenser is very much less than that of the objective, the system is spatially coherent, and lines or sharp edges in the image may suffer from artifacts due to diffraction (Figure 18, solid curve), which are not present with incoherent illumination (short dashes). For truly incoherent illumination, however, the numerical aperture of the condenser must exceed that of the objective by several times. Unfortunately, that is not possible if the objective has even a modest numerical aperture, say, in excess of 0.5. Thus, many researchers set the numerical aperture of the condenser equal to that of the objective, a condition called *full illumination*. The imagery under full illumination, as in most microscopes, is partially coherent, not wholly incoherent, but the artifacts due to diffraction are reduced. Nevertheless, certain quantitative length measurements, such as the width of a stripe on an integrated-circuit mask, are extremely difficult because it is usually unclear where the geometrical images of the edges of the stripe are located when the image is partially coherent.

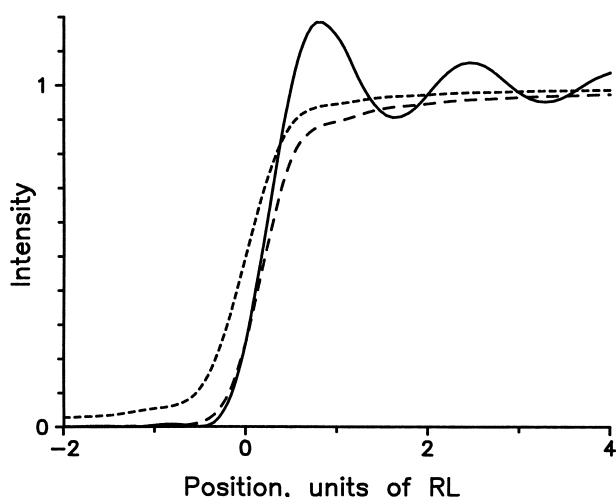


FIGURE 18 Diffraction-limited images of a sharp edge. Solid curve: conventional imaging, coherent light. Short dashes: conventional imaging, incoherent light. Long dashes: scanning confocal microscope. [From Young, M. (2000). *Optics and Lasers, Including Fibers and Optical Waveguides*, 5th ed. Springer, Berlin, Heidelberg, New York. Copyright © 2000 Springer-Verlag.]

C. Video Microscope

In a *video microscope*, the eyepiece and the eye are replaced by a CCD array such as that used in an ordinary closed-circuit television camera. The image may be viewed directly with a video monitor, or digitized and enhanced with a *frame digitizer*.

If the CCD array has the typical dimensions of 8.8×6.6 mm and has 780×480 pixels, then the pixels are approximately $11 \mu\text{m}$ wide. Unfortunately, they are not square and are about twice that value in height. For diffraction-limited performance, it is important to choose the magnification so that two pixels fall within the radius of one Airy disk; this requirement is similar to the Nyquist or sampling theorem in electrical engineering, which states that the sampling rate must exceed two samples for every period of the highest frequency in a signal. The Airy disk radius in the image plane of a typical microscope objective with $g = 160$ mm is equal to the magnification times the resolution limit and is about $20 \mu\text{m}$ for a $40\times$, 0.65-NA objective. This radius unfortunately varies with magnification, so it is hard to generalize further. Let it suffice to say that the standard $40\times$ objective matches the CCD array fairly well, as long as we are concerned with horizontal resolution only, whereas an auxiliary lens, such as a low-power microscope objective, may be necessary for any other magnification. Similarly, if the objective is infinity-corrected, the focal length of the auxiliary lens may have to be chosen carefully to match the width of the pixels in the array.

D. Scanning Confocal Microscope

Figure 19 shows a *scanning confocal microscope*, in which the specimen is illuminated one point at a time through an objective rather than a condensing lens. That point is in turn imaged through a pinhole and onto a photodetector as the specimen is scanned in a raster pattern. The output of the photodetector is typically digitized by a computer and displayed on a monitor in order to render the object visible. The microscope is called confocal because the two objectives share a common conjugate point, loosely called a focal point. If the pinhole is a point (in practice, if it is

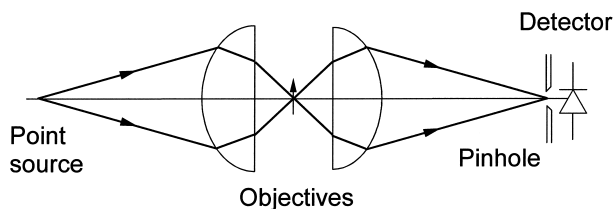


FIGURE 19 A scanning confocal microscope, drawn in transmission for clarity.

smaller than one-half the radius of the Airy disk), then the microscope detects primarily object points that are very near the object plane; other points cast large, out-of-focus images onto the plane of the pinhole and are barely detected. This is one of the main virtues of the scanning confocal microscope: It can image a three-dimensional object one plane at a time, provided that the object is in relief or is at least partly transparent. Object points that do not lie close to that plane are not merely blurred; they are virtually invisible. For this reason, a scanning confocal microscope may be used to scan a great many object planes; if the data are stored in a computer, they may be added plane by plane in order to superimpose those planes and synthesize a final image that shows enhanced depth of field without defocusing. Figure 9 (dashed curve) shows the image of a point in a scanning confocal microscope that is diffraction-limited. The width of the image is about 30% smaller than that of a conventional microscope. This is so because the scanning confocal microscope convolves the images of the two objectives and thereby sharpens the image. The image of an edge is correspondingly sharper as well, as shown by the long dashes in Figure 18.

A more practical version of the scanning confocal microscope uses a *beam splitter*, or partially transmitting mirror, to illuminate the specimen through the objective lens and images the specimen in reflection, rather than transmission. Such a microscope is especially useful in biology and in the study of integrated circuits.

Finally, it is not necessary to move the object to generate an image. Instead, the source may be scanned so that it illuminates the object one point at a time. The image of the source may be scanned optically, as with a pair of mirrors, but a simpler system based on a *Nipkow disk* has more recently been developed. In such a system, a plate (the Nipkow disk) that contains a series of holes is rotated at high speed, and each hole scans one line in the object. The image may be observed by eye or by a video camera, or it may be digitized with a computer.

E. Near-Field Scanning Optical Microscope

The resolution of a conventional microscope or of a scanning confocal microscope is limited by diffraction. The near-field scanning optical microscope (NSOM) gets around the diffraction limit by using a probe to gather an image point by point. The diameter of the probe is much less than one wavelength, and the probe is maintained only a few tens of nanometers from the object. The resolution limit of the system is approximately equal to the diameter of the probe. The probe is manufactured, for example, by heating and stretching an optical fiber or a glass pipette until it breaks. The point is coated with alu-

minum or silver in such a way that the very end of the tip has a small, uncoated region through which light may be transmitted.

Figure 20 shows a typical NSOM. The probe is held at a fixed distance from the specimen by a technique similar to those used in atomic force microscopy. The probe is brought into near contact with the specimen by means of a micropositioner. The probe tip is vibrated parallel to the surface and may be maintained a fixed distance from the surface by using electronic feedback to control the vibration.

The probe is tapered to a diameter well below the wavelength of the light. Most of the power that is directed into the probe cannot be guided by such a small fiber and therefore propagates out of the core of the fiber and into the coating or the surroundings. The transmittance of the probe is typically between 10^{-4} and 10^{-6} . The transmittance back through the probe is similar, so it is hard to measure the power directed back through the probe and toward the source.

An NSOM is more commonly designed similarly to that shown in Figure 20. The probe illuminates a very small area of the specimen, and the light that is reflected by that area is therefore diffracted in all directions. A hemispherical mirror collects the diffracted light and focuses it onto a detector. The probe is scanned in a raster to generate an image of the surface. If the surface is not a plane but has structure, that structure can also be measured by monitoring the height of the probe. If the specimen is transparent, then the hemispherical mirror may be located below the specimen or replaced by a high-aperture lens.

An NSOM may be used to detect fluorescence in a specimen. In such a system, the fluorescence wavelength

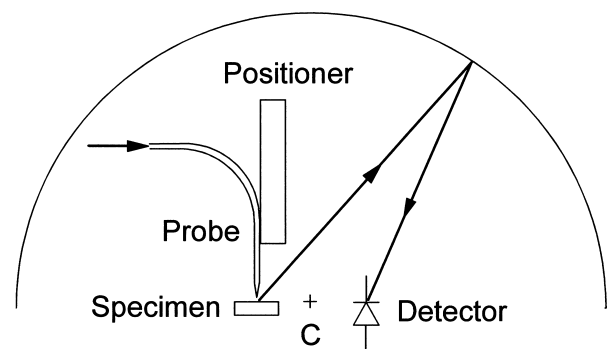


FIGURE 20 A near-field scanning optical microscope. The positioner holds a probe with a subwavelength aperture at a fixed distance from the specimen. C is the center of the hemispherical mirror that collects the light diffracted from the probe. [From Young, M. (2000). *Optics and Lasers, Including Fibers and Optical Waveguides*, 5th ed. Springer, Berlin, Heidelberg, New York. Copyright © 2000 Springer-Verlag.]

differs from the incident wavelength, so the weak fluorescence signal is distinguished from stray incident light with a mirror that selectively reflects one of the wavelengths. Similarly, the specimen itself may be illuminated and, in that case, the light guided by the probe may be detected. The mode structure of semiconductor lasers may be examined in this way. A photoconductor may be illuminated by the probe, and the resulting photocurrent may be detected as a function of position. Finally, the probe may be used to write a pattern onto a semiconductor for integrated circuits.

VII. TELESCOPE

A. Astronomical Telescope

Whereas a microscope is used to examine small objects, a *telescope* is used to magnify large, often distant objects. Specifically, the objective lens projects an image of a distant object into its focal plane, and the image is examined with an eyepiece (Figure 21). If the object is located at ∞ , the magnification of the objective is 0; hence we again use magnifying power instead of magnification to quantify the performance of the telescope.

Suppose that the distant object subtends angle α at the location of the objective. For simplicity, we place the primary focal point of the eyepiece in the image plane and project a virtual image to ∞ . The virtual image subtends angle α' at the eye, and we define the magnifying power MP of the telescope to be α'/α , whence

$$MP = -f'_o/f'_e, \quad (21)$$

where the minus sign indicates that the image is inverted. The focal length f'_e of the eyepiece may be found from the relation $MP_e = d_v/f'_e$ (see Magnifying Glass, above).

Tracing two bundles of rays through the telescope reveals a waist just behind the secondary focal point of the eyepiece. The waist is the image of the aperture stop, that is, of the ring that holds the objective in place. It is called

the *exit pupil* of the telescope. For best viewing, the pupil of the eye should be located in or near the exit pupil; if it is anywhere else, not all the object points will be visible at any one time. Most telescopes are designed so that the exit pupil has a diameter equal to the nominal diameter of the pupil of the eye or 5 mm.

As with the microscope, the magnifying power of the eyepiece need be increased only to the point that the resolved detail in the image plane can be resolved by the eye. Similarly, depending on the assumptions about the angular resolution of the eye, we find that the useful magnifying power of the telescope is

$$MP_u = 5D_{[cm]}, \quad (22)$$

where $D_{[cm]}$ means the diameter of the objective measured in centimeters. The magnifying power of the telescope should not exceed twice the useful magnifying power.

If the telescope is to have a sizable field of view, as measured by the angle α , then the eyepiece could become prohibitively large. To reduce the diameter of the eyepiece, a *field lens* is installed near the secondary focal plane of the objective. The field lens, which may be part of the eyepiece, directs the principal ray toward the axis and allows the diameter of the eyepiece to be reduced. Typically, the focal length of the field lens may be chosen so that the principal ray intersects the eyepiece at its periphery. This shifts the location of the exit pupil somewhat toward the eyepiece and slightly reduces the *eye relief* between the eye and the eyepiece. When the field lens is built into the eyepiece, the second lens in the eyepiece is called the *eye lens*.

B. Terrestrial Telescope

An astronomical telescope yields an inverted image. A pair of *binoculars* uses two reflecting prisms to invert the image and make it erect. What is conventionally called a *terrestrial telescope* uses a *relay lens*, typically at unit magnification, to invert the image. Figure 22 shows a terrestrial

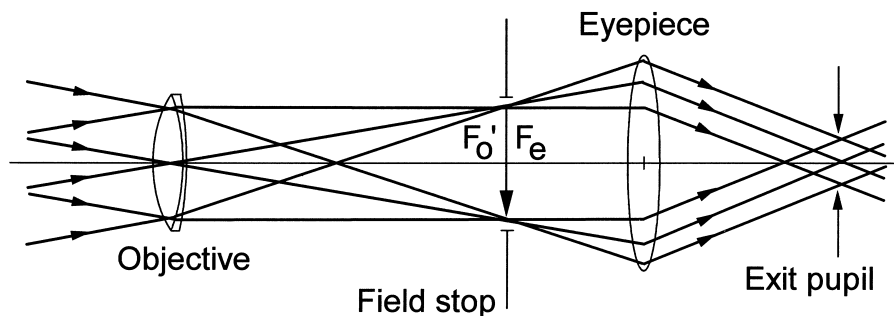


FIGURE 21 An astronomical telescope, showing the field stop and the construction of the exit pupil.

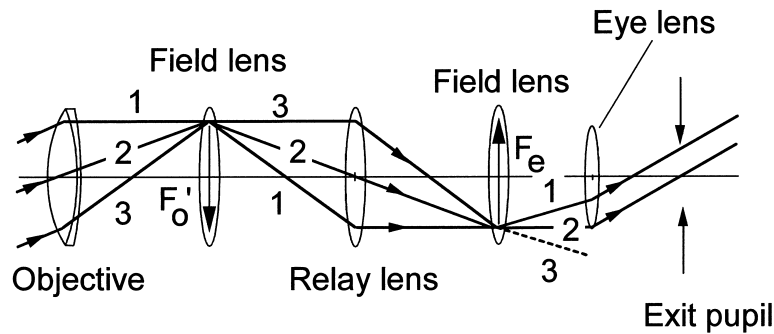


FIGURE 22 A terrestrial telescope, which includes a relay lens for inverting the image and a field lens for reducing the diameter of the relay lens.

telescope that includes a field lens, which is used to reduce the diameter of the relay lens, in the image plane of the objective. The field lens projects an image of the aperture stop into the plane of the relay lens. That image of the aperture stop is called a *pupil* and is analogous to the exit pupil of a telescope or microscope. Figure 1 shows that there is an image everywhere that the ray from O to O' crosses the axis; similarly, in a telescope (or any other optical system), there is a pupil everywhere that the principal ray crosses the axis. The diameter of the relay lens is set equal to the diameter of that pupil.

Unlike most field lenses, the field lens directly before the eyepiece does not project a pupil into the plane of the eye lens because, if it did so, the eye relief would be 0. In almost all other cases, however, the function of a field lens is to project a pupil into the plane of some other lens, and, indeed, the combination of the field lens and the eyepiece projects a pupil into the pupil of the eye. Any system of relay lenses similar to a terrestrial telescope, such as a periscope (which uses mirrors to look along a line of sight displaced from the line of sight of the eye), must necessarily use the principle of projecting a pupil into the apertures of successive lenses. The diameter of any lens is chosen to be equal to the diameter of the appropriate pupil, much as the diameter of a projection lens is matched to the image of the projection lamp.

ACKNOWLEDGMENTS

Thanks to Kent Rochford and Steve Kreger of the National Institute of Standards and Technology and Frank Kowalski of the Colorado School of Mines for their careful and critical reading of the manuscript.

SEE ALSO THE FOLLOWING ARTICLES

HOLOGRAPHY • MICROSCOPY • OPTICAL DIFFRACTION • TELESCOPES, OPTICAL

BIBLIOGRAPHY

- Bass, M., and van Stryland, E. (eds.) (1994). "Handbook of Optics: Devices, Measurements, Properties, 2nd ed," McGraw-Hill, New York.
- Corle, T. R., and Kino, G. S. (1996). "Scanning Microscopy and Related Imaging Systems," Academic Press, San Diego.
- Hecht, E. (1998). "Optics, 3rd ed," Addison Wesley Longman, Reading, MA.
- Inoue, S., and Spring, K. R. (1997). "Video Microscopy: The Fundamentals, 2nd ed," Plenum, New York.
- Smith, W. J. (1990). "Modern Optical Engineering, 2nd ed," McGraw-Hill, New York.
- Young, M. (2000). "Optics and Lasers, Including Fibers and Optical Waveguides, 5th ed," Springer, Berlin.



Liquid Crystal Devices

Michael G. Clark

GEC Hirst Research Centre

- I. Impact of Liquid Crystal Devices
- II. Nematic Devices
- III. Guest–Host Effect
- IV. Megapixel Technologies
- V. Other Devices

GLOSSARY

Active matrix Technique for addressing a matrix display in which each pixel is controlled by an individual semiconductor device incorporated into the structure of the panel.

Birefringence In an optically anisotropic medium with axial symmetry, the difference between the refractive index for light polarized parallel to the symmetry axis (extraordinary index) and that for light polarized perpendicular to the axis (ordinary index).

Cholesteric Chiral (optically active) nematic phase characterized by the director assuming a helical configuration in the absence of external fields and forces.

Dielectric anisotropy In an electrically anisotropic medium with axial symmetry, the difference between the permittivity for fields parallel to the symmetry axis and that for fields perpendicular to the axis.

Director Thermally averaged direction of the molecular orientation in a liquid crystal (especially nematic, cholesteric, or smectic A) phase.

Dyed phase change Device effect that uses an anisotropic dye dissolved in cholesteric liquid crystal to switch between absorbing and nonabsorbing states.

FELCD Ferroelectric liquid crystal device; a device exploiting the ferroelectric property of chiral tilted smectic phases (especially smectic C*).

Frame time Time period for one complete addressing cycle of a display.

Guest–host effect Orientational ordering of a solute molecule by a liquid crystal solvent.

ITO Indium tin oxide; heavily tin-doped indium oxide used for the transparent electrodes in liquid crystal devices.

Line time Time taken to address one line of a matrix display.

Matrix display Display in which the pixels are defined by the intersections of row and column electrodes.

Multiplexing Technique for addressing a matrix display whereby all possible combinations of on and off pixels can be shown.

Nematic Liquid crystal phase in which the constituent molecules show partial orientational ordering, with axial symmetry, but no translational order.

Ordinary/extraordinary ray Ray of light propagating through a birefringent medium in any direction other than parallel to the symmetry axis that splits into an ordinary ray polarized perpendicular to the symmetry

axis and an extraordinary ray polarized perpendicular to the ordinary ray.

Pixel One-addressable picture element of an electronic display.

SBE Supertwist birefringence effect; device effect in a 180° to 270° twisted nematic structure that uses a combination of interference and polarization guiding to achieve an optical effect with sharp threshold behavior, making it particularly suitable for multiplexing.

SLM Spatial light modulator; device using either electrical or optical addressing to impress information onto a light beam.

Smectic A Liquid crystal phase combining nematiclike orientational ordering with one-dimensional partial translational order in the direction parallel to the director.

Thermotropic liquid crystal Substance showing one or more liquid phases appearing as a function of temperature between the solid and the isotropic liquid, in which the constituent molecules show partial orientational ordering and possible partial translational order.

Tilted smectic Liquid crystal phases combining nematiclike orientational ordering with one-dimensional partial translational order in a direction tilted at a nonzero angle to the direction of preferred orientational order.

Twisted nematic Device effect in which the director configuration in the unpowered condition is a 90° twist with the helical axis perpendicular to the plane of the cell.

LIQUID CRYSTALS are widely used in devices for flat-panel electronic displays and other applications. The most usual form of device is one in which a thin film, 5–10 μm thick, of liquid crystal is confined between glass plates bearing transparent electrodes that are used to induce an optical change in selected portions of the film by application of an above-threshold voltage. Liquid crystal displays with color and video display capabilities equal to those of the cathode ray tube are now being developed.

I. IMPACT OF LIQUID CRYSTAL DEVICES

It was realized in the early 1970s that the unusual properties of thermotropic liquid crystals held great promise for use in flat-panel electronic displays and other optical control applications. The advantages particular to liquid crystals of a very large (if not especially fast) electro-optic effect induced by CMOS-compatible voltages and of microwatts per square centimeter power consumption were identified at an early stage. With the discovery of chemically stable nematic liquid crystals, such as the

TABLE I The Worldwide Market for Flat-Panel Electronic Displays

	1990	1991	1992
Total market in millions of U.S. dollars ^a	3826	4579	5565
Distribution by technology as percentage ^a			
Liquid crystal displays	57.8	59.4	60.1
Vacuum fluorescent displays	20.0	18.9	17.7
Plasma display panels	10.8	10.8	11.6
LED displays	8.1	7.2	6.2
Electroluminescent panels	2.7	3.1	3.9
Others	0.6	0.6	0.5

^a Dot-matrix types include electronics.

cyanobiphenyls, it also became clear that an ac-addressed voltage-induced (as opposed to a current-induced) effect in a dielectric liquid was inherently long-lived and offered exceptionally high reliability. Being essentially a light-modulating, rather than light-emitting, device, liquid crystal displays (LCDs) give very good performance in high ambient lighting, with the contrast ratio essentially constant over a wide range of illumination. Thus, LCDs first came to public attention through their use in wristwatches and pocket calculators, both innovative products only made possible through this technology. However, as the advantages of flat format, CMOS drive, and high reliability have become increasingly appreciated, LCDs have been used more and more in situations where auxiliary lighting has to be provided, usually in the form of back-lighting with fluorescent discharge lamps. As a result of its intrinsic simplicity combined with heavy investment in research, liquid crystal technology is already the cheapest per pixel for complex (more than 50,000 pixels) displays, with the exception of the cathode ray tube (CRT), and is predicted to fall below the color CRT in cost after 1992.

As shown in Table I, LCDs now form over half the total flat panel display market, the remainder being fragmented between four emissive technologies. The market for LCDs in 1990 was \$221.3 million and was described as growing at 23% between 1990 and 1991, with equal growth the year after. Articles in this work introduced both the field of flat-panel displays and liquid crystals. In this contribution, the rapidly growing field of liquid crystal devices is reviewed more comprehensively.

II. NEMATIC DEVICES

A. Twisted Nematic

The twisted nematic (TN) is the prototype LCD. Its structure and operation are worth comprehending in detail

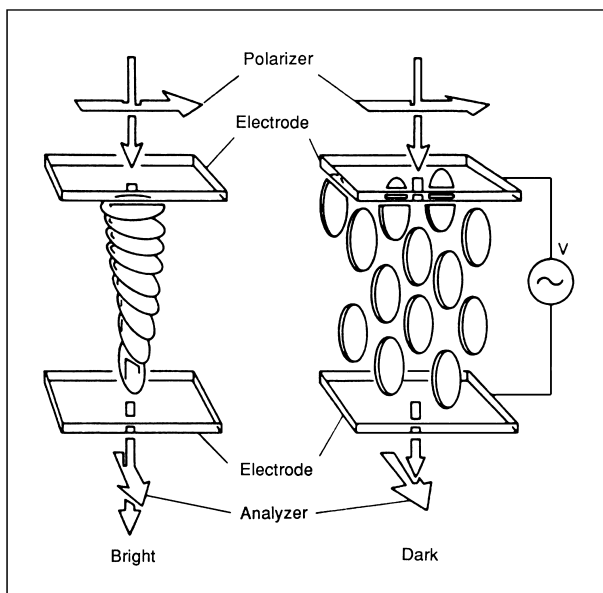


FIGURE 1 Twisted nematic device effect.

both because it presently comprises the vast majority of LCDs in use and because the variety of alternative LCD technologies that are now being developed can readily be understood by comparison and contrast with the TN.

As illustrated in Fig. 1, a TN cell comprises two glass plates between which is confined a thin film of nematic liquid crystal. The plates bear transparent electrodes of indium tin oxide (ITO) that are etched to the desired patterns by conventional photolithographic techniques. Above the electrode layer is a surface alignment layer. The function of this layer is to fix the average orientation of the molecules (called the director) at the surface. The usual technique is to apply a layer of polyimide and to rub this with a roller coated with velvet or nylon cloth. This orients the director parallel to the rubbing direction and at a pretilt angle of approximately 2° to the surface (Fig. 2). The rubbed plates are set at 90° so as to impart a 90° -twisted configuration on the director. The cell is sealed at its edge with epoxy and maintained at a uniform spacing throughout by use of fiber or spherical spacers. For the TN, this spacing is typically 5 to $10\ \mu\text{m}$ with a tolerance of about $\pm 0.5\ \mu\text{m}$. A schematic flowchart for LCD fabrica-

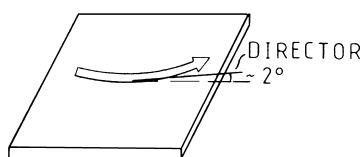


FIGURE 2 Pretilted surface alignment of the director by unidirectional rubbing of a polymer layer.

tion is shown in Fig. 3; note that this chart includes details relevant only to more advanced technologies discussed below.

The basis of the operation of the TN LCD is that both the refractive index and the permittivity vary with the orientation of the electric vector relative to the director. Thus, light polarized parallel to the director experiences the extraordinary refractive index, n_e , whereas the polarization perpendicular to the director experiences the ordinary refractive index, n_o . As a result, provided that $(n_e - n_o)d$ is sufficiently greater than λ , where d is the thickness of the liquid crystal film and λ the wavelength of light *in vacuo*, the twisted configuration shown on the left-hand side of Fig. 1 will guide the plane of plane-polarized light. The device thus transmits when placed between crossed polarizers. Further, provided that the liquid crystal used in the cell has its permittivity for fields parallel to the director, ϵ_{\parallel} , greater than that for fields perpendicular to the director, ϵ_{\perp} , there exists a critical voltage, V_c , given by the following:

$$V_c = \frac{1}{2}\pi[(4K_{11} - 2K_{22} + K_{33})/\epsilon_0(\epsilon_{\parallel} - \epsilon_{\perp})]^{1/2},$$

where K_{11} , K_{22} , and K_{33} are (orientational) elastic constants measuring the stiffness of the director against splay, twist, and bend distortions, respectively. Above this voltage, the 90° -twisted director configuration begins to distort, leading to the configuration shown on the right-hand side of Fig. 1, which does not guide plane-polarized light. Thus, between crossed polarizers, the device goes dark.

This effect has been successfully implemented in seven-segment and similar fixed format displays and in simple dot-matrix displays. Such displays may be employed in a reflective mode, backlit, or in a "transflective" combination of both these modes. Although not outstandingly attractive in appearance, the TN LCD has proved to be versatile, relatively cheap to manufacture and drive, and reliable. Its most obvious limitations are speed (the turnoff time can be tens or hundreds of milliseconds), and the difficulty that it experiences in showing complex information.

The problem of displaying complex information is tied up with that of multiplexing. All multiplexed LCDs are electrically equivalent to the dot-matrix format, in which the picture elements (pixels) form a rectangular array defined by the intersection of row and column electrodes. Thus, N row electrodes and M column electrodes together define $N \times M$ pixels. Even without going into details, it will be evident that if only $N + M$ electrodes are used to address $N \times M$ pixels there are limits on what can be done. Because of its relatively slow response time, the TN responds to the root mean square (rms) of the applied voltage. As a result, it is possible to calculate quite generally, without reference to any specific addressing waveforms, the upper limits attainable by the ratio $V_{\text{rms}}(\text{ON})$:

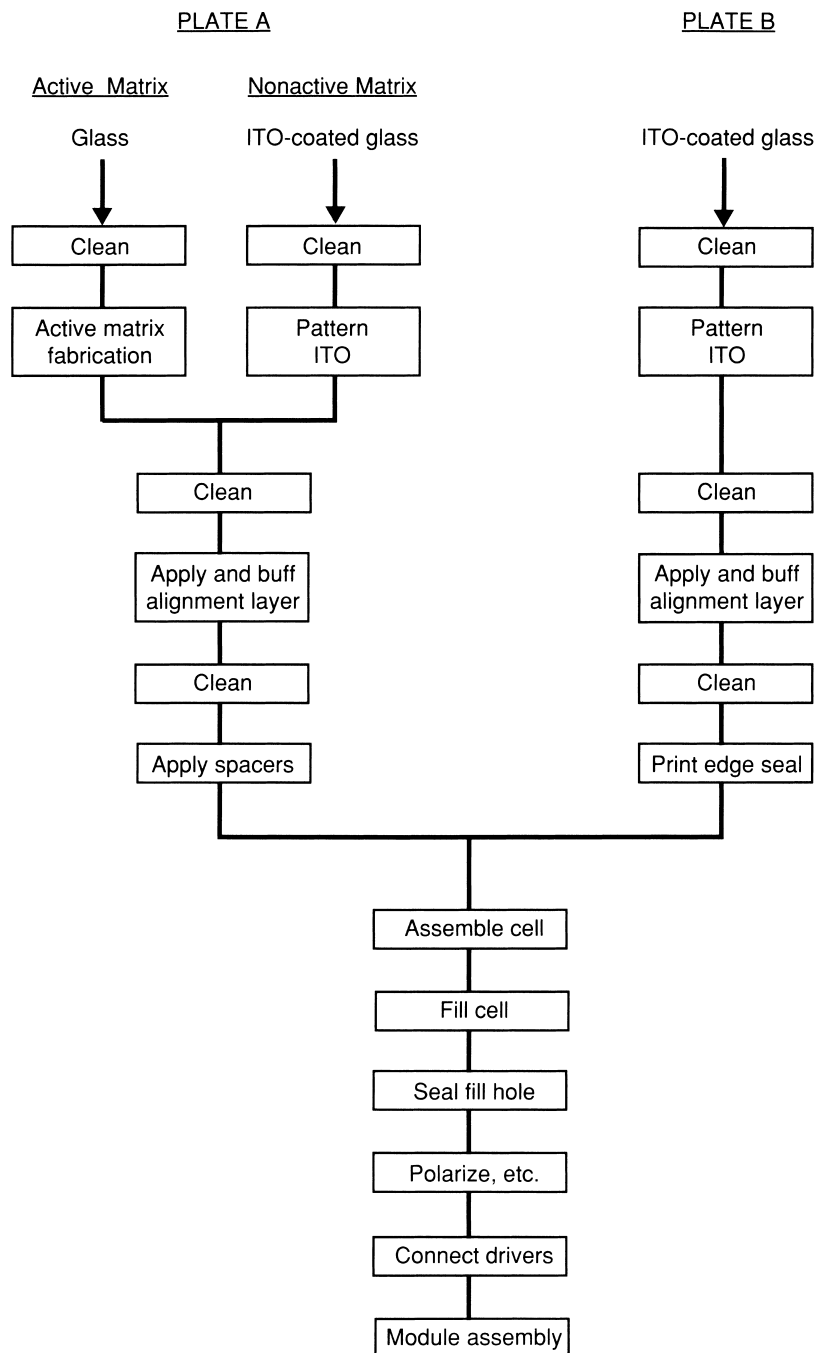


FIGURE 3 Schematic flowchart for the fabrication of a liquid crystal display.

$V_{\text{rms}}(\text{OFF})$ of the voltages that can be delivered to ON and OFF pixels in a multiplexed matrix. This ratio is shown as a function of the number of row electrodes N in Table II.

From Table II, it is evident that the degree of discrimination between ON and OFF pixels decreases rapidly as the level of multiplexing increases, thus demanding an ever sharper threshold for the electro-optic effect in

the liquid crystal. Unfortunately, as Fig. 4 shows, the TN electro-optic effect has a relatively shallow threshold curve that varies strongly with the angle of view, becoming nonmonotonic at extreme angles. For reference in Section II.B, Fig. 4 also illustrates the fact that guiding of plane-polarized light is not lost immediately the director configuration begins to distort at V_c , as evidenced by the

TABLE II Maximum $V_{\text{rms}}(\text{ON}):V_{\text{rms}}(\text{OFF})$ Values Attainable in Multiplexing

Number of rows (N)	Practical limit ^a	Ultimate limit
2	2.414	3
3	1.932	2
4	1.732	1.732
5	1.618	1.633
6	1.543	1.545
7	1.488	1.489
8	1.447	1.450
9	1.414	1.414
10	1.387	1.389
100	1.106	1.106
200	1.073	1.073
500	1.046	1.046
1000	1.032	1.032

^a Practical limit = $[(\sqrt{N} + 1)/(\sqrt{N} - 1)]^{1/2}$ and is exactly equal to the ultimate limit only when \sqrt{N} is an integer.

capacitance–voltage curve, causing the optical threshold to be at a higher voltage than V_c .

As a result of the above considerations, the TN has found only limited acceptance in dot-matrix applications much above $N = 25$ to 50. The size of the matrix can be doubled by juxtaposing two electrically distinct matrices, the columns of one being addressed from the top and of the other from the bottom of the viewing area. Methods for quadrupling the size of the matrix for a fixed multiplexing level have also been suggested, but they are difficult to fabricate in acceptable yield.

B. Supertwisted Nematic

In 1982, over a decade after the invention of the TN, it was found that more tightly twisted configurations

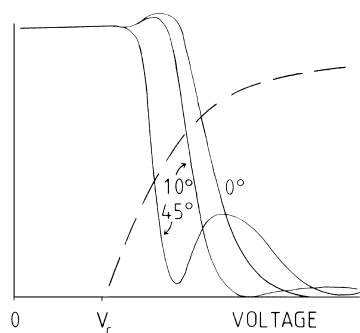


FIGURE 4 Variation with rms voltage of transmission with crossed polarizers for three angles of incidence (solid lines) and cell capacitance (broken line) for a twisted nematic device.

had a sharper threshold. The addition of chiral nematic (cholesteric) material to the normal, optically inactive nematic imparts a natural helical twist to the director configuration. Small amounts of cholesteric dopant are added to the material in a TN cell in order to discriminate between the left-handed and right-handed 90° twist. If larger amounts are added so that the natural pitch of the director's helical configuration is approximately four-thirds of the cell thickness, the director will assume a 270° twist in a TN-like cell. This structure was found to have a much sharper electro-optic threshold curve than the TN.

At first, the device was thought of primarily in the context of the guest–host effect (Section III), but it was quickly realized that the preferred embodiment was a two-polarizer device similar to the TN. Since the twist is much tighter than the TN, guiding is rather imperfect and a planepolarized beam emerges from the 270° structure elliptically polarized. The resultant optical effect is therefore as much an interference effect between ordinary and extraordinary rays as it is a guiding effect. Thus, the liquid crystal film must be constant in thickness to within a fraction of the wavelength of light, the required cell gap tolerance being $\pm 0.1 \mu\text{m}$.

Intensive investigations have shown that, depending on material and cell parameters, the effect can be obtained with twists ranging from 180° to 270° . The highest twist gives the best angle of view but needs a pretilt greater than 5° if formation of scattering defects, because of other director configurations close by in free energy, is to be avoided. Although production-compatible methods for obtaining higher pretilt have been developed, the high-volume manufacturers are reluctant to make any deviation from standard production processes and are at the moment contenting themselves with lower twists, 240° or less.

As mentioned, the supertwist birefringence effect (SBE) is essentially an interference effect and therefore is inherently chromatic. Two modes are possible depending on the relative orientations of the polarizers and the rubbing directions. These are the yellow mode (dark blue on a yellow background) and the blue mode (gray on a blue background). In order to obtain good performance, the cell must be constructed with a gap that is not only uniform, as noted above, but also has the correct thickness for the birefringence and other parameters of the material used. The manufacture of these devices is, therefore, significantly more demanding than the TN. Products with levels of multiplexing up to 200 and higher have been offered, although whether the performance of SBE at greater than 100-way multiplexing will receive widespread user acceptance remains to be determined. Manufacturers attempt to distinguish their products by a variety of names, such as supertwisted nematic (STN) and hypertwisted nematic, for commercial rather than technological reasons.

Several approaches have been proposed to provide a black-on-white SBE-type display. Probably the most effective is to place a second, color-compensating cell between the active display cell and the front polarizer. The compensating cell needs no electrodes because it is not switched, but nevertheless it clearly adds considerably to the cost of the display and its fabrication occupies production capacity. The effect of a color-compensating cell can be approximated by laminating one or more layers of polymer retardation film between the polarizer and the glass of the active cell, but this arrangement gives a poorer angle of view, and it is less certain that the compensation can be adequately maintained over a wide temperature range since the birefringence of the retardation film is likely to have a different temperature dependence to that of the liquid crystal, whereas a compensating cell can be filled with liquid crystal having the same temperature dependence as the active material. Devices of this type are required to satisfy user demand for achromatic displays and so that color displays can be constructed using a backlit configuration with microfilter-equipped pixels dedicated to the primary colors red, blue, and green. Uncompensated, compensated, and to a lesser extent full-color SBE LCDs have found acceptance for a number of applications, most prominently for the display panels of laptop computers.

Finally, it will be evident that since the threshold curve of a supertwist device is much sharper than for the TN (cf. Fig. 4) the whole of the curve lies closer to the appropriate threshold voltage V_c than is the case for the TN. Unfortunately, the region close to V_c is characterized by the phenomenon of “critical slowing down,” which leads to slow device response time. As a result, the SBE device is inherently slower than the TN, although ways to overcome this, for example, low viscosity liquid crystals and thinner cells (dynamics \propto thickness²), has been studied. Alternative device configurations, in which the liquid crystal layer acts as a birefringent layer whose birefringence is electrically controllable, are also being developed in an attempt to combine the optical performance of SBE with faster response.

III. GUEST–HOST EFFECT

A. Anisotropic Dyes

When solute molecules are dissolved in liquid crystal, the molecular interactions between the solute and the partially ordered liquid crystal molecules result in some degree of orientational ordering being imparted to the solute molecules. This is the guest–host effect.

If the solute is a dye, the absorbance A will be different for light polarized parallel (A_{\parallel}) and perpendicular (A_{\perp})

to the nematic director. An optical order parameter S_{op} is defined by the following equation:

$$S_{op} = (A_{\parallel} - A_{\perp}) / (A_{\parallel} + 2A_{\perp}).$$

The value of S_{op} ranges from unity, if the dye transition moment is uniformly parallel to the director, through zero, if the dye orientation is completely randomized to $-\frac{1}{2}$ if the dye transition moment is always perpendicular to the director. Note that S_{op} depends both on the underlying orientational order of the liquid crystal molecules and on the intrinsic molecular anisotropy of the dye. Suitable dyes have been called both pleochroic dyes and dichroic dyes (the latter is more exact); we call them simply anisotropic dyes.

Anisotropic dyes for use in LCDs must satisfy a number of criteria of which the most important are a large S_{op} , a high absorption coefficient, high solubility, stability, purity, high resistivity, and safety.

The interplay of these criteria is well illustrated by the two main classes of compound that have historically dominated the development of anisotropic dyes. Azo dyes tend to offer better order parameters, larger extinction coefficients, and better solubility in liquid crystals, whereas the anthraquinone materials tend to have better photochemical stability. The position now is that following a number of years of intensive development stable anisotropic dyes in a variety of colors are available commercially. Colors may be mixed to give black, high-order-parameter mixtures. Note that such blacks are metamerics and in principle will change hue with changes in the illuminant, although in practice very acceptable mixtures can be obtained.

Although the dyes are used in relatively low concentrations, their effect on the material parameters of the host liquid crystal should not be overlooked. It has been suggested that some anthraquinone dyes can have a disproportionate slowing down effect on the electro-optic response time because of dye-induced changes in the viscosity coefficients.

B. Single-Polarizer Guest–Host Devices

In these devices, the director is electrically switched between configurations in which plan-polarized light propagating through the cell experiences either A_{\parallel} or A_{\perp} . The most obvious configurations satisfying this requirement are either a homogeneously aligned cell containing a material with positive dielectric anisotropy or a homeotropically aligned cell with negative dielectric anisotropy (Fig. 5), although a twisted nematic cell with the polarizer aligned parallel to one of the rubbing directions will also work. Both the homogeneous and the homeotropic alignments should have small pretilts to bias the cell so that all switched areas tilt in the same direction. The need for tilted

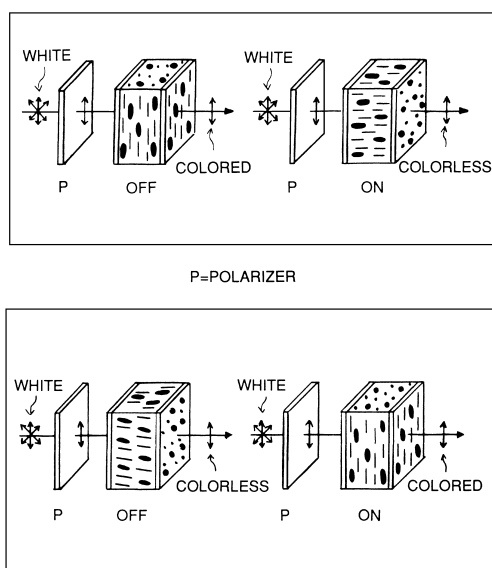


FIGURE 5 The upper panel shows homogeneously aligned, negative contrast, and the lower panel shows homeotropically aligned, positive-contrast, single-polarizer guest-host device effects.

homeotropic alignment and for (less readily available) negative dielectric anisotropy materials militate against the homeotropic configuration.

Single-polarizer guest-host devices have found successful application in public information displays, such as are used in railway stations, airports (departure gate displays at London's Heathrow Airport were an early example), and financial dealing floors. Characters are made up in dot-matrix format, but the dimensions are such that each pixel can have its own drive electrode. Modules comprising one or more characters are constructed so that they can be slotted together to make the complete information board. Each module has its own backlight and microprocessor control. These displays have the advantages of low maintenance costs, high reliability, and computer control with automatic fault detection.

C. Dyed Phase-Change Device

An approach that avoids the use of polarizing film altogether is to dissolve anisotropic dye into a shortish pitch (a few micrometers) cholesteric liquid crystal. The physics of the effects that can be obtained is complex (see also Section IV.C) and is oversimplified for brevity. In the absence of voltage across a cell containing such material, the twisted helical configuration of the director causes the dye to present what is virtually a pseudorandom configuration that is a relatively good absorber of unpolarized light. On application of a sufficient voltage to the material (which should have positive dielectric anisotropy), it

switches to a configuration similar to that shown on the right-hand side of Fig. 1 and which has minimal absorption of light. This switching has been termed (somewhat loosely) a cholesteric-to-nematic phase change; hence, the name of the device.

The device gives bright ON pixels against a colored background. The avoidance of polarizers increases brightness and removes one source of environmental weakness. However, the effect cannot readily be multiplexed and requires a relatively high (although still CMOS-compatible) drive voltage. One advantage shared with single-polarizer guest-host devices is that in the reflective mode the reflector may be incorporated into the inner surface of the back wall of the cell, thus avoiding parallax. However, this advantage does not appear to have been commercially exploited yet.

Black-dyed phase-change devices have been accepted as replacements for mechanical alphanumeric displays in aircraft cockpits. This is because their angle of view is superior to the TN and their appearance more closely resembles the mechanical displays that they are replacing.

D. Fluorescent LCDs

If the anisotropic absorbing dye in a guest-host device is replaced by an anisotropic fluorescent dye, the visual impact of an emissive display may be combined with the virtues of LCD technology. Both single polarizer and dyed phase-change configurations have been demonstrated. The fluorescence is stimulated by use of a UV backlight, although the phase-change device may also be stimulated by the UV contained in ambient sunlight. UV backlights are available in the same sizes and formats as fluorescent backlights (to which they are related by omission of the lamp phosphor), and the backlight drive circuitry and power consumption are essentially identical for both. The best developed anisotropic fluorescent dyes are green fluorophors based on the perylene structure, although red- and blue-emitting fluorophors are also available.

Since the light in a fluorescent LCD is emitted in the liquid crystal film, the geometry of refraction of the light as it leaves the cell causes the display to have a hemispherical angle of view. This technology thus offers the visual impact and good angle of view of an emissive display in applications where a backlit LCD might otherwise be the only alternative.

IV. MEGAPIXEL TECHNOLOGIES

In this section, we review LCD technologies with the capability to yield dot-matrix displays containing over 200 rows of pixels. Such displays are needed for an

enormous range of video uses, information technology, and professional applications. Broadly speaking, these displays are CRT replacements, although it cannot be emphasized too strongly that they also have the capability to innovate new products for which the CRT, with its bulk, fragility, and high-voltage drive, is unsuitable.

A. Active Matrix

The essential principle of an active matrix display is that each pixel has associated with it a semiconductor device that is used to control the operation of that pixel. It is this rectangular array of semiconductor devices (the active matrix) that is addressed by the drive circuitry. The devices, which are fabricated by thin-film techniques on the inner surface of a substrate (usually glass) forming one wall of the LCD cell, may be either two-terminal devices (Fig. 6) or three terminal devices (Fig. 7). Various two-terminal devices have been proposed: ZnO varistors, MIM devices, and several structures involving one or more *a*-Si diodes. Much of the research effort, however, has concentrated on the three-terminal devices, namely thin-film, insulated-gate, field-effect transistors. The subject of thin-film transistors (TFTs) is considered elsewhere in this volume; suffice it to say that of the various materials that have been suggested for the semiconductor, only *a*-Si and poly-Si appear to have serious prospects of commercial exploitation.

Amorphous silicon TFTs are the more developed of the two, research on them having started earlier. Impressive active matrix LCDs have been fabricated using this technology, although longer term polysilicon offers a number

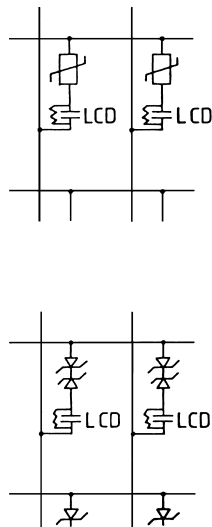


FIGURE 6 Two-terminal active matrix addressing. The pixel is shown by its equivalent circuit of a parallel capacitor-resistor combination.

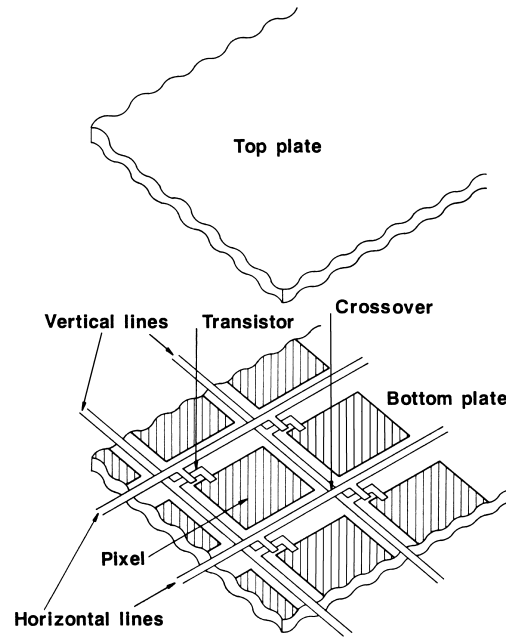


FIGURE 7 Thin-film transistor active matrix addressing. [Courtesy of GEC Hirst Research Centre.]

of advantages: It has a higher mobility and a lack of photo-sensitivity; it gives smaller, self-aligned, transistors; it has better stability and a longer lifetime (the structural hydrogen in *a*-Si is a worry); and there is the possibility of fabricating thin-film drive circuitry on the same substrate as the display. A low-pressure chemical vapor deposition technique for depositing TFT-quality polysilicon at temperatures (630°C or below) compatible with borosilicate glass has been developed, and promising results have been obtained with recrystallization of *a*-Si to yield the very high mobility material required for high-speed driver TFTs.

Active matrix addressing, having been proposed for both LCD and electroluminescence flat-panel technologies, has been the subject of intensive research since the early 1970s. The issue has become not to demonstrate feasibility in the laboratory but to design structures and processes that will give an economically viable yield. One difficulty with the classic structure (Fig. 7) is that a short either at a gate-line/dataline crossover or from a gate to drain or source will cause line defects in the display. Thus, methods have been developed for repairing such faults by using a laser to isolate them, thereby reducing them to point (single-pixel) defects. Such procedures are expensive in production, and an alternative is to use the capacitively coupled transistor (CCT) configuration, shown in Fig. 8, which eliminates crossovers and ensures that gate-drain or gate-source shorts produce only point defects. This configuration requires low parasitic capacitances, as obtained with self-aligned polysilicon TFTs.

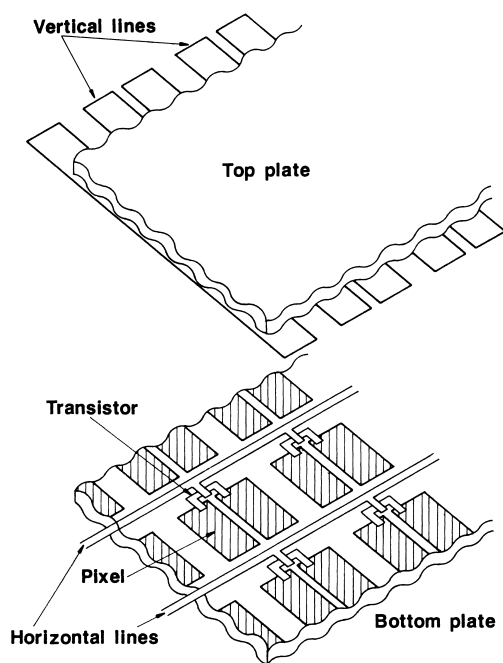


FIGURE 8 CCT fault-tolerant active addressing architecture. [Courtesy of GEC Hirst Research Centre.]

Although an active matrix could in principle drive a variety of liquid crystal effects, most attention has been devoted to driving the twisted nematic effect. Displays are normally used in the transmissive mode, backlit, so that color can be obtained by use of arrays of microfilters aligned with the pixels. The color pixels thus obtained are conventionally triads (red, R; green, G; and blue, B), although (red/green/blue/green) and (red/green/blue/white) arrangements have also been used, the former to improve color balance and the latter to improve brightness (since fully saturated primary colors are relatively infrequently required). The (R/G/B) triads should be arranged in vertical stripes for alphagraphic displays (so that vertical lines of a given color are free of “jaggies”) and in triangles for video displays (since the breaking up of rectangular outlines produces a more attractive result). Several technologies for fabricating color microfilter arrays are available. The filter layer is normally included on the inner surface of the non-active-matrix plate of the cell, either over or under the ITO layer, preferably the latter to avoid voltage drop. A major difficulty is the very high attenuation of the backlight caused by the combination of polarizers and color filters; as little as 5% or less of the luminous energy is transmitted.

In the early 1990s, the largest active matrix displays publicly demonstrated were 15-in.-diagonal devices using 1600×1920 *a*-Si TFTs. Color TV products, also using *a*-Si, with up to 5-in.-diagonal screens, were reportedly

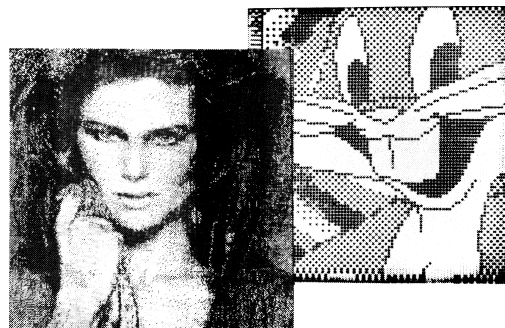


FIGURE 9 The active matrix (left) and ferroelectric (right) LCD technologies are the main contenders for picture-forming display applications. [Courtesy of GEC Hirst Research Centre.]

on sale, although only 2- and 3-in. screens were widely available. The true economics of all these devices remain obscure. Nevertheless, the demonstration that active matrix addressing can offer LCDs with video capability, gray scale, and color has given great impetus to the field. To the extent that one can predict trends, it is likely that efforts will now concentrate on increasing the screen-size capability of *a*-Si technology and, in parallel, developing poly-Si technology with peripheral drive circuitry integrated onto the glass substrate. Figure 9 is a montage of a 5-in.-diagonal polysilicon active matrix display and a ferroelectric LCD (Section IV.B).

B. Ferroelectric LCDs

The limitations on multiplexing any rms-responding monostable liquid crystal effect have been mentioned in Section II.A. Active matrix addressing, described in Section IV.A, is one way of overcoming these limitations. Another is to consider alternative liquid crystal effects that are bistable, or at least non-rms responding. With such effects, the maximum number of rows that can be multiplexed is usually determined by the ratio of the frame time (the time period during which the whole picture must be refreshed or updated) to the line time (the time required to address one row of pixels). This is quite demanding of the line time; a frame time of 40 msec (only 25-Hz frame rate) would require a line time of 40 μ sec for 1000 lines. Bistable behavior is associated with smectic and cholesteric phases, both of which in completely different ways have translational symmetries added to nematolike orientational order. In this section, the ferroelectric tilted smectic devices are reviewed, while (untilted) smectic A and cholesteric devices are described in Section IV.C.

Most liquid crystal phases are centrosymmetric and nonferroelectric. Symmetry low enough for ferroelectricity to occur is found only in smectic phases that are both tilted (as explained below) and optically active (chiral).

Of these, the most important (and which, for simplicity, we shall consider exclusively) is the chiral smectic C, denoted smectic C*. In this phase, like all smectics, partial translational ordering of the molecular centers of gravity is superimposed on the orientational ordering. In the case of smectic C, this translational order can be thought of as a mass density wave as follows:

$$\rho(z) = \rho_0[1 + a \cos(qz - \phi)],$$

where the period $2\pi/q$ of the fluctuations in the number density $\rho(z)$ is of the order of one or two molecular lengths. The direction z is the direction of the number density wave vector and is perpendicular to layers of constant $\rho(z)$. In a tilted smectic phase, it makes a nonzero angle with the director. If the phase is chiral, it also forms the helical axis about which the director spirals. Figure 10 attempts to illustrate these points. In the chiral smectic C, the local point symmetry is a single, twofold axis perpendicular to the mass density wave. This is low enough for spontaneous polarization, that is, ferroelectricity, to develop parallel to this twofold axis. The ferroelectricity is therefore associated with noncentrosymmetric ordering of the transverse component of the molecular electric dipole, not the longitudinal one.

In a bulk sample of smectic C*, the helical ordering due to the chirality causes the spontaneous polarization vector \mathbf{P}_s to average to zero over the sample. Device applications arose from the realization in 1980 that, in a sufficiently thin liquid crystal cell, surface alignment forces might overcome the helical twist. In the form originally envisaged, it was supposed that the directors at the surface could be aligned parallel to the cell walls but free to rotate in the plane of the cell. In this so-called “bookshelf” geometry, the layers were thus perpendicular to the cell walls, allowing two director configurations, as shown in Fig. 11, corresponding to \mathbf{P}_s being perpendicular to the cell

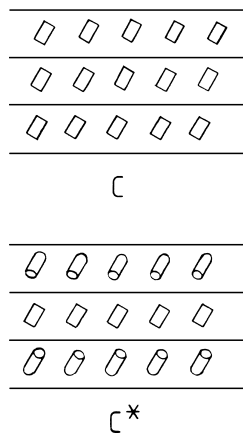


FIGURE 10 Symmetries of the smectic C and C* phases.

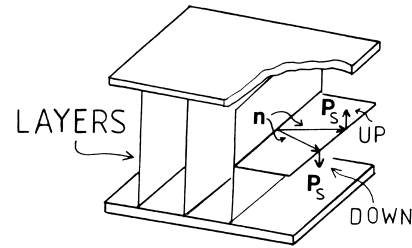


FIGURE 11 Bookshelf geometry in a smectic C* device. The director is represented by \mathbf{n} .

wall and either UP or DOWN. By application of unipolar pulses, it is possible to switch between UP and DOWN states, causing the optic axis of the liquid crystal layer to rotate in the plane of the cell through an angle that on the bookshelf model would be twice the tilt angle. Good contrast can be obtained between crossed polarizers, particularly if the angle through which the optic axis switches is close to the optimum of 45° . Since the optical effect is based on interference of ordinary and extraordinary rays, a cell gap tolerance of $\pm 0.1 \mu\text{m}$ is necessary.

Conditions of surface alignment and cell thickness can be found for which the UP and DOWN states can be bistable, and line address times of the order of $100 \mu\text{sec}$ or less at room temperature have been attained. However, it has become clear that the operation of these devices is not exactly as described by the bookshelf model and that there are several variants depending on the details of device construction. Further, the device mode closest to the original model requires a cell gap of the order of $1.5 \mu\text{m}$. Although impressive laboratory demonstrators have been reported using this technology (the best to date are 14-in.-diagonal 1120×1280 pixels with a $150\text{-}\mu\text{sec}$ line time), the very narrow, tightly tolerated cell gap required is unattractive for production. An alternative device configuration has been reported that has a thicker cell gap (up to $5 \mu\text{m}$) yet can show high contrast over a wide viewing angle and perfect bistability. This is the technology used to construct the $3\frac{1}{2}$ -in.-diagonal ferroelectric LCD shown in Fig. 9.

Color can be introduced into ferroelectric LCD (FELCD) technology by use of color microfilters aligned with the pixels. Although several suggestions for gray scale in FELCDs have been made, including combining active matrix addressing with FELCD technology, the issue of whether a practical method of introducing gray scale is possible remains to be resolved. Although significant development work is still required, FELCDs are likely to supersede the SBE by offering higher levels of multiplexing (more rows of pixels), a better viewing angle, and rapid line updating, which can be used to display smoothly the operation of a cursor or mouse (cf. the intrinsically slow

response of the SBE). If gray scale can be successfully introduced, video applications will also be accessible.

C. Smectic A and Cholesteric

The smectic A is an untilted phase in which the mass density wave is parallel to the director. The cost in free energy of buckling the layers into saddle-shaped deformations is low, with the result that it is relatively easy to construct devices that show bistability between a scattering focal conic director configuration in which the layers are buckled and a clear homeotropic configuration in which the director is perpendicular to the cell walls and the layers parallel to the walls. Transitions between these two textures have been exploited in laser-written projection displays and in both thermo-optic and electrooptic matrix displays. The various mechanisms employed are summarized in Fig. 12.

In the laser-written projection display, information is written as scattering lines on a clear homeotropic texture by using local heating with a laser. Dye may be dissolved in the liquid crystal to aid absorption of the laser energy. Lines are selectively erased by writing over them with the laser in the presence of a moderate electric field applied with ITO electrodes extending over the whole active area of the cell. This field is insufficient to cause erasure in unheated areas. Bulk erasure is achieved by application of a higher field, possibly assisted by heating of the whole cell. The cell is viewed in projection to give a large-area display useful in applications such as map overlay (for example, for command and control) and computer-aided design. The speed of writing in the vector mode is sufficient for such applications, but the raster-mode writing required for an alphagraphic display would be too slow.

The same physics have been demonstrated in a thermo-optic matrix display that uses row electrodes as heating

bars. Each row is addressed in turn by heating it into the isotropic phase and then applying appropriate voltages to the column (data) electrodes to ensure that those pixels in the row intended to be homeotropic cool in sufficient field, while those intended to be focal conic cool in zero field.

An electro-optic smectic A matrix display of significant size and complexity ($12\frac{1}{2}$ -in.-diagonal 780×420 pixels) has been brought to small-scale production. This used the electrohydrodynamic instability induced by passage of a low-frequency current (dynamic scattering) to effect the homeotropic-to-focal conic transition and a high ac field to induce the reverse transition. The display had impressive performance, but it was limited by the cost and reliability questions associated with the high voltages and dynamic scattering used in the addressing. In both of these types of matrix display, the clear and scattering textures may be used either directly to give optical contrast or via the guest-host effect (Section III) using dissolved anisotropic dye.

Cholesteric materials, as mentioned above, have a natural helical twisting of the director, with the director perpendicular to the helical axis. This gives a stratified structure with layers of constant director orientation. Although the pitch (micrometers) is 3 orders of magnitude greater than the layer repeat in smectics, a surprisingly similar focal conic texture can be formed. As mentioned in Section III.C, this scattering texture can be transformed into a clear nematiclike texture by application of an electric field (assuming positive dielectric anisotropy). Since the transition back to focal conic has to be nucleated, if the voltage is then lowered to a sustaining value, the clear texture can be maintained (at least for a frame period). This can be described as a hysteresis curve, as shown in Fig. 13. Various schemes have been devised whereby pixels may be set in either clear or scattering states with voltages greater than V_c or less than V_0 , respectively (see Fig. 13),

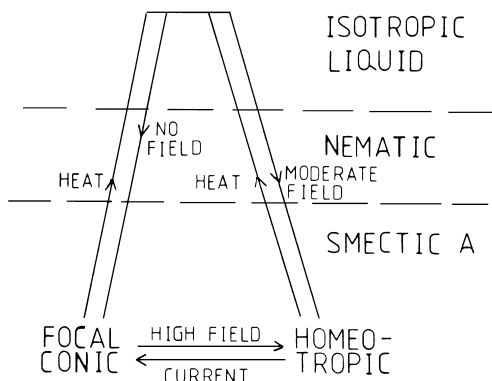


FIGURE 12 Mechanisms for inducing transitions between scattering (focal conic) and clear (homeotropic) textures in a smectic A device. The material is assumed to have positive dielectric anisotropy.

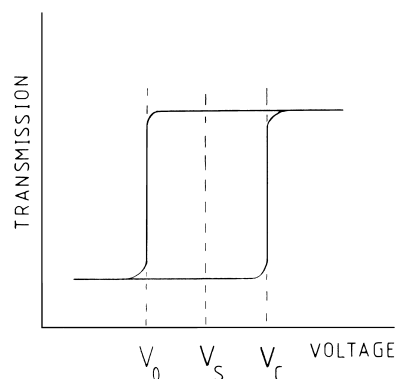


FIGURE 13 Schematic plot of hysteresis in the electrooptic response of a cholesteric phase-change device.

and then sustained in the desired state with a voltage, V_s . Although laboratory demonstrations have been fabricated, there appears to be some reluctance to produce, perhaps because of the complexity of the effects involved.

V. OTHER DEVICES

A. Thermochromic Devices

Cholesteric materials can be made that have the natural pitch, P , comparable with the wavelength of visible light. Such materials show selective reflection of light, which satisfies a Bragg-like condition between wavelength and pitch. The reflection, which is circularly polarized, is centered (for normal incidence on a thick sample) at a wavelength given by the following:

$$\lambda_0 = nP,$$

where

$$n = \frac{1}{2}(n_e + n_o),$$

n_e and n_o being the extraordinary and ordinary refractive indices of the liquid crystal, and has a width given by the following:

$$\Delta\lambda = P\Delta n,$$

where

$$\Delta n = n_e - n_o.$$

The pitch of a cholesteric may vary quite rapidly with temperature, particularly if the phase diagram of the material contains a cholesteric to smectic A transition just below the temperature range of interest since the pitch in this case diverges to infinity as the smectic A phase is approached. Thus, the color of the selective reflection can be used as a sensitive indicator of temperature. Material for this purpose is incorporated into plastic films with a black backing to absorb the unreflected polarization, or it is encapsulated to apply as a paint and used for a variety of medical, engineering, and consumer applications.

The thermochromic effect could also be used for thermal imaging, although better results are obtained with a different approach based directly on the temperature dependence of the birefringence of a liquid crystal. Unfortunately, the large investment already made in other approaches to thermal imaging militates against development of liquid crystal thermal imagers.

B. Polymer Dispersed Nematics

Dispersions of micron-sized droplets of nematic liquid crystal in a polymer matrix form the basis of a potentially

important class of electrooptic devices. The unpowered nematic/polymer film scatters light very efficiently and is milky (translucent) in appearance; the film clears to achieve a high degree of transparency where an electric field is applied across it. Electrically controllable absorption may be obtained by incorporating an anisotropic dye (Section III.A) into the nematic. The use of a polymer film as the matrix means that it is easy to make curved and large (3-m^2) devices. Although a number of display applications of this material have been demonstrated, perhaps the most promising application is as an architectural material for windows and screens which can be switched between opaque and clear.

Polymer dispersed nematic films are made by one of two distinct processes. In one, the nematic is emulsified in either an aqueous solution of a film-forming polymer (for example, poly vinyl alcohol) or an aqueous colloidal dispersion (for example, a latex). This emulsion is coated onto a conductive substrate and allowed to dry, during which time the polymer coalesces around the nematic droplets. Laminating a second conductive substrate to the dried film completes the device. Alternatively, the nematic is mixed with a precursor to the polymer to form an isotropic solution. When polymerization is initiated, typically with heat or light, nematic droplets nucleate *in situ* as the polymer chains grow.

A typical film consists of droplets with dimensions in the range of 1 to $3\text{ }\mu\text{m}$ dispersed in a film 10 to $30\text{ }\mu\text{m}$ thick. The ordinary refractive index (n_o) of the nematic approximately matches the refractive index of the polymer, but the extraordinary one does not, giving rise to scattering of light at the nematic-polymer interface when the droplets are randomly aligned. Application of a voltage sufficient to align the droplets (a saturation voltage of 100 V dc is typical) removes the refractive index mismatch that causes this scattering.

C. Fast Shutters

A shutter is by definition an unmultiplexed single-pixel device to which only two (possibly rms) drive voltages, V_{ON} and V_{OFF} , are applied, one of which may be zero if required. The twisted nematic device itself is a shutter, but there exist applications for which its natural turnoff time is too slow.

The three applications presently attracting greatest attention are linear shutter arrays for nonimpact printers, field sequential color displays, and stereoscopic displays. In the first application, the array of shutters is used to impart information onto a photosensitive print drum. In the second, a frame of information displayed on an emissive display device (CRTs and vacuum fluorescent dot-matrix panels have both been used) is divided into two (red/green)

or three (red/green/blue) successive fields according to whether limited or full color is desired, and either a single shutter (for R/G) or a stack of two shutters (for R/G/B) is placed in front of the display and switched so as to transmit only the required color during each field. Colored polarizers are used in place of the normal neutral density type. The field rate must be sufficiently rapid for persistence of vision to integrate the fields so that information written, for example, with equal intensity in both red and green fields, is seen as flicker-free yellow. In the third application, each frame on display is divided into successive left-eye and right-eye fields that are selected either by a single polarization-switching shutter covering the display screen with the viewer wearing passive polarizing eyeglasses or by the viewer wearing active eyeglasses in which the left- and right-eye shutters switch 180° out of phase in synchronization with the fields. Again, persistence of vision is employed to integrate left- and right-eye fields to give a flicker-free stereoscopic display.

Table III summarizes the available shutter technologies. The cholesteric phase-change effect has already been mentioned briefly in Sections III.C and IV.C. A facet of its behavior not mentioned there was that if the voltage applied to obtain the clear, field-induced nematic is then dropped abruptly to zero the device switches rapidly (more rapidly with shorter pitch) to a scattering texture (different to the focal conic). The dual-frequency nematic is a derivative of the basic TN device (Section II.A), in which the liquid crystal used has positive dielectric anisotropy for low-drive frequencies (10^2 Hz) and negative dielectric anisotropy for higher ones (10^4 Hz). Thus, its turnoff, instead of being reliant on viscoelastic forces alone, can be sped up by application of a high-frequency pulse. The pi-cell is a homogeneously aligned nematic cell with the pretilts arranged so that the zero field director configuration is splayed (H in Fig. 14). It is switched by application of a field into a metastable state (V in Fig. 14) in which a modest and rapid change in director configuration (to V') is sufficient to produce a change of 180° in the relative phase

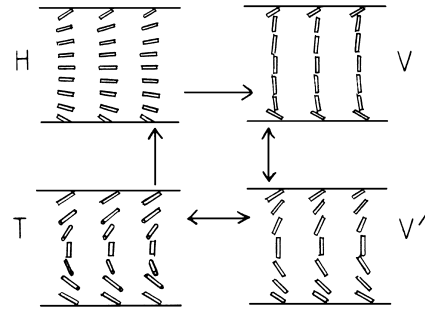


FIGURE 14 Director configurations occurring during operation of a pi-cell. The fully relaxed splayed configuration H is switched into a fully ON state V, which can then be switched rapidly to V' and back. On removal of the switching waveform, the cell relaxes first to form a 180° twisted structure T and eventually back to H.

retardation of the extraordinary ray. Thus, oriented at 45° to the polarization of plane-polarized incident light, it acts as a polarization switch. By placing two or more pi-cells in series with their alignment directions parallel, a more than proportional decrease in switch-off time is obtained. The ferroelectric device was described in Section IV.B. The electroclinic effect can occur in the smectic A phase of a material showing a smectic to smectic C* transition with decreasing temperature. In this pretransitional effect, tilt is induced in the smectic A phase by application of a unipolar voltage, the direction of tilt depending on the sign of the voltage. Although both dual-frequency and pi-cell shutters have achieved limited commercial use, the rapid development of ferroelectric devices indicates that they are likely to dominate shutter applications in the immediate future.

D. Spatial Light Modulators

Rather confusingly, the term spatial light modulator (SLM) is used for two entirely different types of devices, electrically addressed SLMs and optically addressed SLMs. The inference to be drawn from the use

TABLE III Liquid Crystal Fast Shutter Technologies

Technology	Status	Typical switch-off time at room temperature	Remarks
Cholesteric phase change	Production	3 msec	Leaky; suitable for stereoviewer spectacles only
Dual-frequency nematic	Production	2 msec	Have been used in commercial products
Pi-cell	Production	2 msec	Have been used in commercial products
Multiple pi-cell	Development	300 μ sec	Superseded by ferroelectric LCD
Ferroelectric	Production	100 μ sec	$\tau(\text{ON}) \approx \tau(\text{OFF})$
Electroclinic	Research	10 μ sec	Limited temperature range

of the term SLM is that the device is intended for use in applications such as optical processing and digital optics.

An electrically addressed SLM is essentially a dot-matrix display device made to optical standards. The salient points are that the cell must be optically uniform and that the pixels must be small in order to keep the aperture down to a convenient size. Evidently, such a device is equally valuable as the display element in a projection display. The only SLM requirement not essential in the display application is the requirement for uniform optical thickness of the cell walls, although the optical quality of materials used in dot-matrix display construction is in practice high. All of the dot-matrix technologies are in principle candidates for application as electrically addressed SLMs; most attention has focused on the active matrix and ferroelectric technologies (Sections IV.A and IV.B). Note that in this application the active matrix technology is not effectively restricted to glass substrates; both quartz and single-crystal silicon are substrates on which the active matrix can be fabricated and which may have advantages in the SLM application.

An optically addressed SLM is a shutter-type device consisting essentially of the following layers: a transparent electrode, a photoconductor, a dielectric mirror, a liquid crystal, and a transparent electrode, together with various alignment and barrier layers that are omitted from this list for simplicity. The voltage applied between the transparent electrodes is insufficient to switch the liquid crystal layer when the voltage dropped across the dark-state photoconductor is taken into account. However, in regions where the photoconductor is illuminated by a write beam, the fraction of the applied voltage that appears across the liquid crystal layer is increased, causing the liquid crystal to switch in those regions. Thus, information written on the device is impressed on a read beam reflected off the liquid crystal layer. The function of the dielectric mirror is to prevent this read beam from interacting with the photoconductor.

Both CdS and *a*-Si have been successfully used as the photoconductor; 45°-twisted nematic layers and, on an experimental basis, ferroelectric layers have been used for the liquid crystal. CCD structures and silicon vidicon microdiode arrays have been used in place of the photoconductive layer. The device is useful both when the write beam is coherent (for example, a scanned laser) and when it is incoherent (for example, a CRT). In the latter case, the SLM can be used as an incoherent-to-coherent converter. The CRT-written device has also found application as a projection display. There exists a very large potential market for optically addressed SLMs in a variety of optical processing applications and for projection displays.

At the moment, the growth of this market is limited by the high cost of these devices.

E. Optical Control Devices

All liquid crystal devices, excepting possibly the fluorescent LCDs described in Section III.D, could be called optical control devices. The purpose of this section is to emphasize that the liquid crystal devices described previously can be used in applications whose range and variety is limited only by human ingenuity. To do this, we cite briefly a few selected applications in addition to those already discussed, classifying the applications according to the liquid crystal optical effect used: scattering, absorption, polarization guiding, variable birefringence, or variable refractive index.

Electrically controllable scattering can be obtained with cholesteric and smectic A devices (Section IV.C), with polymer dispersed nematic films (Section V.B), or by use of the dynamic scattering electrohydrodynamic instability induced by current flow in nematics. It has application, for example, in antidazzle rearview mirrors for automobiles and an ingenious device for training student pilots to fly in fog. Absorption via the guest–host effect (Section III) can be applied to give electronically controllable filters and attenuators; both neutral density and colored types are possible. The birefringence of a liquid crystal is exploited via polarization guiding in the twisted nematic (Section II.A). This device makes an excellent electro-optic shutter, provided its relatively slow switching time is acceptable. Fast twisted nematic shutters (Section V.C) have been used for automatic welding goggles, in which detection of the welding arc by a photodiode initiates dimming of the goggle lenses. The birefringence can also be employed directly by using an untwisted device as an electrically controllable birefringent plate. A stack of such plates can be used as a state of polarization (SOP) converter, which will change any elliptical SOP into any other chosen SOP. Such converters are needed in coherent optical communications systems. The possibility of variable optical phase retardation offered by LCDs has also been exploited in several proposals for electrically controllable lenses and prisms. Switching the external refractive index seen at a glass–LC interface by a (polarized) ray in the glass can be used to switch between total internal reflection of the ray and its passage across the interface. In this way, optical switches needed to reconfigure optical fiber networks have been demonstrated.

F. Liquid Crystal Polymers

Molecular fragments that confer liquid crystalline properties can be incorporated into polymer chains in a variety

of ways. The two basic strategies are to incorporate such fragments into the backbone of the polymer or to hang them, usually by use of alkyl chains, as side chains off the backbone. Backbone liquid crystal polymers (LCPs) have found success as structural materials since the greater ordering can confer greater strength (for example, Kevlar). Side-chain LCPs are of greater interest here since the rationale of incorporating the side-chain moiety is to transfer some desired property (absorption, fluorescence, optical nonlinearity, ferroelectricity, etc.) into the polymeric medium. The perceived advantages relative to the corresponding monomeric systems are (1) that the polymeric nature allows new device embodiments, for example, plastic sheet or film, not possible with a fluid that must be contained in a cell; (2) that an imposed alignment of the molecular fragments can be preserved even against mechanical disturbance by storing it is the highly viscous polymer smectic or glass phases; and (3) that greater concentrations of an active guest molecule might be attainable than would be possible with a solution of monomeric components.

Promising nonlinear optical effects can be obtained by poling thin films of such materials with electric fields. Another promising application is to employ the physics described for the laser-written projection display in Section IV.C to write, erase, or update information on a smectic LCP film. The great advantage of an LCP relative to a monomer is that the information is permanently stored unless deliberately erased. Since these media are autodeveloping (that is, no postwriting processing is needed), they are competitive with conventional micrographic media (microfiche, microfilm) because they offer a shorter total cycle time, freedom from chemical processing, updatability, and erasability. Like microfiche, they store information in an instantly human-readable form that can be turned into paper copy by a conventional microfiche reader/printer. Figure 15 shows an image written in microfiche format and resolution onto an LCP film using a low-power laser.

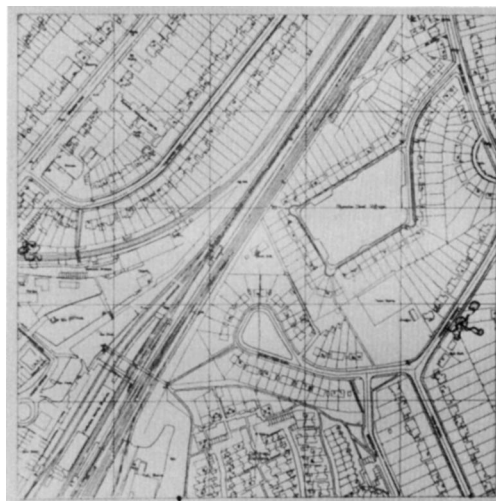


FIGURE 15 Image written in microfiche format and resolution on a liquid crystal polymer film using a low-power HeNe laser. [Courtesy of GEC Hirt Research Centre.]

SEE ALSO THE FOLLOWING ARTICLES

IMAGE PROCESSING • LIQUID CRYSTALS (PHYSICS) • THIN-FILM TRANSISTORS

BIBLIOGRAPHY

- Chigrinov, V. G. (1999). "Liquid Crystal Devices: Physics and Applications," Artech House, Norwood, MA.
- Ivashchenko, A. V. (1995). "Dichroic Dyes for Liquid Crystal Displays," CRC Press, Boca Raton, FL.
- Kramer, L., and Buka, A. (1996). "Pattern Formation in Liquid Crystals," Berlin, Springer-Verlag.
- Lueder, E. (2001). "Liquid Crystal Displays: Active and Passive Matrix Addressing Techniques," Wiley, New York.
- Sonin, A. A. (1998). "Freely Suspended Liquid Crystalline Films," Wiley, New York.
- Yeh, P. (1999). "Optics of Liquid Crystal Displays," Wiley, New York.



Microoptics

Kenichi Iga

Tokyo Institute of Technology

- I. Microoptics: Its Roles in Optoelectronics
- II. Basic Theory For Microoptics
- III. Fundamental Microoptic Elements
- IV. Microoptic Components
- V. Microoptic Systems
- VI. Characterization of Microoptic Components and Systems
- VII. Stacked Planar Optics

GLOSSARY

Distributed-index microlens Microlenses that utilize the lens effect due to refractive index distribution and this provides a flat surface.

Microlens Lens for microoptic components with dimensions small compared to classical optics; large numerical aperture (NA) and small aberration are required.

Microoptic components Optical components used in microoptics such as microlenses, prisms, filters, and gratings, for use in constructing microoptic systems.

Microoptic systems Optical systems for lightwave communications, laser disks, copy machines, lightwave sensors, etc., that utilize a concept of microoptics.

Planar microlens Distributed-index microlens with a three-dimensional index profile inside the planar substrate; planar technology is applied to its construction and two-dimensional arrays can be formed by the photolithographic process.

Stacked planar optics Microoptic configuration composed of two-dimensional optical devices such as

planar microlenses and other passive as well as active devices.

MICROOPTICS utilizes a number of tiny optical components for use with electrooptics such as lightwave communications, laser disks, and copying machines.

Since optical fiber communication began to be considered as a real communication system that provides many possibilities, a new class of optical components has become necessary. It must be compact and lightweight, have high performance, and contain various functions that are different from classical optics. Integrated optics, which utilizes a concept of a planar dielectric waveguide, is thought to be the optics of the future, but presently it is difficult to find systems in use that consist of components with the *guided-wave* configuration. On the other hand, microoptic components made of *microlenses* and other tiny optical elements are pragmatically used in real optical systems such as lightwave communications and *laser disk* systems. They make use of all concepts of classical

optics as well as of beam optics and partly even of guided-wave optics. One of the new important devices introduced in the course of research is a *distributed index* (DI) or gradient index (GI) *lens*. It uses the refraction of light rays by the index gradient coming from the nonuniform index existing inside the medium. The point is that we can make a lens with flat surfaces, which is essential for optical fiber applications since we can put fibers directly in contact with the lens, while the classical lens immediately loses its lens action when some other materials touch its surface.

Therefore, the distributed index lens plays an important role in microoptics. We can thus define microoptics as an optics that utilizes microlenses and other tiny optical elements, with the high performance and reliability required in heavy-duty electrooptic systems.

I. MICROOPTICS: ITS ROLES IN OPTOELECTRONICS

A. Comparison of Optical Components

Great progress in optical fiber communication has been made and many working systems have been installed. Three types of optical components used in optical fiber communication systems have been considered:

1. Microoptics, which consists of microlenses such as gradient-index lenses or tiny spherical lenses
2. Optical fiber circuits made from manufactured fibers
3. Integrated optics

There have been many problems in the first two schemes, such as optical alignment and manufacturing process, and integrated optics devices are still far from the usable level, as shown in Table I.

TABLE I Advantages of Stacked Planar Optics^a

Classification of encountered problem	Optical system		
	Discrete	Waveguide	Stacked
Fabrication process	Yes	No	No
Surface preparation	Yes	Yes	No
Optical alignment	Yes	Yes	No
Coupling	No	Yes	No
Integration of different materials	No	Yes	No
Single-multicompatibility	No	Yes	No
Polarization preference	No	Yes	No
Mass-production	Yes	No	No
2D array	Yes	Yes	No
Large-scale optics ←			

^a From Iga, K., Kokubun, Y., and Oikawa, M. (1984). "Fundamentals of Microoptics," Academic Press, Orlando, FL.

B. Roles of Microoptics

The role of microoptics is then believed to be not as a substitute for other components such as guided-wave optic components or fiber optic circuits, but for the purpose of fully utilizing these optical systems more effectively by cooperating with them. It is hoped that some more modern concepts will evolve to integrate microoptic devices without leaving them as old discrete optics. One of the ideas may be *stacked planar optics*, which will be discussed in detail in Section VII.

II. BASIC THEORY FOR MICROOPTICS

Light propagation is described simply by a ray that indicates a path of light energy. Actually, sometimes a very thin beam from a laser or collimated light through a pin-hole appears as a "ray" that certainly indicates the path of light energy. In a conventional optical system, which is constructed with lenses, mirrors, prisms, and so on, a ray is represented by a straight line that is refracted or reflected at a surface where the refractive index changes. On the other hand, in a distributed index medium, a ray does not follow a straight line; rather, it takes a curved trajectory as if it were affected by a force toward the higher refractive index, as shown in Fig. 1. The ray trajectory in a distributed index medium is calculated from ray equations that are second-order partial differential in nature.

If we give certain meaning to the parameter $d\tau = ds/n$, we obtain differential ray equations in Cartesian coordinates:

$$d^2x/d\tau^2 = \partial(\frac{1}{2}n^2)/\partial x, \quad (1a)$$

$$d^2y/d\tau^2 = \partial(\frac{1}{2}n^2)/\partial y, \quad (1b)$$

$$d^2z/d\tau^2 = \partial(\frac{1}{2}n^2)/\partial z, \quad (1c)$$

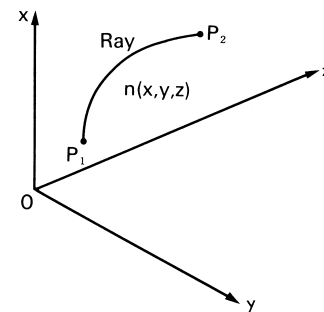


FIGURE 1 Ray propagation in a distributed-index medium. [From Iga, K., Kokubun, Y., and Oikawa, M. (1984). "Fundamentals of Microoptics," Academic Press, Orlando, FL.]

where s is the optical path length. When we treat a GI rod lens, the partial derivative of n^2 with respect to z is zero when we choose the z axis as the optical axis. Equation (1c) is integrated and we have

$$dz/d\tau = C_i. \quad (2)$$

For the ray incident at $x = x_i$, $y = y_i$, and $z_i = 0$, we have

$$C_i = n(x_i, y_i, 0) \cos \gamma_i. \quad (3)$$

Since $d\tau = ds/n$, here $\cos \gamma_i$ is a directional cosine of the ray at $z = 0$ with respect to the z axis.

Eq. (1) using Eqs. (2) and (3) gives

$$\frac{d^2x}{dz^2} = \frac{1}{2n^2(r_i) \cos^2 \gamma_i} \frac{\partial n^2(x, y)}{\partial x}, \quad (4)$$

$$\frac{d^2y}{dz^2} = \frac{1}{2n^2(r_i) \cos^2 \gamma_i} \frac{\partial n^2(x, y)}{\partial y}, \quad (5)$$

where z denotes the axial distance, r_i the incident ray position, and i the direction cosine of the incident ray.

As will be discussed later, optical components with a distributed index may play an important role in the microoptics area. There have been various ways of expressing the index distribution of such devices. One method is a power series expansion of the refractive index with respect to coordinates. The expression of the refractive index medium that we present is of the form

$$n^2(r) = n^2(0) [1 - (gr)^2 + h_4(gr)^4 + h_6(gr)^6 + \dots] \quad (6)$$

for a circularly symmetric fiber or rod, where g is a focusing constant expressing the gradient index, and h_4 and h_6 represent higher-order terms of the index distribution and are closely related to aberration. In Eq. (6) $n(0)$ expresses the index at the center axis when $r = 0$.

When we attempt to solve a diffraction problem, however, we should treat the light as a wave. By using a wavefront having a certain relation to the light ray, we can accomplish this.

The wavefront is defined as a surface constructed with a set of equioptical path-length points from a light source. On a wavefront the light phase is constant.

If the wavefront is nonspherical and the displacement of the actual wavefront from a spherical one is, we call it a wave aberration. Ray aberration is also reduced from wave aberration.

Spatial information, which is described by an electric field $f(x, y, 0)$ at $z = 0$ in Cartesian coordinates, is transformed while propagating in free space by a relation obeying the well-known Fresnel–Kirchhoff integral,

$$f(x, y, z) = \frac{j}{kz} \exp(-jkz) \iint f(x', y', 0) \times \exp \left[-\frac{jk}{2z} ((x - x')^2 + (y - y')^2) \right] dx' dy', \quad (7)$$

where z is the propagation distance and k the propagation constant defined by $2\pi/\lambda$. In polar coordinates, the expression is given in the same way as

$$f(r, \theta, z) = \frac{j}{kz} \exp(-jkz) \iint f(r', \theta', 0) \times \exp \left[-\frac{jk}{2z} (r^2 - 2rr' \cos(\theta - \theta') + r'^2) \right] r' dr' d\theta'. \quad (8)$$

A simple example is diffraction from a circular aperture, where $f(r, \theta, 0)$ is given by a constant within the region $r < A$ with A the radius of the aperture. The result is given by

$$|f(r, \theta, z)/f(0, \theta, z)|^2 = [2J_1(\rho)/\rho]^2, \quad (9)$$

where $J_1(\rho)$ is the first-order Bessel function and $\rho = ka(r/z)$. The profile is known as an Airy pattern as shown in Fig. 2.

The Fresnel–Kirchhoff (FK) integral for transformation of a light beam through free space and a GI medium is discussed. Here we shall describe how a Gaussian beam changes when propagating in free space and in a GI medium.

We shall assume a Gaussian beam at $z = 0$ given by

$$f(x', y', 0) = E_0 \exp[-\frac{1}{2}((x'^2 + y'^2)/s^2)]. \quad (10)$$

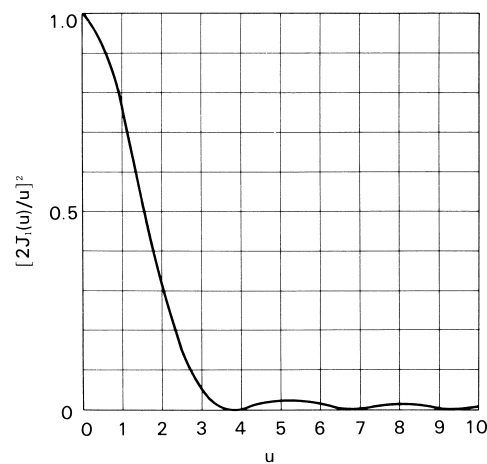


FIGURE 2 Airy function. [From Iga, K., Kokubun, Y., and Oikawa, M. (1984). "Fundamentals of Microoptics," Academic Press, Orlando, FL.]

Positional change of the beam can be calculated by Eq. (7) and is expressed as

$$f(x, y, z) = E_0 \exp(-jkz)(s/w) \times \exp\left[-\frac{1}{2}p(x^2 + y^2) + j\phi\right]. \quad (11)$$

Here the spot size w and radius R of the phase front are given by

$$w = s\sqrt{1 + (z/ks^2)^2}, \quad (12)$$

$$R = z[1 + (ks^2/z)^2].$$

Parameters P and ϕ are defined by

$$P = 1/w^2 + j(k/R), \quad (13)$$

$$\phi = \tan^{-1}(z/ks^2).$$

From Eq. (11) it can be seen that the transformed beam is still Gaussian, although there are spot-size and phase-front changes. It is clear that R expresses the phase front if we take into consideration the phase condition:

$$kz + (k/2R)r^2 = \text{const.} \quad (14)$$

With this equation, the functional dependence of the phase front $z = -(1/2R)r^2$ can be reduced. When R is positive, the phase front is convex, as seen from $z = +\infty$.

Let us next examine the parameter z/ks^2 that appears in Eqs. (12) and (13). When this parameter is rewritten as

$$z/ks^2 = (1/2\pi)(s^2/\lambda z)^{-1} \quad (15)$$

and the Fresnel number N is defined as

$$N = s^2/\lambda z, \quad (16)$$

the *Fresnel number* is a function of wavelength, distance, and spot size, and expresses normalized distance. Regions can be characterized according to N such that

$$N \ll 1 \quad (\text{Fraunhofer region}),$$

$$N > 1 \quad (\text{Fresnel region}).$$

When the point of observation is located at a point some distance from the origin ($N \ll 1$), the spot size w can be approximated from Eq. (13) as

$$w \cong z/ks. \quad (17)$$

The spreading angle of the beam is, therefore,

$$2 \Delta\theta = 2w/z = 0.64(\lambda/2s). \quad (18)$$

This is analogous to the spreading angle of a main lobe of a diffracted plane wave from a circular aperture, given by

$$2 \Delta\theta = 1.22(\lambda/D). \quad (19)$$

Figure 3 presents the waveform coefficients P_0 , P_1 , and P_2 at $z=0$, z_1 , and z_2 , respectively. If the spot sizes and curvature radii of the wavefront are given by s , w_1 ,

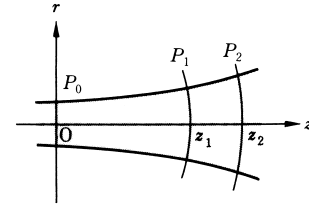


FIGURE 3 Transformation of waveform coefficients. [From Iga, K., Kokubun, Y., and Oikawa, M. (1984). "Fundamentals of Microoptics," Academic Press, Orlando, FL.]

and w_2 and ∞ , R_1 , and R_2 , respectively, the coefficients can be expressed as

$$P_0 = 1/s^2,$$

$$P_1 = 1/w_1^2 + jk/R_1, \quad (20)$$

$$P_2 = 1/w_2^2 + jk/R_2.$$

From Eqs. (10)–(12),

$$1/P_0 = 1/P_1 + jz_1/k, \quad (21)$$

$$1/P_0 = 1/P_2 + jz_2/k.$$

When P_0 is eliminated, the relationship between P_1 and P_2 is reduced to

$$P_1 = P_2/[1 + (j/k)(z_2 - z_1)P_2]. \quad (22)$$

This is a special case of the linear transform:

$$P_1 = (AP_2 + B)/(CP_2 + D) \quad (23)$$

It is very convenient to use the matrix form

$$\tilde{F} = \begin{bmatrix} A & B \\ C & D \end{bmatrix} \quad (24)$$

to calculate the transform for a system composed of many tandem components. It is then possible to obtain a total F matrix with the product of the matrices expressed as

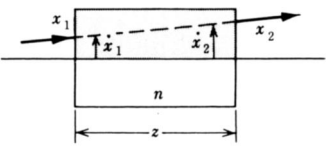
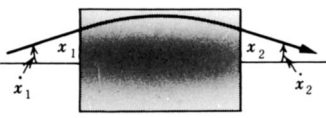
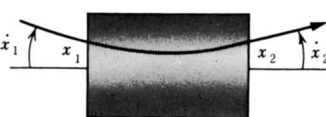
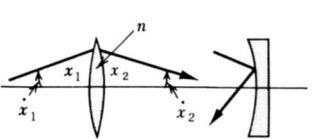
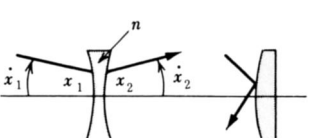
$$\tilde{F} = \tilde{F}_1 \times \tilde{F}_2 \times \tilde{F}_3 \times \cdots \quad (25)$$

Table II presents a tabulation of the waveform matrices associated with some optical components. It is not difficult to obtain these matrix forms by calculating the change of a Gaussian beam when it passes through these optical components. Kogelnik also proposed a matrix form for the same purpose, but it is somewhat different from the definition introduced here.

Figure 4 shows that ray position x_1 and ray slope \dot{x}_1 at the incident position are related to x_2 and \dot{x}_2 by the same matrix representation; that is,

$$\begin{bmatrix} jk\dot{x}_1 \\ x_1 \end{bmatrix} = \begin{bmatrix} A & B \\ C & D \end{bmatrix} \begin{bmatrix} jk\dot{x}_2 \\ x_2 \end{bmatrix}. \quad (26)$$

TABLE II Waveform Matrices for Various Optical Systems^a

Optical system	Waveform matrix, F
 <p>FREE SPACE</p>	$\begin{bmatrix} 1 & 0 \\ j\frac{z}{k} & 1 \end{bmatrix}, k = k_0 n$
 <p>POSITIVE DI LENS</p>	$\begin{bmatrix} \cos gz & jkg \sin gz \\ j\frac{1}{kg} \sin gz & \cos gz \end{bmatrix}, k = k_0 n(0)$
 <p>NEGATIVE DI LENS</p>	$\begin{bmatrix} \cosh gz & \sim jkg \sinh gz \\ j\frac{1}{kg} \sinh gz & \cosh gz \end{bmatrix}, k = k_0 n(0)$
 <p>CONVEX LENS CONCAVE MIRROR</p>	$\begin{bmatrix} 1 & \frac{jk}{f} \\ 0 & 1 \end{bmatrix}, k = k_0 n$
 <p>CONCAVE LENS CONVEX MIRROR</p>	$\begin{bmatrix} 1 & -\frac{jk}{f} \\ 0 & 1 \end{bmatrix}, k = k_0 n$

^a From Iga, K., Kokubun, Y., and Oikawa, M. (1984). "Fundamentals of Microoptics," Academic Press, Orlando, FL.

The propagation constant k is included in Eq. (26) to make it possible to treat a tandem connection of optical components having different refractive indices.

A. Guided-Wave Theory

In order to obtain a basic concept of guided-wave components including optical fibers and planar dielectric waveguides, we summarize the treatment of a simple planar waveguide consisting of three layers, as shown in Fig. 5. The tangential field components of TE modes are E_y and H_z . Since these two components are directly related to the boundary condition, it is convenient to deal with the E_y component instead of the H_z component. From Maxwell's equation, E_y must satisfy

$$d^2 E_y / dx^2 + (k_0^2 n^2 - \beta^2) E_y = 0. \quad (27)$$

Special solutions of Eq. (27) in the core are cosine and sine functions. In the cladding, the solutions are classified into two types—namely, the evanescent (exponentially decaying) solution for $n_2 k < \beta < n_1 k$ and the sinusoidal oscillating solution for $\beta < n_2 k$. The former is called the guided mode. Some amount of optical power of a guided mode is confined in the core, and the remainder permeates from the cladding. The latter is called a set of radiation modes and the power is not confined in the core. The group of all guided and radiation modes constitutes a complete orthogonal set and any field can be expanded in terms of these guided and radiation modes.

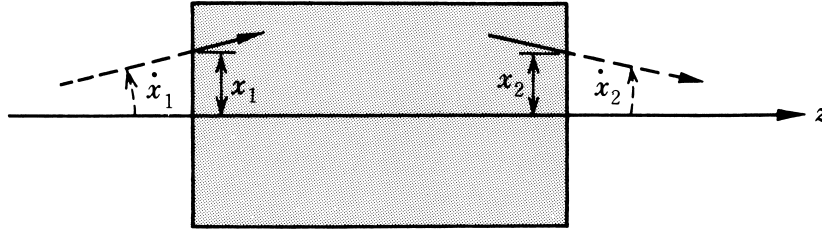


FIGURE 4 Relationship of ray positions and ray slopes. [From Iga, K., Kokubun, Y., and Oikawa, M. (1984). "Fundamentals of Microoptics," Academic Press, Orlando, FL.]

The solutions of eigenvalue equations can be normalized by introducing new parameters b and V defined by

$$b = (\beta/k_0 - n_2)/(n_1 - n_2) \quad (28)$$

where

$$V = k_0 a n_1 \sqrt{2 \Delta} \quad (29)$$

and giving a simple closed form,

$$V = (\pi/2)/\sqrt{(1-b)}[(2/\pi) \tan^{-1}(\sqrt{b/(1-b)} + N)]. \quad (30)$$

From this expression we obtain a dispersion curve that relates V and b , as shown in Fig. 6. When the waveguide parameters n_1 , n_2 , and a and wavelength λ are given, the propagation constant β_N of any mode can be obtained from (Eq. 30). The mode number is labeled in the order of increasing N , and TE₀ is the fundamental mode. This mode number N corresponds to the number of nodes in the field distribution.

When the propagation constant of one guided mode reaches $n_2 k$ ($b \rightarrow 0$), the mode is cut off, and the V value is then called the cutoff V value. By putting $\gamma = 0$ and $\kappa a = V$, the cutoff V value of TE modes is easily obtained from Eqs. (28–30) as

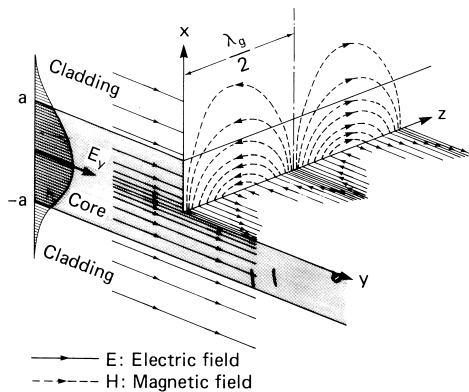


FIGURE 5 Planar waveguide. [From Iga, K., and Kokubun, Y. (1986). "Optical Fibers," Ohm, Tokyo.]

$$V = (\pi/2)N \quad (N = 0, 1, 2, \dots). \quad (31)$$

The cutoff V value of TE₁ gives a single mode condition, because, when V is smaller than $\pi/2$, only the TE mode can propagate. The single mode condition is important for designing single mode waveguides and single mode optical fibers with an arbitrary refractive index profile.

From Eq. (30) we can obtain the group velocity v_g , which is the velocity of a light pulse through the waveguide.

III. FUNDAMENTAL MICROOPTIC ELEMENTS

A. Microlens

In microoptics, several types of microlenses have been developed. A spherical microlens is used mostly to gather light from a laser diode by taking advantage of the high numerical aperture and the small Fresnel reflection of light from its spherical surface. Moreover, the precise construction of the sphere is very easy but its mount to a definite position must be considered specially. Spherical aberration is of course large.

A distributed-index GI rod lens with a nearly parabolic radial index gradient related to light focusing fiber was

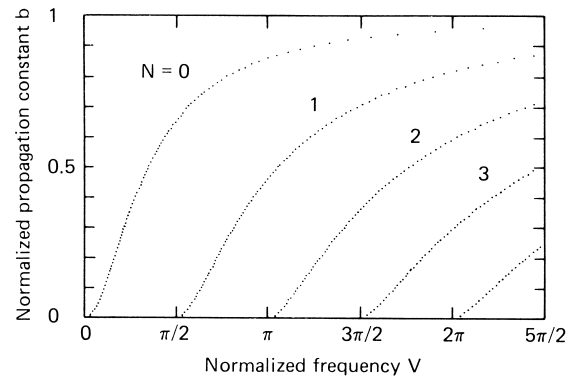


FIGURE 6 Dispersion curve for a TE mode of a planar waveguide.

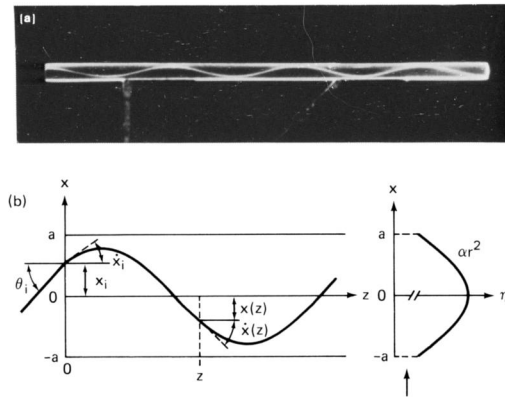


FIGURE 7 Sinusoidal light ray trajectory in a GI rod. [From Iga, K., Kokubun, Y., and Oikawa, M. (1984). "Fundamentals of Microoptics," Academic Press, Orlando, FL.]

developed in 1968. The ray trace in the GI rod lens is expressed by

$$x = x_i \cos(gz) + (\dot{x}_i/g) \sin(gz), \quad (32)$$

where x_i and \dot{x}_i are the incident ray position and slope.

It is readily known that the ray is transmitted with a sinusoidal trace, as shown in Fig. 7, and the periodicity pitch of the trace is

$$L_p = 2\pi/g. \quad (33)$$

If we cut the rod into lengths $L_p/4$, the lens acts as a single piece of positive (focusing) lens. On the other hand, if the rod length is $\frac{3}{4}L_p$ an elect image can be formed by a single piece of lens. In addition, it has the merit of having a flat surface, and other optical elements, such as dielectric multilayer mirrors, gratings, and optical fibers, can be cemented directly without any space between them. Another feature is its ability to constitute the conjugate image device (i.e., real images with unity magnification can be formed by a single piece of rod microlens). This is illustrated in Fig. 8.

A planar microlens was invented for the purpose of constructing two-dimensional lightwave components. A huge number of microlenses of 0.1–2 mm in diameter can be arranged two-dimensionally in one substrate, and their position is determined by the precision of the photomasks. The applications to lightwave components and multi-image forming devices are considered.

The importance of a microlens for focusing the light from a diode laser onto a laser disk is increasing in potential. Some new types of microlenses are being developed. One of them is a mold lens with aspheric surfaces that can be permitted for diffraction-limited focusing.

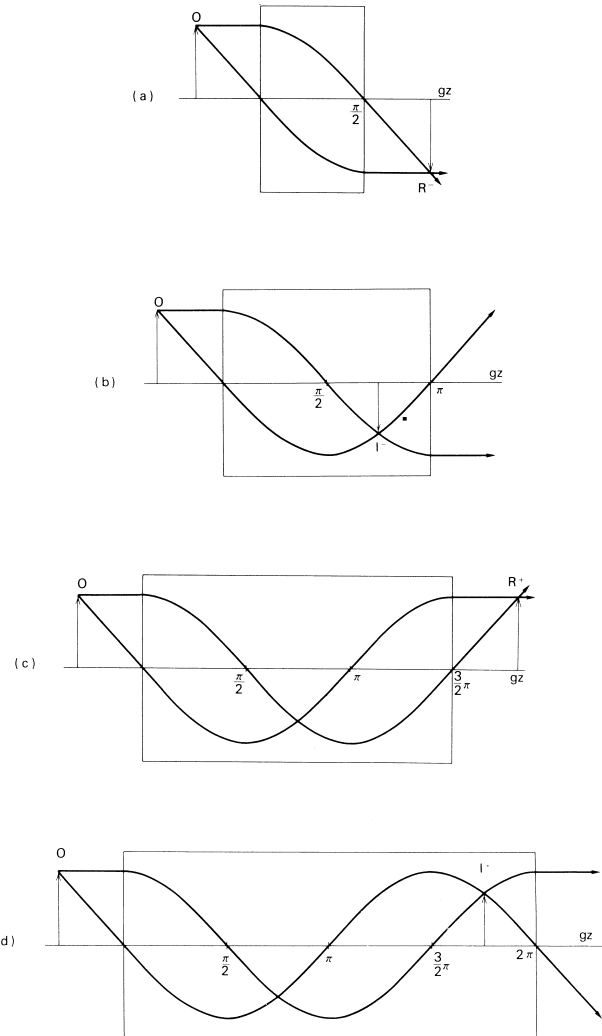


FIGURE 8 Imaging configurations with various length GI rods. [From Iga, K., Kokubun, Y., and Oikawa, M. (1984). "Fundamentals of Microoptics," Academic Press, Orlando, FL.]

B. Grating

A grating normally used in monochromators is used for multiplexing or demultiplexing different wavelengths. The grating fundamentals are described in standard optics textbooks. The band elimination property is sharp, and reliable components can be constructed. One slight problem is that the angle of the exit ray is changed if the wavelength varies. This is serious in wavelength demultiplexers because the wavelength of currently used semiconductor lasers varies considerably with temperature.

C. Multilayer Mirror and Filter

The dielectric-coated multilayer mirror or filter does not have the problem that grating filters have. The basic theory

of multilayer mirrors and filters can be found in classical optics references. The reflectivity and transmittance can be designed. The change in optical properties with age must be especially examined in conventional applications (e.g., resistivity versus moisture and temperature change).

D. Aperture Stop and Spatial Frequency Filter

An aperture stop is used to eliminate light from an unwanted direction. In special cases various types of spatial filters are inserted in order to cut out unwanted modes.

E. Optical Fiber

A fiber component consists of manufactured optical fibers. A branch is obtained by polishing fibers to bare the core region for coupling. A directional coupler, polarizer, Faraday rotator, and so on can be constituted only of optical fibers. The merit of the fiber component is that the mode of the component is identical to that of the fibers used.

IV. MICROOPTIC COMPONENTS

A. Focuser

A light-focusing component (focuser) is the simplest and still-important component. We first discuss a simple focuser composed of a single lens. The diffraction limit of the focused spot D_s is given by

$$D_s = 1.22\lambda/\text{NA} \cong 1.22f\lambda/a, \quad (34)$$

where λ is the wavelength, $2a$ the aperture stop diameter, f the focal length of the employed lens, and NA the numerical aperture of the lens. This formula can be applied to a GI lens.

The focuser for the laser disk system must have the smallest aberration since the lens is used in the diffraction limit. Another important point is the working distance of the lens. In the laser disk application some space must be considered between the objective lens and the focal point because we have to have a clearance of about 2 mm and must also consider the disk thickness. Therefore, we need a lens with a large diameter (~ 5 mm) as well as a large NA (~ 0.5). First, the microscope objective was considered and several types of microlenses have been developed, such as the mold plastic aspheric lens and the GI lens with a spherical surface.

A microlens is employed for the purpose of focusing light from a laser into an optical fiber. A spherical lens is a simple one with an NA large enough to gather light from a semiconductor laser emitting light with 0.5–0.6 NA. In the

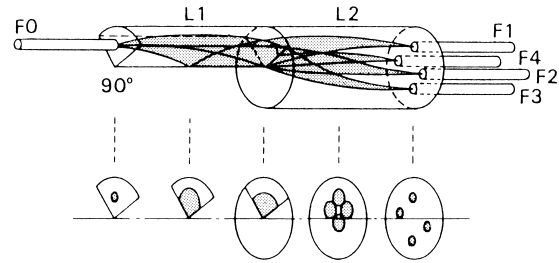


FIGURE 9 Optical branch made of GI rod lenses. [From Kobayashi, K., Ishikawa, R., Minemura, K., and Sugimoto, S. (1979). *Fibers Integr. Opt.* 2, 1.]

single mode fiber application the combination of spherical and GI lenses is considered.

B. Branch

High-grade optical communication systems and electrooptic systems need a component that serves to divide light from one port to two or more ports. This is a branching component not necessarily used for separating power equally into two branches. A simple branch consists of a lens and prisms. In Fig. 9 we show a branch made of a manufactured distribution index lens. A microoptic branch is used inherently both for single mode and multimode fibers and no polarization preference exists. Another possibility is to utilize a planar waveguide as shown in Fig. 10. We have to design separately for single or multimode fibers.

C. Coupler

A power combiner, or simply a coupler, is a component for combining light from many ports into one port. In general, we can use a branch as a combiner if it is illuminated from the rear. A directional coupler is a component consisting of a branch and coupler, as shown in Fig. 11. The light from port 1 is divided into ports 2 and 3, and the light from port 3 exits from ports 1 and 4. A component consisting of a

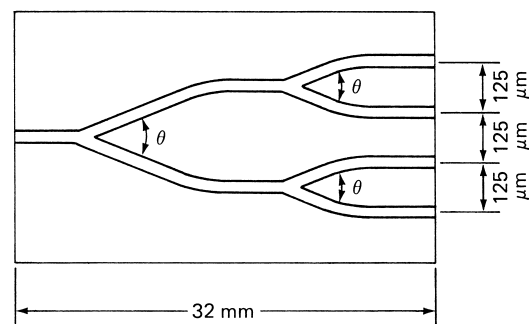


FIGURE 10 Waveguide branch. [From Okuda, E. Tanaka, I., and Yamasaki, T. (1984). *Appl. Opt.* 23, 1745.]

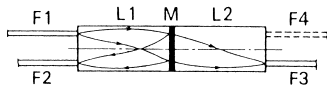


FIGURE 11 Directional coupler. [From Uchida, T., and Kobayashi, K. (1982). *Jpn. Annu. Rev. Electron. Comput. Telecommun. Opt. Devices Fibers* 4, 179.]

GI lens and a half-mirror is shown in Fig. 11. There is one component made of coupled fibers in which two fibers are placed so that the light in the two fibers can be coupled with each other (i.e., light propagating in one waveguide couples into the other while propagating along the guide).

A star coupler or optical mixer branches m ports into n ports, which serves to send light to many customers as in data highway or local area networks (LAN). Figure 12 shows a mixer made of manufactured fibers and Fig. 13 utilizes a planar waveguide configuration.

D. Wavelength Multiplexer/Demultiplexer

A wavelength *multiplexer* (MX)/*demultiplexer* (DMX) is a device that combines/separates light of different wavelengths at the transmitter receiver, which is needed inevitably for a communication system using many wavelengths at the same time. There exists a device consisting of a multilayer filter and GI lenses, as shown in Fig. 14, which is good for several wavelengths and one with a grating and lenses as in Fig. 15. The grating DMX is available for many wavelengths but the remaining problem is that the beam direction changes when the wavelength varies (e.g., as a result of the change in source temperature). Therefore in this case the wavelength of the utilized light must be stabilized.

E. Optical Isolator

An optical isolator is used in sophisticated lightwave systems where one cannot allow light reflection that might perturb the laser oscillator. Of course, to maintain system performance one wishes to introduce an optical isolator, but the cost sometimes prevents this. The principle of the isolator is illustrated in Fig. 16. The Faraday rotator rotates the polarization of incident light by 45° . Then the

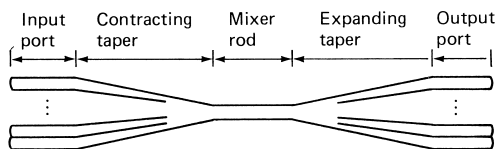


FIGURE 12 Optical mixer made of fabricated fibers. [From Ohshima, S., Ito, T., Donuma, K., and Fujii, Y. (1984). *Electron. Lett.* 20, 976.]

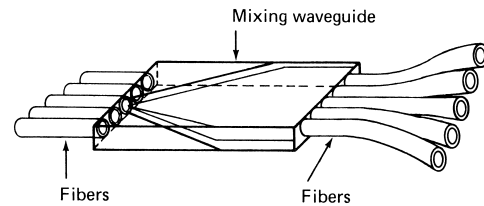


FIGURE 13 Waveguide star coupler. [From Minowa, J., Tokura, N., and Nosu, K. (1985). *IEEE J. Lightwave Technology* LT-3, 3, 438.]

reflected light can be cut off by the polarizer at the input end while the transmitted light passes through the analyzer at the exit end. Lead glass is used for the short-wavelength region and YIG ($\text{Y}_2\text{Fe}_5\text{O}_{12}$) for the long-wavelength region. A device with 1 dB of insertion loss and 30 dB of isolation is developed for wavelength multiplexing communications.

F. Functional Components

Functional components such as a switch and a light modulator are important for electrooptic systems. Several types of optical switches have been considered:

1. Mechanical switch
2. Electrooptic switch
3. Magneto optic switch

A switch as shown in Fig. 17, using the same idea as that of the isolator, will be introduced to switch laser sources in undersea communications for the purpose of maintaining the system when one of the laser devices fails.

A beam deflector is important in the field of laser printers. A rotating mirror is used now, but some kind of electro-optically controlled device is required to simplify the system.

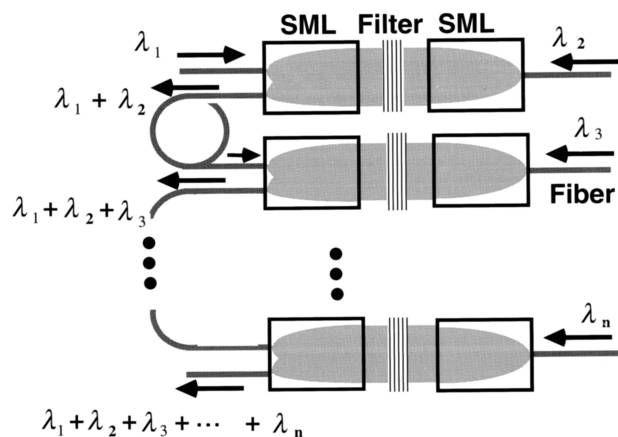


FIGURE 14 DMX using multilayer filter.

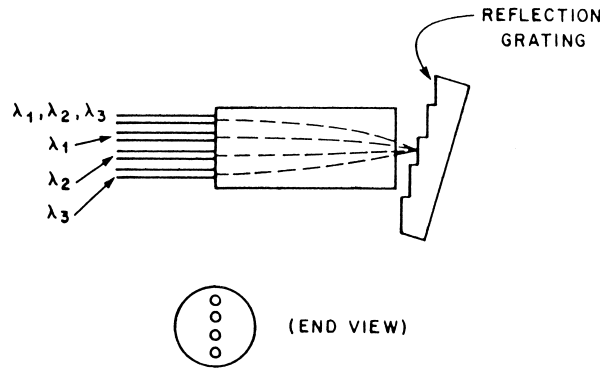


FIGURE 15 DMX using a grating. [From Tomlinson, W. J. (1980). *Appl. Opt.* **19**, 1117.]

G. Imager

In this section we deal with some imaging components (imagers), especially those consisting of distribution-index (GI) lenses, including a simple imager, conjugate imager, and multiple imager. In the early history of light-wave transmission, a GI medium was considered a promising device among continuously focusing light guides having geometries of fibers (e.g., by Kapany) or slabs (e.g., by Suematsu).

Some types of lenslike and square-law media have been proposed and studied for light beam waveguides for laser communication systems. Various types of gas lenses and focusing glass fibers provide examples of such media, whose dielectric constants are gradually graded according to a square law with regard to the distance from the center axis. As for the optical characteristics of such media, it is known that a Hermite–Gaussian light beam is guided along the axis of the lenslike medium and, moreover, images are transformed according to a definite transform law that not only maintains the information concerning their intensity distribution but also that concerning their phase relation. This is thought to be one of the significant characteristics of gas lenses and focusing glass fibers, a characteristic different from that of a step-index glass fiber.

Various authors have reported on imaging properties of GI media. The imaging property of a gas lens was investi-

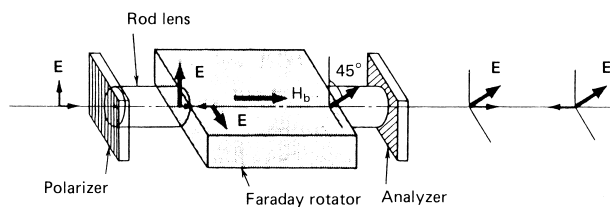


FIGURE 16 Optical isolator. [From Suematsu, Y., and Iga, K. (1976). "Introduction to Optical Fiber Communications," Ohm, Tokyo (Wiley, New York, 1980).]

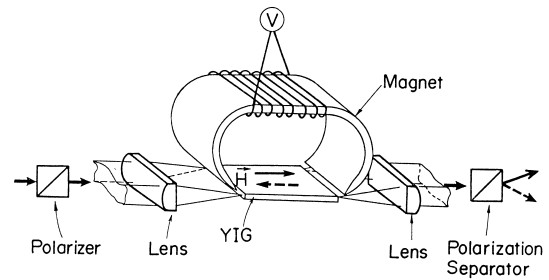


FIGURE 17 Optical switch. Arrow indicates a vector quantity. [From Shirasaki, M. (1984). *Jpn. Annu. Rev. Electron. Comput. Telecommun. Opt. Devices Fibers* **11**, 152.]

gated, in which the lens formula and optical transfer function were obtained on the basis of geometrical optics, and some imaging experiments were made using a flow-type gas lens. In a paper describing the optical characteristics of a light-focusing fiber guide (SELFOC), Uchida *et al.* mentioned the experimental imaging property and measured the resolving power of the SELFOC lens. In each of these papers, imaging properties of GI media are interpreted in terms of geometrical optics.

When a GI medium is applied to coherent optics such as in interferometry and holography, however, it is important that a two-dimensional system theory based on wave optics be introduced into the treatment of transforms by an optical system with a GI medium. In this article we introduce such a theory, which applies an integral transform associated with a GI medium into the system theory of optics. This will enable us to learn not only about the imaging condition but also about some types of transform representations.

We express the index profile of a GI medium by Eq. (6). If we express the transverse field component by $\exp(j\omega t - j\beta z)$, the function ψ for the index profile is given approximately by the scalar wave equation,

$$\frac{1}{r} \frac{d}{dr} \left(r \frac{d\psi}{dr} \right) + \frac{d^2\psi}{d\theta^2} + k_0^2 n^2(r) \psi = \beta^2 \psi, \quad (35)$$

where $k_0 = 2\pi/\lambda$.

The normal modes associated with a square-law medium, using the first two terms of Eq. (35), are known to be Hermite–Gaussian functions, as shown in Fig. 18.

The characteristic spot size w_0 of the fundamental mode is given by

$$w_0 = a/\sqrt{V}, \quad (36)$$

where the normalized frequency V is written as

$$V = kn(0)a\sqrt{2\Delta}, \quad (37)$$

with $\Delta = [n(0) - n(a)]/n(0)$. In the usual DI lenses, V is larger than 3000, since we have $\lambda = 0.5 \mu\text{m}$, $n(0) = 1.5$,

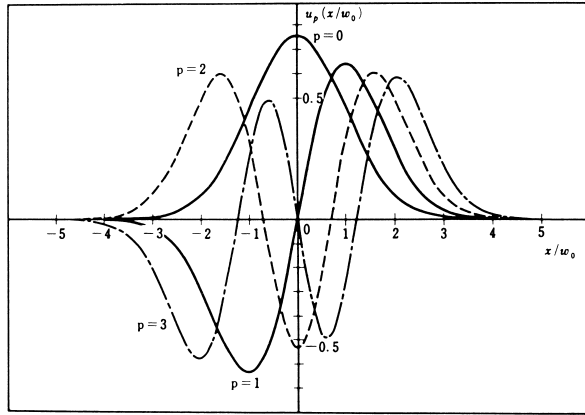


FIGURE 18 Hermite-Gaussian functions.

$a = 0.5$ mm, and $\Delta = 5\%$. The characteristic spot size w_0 is therefore smaller than the core radius a by a factor of 50–100. The ratio $w_0/a = \sqrt{V}$ is a measure to indicate whether we can use this gradient medium as an imaging lens, because a sinusoidal ray trace is distorted if w_0/a is not small. The propagation constant associated with the index profile given by Eq. (35), including higher-order terms, is obtained by a perturbation method and is expressed in terms of a series expansion in powers of $g/k(0)$, where $k(0) = k_0 n(0)$, as

$$\begin{aligned} \beta_{lm}/k_0 = & 1 - (2l + m + 1)(g/k(0)) \\ & + \frac{1}{2} \left\{ h_4 \left[\frac{3}{2}(2l + m + 1)^2 \right. \right. \\ & \left. \left. + \frac{1}{2}(1 - m^2) \right] - (2l + m + 1)^2 \right\} \\ & \times (g/k(0))^2 + O[(g/k(0))^3] + \dots \quad (38) \end{aligned}$$

We should note that $g/k(0)$ is of the order of 10^{-4} to 10^{-3} . If $m = 0$, Eq. (38) is associated with the radially symmetric mode that corresponds to meridional rays.

If we calculate the group velocity v_g by differentiating $n(0)$ with respect to ω , we see that the minimum dispersion condition is $h_4 = \frac{2}{3}$, the same result obtained from the WKB method.

V. MICROOPTIC SYSTEMS

A. Lightwave Communications

There is a wide variety of microoptic components used in lightwave communication systems. The simplest system consists of a light source such as a laser diode, an optical fiber, and a detector at the receiving end, as shown in Fig. 19. The component employed is said to focus the laser light into the fiber. The wavelength multiplexing (WDM) system needs a more sophisticated combination of devices. At the output end, the MX combines multiple wave-

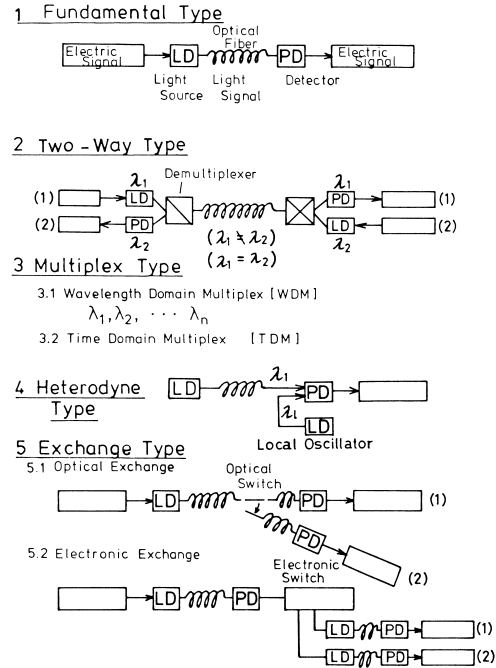


FIGURE 19 Some lightwave communication systems. [From Suematsu, Y. (1983). *Proc. IEEE* 71, 692–721.]

lengths into a single piece of fiber, and on the contrary, the DMX is utilized for the purpose of dividing different signals on different wavelengths. The overall bandwidth can be 10 THz. The important points have been introduced in the previous section, and we have to pay attention to the near-end reflection, which affects the perturbation of the laser oscillator to obtain stable operation. Many long-haul transmission systems have been installed in many countries. Transatlantic and transpacific underseas cables will be in use for international communications having a very wide bandwidth of about several thousand voice channels. The LAN will be a most popular lightwave communication. A video transmission system is an attractive medium for education and commercials. At the Tokyo Institute of Technology television classrooms were introduced to connect two campuses separated by 27 km by eight-piece single mode fibers. Every minute 400-Mbit/sec signals are transmitted back and forth. The quantity of the microoptic components employed will increase more and more as higher speeds and more complex designs are considered.

B. Laser Disks

One of the most popular systems is a laser disk for audio known as a compact disk (CD), in which PCM signals are recorded on the rear side of a plastic transparent disk and laser light reads them as shown in Fig. 20. A digital video

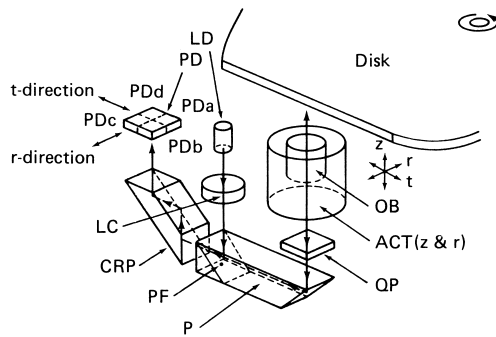


FIGURE 20 Compact audio disk system. [From Musha, T., and Morokuma, T. (1984). *Jpn. Annu. Rev. Electron. Comput. Telecommun. Opt. Devices Fibers* **11**, 108.]

disk and optical disk file for computer memory will be the successor to the CD. The optics used there consist of a combination of microoptic components. The light from a semiconductor laser is collimated by a collimating lens with large NA, passes through a beam splitter, and is focused on the rear surface of the disk by a focusing objective lens. The light reflected from the disk changes its radiation pattern as a result of the depth of the pits recorded on the disk, and the variation of light intensity is detected by a matrix light detector. The most important element is the focusing lens, since it must be designed with a diffraction limit while remaining lightweight because motional feedback is employed to compass the disk deformation.

C. Copiers

A lens array consisting of distributed index lenses is introduced in a desktop copying machine in which conjugate images (erect images with unit magnification) of a document can be transferred to the image plane as illustrated in Fig. 21. This configuration is effective for eliminating a

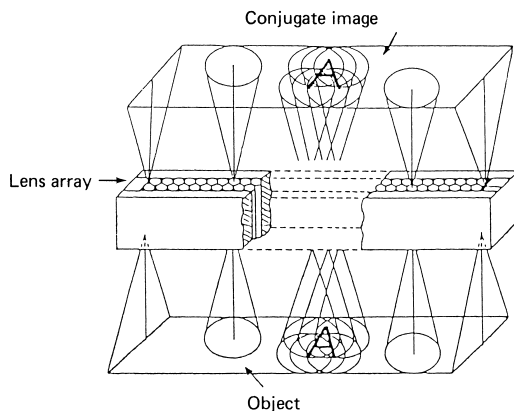


FIGURE 21 Conjugate image made by DI lens array for a copying machine. [From Kitano, I. (1983). *Jpn. Annu. Rev. Electron. Comput. Telecommun. Opt. Devices Fibers* **5**, 151.]

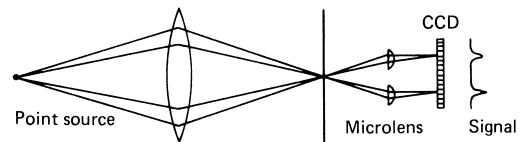


FIGURE 22 Autofocus microoptic system. [From the Minolta catalogue, Nagoya, Japan.]

long focal length lens and the total volume of the copying machine can be drastically reduced. The number of lenses produced is more than 10^9 .

D. Autofocusers

An autofocusing element to provide easy focusing for a steel camera is now becoming popular. Multiple images are formed by a lens array, and a CCD detector generates an error-defocus signal until the main lens is automatically moved to the correct position, as in Fig. 22. Some types of microlens arrays have now been developed. A planar microlens array will be one of such arrays since it has the merit of being easy to mask by a photolithographic technique.

E. Fiber Sensors

A fiber gyro and other lightwave sensing systems are considered in various measurement demands. The microoptic element employed is some type of interferometer consisting of a beam splitter and half-mirrors, as shown in Fig. 23. Single mode components that match the mode of the single mode employed and polarization-maintaining fiber are necessary. An optical circuit made of manufactured fiber is an interesting method.

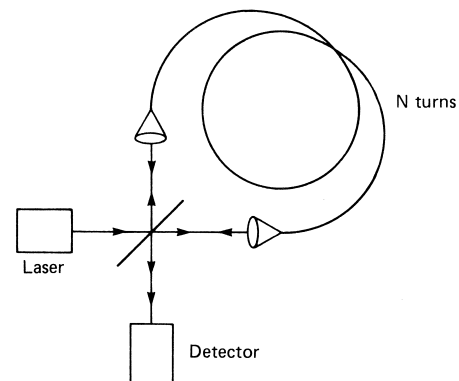


FIGURE 23 Fiber optic gyro. [From Ezekiel, S., and Arditty, H. J. (1982). "Fiber-Optic Rotation Sensors," Springer-Verlag, Berlin and New York.]

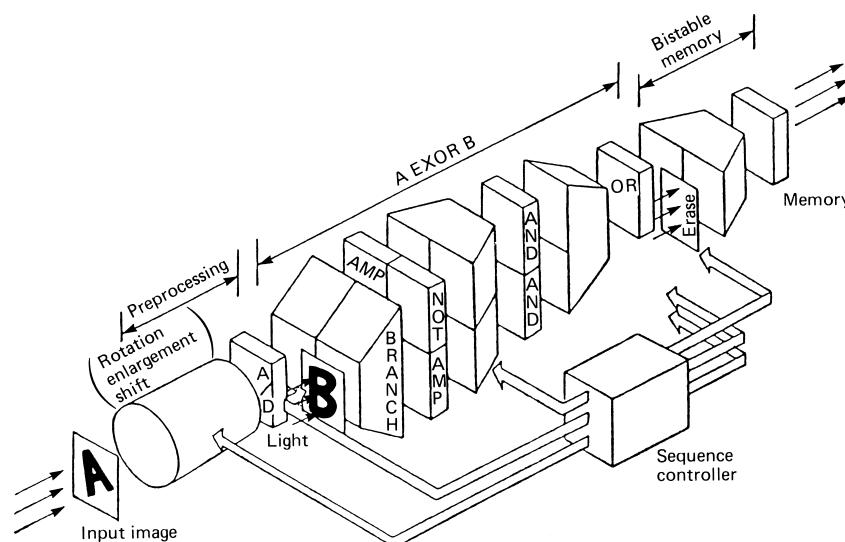


FIGURE 24 Idea of an optical computer. [From Seko, J. (1984). *Oyo Butsuri* **53**, 409.]

F. Optical Computers

A future technique may be an optical parallel processor such as a TSE computer; this idea is illustrated in Fig. 24. Some functional devices such as optical AND and OR elements based on semiconductor materials must be developed to obtain this sophisticated system. The two-dimensional configuration made of planar microlenses and a surface-emitting laser array will be very helpful.

VI. CHARACTERIZATION OF MICROOPTIC COMPONENTS AND SYSTEMS

A. Measurement of Refractive Index Distribution

1. Longitudinal Interference Method

The sample (fiber or preform rod sample) is cut into a thin round slice and its surfaces are polished to be optically flat. This thin sample is examined under an interference microscope, as shown in Table III. The refractive index profile is calculated from the fringe shift. An automatic measuring system has been developed by incorporating an image processing apparatus, such as a vidicon camera, and a computer. The spatial resolution limit is about $0.7\ \mu\text{m}$. If the sliced sample is too thick, the incident ray is refracted through the sample and causes an error. Therefore, the sample must usually be polished to a thickness of less than $100\ \mu\text{m}$. This takes much time, which prevents this method from being introduced into the fiber manufacturing process as a testing system. Accuracy of index is limited to about 0.0005 because of the roughness of the polished surfaces.

2. Transverse Interference Method

The sample is immersed in index-matching oil and is observed in its transverse direction by using an interference microscope (Table III). The index profile is calculated from the fringe shift. Before the author began this study, analysis based on the straight-ray trajectory had always been used to calculate the index profile. However, it is now known that accuracy can be increased by using an analysis that includes ray refraction. There also exists another method that uses the ray refraction angle to calculate the index profile, but the accuracy is not very good.


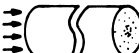
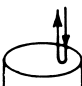
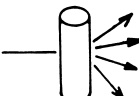



3. Transverse Differential Interference Method

The transverse differential method is an interference method modified to apply to thick samples, such as focusing rod lenses and optical fiber preform rods. Instead of a transverse interference pattern, a transverse differential interference pattern, differentiated with respect to the transverse distance, is used to calculate the index profile.

4. Focusing Method

When an optical fiber with an axillary symmetric index distribution is illuminated in its transverse direction by an incoherent light source, the fiber acts as a lens so that the incident light is focused on a plane placed behind the fiber as seen from Table III. If the light intensity distribution is uniform with respect to the incident plane, the index profile can be calculated from the focused light intensity distribution. This method can be applied to preform rods as well as to fibers, and is one of the promising methods,

TABLE III Various Measuring Methods of Index Profile^a

		Measurement time	Accuracy	Sample preparation	Correction of elliptical deformation
Longitudinal interference method		Long	Good	Mirror polish	Easy
Near-field pattern method		Short	Fairly good	Cleave	Possible
Reflection method		Medium	Fairly good	Cleave	Possible
Scattering pattern method		Short	Fairly good	Not necessary	Difficult
Transverse interference method		Short	Good	Not necessary	Possible
Focusing method		Short	Good	Not necessary	Not practical
Spatial filtering method		Short	Good	Not necessary	Possible

^a From Iga, K., Kokubun, Y., and Oikawa, M. (1984). "Fundamentals of Microoptics," Academic Press, Orlando, FL.

along with the transverse and transverse differential interference methods. This method can be applied to axially nonsymmetric preforms.

B. Measurement of Composition

1. X-Ray Microanalyzer (XMA) Method

The XMA method measures the dopant concentration profile, which is related to the index profile, by means of an XMA (X-ray microanalyzer). The contribution of dopants such as P_2O_5 , GeO_2 , and B_2O_5 to the refractive index can be obtained separately, but accuracy is not good because of the low signal-to-noise ratio.

2. Scanning Electron-Beam Microscope (Etching) Method

When the end surface of a distributed index sample is chemically etched, the etching speed depends on the dopant concentration. Therefore, the unevenness can

be observed by a scanning electron-beam microscope (SEM).

C. Reflection Pattern

The reflection coefficient of a dielectric material is related to the refractive index at the incident surface. The refractive index profile of an optical fiber can be measured by utilizing this principle. A laser light beam with a small spot size is focused into the end surface of a sample, and the reflection coefficient is measured by comparing the incident and reflected light intensity, as shown in Table III. The refractive index profile is obtained from the reflection coefficient profile by shifting the reference point. Accuracy is strongly affected by the flatness of the end surface. A fractured end surface gives better results than does a polished end surface. For borosilicate fibers, the result changes rapidly with time because of atmospheric exposure of the dopant. Spatial resolution is usually limited to about 1–2 μm by the spot size of the incident beam. This

effect of the finite beam spot size can be corrected by numerical calculation. A 0.3- μm spatial resolution and 5% total accuracy of the refractive index has been obtained by this correction.

D. Scattering Pattern

The scattering pattern is classified into both a forward scattering pattern method and a backward scattering pattern method. In the case of the forward scattering pattern method, the sample is immersed in an index-matching oil and the forward scattering pattern is observed. The refractive index profile is calculated from the scattering pattern by using a computer. The error in this method is determined from the product of the core radius a and the index difference Δn , and it increases with an increase of $a \Delta n$; as a numerical example, when $a \Delta n = 0.04 \text{ mm}$, the error is 5%. Therefore, this method is applicable only to single mode fibers. Since this method requires many sampling points (500–1000), it is necessary to collect the data automatically.

On the other hand, the index profile can also be obtained from the backward scattering pattern. This method does not require index-matching oil and is applicable to thick samples such as preform rods. However, since the backward scattering pattern is tainted by externally reflected light, it is not suitable for precise measurements. Furthermore, the accuracy of this method is very sensitive to the elliptical deformation of the core cross section.

E. Near-Field Pattern

When all the guided modes of a multimode waveguide are excited uniformly by using an incoherent light source, the near-field pattern of output optical power is similar to the refractive index profile. Since it is difficult to satisfy the incident condition strictly and the near-field pattern is affected by leaky modes and the absorption loss difference of guided modes, this method cannot provide accurate measurements. Although several improvements, such as a correction factor for leaky modes, a refracting ray method that is free from the leaky mode effect, and a spot scanning method, have been made to increase accuracy, this method is being used only as an auxiliary technique.

F. Far-Field Pattern

This method utilizes the far-field pattern of output optical power instead of the near-field pattern. This method is applicable only to single mode waveguides. The former method is not very accurate because of modal interference within the far field. The latter requires an optical detector with a large dynamic range and the error is more than 5%.

G. Spot Diagram

The spot diagram is a standard characterization of classical optical systems. This is also applied to dielectric optical waveguides and fibers as well as microoptic components. In Fig. 25, we show a typical spot diagram for a GI rod lens that exhibits some aberrations due to the index distribution.

H. Optical Transfer Function

The optical transfer function (OTF) presents the fineness with which we can transmit spatial information in the spatial frequency domain. The OTF $H(s)$ is defined with s the spatial frequency,

$$H(s) = \frac{1}{N} \sum_{i=1}^N \exp(jsx_i), \quad (39)$$

where x_i is the ray position given by the spot diagram and N is the total number of spots. Figure 26 gives one example of OTF obtained from the spot diagram of a DI lens.

I. Phase Space Expression

It is very convenient to utilize a phase space consisting of the ray position and ray slope to express a mode spectrum of a multimode waveguide. Figure 27 shows a phase space plot of the output end of a branching waveguide.

VII. STACKED PLANAR OPTICS

A. Concept of Stacked Planar Optics

Stacked planar optics consists of planar optical components in a stack, as shown in Fig. 28. All components must have the same two-dimensional spatial relationship, which can be achieved from planar technology with the help of photolithographic fabrication, as used in electronics. Once we align the optical axis and adhere all of the stacked components, two-dimensionally arrayed components are realized; the mass production of axially aligned discrete components is also possible if we separate individual components. This is the fundamental concept of stacked planar optics, which may be a new type of integrated optics.

B. Planar Microlens Array

To have stacked planar optics, all optical devices must have a planar structure. The array of microlenses on a planar substrate is required in order to focus and collimate the

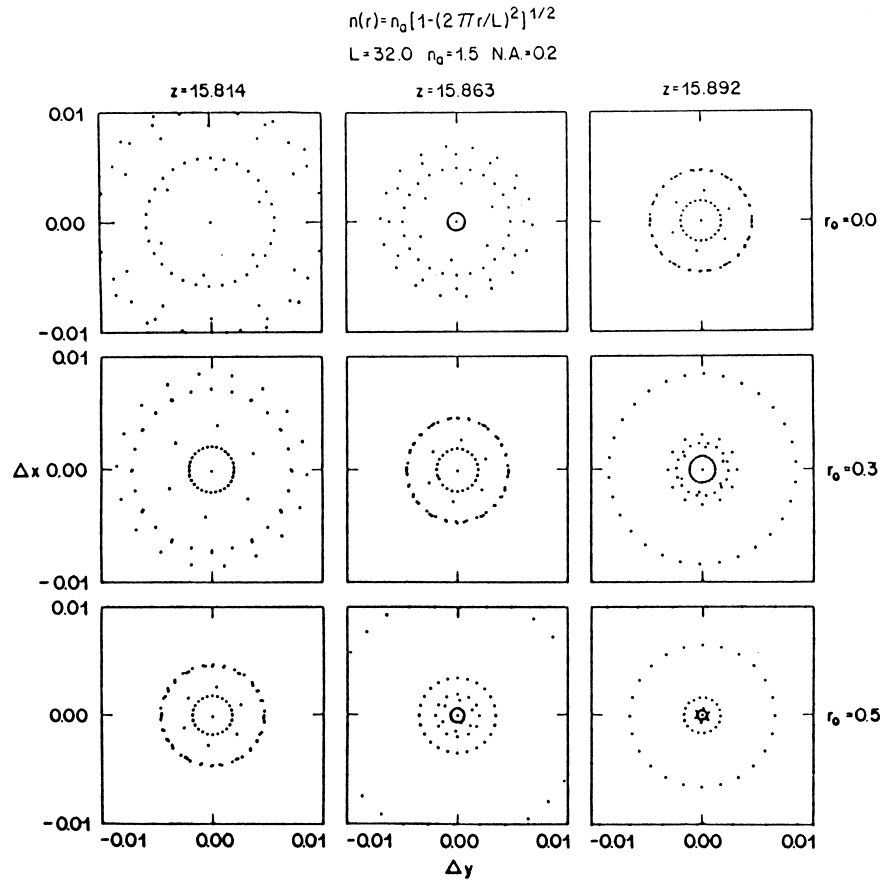


FIGURE 25 Spot diagram for a GI lens. [From Tomlinson, W. J. (1980). *Appl. Opt.* **19**, 1117.]

light in optical circuits. A planar microlens is fabricated by selective diffusion of a dopant into a planar substrate through a mask, as shown in Fig. 29. We can have an array of planar microlenses with a 1.6- to 2.0-mm focal length and a numerical aperture (NA) of 0.34. We have confirmed that the substrate NA can be increased to as high as 0.54 by

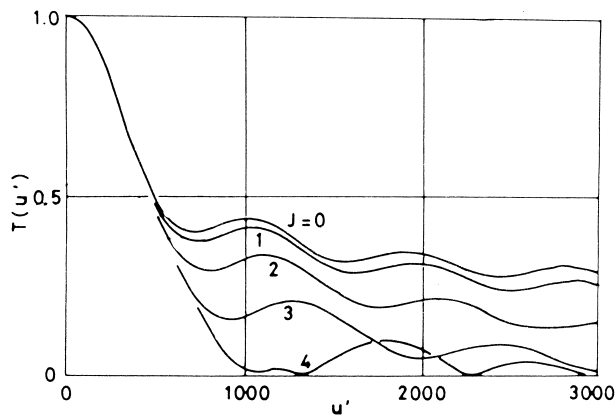


FIGURE 26 OTF for a DI lens. [From Iga, K., Hata, S., Kato, Y., and Fukuyo, H. (1974). *Jpn. Appl. Phys.* **13**, 79.]

stacking two microlenses. This value is considered to be large enough for use as a focusing light from laser diodes.

A planar microlens is fabricated by using an electromigration technique, which was described by Iga, Kokubun, and Oikawa. The substrate is a planar glass $40 \times 40 \times 3$ mm, where planar microlenses were formed as a 40×40 matrix with a 1-mm pitch. The radius of the mask is about $50 \mu\text{m}$ and the radius of the resultant lens is 0.45 mm. The focused spot of the collimated He-Ne laser beam ($\lambda = 0.63 \mu\text{m}$) was measured with the planar microlens. We could observe an Airy-like disk originating from diffraction and aberration, as shown in Fig. 30. The spot diameter is $3.8 \mu\text{m}$, small enough in comparison with the $50\text{-}\mu\text{m}$ core diameter of a multimode fiber, even when we use it in the long-wavelength region $1.3\text{--}1.6 \mu\text{m}$. The data for available planar microlenses are tabulated in Table IV.

C. Design Rule of Stacked Planar Optics

A proposed possible fabrication procedure for stacked planar optics is as follows:

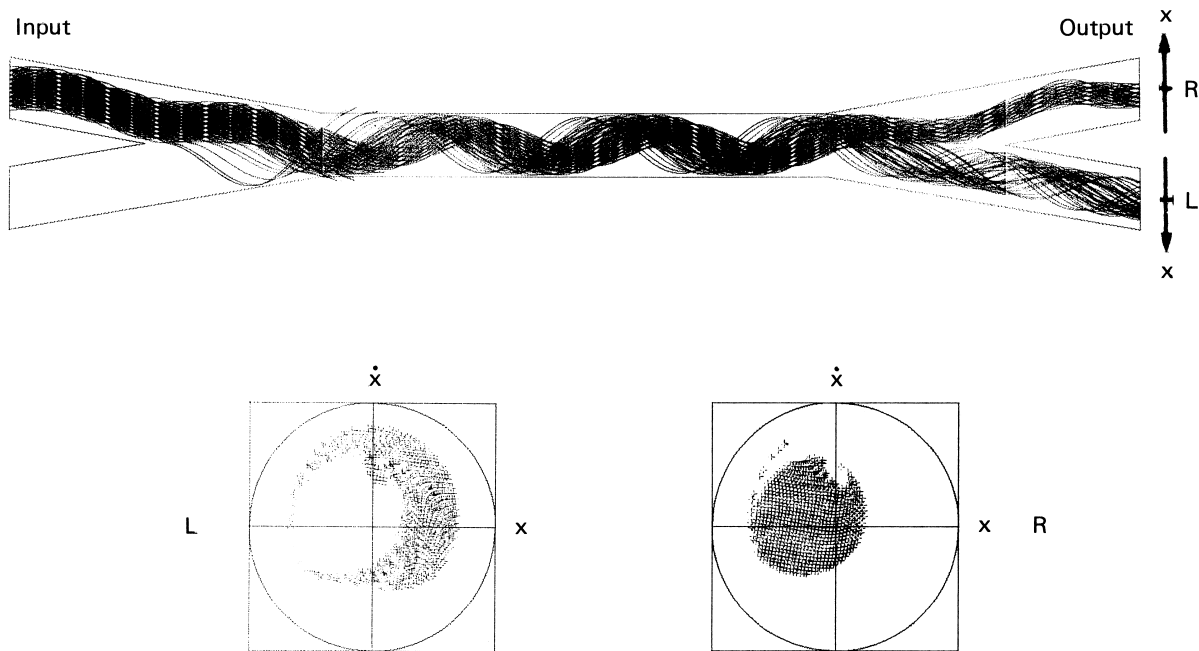


FIGURE 27 Phase space plot of a waveguide branch. [From Kokubun, Y., Suzuki, S., Fuse, T. and Iga, K. (1986). *Appl. Opt.*]

1. Design of planar optical devices (determination of thickness, design of mask shape, etc.)
2. Construction of planar optical devices
3. Optical alignment
4. Adhesion
5. Connection of optical fibers in the case of arrayed components
6. Separation of individual components in the case of discrete components and connection of optical fibers

Features of stacked planar optics include

1. Mass production of standardized optical components of circuits, since the planar devices are fabricated by planar technology
2. Optical alignments, and
3. Connection in tandem optical components of different materials such as glass, semiconductors, and electrooptical crystals. This had been thought difficult in integrated optics consisting of planar substrates in which the connection of different components requires optical adjustment, since light is transmitted through a thin waveguide of only a few microns in thickness and width.

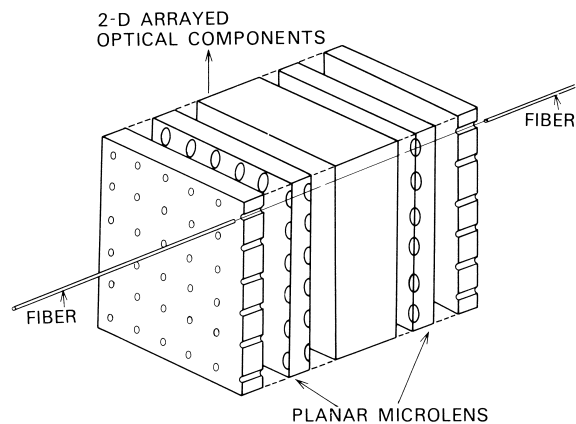


FIGURE 28 Stacked planar optics. [From Iga, K., Oikawa, M., Misawa, S., Banno, J., and Kokubun, Y. (1982). *Appl. Opt.* **21**, 3456.]

D. Applications of Stacked Planar Optics

Many kinds of optical circuits can be integrated in the form of stacked planar optics, as is summarized in [Table V](#).

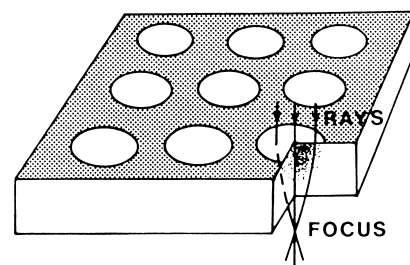


FIGURE 29 Distributed-index planar microlens. [From Iga, K., Oikawa, M., Misawa, S., Banno, J., and Kokubun, Y. (1982). *Appl. Opt.* **21**, 3456.]

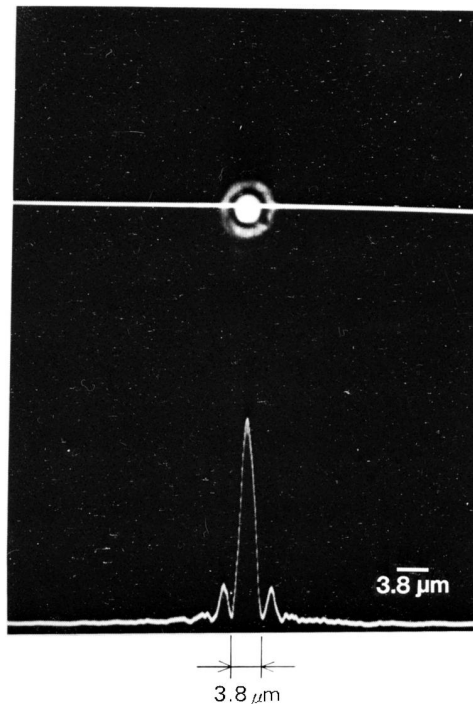


FIGURE 30 Focused spot of a planar microlens. [From Iga, K., Oikawa, M., Misawa, S., Banno, J., and Kokubun, Y. (1982). *Appl. Opt.* **21**, 3456.]

We introduce some components as examples of stacked planar optics. The optical tap in Fig. 31 is the component for monitoring the part of the light being transmitted through an optical fiber. The problem of optical tap is that of reducing the scattering and diffraction loss at the component. The light from the input fiber is focused by the use of a partially transparent mirror placed at the back surface of the device. Some of the light is monitored by the detector, which is placed on the back of the mirror.

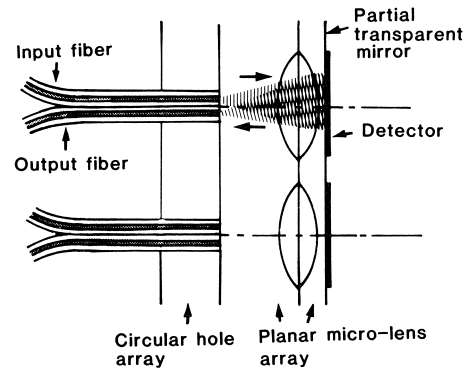


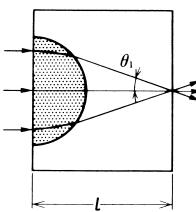
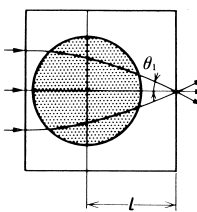
FIGURE 31 Optical tap consisting of stacked planar optics. [From Iga, K., Oikawa, M., Misawa, S., Banno, J., and Kokubun, Y. (1982). *Appl. Opt.* **21**, 3456.]

The main beam is again focused by the same lens on the front surface of the output fiber. With this configuration we can fabricate many optical taps on the same substrate.

The 2×3 branching component has been produced with two pieces of stacked planar microlenses and a half-mirror, as shown in Fig. 32. The light from the input fiber was collimated by the first stacked planar microlens, and a part of the light was reflected by the half-mirror and focused again to output fiber 2. Output fiber 2 was put in contact with input fiber 1. The off-axial value was then $62.5 \mu\text{m}$, the fiber radius. The collimated light through the half-mirror was also focused to output fiber 3. In order to put the fibers on each surface, the thickness of the planar microlens was carefully designed and adjusted by using the ray matrix method.

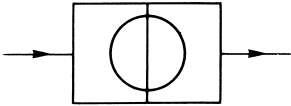
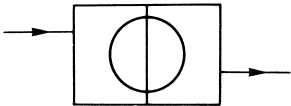
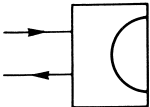
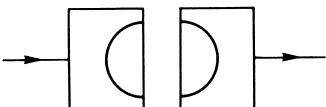
We show an example of optical component using a planar microlens array and micro-optical-bench (MOB) as in Fig. 33. This is based on vertical cavity surface emitting laser array and alignment free MOB.

TABLE IV Design Data of Planar Microlens^a

	Required	Simple	Pained
			
NA	0.2–0.5	0.24	0.38
NA _{eff} (aberration free)	0.2	0.18	0.2
Diameter $2a$ (mm)	0.2–1.0	0.86	0.86
Focal distance l (in glass)	0.3–3.8	3.0	1.76
Lens pitch L_p	0.13–1.0	1.0	1.0

^a From Iga, K., Kokubun, Y., and Oikawa, M. (1984). "Fundamentals of Microoptics," Academic Press, Orlando, FL.

TABLE V Basic Optical Components for Stacked Planar Optics and Optical Circuits^a

Basic components	Application	Reference
Coaxial imaging components	Coupler ^b	4, 11, 12
		
Noncoaxial imaging components (transmission-type)	Branching circuit ^b	3, 4, 12, 13
	Directional coupler ^b	13
	Star coupler ^c	12
	Wavelength demultiplexer ^b	4, 12, 13
		
Noncoaxial imaging components (reflection-type)	Wavelength demultiplexor ^b	12
	Optical tap ^b	12, 13
		
Collimating components	Branching insertion circuit ^b	3, 4, 12, 13
	Optical switch ^c	12, 13
	Directional coupler ^b	3, 4, 11, 13
	Attenuator ^b	11
		

^a From Iga, K., Kokubun, Y., and Oikawa, M. (1984). "Fundamentals of Microoptics," Academic Press, Orlando, FL.

^b Circuit integrated in a two-dimensional array.

^c Circuit integrated in a one-dimensional array.

For the moment we are not concerned with the coupling effect among optical components in the stacked planar optical circuit. But we can construct a three-dimensional optical circuit that structures the network by allowing coupling among adjacent components.

Since the accumulation of lens aberration may bring about coupling loss, the number of stackings is limited by the aberration of the planar microlenses. The reduction of aberration in the planar microlens is important, therefore, if we apply stacked planar optics to more complex components with a large number of stacks.

Stacked planar optics, a new concept in integrating optical circuits, has been proposed. By using stacked planar optics, we not only make possible the monolithic fabrication of optical circuits, such as the directional coupler and wavelength demultiplexer, but we can also construct three-dimensional optical circuits by allowing coupling among individual components in the array with a suitable design.

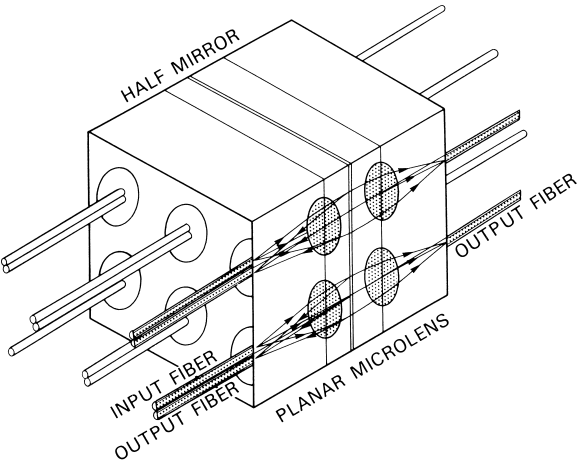


FIGURE 32 Branching component made of planar microlenses. [From Oikawa, M., Iga, K., and Misawa, S. (1984). *Dig. Tech. Pap. Top. Meet. Integr. Guided Wave Opt.*, 1984.]

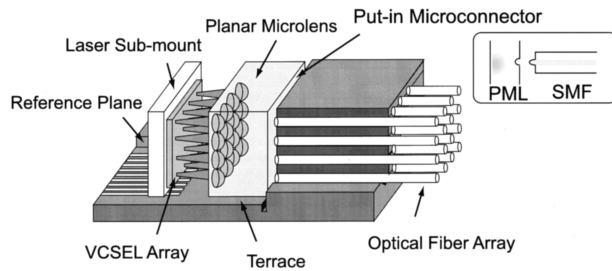


FIGURE 33 Micro-optical-bench. [From Aoki, Y., Shimada, Y., Mizuno, R. J., and Iga, K. (1999). *Opt. Rev.* **7**, 54.]

SEE ALSO THE FOLLOWING ARTICLES

DIFFRACTIVE OPTICAL COMPONENTS • LASERS, OPTICAL FIBER • OPTICAL FIBER COMMUNICATIONS • OPTICAL FIBER TECHNIQUES FOR MEDICAL APPLICATIONS • OPTI-

CAL INTERFEROMETRY • OPTICAL ISOLATORS, CIRCULATORS • OPTICAL WAVEGUIDES AND WAVEGUIDE DEVICES

BIBLIOGRAPHY

- Ezekiel, S., and Arditty, H. J., eds. (1982). "Fiber-Optic Rotation Sensors," Springer-Verlag, Berlin.
- Iga, K., and Kokubun, Y. (1986). "Optical Fibers," Ohm, Tokyo.
- Iga, K., Kokubun, Y., and Oikawa, M. (1984). "Fundamentals of Microoptics," Academic Press, New York.
- Suematsu, Y. (1983). *Proc. IEEE* **71**, 692–721.
- Suematsu, Y., and Iga, K. (1976). "Introduction to Optical Fiber Communication," Ohm, Tokyo (Wiley, New York, 1980).
- Tomlinson, W. J. (1980). *Appl. Opt.* **19**, 1117.
- Uchida, T., and Kobayashi, K. (1982). *Jpn. Annu. Rev. Electron. Comput. Telecommun. Opt. Devices & Fibers* **4**, 172.
- Uchida, T., Furukawa, M., Kitano, I., Koizumi, K., and Matsumura, H. (1970). *IEEE J. Quant. Electron.* **QE-6**, 606.



Nonimaging Concentrators (Optics)

R. Winston
J. O’Gallagher

University of Chicago

- I. Introduction
- II. We Do It with String
- III. Geometrical Vector Flux
- IV. Designs that Are Functionals
of the Acceptance Angle
- V. Solar Thermal Application
- VI. Two-Stage Maximally Concentrating Systems
- VII. Ultra High Flux and Its Applications
- VIII. Conclusion

GLOSSARY

Čerenkov radiation Faint light produced by charged particles moving in a medium at a velocity greater than that of light in the medium.

Compound parabolic concentrator Name given generically to a class of nonimaging collectors with reflecting walls (not necessarily parabolic) that concentrate flux by the theoretical limit.

Dielectric compound parabolic concentrator Nonimaging collector that operates by total internal reflection.

Edge-ray principle (maximum slope principle) Method for designing optical systems with maximum collecting power.

Étendue Product of area times projected solid angle.

Flow line concentrator (trumpet) Nonimaging collector in which the reflecting wall follows the lines of vector flux from a Lambertian source.

Nonimaging optics Optical theory and design that departs from traditional methods and develops techniques for maximizing the collecting power of concentrating elements and systems.

Phase space Abstract space wherein half the coordinates specify locations and half the direction cosines of light rays.

NONIMAGING OPTICS departs from the methods of traditional optical design to develop techniques for maximizing the collecting power of concentrating

elements and systems. Designs that exceed the concentration attainable with focusing techniques by factors of 4 or more and approach the theoretical limits are possible. This is accomplished by applying the concepts of Hamiltonian optics, phase space conservation, thermodynamic arguments, and radiative transfer methods.

I. INTRODUCTION

Geometrical optics is arguably the most classical and traditional of the branches of physical science. Optical design of instruments and devices have been worked out and improved over centuries. From the telescopes of Galileo to the contemporary camera lens, progress, although impressive, has been largely evolutionary, with modern design benefiting enormously from the availability of fast, relatively inexpensive digital computers. In one important respect, however, conventional image-forming optical design is quite inefficient, that is, in merely concentrating and collecting light. This is well illustrated by an example taken from solar energy concentration.

The flux at the surface of the sun ($\approx 63 \text{ W/mm}^2$) falls off inversely with the square of the distance to a value $\approx 1.37 \text{ mW/mm}^2$ above the earth's atmosphere, or typically 0.8 to 1 mW/mm^2 on the ground. The Second Law of Thermodynamics permits an optical device (in principle) to concentrate the dilute solar flux at earth in order to attain temperatures up to but not exceeding that of the sun's surface. This places an upper limit on the solar flux density achievable on earth and correspondingly on the concentration ratio of any optical device. From simple geometry, this limiting concentration ratio is related to the sun's angular size (2θ) by

$$C_{\max} = 1/\sin^2 \theta \approx 1/\theta^2 \text{ (small angle approximation).} \quad (1)$$

Because $\theta = 0.27^\circ$, or 4.67 mrad , $C_{\max} \approx 46,000$. When the target is immersed in a medium of refractive index n , this limit is increased by a factor n^2 :

$$C_{\max} = n^2/\sin^2 \theta. \quad (2)$$

This means that a concentration of about 100,000 will be the upper limit for ordinary ($n \approx 1.5$) refractive materials.

However, conventional means for concentrating sunlight will fall substantially short of this limit, not for any fundamental reason but because imaging optical design is quite inefficient for delivering maximum concentration. For example, consider a paraboloidal mirror used to concentrate sunlight at its focus (Fig. 1). We can relate the concentration ratio to the angle 2Φ subtended by the paraboloid at its focus and the sun's angular size (2θ)

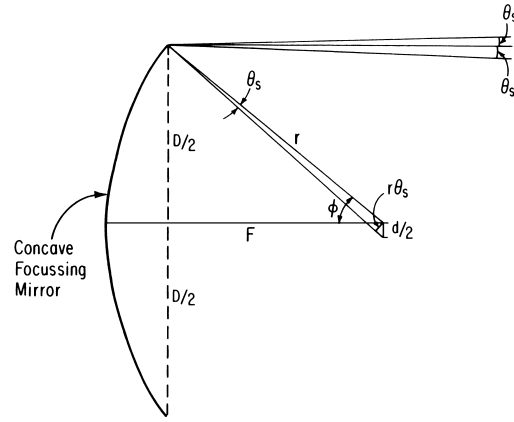


FIGURE 1 Ray diagram for analyzing two-stage concentrators.

$$C_{\text{para}} = (\sin \Phi \cos \Phi / \theta)^2 = \left(\frac{1}{4}\right) \sin^2 2\Phi / \theta^2, \quad (3)$$

where we have used the small angle approximation for θ .

In Eq. (3) C_{para} is maximized at $\Phi = \pi/4$, or

$$C_{\text{para, max}} = 1/(4\theta^2) = \left(\frac{1}{4}\right) C_{\max}. \quad (4)$$

In fact, this result does not depend on the detailed shape of the paraboloid and would hold for any focusing mirror. One fares no better (and probably worse) with a lens because the optimum paraboloid in the above example is equivalent in concentrating performance to a lens with focal ratio $f = 1$ that has been corrected for spherical aberration and coma. Such high aperture lenses are typically complex structures with many components. The limit given by Eq. (1) would require a lens with focal ratio $f = 0.5$, which, as every optical designer knows, is unattainable.

The reason for this large shortfall is not hard to find. The paraboloid images perfectly on axis but has severe off-axis aberration (coma), which produces substantial image blurring and broadening. The essential point is that requiring an image is unnecessarily restrictive when only concentration is desired. Recognition of this restriction and relaxation of the associated constraints led to the development of nonimaging optics. A nonimaging concentrator is essentially a "funnel" for light. Nonimaging optics departs from the methods of traditional optical design to develop instead techniques for maximizing the collecting power of concentrating elements and systems. Nonimaging designs exceed the concentration attainable with focusing techniques by factors of 4 or more and approach the theoretical limit (ideal concentrators). The key is simply to dispense with image-forming requirements in applications where no image is required. The traditional approaches of aberration theory are replaced by a few generic ideas. In the "edge ray" method, maximum concentration is achieved by ensuring that rays collected at the extreme angle for

which the concentrator is designed are redirected, after at most one reflection, to form a caustic on the absorber. Alternatively, in the “flow-line” method, reflective surfaces that follow the lines of net flux, combined with refractive surfaces can lead to “ideal” concentrating elements and systems in three dimensions.

II. WE DO IT WITH STRING

One way to design nonimaging concentrators is to reflect the extreme input rays into the extreme output rays. We call this the “edge-ray method.” An intuitive realization of this method is to wrap a string about both the light source and the light receiver, then allow the string to unwrap from the source and wrap around the receiver. The locus traced out turns out to be the correct reflecting surface. Let’s see how this works in the simplest kind of problem nonimaging optics addresses: collecting light over an entrance aperture AA' with angular divergence $\pm\theta$ and concentrating the light onto an exit aperture BB' (Fig. 2). We attach one end of the string to the edge of the exit aperture B and loop the other end over a line $A'C$ inclined at angle θ to the entrance aperture (this is the same as attaching to a “point at infinity”). We now unwrap the string and trace out the locus of the reflector, taking care that the string is taught and perpendicular to $A'C$. Then we trace the locus of the other side of the reflector. We can see with a little algebra that when we are done, the condition for maximum concentration has been met. When we start unwrapping the string, the length is $A'B + BB'$. When we finish, the same length is $AC + AB'$. But $AC = AA' \sin \theta$, while $AB' = A'B$. So we have achieved $AA'/BB' = 1/\sin \theta$, which is maximum concentration. To see why this works, we notice that the reflector directs all rays at $\pm\theta$ to the edges BB' so that rays at angles $>\pm\theta$ are reflected out of the system and rejected. There is a conservation theorem for light rays called conservation of phase space, or “étendue,” which implies that

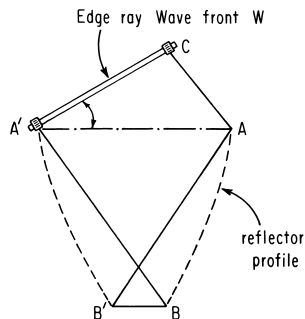


FIGURE 2 $\int_W^{B'} n dl = \text{Constant}$. $AC + AB' = A'B + BB'$
 $AC = AA' \sin \theta$
 $AB' = A'B$
 $\Rightarrow AA' \sin \theta = BB'$

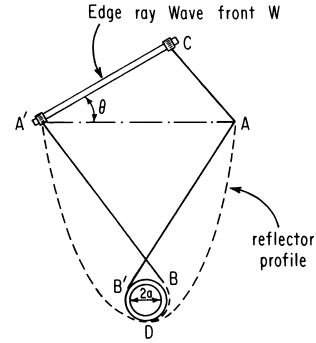


FIGURE 3 String method: $n dl = \text{Constant}$; $AC + AB' + B'D = A'B + BD + 2a$; $AA' \sin = 2a$.

the rays that have angles $<\theta$ are all collected. Next we can try a more challenging problem, where the “exit aperture” is a cylinder of radius a (Fig. 3). Now we attach the string to a point on the cylinder and wrap it around the cylinder. When the string is unwrapped, we find that $AA'/2\pi a = 1/\sin \theta$, which is maximum concentration on the surface of the cylinder. There are, of course, many variations of this method (Fig. 4), which are well illustrated by the following rather extreme example (Fig. 5). Suppose we want to transfer all energy between identical cylinders. The locus is traced out as the string unwraps from one cylinder and wraps around the other. It is interesting to note that this configuration bears no relationship to the elliptical cavity commonly used for this purpose. In fact, as the radii of the cylinders go to zero, our design does not go over into an ellipse but degenerates into a figure with zero area. The point is that the ellipse is never correct for transferring energy between cylinders. This dispels a variety of paradoxes based on using the ellipse for energy transfer.

Although the construction of Fig. 4a gives the two-dimensional compound parabolic concentrator, or CPC, rotating the profile about the axis of symmetry gives the three-dimensional CPC with diameter A_1 at the entrance and A_2 at the exit. The two-dimensional CPC is an ideal concentrator, that is, it works perfectly for all rays within the acceptance angle θ_0 . The three-dimensional CPC is very close to ideal. The flat absorber case is a natural candidate for rotating about the axis because the ratio of diameters ($\sin \theta$) agrees with the ratio of maximum skew. Other absorber shapes, such as circular cross-sections (Fig. 4d) (cylinders in two-dimensional, spheres in three-dimensional), do not have this correspondence because the area of the sphere is πA^2 , whereas the entrance aperture area is $\pi A_1^2/4$.

III. GEOMETRICAL VECTOR FLUX

There is an alternative method for designing “ideal” optical systems that bears little resemblance to the “string

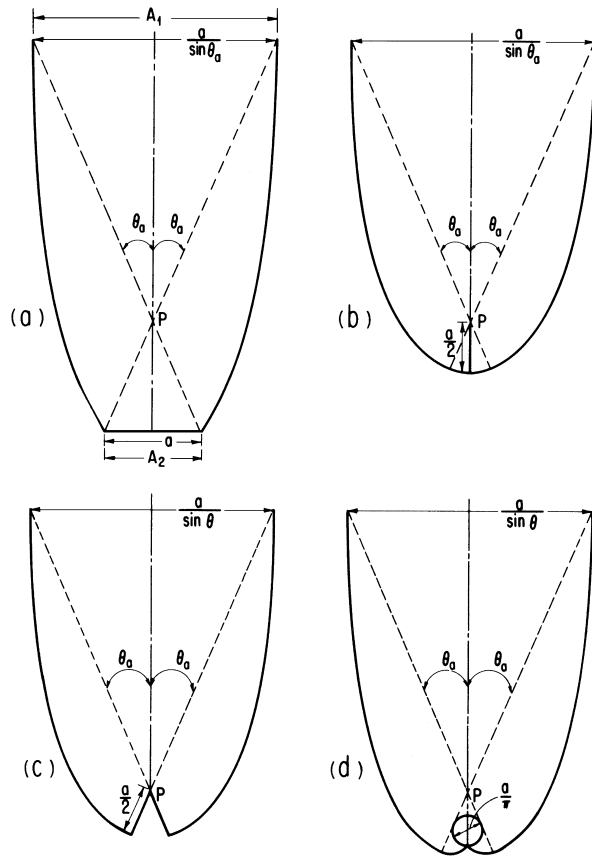


FIGURE 4 Cross-sectional profiles of ideal trough concentrators generalized for absorbers of different shapes. In practice the reflectors are usually truncated to about half their full height to save reflector material with only negligible loss of concentration. Such designs have come to be called compound parabolic concentrators, or CPCs.

method” already described. We picture the aggregate of light rays traversing an optical system as a fluid flow in phase space (see Fig. 4). We can construct a vector field in this space with a surface that is integral through any cross section of the optical system is always the same. Such a field is referred to as a conserved flux, and we have called it geometrical vector flux. Such quantities have proved extremely useful in mechanics, where they were introduced in the early part of this century by Poincaré. We can place reflectors along the lines of flow of this vector field and ask how the field is changed. Because energy cannot flow into a reflector (it all gets reflected), we might hope that the effect is not too severe. In fact, in some geometries the field is not changed at all. This can lead to an important class of nonimaging concentrators. Consider a flat circular Lambertian source. In this case, the lines of flow are hyperbole with foci at the edge of the circle (confocal hyperbole). So a hyperboloid of revolution with its small aperture covered by a circular cap will appear to be a circle of larger radius.

Place this at the focus of a telescope of a solar furnace, and the instrument is “fooled” into having a larger target for the light. In other words, the light is concentrated. Place an appropriate lens over the large aperture, and you obtain an ideal concentrator, which concentrates by the $1/\sin^2 \theta$ limit. The local length of the lens should match the length of the hyperboloid. A variant of this places a lens at both apertures to produce a device for perfectly changing the angular aperture of a beam, say, from θ_1 to θ_2 , an impossible task for imaging optics. The point is that the flow-line designs are perfect in three dimensions, while the string designs rotated about an axis are not. On the other hand, the number of flow-line designs is much more restricted because in very few geometries is one lucky enough to leave the flow-lines undisturbed after adding reflectors.

IV. DESIGNS THAT ARE FUNCTIONALS OF THE ACCEPTANCE ANGLE

A. Tailored Edge-Ray Design Method for Nonimaging Concentrators

The edge-ray design method, which has produced a variety of useful solar concentrators, was further enhanced in the past decade by allowing the acceptance angle to be a function of another parameter, causing the design itself to be a *functional*. The development of new techniques to “tailor” the design of concentrators has enabled the solution of many new types of optical design problems. The basic advance is that nonimaging concentrators can be generalized beyond simple designs that accept only a fixed acceptance angular cone of rays. Using tailoring, edge-ray approaches yield designs that are more general. For instance, in some cases, this will allow two-stage systems with short focal lengths to increase concentration above what had been previously thought possible, as shown for parabolic dish systems by Friedman, Gordon, Rabl, and Ries. Further developments by Jenkins and Winston show how to use a general design method that involves numerical integration of a simple differential equation and can be used to generate most types of nonimaging concentrators including already known solutions, such as CPCs and flow-line concentrators or “trumpets.” The design geometry for the reflector uses numerical integration of a simple differential equation. This geometry is illustrated in Fig. 7, and the differential equation given in polar coordinates by

$$\frac{d(\ln R)}{d\phi} = \tan \alpha(R, \phi), \quad (5)$$

is completely specified if α (the slope angle of the reflector) is known for all points (R, ϕ) in the plane of the reflector. The value of α depends on the incident flux distribution on the aperture of the concentrator and the

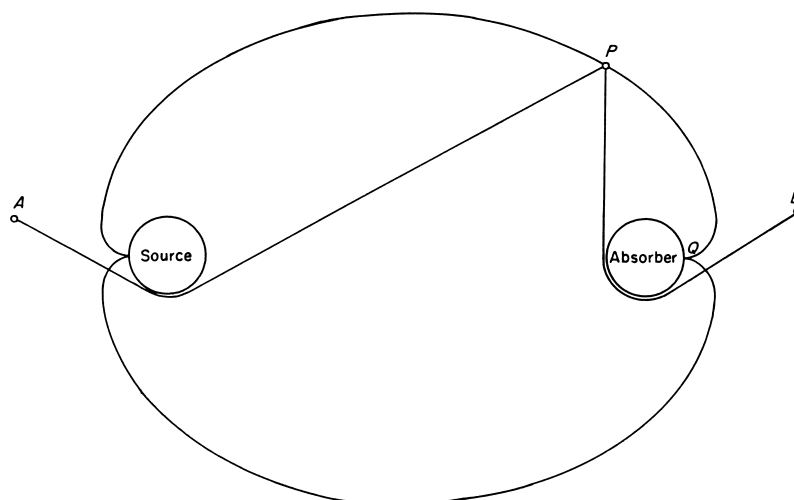


FIGURE 5 The case of a convex source and a convex absorber treated by the string construction.

target absorber's shape. For a flat absorber, shapes similar to CPCs are attained. The geometry of a "tailored" CPC with flat target absorber is shown in Fig. 8. The coordinate origin of the differential equation is placed on the edge of the target and the reflector slope depends only on the variable acceptance angle at the point R , φ . All the concentrator designs used in the experiments described in this article are based on the well established long focal length configurations. However, it should be kept in mind that the use of tailored designs may permit the design of solar furnaces with as yet undeveloped nonimaging secondaries to approach the ideal limits of concentration in a short focal length configuration, thus significantly reducing the size of these systems. Currently, long focal length furnaces with nonimaging secondaries are needed to produce the highest fluxes. The use of fixed angle acceptance secondary concentrators in existing dishes and furnaces may not be the optimum way to compensate for the aberrations induced by the imaging primary mirror. Optimizing the compensation for these aberrations by tailoring both the primary and secondary reflectors should result in more compact two-stage concentrating systems. The future potential for solar furnaces configurations using these new design techniques is just beginning to be explored.

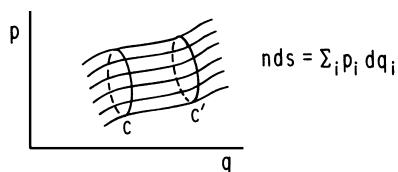


FIGURE 6 Flow lines in phase space. $\int n ds$ along curves C and C' is the same.

Finally, we mention a numerical optimization approach that is sophisticated enough to solve complex problems that do not readily lend themselves to analytical algorithms. An example would be concentrating onto a spherical receiver, where as already noted in Section II, there is a factor of 4 mismatch between étendue and skew invariants. This approach has been called "inverse engineering" by its authors Shatz and Bortz.

V. SOLAR THERMAL APPLICATION

The preferred configurations are those that eliminate heat losses through the back of the absorber, as shown in Fig. 4b and 4d. Note that the geometric concentration is defined relative to the full surface area of the fin absorber (both

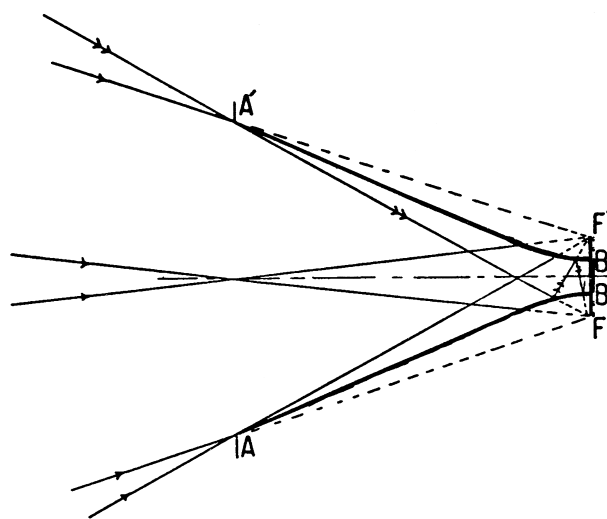


FIGURE 7

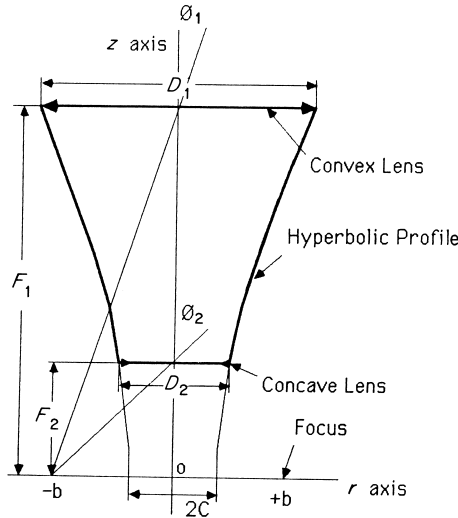


FIGURE 8 The geometry of a "tailored" CPC with a flat target absorber.

sides) or tube (full circumference) and that these designs effectively have no back.

In recent years work has centered on the development of CPCs for use with spectrally selective absorbers enclosed in vacuum. Work at Argonne National Laboratory after 1975 was strongly influenced by the emergence of the Dewar-type evacuated absorber. This thermally efficient device, developed by two major U.S. glass manufacturers (General Electric Co. and Owens Illinois) could, when coupled to CPC reflectors, supply heat at higher temperatures than flat-plate collectors while retaining the advantage of a fixed mount. The heat losses associated with such absorbers are so low that the moderate levels of concentration associated with CPCs provide a dramatic improvement in thermal performance and achieve excellent efficiencies up to 300°C. In fact, the gains associated with even higher concentration ratios are marginal and probably negligible when the added complication and expense of active tracking are considered. This is illustrated in Fig. 9, which shows the calculated efficiency relative to total insolation of a collector consisting of evacuated tubes with a selective absorber coating under increasing levels of concentration.

In characterizing the thermal performance of CPCs in general, we describe the thermal collection efficiency η as a function of operating temperature T by

$$\eta(T) = \eta_0 - \frac{U(T - T_a)}{I} - \frac{\sigma \varepsilon (T^4 - T_a^4)}{CI}, \quad (5)$$

where C is the geometric concentration ratio, T_a the ambient temperature, I the total insolation, σ the Stefan-Boltzmann constant, ε the absorber surface emittance, and U the linear heat loss coefficient.

The optical efficiency η_0 is the fraction of the incident solar radiation (insolation) actually absorbed by the receiver surface after transmission and reflection losses. The other terms represent parasitic conduction losses and radiation losses from the absorber surface, both of which increase substantially as the collector working fluid temperature is increased. The most important point to be seen from Fig. 9 is that the dramatic reduction in relative thermal losses produced by adding a 1.5× reflector or again effectively tripling the concentration from 1.5× to 5× is not continued as one increases the concentration much beyond 5×. For example, at approximately 300°C ($\Delta T/I \cong 0.3$), the thermal losses for a 5× are already so low that a further factor of 5 reduction corresponds to only a negligible fraction of the operating efficiency. It is not necessary, or even desirable, to increase the concentration much beyond 5× when using an evacuated selective absorber at these temperatures. Thus the combination of CPCs up to about 5× with these thermally efficient absorber tubes represents a nontracking strategy for practical solar thermal collection up to power generation temperatures with many unique advantages.

The most developed evacuated CPC design is that optimized for totally stationary collection throughout the year, having an acceptance half-angle $\theta = \pm 35^\circ$ corresponding to a maximum ideal concentration of 1.7×. Commercial versions, truncated to net concentrations of about 1.1–1.4×, are now available with typical optical efficiencies between 0.52 and 0.62, depending on whether a cover glass is used. The thermal performance is excellent, with values of $\varepsilon \cong 0.05$ and $U \cong 0.5 \text{ W/m}^2 \text{ K}$ in Eq. (5). Other versions have opened the acceptance angle to $\pm 50^\circ$ with $C \cong 1.1 \times$ to allow polar orientation and are characterized by $U = 1.3 \text{ W/m}^2 \text{ K}$ (lumping the radiative losses in the

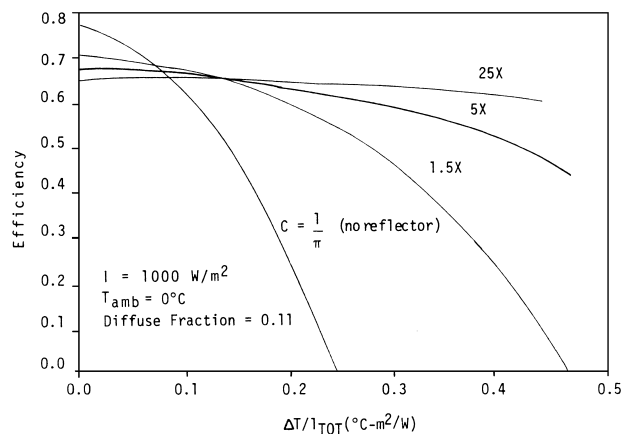


FIGURE 9 Calculated thermal performance curves for evacuated tubular absorbers under increasing levels of concentration. Note that the improvement in increasing the concentration above 5× is marginal.

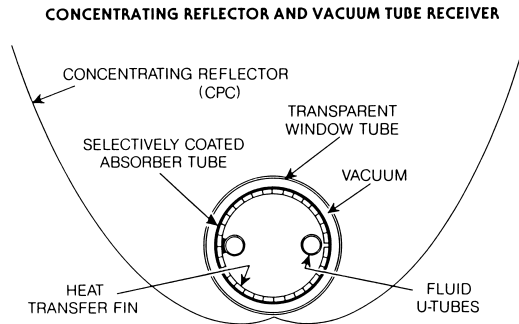


FIGURE 10 Cross section of a contemporary commercial evacuated-tube CPC according to the basic design developed by Argonne National Laboratories.

linear term). An illustration of the basic configuration is shown in Fig. 10. Several manufacturers introduced collectors of these types for applications ranging from heating to absorption cooling to driving Rankine cycle engines. A number of installations with areas greater than 1000 m² have been deployed successfully.

Experimental prototype CPCs have been built in higher concentrations for use with evacuated absorbers. One version with $C = 5.25 \times$ is shown in Fig. 11a. This CPC, studied at the University of Chicago, is a large trough coupled to the same glass Dewar-type evacuated tube as used in the $1.5 \times$ shown previously. Performance measurements for two modules with different reflecting surfaces are shown in Fig. 11b. The upper curve is for a module with a silver foil reflector. It has been operated at 60% efficiency (relative to a direct beam) at 220°C above ambient. This is to be compared with the measured performance of a fully tracking parabolic trough tested by Sandia Laboratories, as shown by the dashed line in Fig. 11b. The performance of the CPC is comparable to that of the parabolic trough at all temperatures tested. The lower curve is for a module with aluminized Mylar reflectors, and even with poorer reflectors, it exhibits quite respectable performance. The angular acceptance properties of the module are in excellent agreement with the design value of $\pm 8^\circ$, which allows collection with 12–14 annual tilt adjustments.

An experimental CPC collector under development at the University of Chicago and Argonne National Laboratory that should ultimately lead to the most practical general-purpose solar thermal collector is the integrated stationary evacuated concentrator (ISEC). The optical efficiency of evacuated CPC solar collectors can be significantly improved over that of contemporary commercial versions discussed above by shaping the outer glass envelope of the evacuated tube into the concentrator profile. Improved performance results directly from integrating the reflecting surface and vacuum enclosure into a single unit. This concept is the basis of a new evacuated CPC

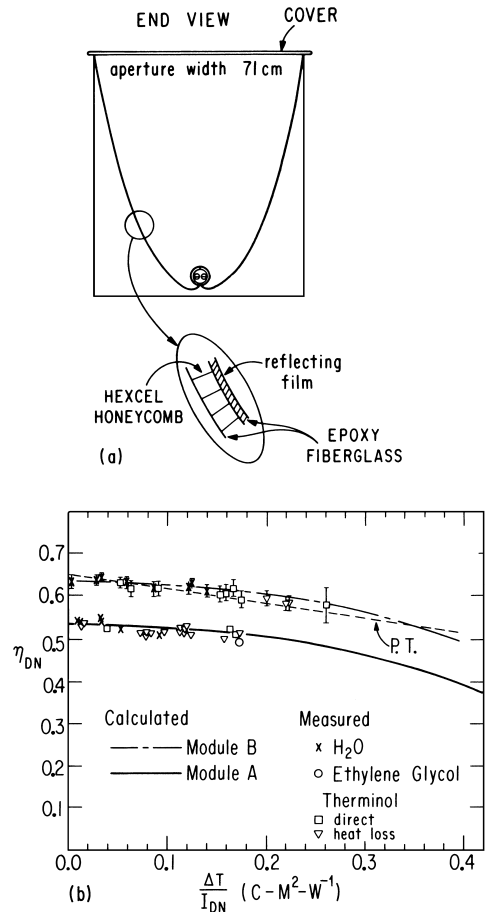


FIGURE 11 (a) Profile of an experimental $5 \times$ CPC for an evacuated tubular absorber built and tested at the University of Chicago. (b) Measured performance curves for two $5.25 \times$ CPC prototype modules. The performance is comparable to that of a commercial parabolic trough, shown by the dashed line.

collector tube that has a substantially higher optical efficiency and a significantly lower rate of exposure-induced degradation than external reflector versions. These performance gains are a consequence of two obvious advantages of the integrated design.

1. Placing the reflecting surface in vacuum eliminates degradation of the mirror's reflectance, thus permitting high-quality (silver or aluminum with reflectance $\rho = 0.91\text{--}0.96$) first-surface mirrors to be used instead of anodized aluminum sheet metal or thin-film reflectors ($\rho = 0.80\text{--}0.85$) typical of the external reflector designs.
2. The transparent part of the glass vacuum enclosure also functions as an entrance window and thus eliminates the need for an external cover glazing. This increases the initial optical efficiency by a factor of $1/\tau$, where typical transmittances $\tau = 0.88\text{--}0.92$.

For the past several years, the solar energy group at the University of Chicago has been developing this concept in collaboration with GTE Laboratories, which fabricated the tubes. Of the 80 prototype built, 45 were assembled into a panel with a net collecting area of about 2 m².

The ISEC shown in Fig. 12 is an extended cusp tube CPC matched to a circular absorber of diameter 9.5 mm. The design acceptance half-angle $\theta_a = \pm 35^\circ$ was chosen to permit stationary operation throughout the year. After truncating the CPC, the net concentration was 1.64 \times . This collector was tested at Chicago for three years and routinely achieved the highest high-temperature performance yet measured for a fixed stationary mount collector. Performance curves based on these tests are shown in Fig. 13,

along with curves for three other collector types. Note in particular that for temperatures up to about 200°C, the ISEC is comparable to a fully tracking trough and remains respectable up to temperatures approaching 300°C. The relative performance advantages are similar when the comparison is made on an annual energy delivery basis at a variety of locations, as shown for one location in Fig. 14.

The design problems for nonevacuated CPC collectors are entirely different from those for CPCs with evacuated absorbers. The nonevacuated collectors are particularly vulnerable to high heat losses if an improper design is used. One must be careful to minimize or eliminate heat loss via conduction through the reflectors. This can be accomplished by using reflectors with a thickness that is

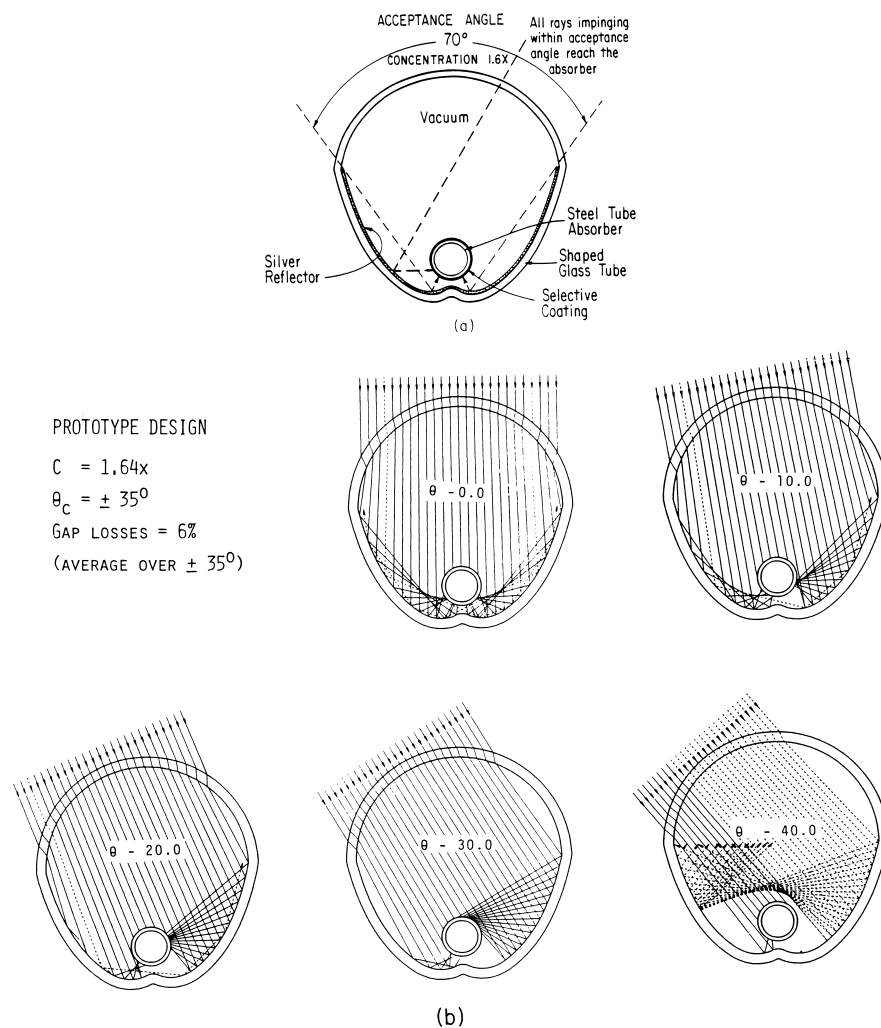


FIGURE 12 (a) Details of the actual profile shape and collector design for the integrated stationary evacuated concentrator (ISEC) tube, which has achieved a thermal efficiency of 50% at 200°C. (b) Ray trace diagram showing how essentially all the solar energy incident within $\pm 35^\circ$ is directed onto the absorber tube. Because the reflector cannot physically touch the absorber, as required for an ideal concentrator, a small fraction is lost in the gap between the reflectors and the absorber.

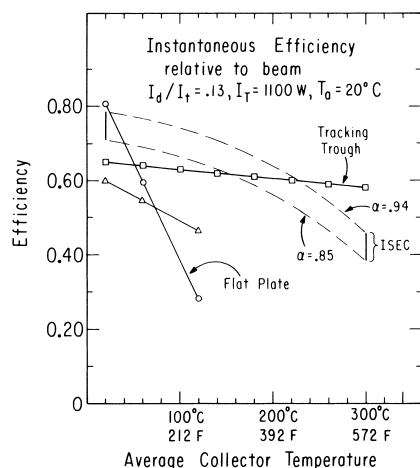


FIGURE 13 Comparative peak performance for a tracking parabolic trough, an ISEC, contemporary (external reflector) evacuated CPCs (triangles) and flat plates. The ISEC's superior performance up to temperatures above 200°C, achieved with no moving parts, makes it an extremely flexible solar thermal collector.

negligible compared to the overall dimension (e.g., height, aperture) of the trough, such as metallized plastics or films, or by thermally decoupling the absorber from the reflectors by a small gap maintained by insulating standoffs. Two prototype CPCs with nonevacuated absorbers, a 3× and a 6×, were built and tested extensively as part of the early program at Chicago. The features and performance of these collectors are summarized here. The optical efficiencies and total heat loss coefficients were $\eta_0 = 0.68 \pm 0.01$ and $U = 1.85 \pm 0.1 \text{ W/m}^2 \text{ } ^\circ\text{C}$ for the 6×, and $\eta_0 = 0.61 \pm 0.03$ and $U = 2.7 \pm 0.02 \text{ W/m}^2 \text{ } ^\circ\text{C}$ for the 3×.

The efficiencies of these nonevacuated CPCs are to be compared with typical values for flat plates of $\eta_0 = 0.70$ –

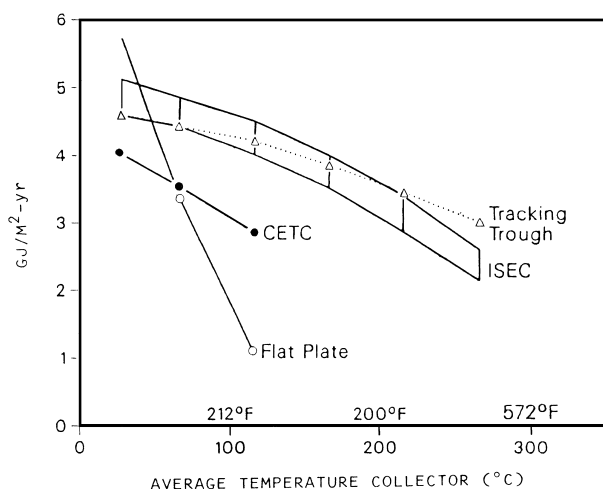


FIGURE 14 Comparative annual energy delivery at Phoenix, Arizona, for the collector types in Fig. 13.

0.78 and $U = 4.5\text{--}7 \text{ W/m}^2 \text{ K}$. Despite lower optical efficiencies, the CPCs outperform typical flat-plate collectors above temperatures as low as 10°C above ambient (for the 6×) to about 35°C above ambient (for the 3×). This is particularly important because the 3× should represent a relatively inexpensive collector design. Although a detailed economic analysis cannot be based on the prototype construction methods used here, several unique features contribute to its low cost potential, among them the relatively small absorber cost and very limited insulation requirements.

VI. TWO-STAGE MAXIMALLY CONCENTRATING SYSTEMS

The principal motivation for employing optical concentration with photovoltaic cells in solar energy is economic. By using what one hopes are relatively inexpensive lenses or mirrors to collect the sun's energy over a large area and redirect it to the expensive but much smaller energy conversion device, the net cost per unit total area of collection can be reduced substantially. Alternatively, to generate electricity through the thermodynamic conversion of solar heat to mechanical energy, concentration is required to achieve the high temperatures necessary to drive a heat engine with reasonable efficiency. In this case, the solar flux is directed to an absorber (often a cavity) small enough that the heat losses, even at high temperatures, remain relatively small. It often turns out in both the photovoltaic and thermal conversion cases that the desired concentration is much higher than can be achieved with non-tracking CPC-type devices. Conventionally, these higher concentrations are achieved by means of some kind of focusing lens or paraboloidal mirror that is not maximally concentrating.

It is not widely recognized that the nonimaging techniques described in the previous section can be used to design secondary elements that can augment the concentrations of more conventional focusing elements used as primaries, and that such a hybrid optical system can also approach the allowable limit. Applications of such two-stage designs lie in the regime of higher concentration and small angular acceptance, where the geometry of a single-stage CPC becomes impractical. The fundamental advantage is the same as in lower-concentration applications and may be expressed in complementary ways: either significant additional system concentration can be attained (i.e., a smaller, lower-cost absorber) or the angular tolerances and precision can be relaxed while maintaining the same level of concentration.

The limits of achievable levels of solar concentration are represented in Fig. 15 for both line focus and point focus

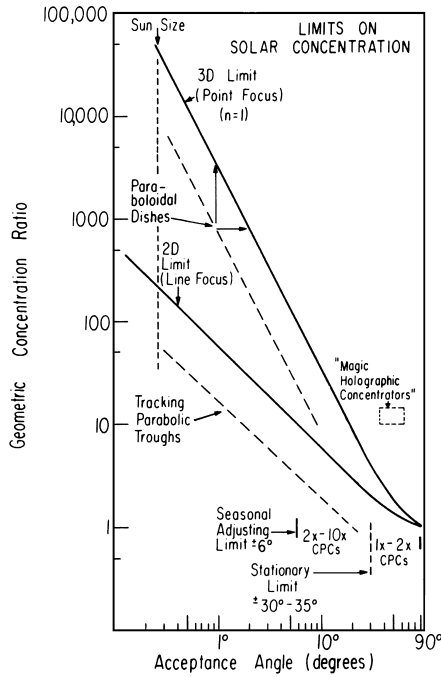


FIGURE 15 Maximum geometric concentration ratio corresponding to a given half-angle of incidence permitted by physical conservation laws (the thermodynamic limit) (solid lines). Traditional focusing concentrators fall about a factor of 4 below this limit (dashed lines). Proposed concentrator designs purporting to achieve 10:1 ratios with no tracking based on holographic techniques cannot work because they violate this limit.

geometries across the range of possible desired angular acceptances from wide angles permitting stationary or seasonal adjusting CPCs down to the angular subtense of the sun (± 4.6 mrad). If one could achieve the thermodynamic limit in a point focus geometry (no configuration solution that could accomplish this is known) with no slope or alignment errors, one could reach the thermodynamic limit of 46,000 suns and, in principle, reach the sun's surface temperature of 6000 K. In practice, slope and alignment tolerances, typically about $\pm 0.5^\circ$, and the aberrations associated with focusing designs limit the actual values to 30–70 \times in line focus and 1000–5000 \times in point focus geometries. Use of a nonimaging secondary can increase these limits (or tolerances) as indicated.

It has been proposed that holographic optical elements (HOEs) could achieve concentrations in the range 10–20 \times without tracking by stacking individual elements designed to be effective at different times of the day. This is impossible, as indicated by the dashed box in Fig. 15, because it would violate the thermodynamic limit.

In this section, we describe both photovoltaic applications, where the primary is a lens, and thermal electric applications, where the primary is a paraboloidal mirror. In each case, we discuss only the point focus configura-

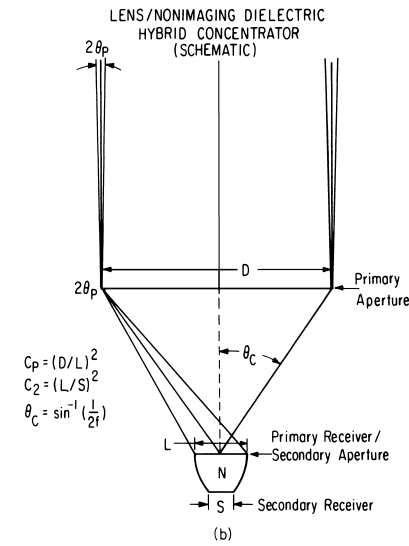
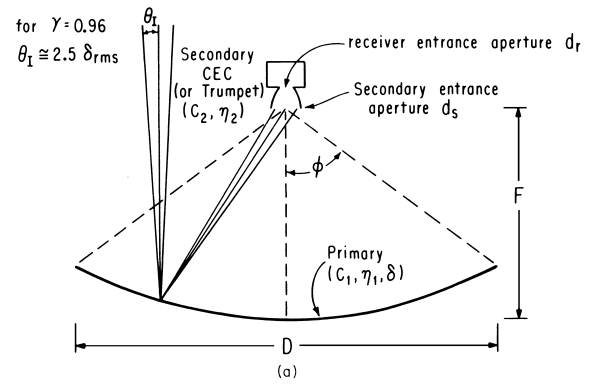


FIGURE 16 Geometric and optical parameters characterizing (a) the two-stage point focus reflecting dish and (b) lens concentrators, discussed in the text.

tion. Schematic drawings of the basic elements for the two cases are shown in Fig. 16.

For the thermal application (Fig. 16a), the primary is characterized by its focal length F and aperture diameter D , which define the rim angle

$$\tan \phi = |2f - 1/8f|^{-1}, \quad (6)$$

where $f \equiv F/D$. The secondary is a nonimaging concentrator of either the compound elliptical concentrator (CEC) type, a variant of the more familiar compound parabolic concentrator, or the hyperbolic trumpet type. It is convenient to simplify the analysis by characterizing the primary as having a conical angular field of view of half-angle $\pm \theta_1$. This is chosen to accommodate the angular tolerance budget of the primary, including concentrator slope errors, specular spread, pointing error, and incoming direct sunlight.

The thermodynamic limit in a point focus reflecting geometry is given by Eq. (4b), while geometric arguments show that the geometric concentration of the primary alone must be

$$C_1 < \sin^2 \phi \cos^2 \phi / \sin^2 \phi_1 \quad (7)$$

to intercept all the energy incident within $\pm\theta_1$. The limiting concentration for the secondary is

$$C_2 < 1 / \sin^2 \phi. \quad (8)$$

Therefore, the maximum combined concentration is

$$C_1 C_2 = \cos^2 \phi / \sin^2 \theta_1, \quad (9)$$

which approaches the maximum limit for small ϕ (large f).

For the photovoltaic concentrator (Fig. 16b), the nonimaging secondary is formed from a transparent dielectric material with an index of refraction $n = 1.3$ – 1.5 in contact with the solar cell. It is usually possible to ensure that total internal reflection (TIR) occurs for all rays accepted by the secondary, thus saving the cost of applying and protecting metallic reflecting surfaces. In such a totally internally reflecting dielectric CPC (DCPC), the refractive power of the dielectric operates in combination with the reflecting profile shape so that secondary concentration ratios are achieved that are a factor n^2 larger than is possible with conventional reflecting secondaries. The primary is a lens (usually a Fresnel lens) of aperture diameter D that focuses normally incident rays at a point that defines the center of the primary receiver. Rays from outer edges of the primary are brought together with a convergence angle $2\theta_c$. If the incident rays subtend an angle of $\pm\theta_p$ on either side of the aperture normal, one can show that the maximum possible geometric concentration C_{\max} attainable is

$$C_{\max} = n^2 / \sin^2 \theta_p. \quad (10)$$

Here n is the index of refraction at the final absorber surface relative to that just outside the collecting aperture. Equation (10) defines the “ideal” limiting concentration allowed by physical conservation laws for a dielectric secondary. For a focusing lens of focal ratio f' , where

$$f' \equiv 1/2 \sin \theta_c, \quad (11)$$

such that it corresponds to the generalized focal ratio used to express the Abbe sine condition for off-axis imaging, one can show that the *actual* concentration achieved by the lens alone is

$$C_p = 1/4 f'^2 \sin^2 \phi_p. \quad (12)$$

Comparing Eqs. (12) and (10), we see that for practical lens systems where $f' \gtrsim 1$, C_p falls short of the limit by a factor of $(4 f'^2 n^2)^{-1}$.

If a DCPC-type secondary with entrance aperture diameter L and exit aperture diameter S is placed at the focal spot of the imaging primary lens, it can achieve an additional geometric concentration $C = L/S$, which is given by Eq. (10) with θ_c replacing θ_p , or

$$C_2 = n^2 / \sin^2 \theta_c = 4 f'^2 n^2. \quad (13)$$

Combining Eqs. (12) and (13) shows that, in principle, the two-stage system can attain an overall geometric concentration equal to the thermodynamic limit. In practice, certain compromises are needed that reduce this somewhat; but typically secondary concentrations in the range 7 – $10\times$ are readily achievable.

The limits of concentration for both reflecting and dielectric refracting systems with and without optimized nonimaging concentrators are shown in Fig. 17 as a function of the focal ratio of the primary.

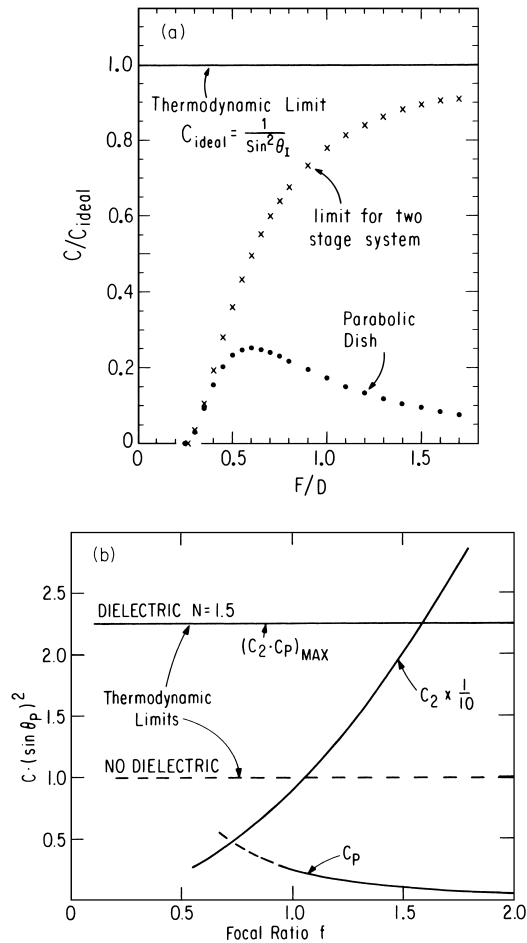


FIGURE 17 Maximum geometric concentration achieved by (a) reflecting dish primaries with nonimaging CPC-type secondaries and (b) lens primaries combined with refracting dielectric CPC (DCPC) secondaries.

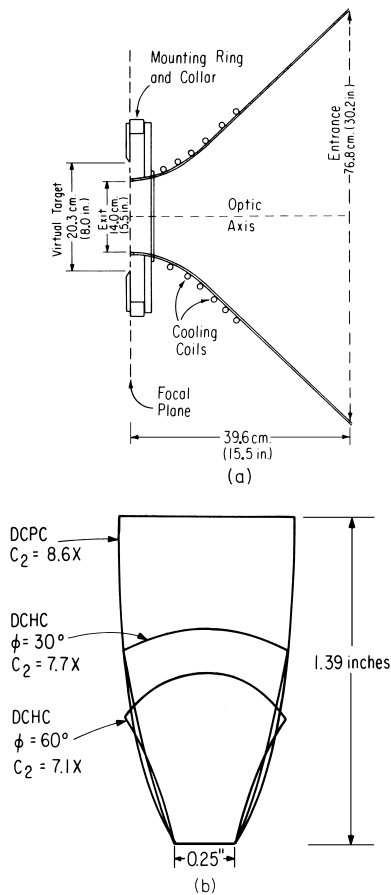


FIGURE 18 (a) Practical thermal secondary concentrator referred to as the “trumpet,” which has been built and tested at the University of Chicago. (b) Profiles for a DCPC and two dielectric compound hyperbolic concentrators (DCHCs) used for photovoltaic secondaries.

A. Applications

Figure 18 shows examples of practical nonimaging secondaries used to date for both kinds of applications. The flow line or “trumpet”-shaped secondary is a recent development with particular advantages in a retrofit mode, that is, to increase the concentration of a dish that is already designed and built. Figure 18b shows how introducing a small amount of curvature into the front surface of the secondary provides some of the concentration so that the overall height of the side-walls can be reduced, yielding substantial savings in material.

For thermal applications with a CPC-type secondary, a particularly attractive option is the use of a primary of longer focal ratio (f -number $\gtrsim 1.0$) with flat mirror facets, in which case the number of required facets can be reduced by a factor approximately equal to the secondary concentration ratio. For example, with a secondary with $C^2 = 5\times$ the number of flat facets required to achieve $150\times$ can be

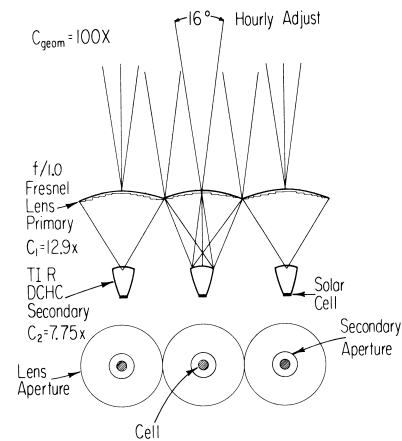


FIGURE 19 Illustration of a lens/DCHC combination that could provide 100:1 concentration while requiring only approximately hourly adjustment. Prototypes of actual devices based on this concept are under development.

reduced from nearly 200 to less than 40. Details of such a design are being developed.

A two-stage photovoltaic concentrator presently being developed is shown in Fig. 19. This system has no real analog in contemporary solar concentrator configurations because conventional $100\times$ concentrators must track quite accurately ($\lesssim \pm 1^\circ$) and crude tracking devices have much lower concentrations ($\sim 5\text{--}15\times$). Here a good primary is designed to provide a concentration of about $13\times$ with an acceptance angle wide enough to accommodate the sun’s movement for close to 1 hour. A TIR secondary provides the additional concentration (here $7.7\times$) required to make the economic savings associated with reduced cell area really worthwhile. No scale is shown in Fig. 19 because a study of the trade-offs associated with size is one of the objectives of the present work. For reference, note that an 0.5-in. cell would correspond to a 5-in.-diameter circular collecting lens. These individual elements could be arranged in a hexagonal close-packed geometry if a high packing density is required, or on a square lattice if not. Finally, note that in this geometry the function of the secondary in redistributing the concentrated sunlight more uniformly on the cell is especially useful.

Several years ago, our group at The University of Chicago set out to demonstrate the ability of nonimaging optics to concentrate sunlight at or near the theoretical limit implied by the sine law. Actually, we chose to immerse the “hot zone” in a refractive medium to further boost the limit by the square of the index of refraction. These have been small-scale “roof-top” experiments carried out by talented and dedicated graduate students. In the first round of experiments we attained “56,000 suns.” This was a fairly conservative design where the nonimaging optical element was a hollow silver funnel filled with

an oil (refractive index 1.53). In the latest round of experiments we reached “84,000 suns.” This design was more daring and used a solid sapphire nonimaging optical element (refractive index 1.78). At such huge concentrations, the power flux measured is actually some 15% greater than it is at the surface of the sun. So the answer to the question: Where would one go to find the highest solar flux in the solar system? is—to the roof-top of a building at The University of Chicago!

VII. ULTRA HIGH FLUX AND ITS APPLICATIONS

There are three different regimes for applications. These are ultra-high flux applications in air, ultra-high flux applications in a refractive medium, and conventional dish-thermal retrofit applications.

Corresponding secondaries of choice in each of these regimes respectively are reflecting CPCs (or CECs), refracting hi-index nonimaging concentrators (either dielectric filled CPCs (DCPCs), or dielectric totally internally reflecting concentrators (DTIRC)s), and flow-line or “trumpet” concentrators. Each type of secondary has advantages and disadvantages, and particular features that must be borne in mind in optimizing its design.

Following through with a long-standing desire to explore the development of these techniques for larger scale, higher power applications, the National Renewable Energy Laboratory (NREL), designed and constructed a scaled-up solar concentrating furnace facility specifically intended to make use of nonimaging optics. This high flux solar furnace (HFSR) concept uses a modified long focal length design and is capable of delivering up to 10 kilowatts to the focal zone and, in particular, to the entrance aperture of a secondary concentrator.

Recent experiments and others performed at the National Renewable Energy Laboratory’s (NREL) High Flux Solar Furnace (HFSF) have demonstrated the effectiveness of using concentrated sunlight and advanced nonimaging secondaries to pump lasers and produce fullerenes (potentially useful new forms of molecular carbon). The recently developed techniques that allow more flexibility in design have been used to develop two new configurations, each of which couples the high solar flux available at the HFSF to unusually shaped targets which impose unusual constraints.

A. Solar-Pumped Lasers

There are two methods for pumping lasers using sunlight, end-pumping and side-pumping. There are advantages and disadvantages for each. The end-pumping scheme can be used on very small scales, but the pump light enters the

laser crystal entirely from one end. Previous high-flux measurements show that this system can not be scaled up indefinitely as the optical coupling between dielectric surfaces degrades when exposed to ultrahigh solar fluxes at the kW power scale. An end-pumping scheme has a maximum output power to about 5 W. The side-pumping scheme does not face this limitation because the lasing medium is excited sunlight entering over a much larger aperture (the side walls of the laser crystal). Current research suggests that space-based applications hold the most promise for solar lasers. In space, the solar insolation level, the input pump source, is more intense and much more stable. This in turn allows more stable lasing configurations that increase laser output efficiency and brightness (and building an electrically pumped laser is difficult).

B. Fullerene Production

Researchers at Rice University announced the creation and isolation of macroscopic quantities of C_{60} and C_{70} molecules and other less common fullerenes by placing a carbon arc inside a partially evacuated inert atmosphere (helium gas at around $\frac{1}{6}$ of an atm) and supplying a high voltage. However, there turned out to be a problem with using arc lamps to produce fullerenes in that as the size of the arcs increased, the yield relative to total soot vaporization decreased. This led to the consideration of using highly concentrated sunlight. Sunlight on earth has a relatively low UV content and mirrors can be made less reflective at lower wavelengths, making it a promising source to produce fullerenes. Using solar flux to produce small amounts of fullerenes was demonstrated simultaneously by both researchers at Rice University and researchers at NREL. The method at NREL used a nonimaging secondary CPC in conjunction with the high flux solar furnace. More recent work is underway to develop mass-production techniques. The scalability of solar as opposed to small-scale arc lamp systems may make this the preferred method of fullerene production.

C. Solar Processing of Materials

Ideally, concentrated sunlight can cause materials to reach equilibrium temperatures approaching those found on the sun’s surface (5800 K). This allows solar to cover a wide range of applications at all temperatures below this level. The high fluxes also lead to extremely high heating rates in nonequilibrium setups because the heating rate is proportional to the incident flux minus the reradiation losses, which are small for low temperatures. Surfaces of materials can be superheated to induce chemical reactions that modify the properties and composition on a material’s surface while leaving the bulk material unchanged.

D. Dish Thermal Applications (Reflecting Trumpet Secondary)

A project is being carried out by the University of Chicago to design a practical trumpet secondary concentrator to be used in combination with a faceted membrane primary concentrator for dish-stirling applications. This is a retrofit design for a dish that was originally designed as a single stage and so it cannot attain the full power of a fully optimized two-stage concentrator. Nevertheless, preliminary ray trace studies show that the addition of a small trumpet to the receiver aperture will allow a reduced aperture size and corresponding lower thermal losses at the operating temperature of 675°C.

E. Solar Thermal Applications of High-Index Secondaries

There is interesting research, at the Weizmann Institute of Science (WIS) in Israel, studying the use of large-scale dielectric secondary concentrators and an efficient extractor tip to create high-temperature gas turbine engines. The use of high-index concentrators with TIR reflection conditions (no losses from reflection that occur on metal surfaces) minimize optical losses in the system and increases the concentration limit by n^2 . In general, the use of a higher-index material to form the aperture of a gas turbine will not give higher operating temperatures. This is because the reradiation is proportional to n^2 times the area of the secondary exit aperture, which normally cancels the n^2 gain in concentration obtained by using higher-index materials. All gas turbines operate at temperatures much less than the sun's (6000 K) and, therefore, the spectrum of reradiation is considerably red-shifted into the infrared (IR) region. By applying selective coatings onto the dielectric concentrator and extractor optics, one may prevent IR radiation from reradiating and it may be possible to reduce heat losses inside high-temperature electrical generation systems. Research is in progress at WIS to study the coupling of dielectric secondaries with gas turbines. They use a six-sided extractor tip to let light out of the high-index secondary and into the high-temperature gas environment of the turbine ($n = 1$). This is another possible use of high-index secondaries other than end-pumping solid-state laser crystals.

F. Using Highly Concentrated Sunlight in Space

The techniques of nonimaging optics are particularly valuable in space or lunar environments where the use of solar thermal energy has obvious advantages. Earlier preliminary studies have explored this concept for the produc-

tion of cement from lunar regolith and for solar thermal propulsion in space. For example, extremely high temperatures, in the range 1700–1900°C, are necessary for the production of cement from lunar minerals. Such temperatures will, in turn, require very high levels of solar flux concentration. Energy budgets for the support of permanent staffed operations on the lunar surface are expected to be limited. For high temperature thermal (i.e., >300°C) end uses, direct solar energy has obvious advantages over most other practical power sources. Conventional combustion processes are clearly impractical and conversion of electricity (either solar or nuclear generated) to high-temperature heat represents a very wasteful use of high-quality energy. On the other hand, solar radiation is abundant and nondepletable. Most importantly, it is readily converted to heat with high efficiency, although at high temperatures this requires high concentration.

Earlier work illustrated the feasibility of some particular two-stage configurations and indicated that the corresponding solar thermal conversion efficiency can be about 2.5 times that of the corresponding conventional design at 1500°C. A preliminary design configuration for such a high-flux nonimaging solar concentrating furnace for lunar applications was proposed. It employs a tracking heliostat, a fixed, off-axis, two-stage concentrator with a long focal length utilizing a nonimaging trumpet or CPC type secondary deployed in the focal zone of the primary. An analysis of the benefits associated with this configuration used as a solar furnace in the lunar environment shows that the thermal conversion efficiency can be about 3–5 times that of the corresponding conventional design at 2000°C. Furthermore, this configuration allows several other advantageous practical design options in addition to high performance. For instance, the furnace itself and associated support structure and equipment need not shade the primary collecting aperture and spherical or faceted primaries may be able to be used.

Solar thermal propulsion systems in space will require very high temperatures to generate necessary levels of thrust by the direct solar heating and resulting expansion and expulsion of the propellant material. The generation of such temperatures, in the range 1400–2200°C, will require very high levels of solar flux concentration. In practice, to attain such levels it may be useful, and perhaps even necessary, to incorporate some form of ideal or near ideal nonimaging concentrator.

Applications of the techniques of nonimaging optics in the design of solar thermal concentrators has resulted in significant performance improvements in conventional short focal length dish concentrators, dramatic changes in the approach to solar furnace design, and the establishment of new solar flux concentration records in air and in refractive media. These accomplishments in just over ten

years of active experimental development are indications of what can be achieved with these powerful methods. The wide variety and range of applications make this a very exciting field. It extends the bounds of solar energy research into solar manufacturing and even moves us into outer space.

VIII. CONCLUSION

Nonimaging optics departs from the methods of traditional optical design in order to develop techniques for maximizing the collecting power of concentrating elements and systems. Designs that exceed the concentration attainable with focusing techniques by factors of 4 or more and approach the thermodynamic limit are possible. This is accomplished by applying the concepts of Hamiltonian optics, phase space conservation, thermodynamic arguments, and radiative transfer methods. Concentrators based on this approach are finding increasing use in solar energy and a variety of other applications and include the now well-known nontracking compound parabolic concentrator (CPC).

Compound parabolic concentrators permit the use of low to moderate levels of concentration for solar thermal collectors without the requirement of diurnal tracking. When used in conjunction with an evacuated absorber with selective surfaces, a fully stationary CPC (designs with concentrations of $1.1\text{--}1.4\times$ are now commercially available) has a typical efficiency of about 40% at 150°C (270°F) above ambient with available conventional materials. An experimental $5\times$ CPC requiring approximately 12 tilt adjustments annually, when used with a similar available evacuated absorber, has a measured efficiency of 60% at 220°C above ambient and is capable of efficiencies near 50% at 300°C . With such thermally efficient absorbers, higher concentrations are not necessary or desirable.

Argonne National Laboratory and the University of Chicago are working on a research and development program to develop an advanced evacuated-tube collector that will be suitable for mass production by industry, will compete successfully with conventional flat-plate collectors at domestic hot water (DHW) temperatures, and will be suitable for industrial process heat (IPH) and/or cooling applications. The essence of the design concept for these new collectors is the integration of moderate levels of nonimaging concentration inside the evacuated tube itself. This permanently protects the reflecting surfaces and allows the use of highly reflecting front-surface mirrors with reflectances greater than 95%. Fabrication and long-term testing of a proof-of-concept prototype have established the technical success of the concept. Present work

is directed toward the development of a manufacturable unit that will be suitable for the widest possible range of applications.

The temperature capabilities of CPCs with nonevacuated absorbers are somewhat more limited. However, with proper design, taking care to reduce parasitic thermal losses through the reflectors, nonevacuated CPCs will outperform the best available flat-plate collectors at temperatures above about $50\text{--}70^\circ\text{C}$.

Nonimaging, near-ideal secondary or terminal concentrators have advantages for any solar concentrating application. In the near term, they are being applied to increase the angular tracking, alignment, and slope definition requirements for the primaries in an effort to reduce system cost and complexity. In future applications, very high geometric concentrations that cannot otherwise be attained can be achieved through the use of these devices.

ACKNOWLEDGMENT

Work described in this review was supported in large part by the U.S. Department of Energy, in particular the Engineering Research Program of the Office of Basic Energy Science and the National Renewable Energy Laboratory.

SEE ALSO THE FOLLOWING ARTICLES

ENERGY FLOWS IN ECOLOGY AND IN THE ECONOMY • FULLERENES AND CARBON NANOTUBES • HOLOGRAPHY • LASERS • PHOTOVOLTAIC SYSTEM DESIGN • SOLAR THERMAL POWER STATIONS

BIBLIOGRAPHY

- Cooke, D., Gleckman, P., Krebs, H., O'Gallagher, J., Sagie, D., and Winston, R. (1990). "Brighter than the sun," *Nature* **346**, 802.
- Friedman, R., Gordon, J., and Ries, H. (1993). "New high-flux two-stage optical designs for parabolic solar concentrators," *Solar Energy* **51**, 317–325.
- Gleckman, P., O'Gallagher, J., and Winston, R. (1989). "Concentration of sunlight to solar-surface levels using non-imaging optics," *Nature* **339**, 198.
- Gordon, J., and Rabl, A. (1992). "Nonimaging CPC-type reflectors with variable extreme direction," *Appl. Opt.* **31**, 7332–7338.
- Jenkins, D., and Winston, R. (1996). "Integral design method of non-imaging Optics," *J. Opt. Soc. Am. A*, **13**, 2106–2116.
- Jenkins, D., O'Gallagher, J., and Winston, R. (1997). "Attaining and using extremely high intensities of solar energy with nonimaging concentrators, advances in solar energy," **11**, 43–108, *Am. Solar Energy Soc.*, Boulder, CO.
- Jenkins, D., Winston, R., Bliss, J., O'Gallagher, J., Lewandowski, A., and Bingham, C. (1996a). "Solar concentration of 50,000 achieved with output power approaching 1 kW," *J. Sol. Energy Eng.*

- Lewandowski, A., Bingham, C., O'Gallagher, J., Winston, R., and Sagie, D. (1991). "Performance characterization of the SERI high flux solar furnace," *Solar Energy Materials*, **24**, 550–563.
- Ning, X., O'Gallagher, J., and Winston, R. (1987). "The optics of two-stage photovoltaic concentrators with dielectric second stages," *Appl. Opt.* **26**, 1207.
- O'Gallagher, J., and Winston, R. (1983). "Development of compound parabolic concentrators for solar energy," *J. Ambient. Energy* **4**, 171.
- O'Gallagher, J., and Winston, R. (1986). "Test of a "trumpet" secondary concentrator with a paraboloidal dish primary," *Solar Energy* **36**, 37.
- O'Gallagher, J., Welford, W. T., and Winston, R. (1986). "Axially symmetric nonimaging flux concentrators with the maximum theoretical concentration ratio," *J. Opt. Soc. Am.* **4**, 66.
- Rabl, A., O'Gallagher, J., and Winston, R. (1980). "Design and test of non-evacuated solar collectors with compound parabolic concentrators," *Solar Energy* **25**, 335.
- Ries, H., and Winston, R. (1994). "Tailored edge-ray reflectors for illumination," *J. Opt. Soc. Am. A*, **11**, 1260–1264.
- Shatz, N., and Bortz, J. (1995). "Inverse engineering perspective on non-imaging optical designs," *Nonimaging Optics: Maximum Efficiency Light Transfer III*, *Proc. SPIE*, **2538**, 136–156.
- Snail, K. A., O'Gallagher, J. J., and Winston, R. (1984). "A stationary evacuated collector with integrated concentrator," *Solar Energy* **33**, 441.
- Welford, W. T., and Winston, R. (1982). "Conventional optical systems and the brightness theorem," *Appl. Opt.* **21**, 1531.
- Welford, W. T., and Winston, R. (1989). "High Collection Nonimaging Optics," Academic Press, San Diego.
- Winston, R. (ed.), (1995). "Selected papers on nonimaging optics," *SPIE Milestone Series*, **MS 106**, SPIE Opt. Engineering Press, Bellingham, WA.
- Winston, R. (ed.) (1999). "Nonimaging optics: Maximum efficiency light transfer," **V**, *Proc. SPIE*, **3781**. Previous sessions are I, Proceedings SPIE 1528 (1991); II, Proceedings SPIE 2016 (1993); III Proceedings SPIE 2538 (1993); IV Proceedings SPIE 3139 (1997).
- Winston, R., and Ries, H. (1993), "Nonimaging reflectors as functionals of the desired irradiance," *J. Opt. Soc. Am.* **10**, 1902–1908.



Nonlinear Optical Processes

John F. Reintjes

Naval Research Laboratory

- I. Optical Nonlinearities—Physical Origins
- II. Optical Nonlinearities—Mathematical Description
- III. Specific Nonlinear Processes
- IV. Applications

GLOSSARY

Anti-Stokes shift Difference in frequency between the pump wave and the generated wave in a stimulated scattering interaction when the generated wave is at a higher frequency than the pump wave.

Brillouin shift Difference in frequency between the incident wave and the scattered wave in a stimulated Brillouin interaction. It is equal to the sound wave frequency.

Coherence length Distance required for the phase of a light wave to change by $\pi/2$ relative to that of its driving polarization. Maximum conversion in parametric processes is obtained at L_{coh} when $\Delta k \neq 0$.

Constitutive relations Set of equations that specify the dependence of induced magnetizations or polarizations on optical fields.

Depletion length Distance required for significant pump depletion in a frequency-conversion interaction.

Mode-locked laser Laser that operates on many longitudinal modes, each of which is constrained to be in phase with the others. Mode-locked lasers produce pulses with durations of the order of 0.5–30 psec.

Nonlinear polarization Electric dipole moment per unit volume that is induced in a material by a light wave

and depends on intensity of the light wave raised to a power greater than unity.

Parametric interactions Interactions between light waves and matter that do not involve transfer of energy to or from the medium.

Phase conjugation Production of a light wave with the phase variations of its wave front reversed relative to those of a probe wave.

Phase matching Act of making the wave vector of an optical wave equal to that of the nonlinear polarization that drives it.

Population inversion Situation in which a higher-energy level in a medium has more population than a lower-energy one.

Q-switched laser Laser in which energy is stored while the cavity is constrained to have large loss (low Q) and is emitted in an intense short pulse after the cavity Q is switched to a high value reducing the loss to a low value.

Raman shift Difference in frequency between the incident wave and the scattered wave in a stimulated Raman interaction. It is equal to the frequency of the material excitation involved in the interaction.

Slowly varying envelope approximation Approximation in which it is assumed that the amplitude of an

optical wave varies slowly in space compared to a wavelength and slowly in time compared to the optical frequency.

Stimulated scattering Nonlinear frequency conversion interaction in which the generated wave has exponential gain and energy is transferred to or from the nonlinear medium.

Stokes shift Difference in frequency between the pump wave and the generated wave in a stimulated scattering interaction when the generated wave is at lower frequency than is the pump wave.

Susceptibility, n th order Coefficient in a perturbation expansion that gives the ratio of the induced polarization of order n to the n th power of the electric or magnetic fields.

Wave number (cm^{-1}) Number of optical waves contained in 1 cm. It is equal to the reciprocal of the wavelength in centimeters and has the symbol $\bar{\nu}$. It is commonly used as a measure of the energy between quantum levels of a material system or the energy of the photons in an optical wave and is related to the actual energy by $E = h\bar{\nu}c$, where h is Planck's constant and c is the speed of light.

Wave vector Vector in the direction of propagation of an optical wave and having magnitude equal to 2π divided by the wavelength.

Wave-vector mismatch Difference between the wave vector of an optical wave and that of the nonlinear polarization that drives it. It is also called the phase mismatch.

Zero-point radiation Minimum energy allowed in optical fields due to quantum-mechanical effects.

NONLINEAR OPTICS is a field of study that involves a nonlinear response of a medium to intense electromagnetic radiation. Nonlinear optical interactions can change the propagation or polarization characteristics of the incident waves, or they can involve the generation of new electromagnetic waves, either at frequencies different from those contained in the incident fields or at frequencies that are the same but with the waves otherwise distinguishable from the incident waves—for example, by their direction of polarization or propagation. Nonlinear optical interactions can be used for changing or controlling certain properties of laser radiation, such as wavelength, bandwidth, pulse duration, and beam quality, as well as for high-resolution molecular and atomic spectroscopy, materials studies, modulation of optical beams, information processing, and compensation of distortions caused by imperfect optical materials. Nonlinear optical interactions are generally observed in the spectral range covered by lasers, between the far infrared and the extreme

ultraviolet (XUV) (Fig. 1), but some nonlinear interactions have been observed at wavelengths ranging from X rays to microwaves. They generally require very intense optical fields and as a result are usually observed only with radiation from lasers, although some nonlinear interactions, such as saturation of optical transitions in an atomic medium, can occur at lower intensities and were observed before the invention of the laser. Nonlinear optical interactions of various kinds can occur in all types of materials, although some types of interactions are observable only in certain types of materials, and some materials are better suited to specific nonlinear interactions than others.

I. OPTICAL NONLINEARITIES— PHYSICAL ORIGINS

When an electromagnetic wave propagates through a medium, it induces a polarization (electric dipole moment per unit volume) and magnetization (magnetic dipole moment per unit volume) in the medium as a result of the motion of the electrons and nuclei in response to the fields in the incident waves. These induced polarizations and magnetizations oscillate at frequencies determined by a combination of the properties of the material and the frequencies contained in the incident light waves. The optical properties of the medium and the characteristics of the radiation that is transmitted through it result from interference among the fields radiated by the induced polarizations or magnetizations and the incident fields.

At low optical intensities, the induced polarizations and magnetizations are proportional to the electric or magnetic fields in the incident wave, and the response of the medium is termed linear. Various linear optical interactions can occur, depending on the specific properties of the induced polarizations. Some of the more familiar linear optical effects that result from induced polarizations that oscillate at the same frequency as the incident radiation are refraction, absorption, and elastic scattering (Rayleigh scattering from static density variations, or Tyndall or Mie scattering). Other linear optical processes involve inelastic scattering, in which part of the energy in the incident wave excites an internal motion of the material and the rest is radiated in a new electromagnetic wave at a different frequency. Examples of inelastic scattering processes are Raman scattering, which involves molecular vibrations or rotations, electronic states, lattice vibrations, or electron plasma oscillations; Brillouin scattering, which involves sound waves or ion-acoustic waves in plasmas; and Rayleigh scattering, involving diffusion, orientation, or density variations of molecules. Although these inelastic scattering processes produce electromagnetic waves at

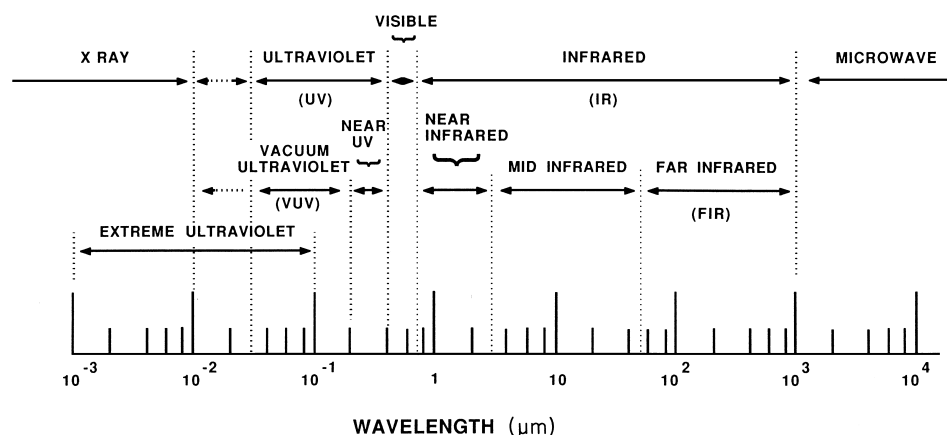


FIGURE 1 Designations of various regions of the electromagnetic spectrum between microwaves and X rays. Dashed arrows indicate inexact boundaries between spectral regions.

frequencies different from those in the incident waves, they are linear optical processes when the intensity of the scattered wave is proportional to the intensity of the incident wave.

When the intensity of the incident radiation is high enough, the response of the medium changes qualitatively from its behavior at low intensities, giving rise to the nonlinear optical effects. Some nonlinear optical interactions arise from the larger motion of the electrons and ions in response to the stronger optical fields. In most materials, the electrons and ions are bound in potential wells that, for small displacements from equilibrium, are approximately harmonic (i.e., have a potential energy that depends on the square of the displacement from equilibrium), but are anharmonic (i.e., have a potential energy that has terms that vary as the third or higher power of the displacement from equilibrium) for larger displacements, as shown in Fig. 2.

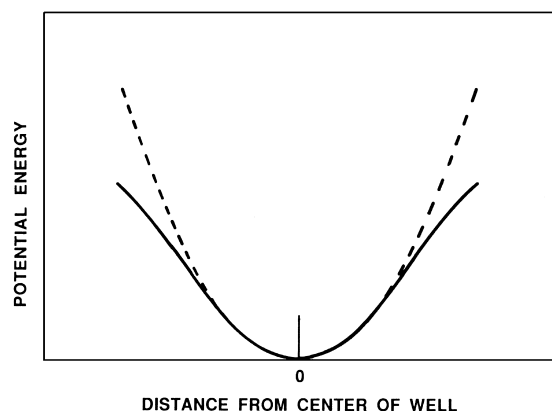


FIGURE 2 Schematic illustration of an anharmonic potential well (solid line) that can be responsible for nonlinear optical interactions. A harmonic potential well that does not result in nonlinear interactions is shown for comparison (dashed line).

As long as the optical intensity is low, the electron or ion moves in the harmonic part of the well. In this regime, the induced polarizations can oscillate only at frequencies that are contained in the incident waves and the response is linear. When the incident intensity is high enough, the charges can be driven into the anharmonic portion of the potential well. The additional terms in the potential introduce terms in the induced polarization that depend on the second, third, or higher powers of the incident fields, giving rise to nonlinear responses. Examples of nonlinear optical processes that occur in this manner are various forms of harmonic generation and parametric frequency mixing.

A second type of nonlinear response results from a change in some property of the medium caused by the optical wave, which in turn affects the propagation of the wave. The optical properties of a material are usually described in the limit of extremely weak optical fields and are therefore considered to be intrinsic properties of the medium. When the optical field is strong enough, however, it can change certain characteristics of the medium, which in turn changes the way the medium affects the optical wave, resulting in a nonlinear optical response. An example of such a response is a change in the refractive index of a medium induced by the optical wave. Such changes can occur, for example, because of orientation of anisotropic molecules along the incident fields, or because of changes in the density of the medium as a result of electrostriction or as a result of a temperature change following absorption of the incident wave. Many of the propagation characteristics of optical waves are determined by the refractive index of the medium. In the situations just described, the refractive index depends on the intensity of the optical wave, and, as a result, many of the propagation characteristics depend on the optical intensity as well.

II. OPTICAL NONLINEARITIES— MATHEMATICAL DESCRIPTION

A. Maxwell's Equations for Nonlinear Materials

The propagation of an optical wave in a medium is described by the Maxwell equation for the optical field, including the effects of the induced polarizations and magnetizations, and a specification of the dependence of the induced polarizations and magnetizations on the optical fields. The Maxwell equation for the electric field of an optical wave in a medium is

$$\nabla^2 E(r, t) - (1/c^2) \partial^2 E(r, t) / \partial t^2 = \mu_0 \partial^2 P(r, t) / \partial t^2 + \partial[\nabla \times M(r, t)] / \partial t. \quad (1)$$

Here $E(r, t)$ is the electric field of the optical wave, and the terms on the right are the induced polarizations (P) and magnetizations (M) that describe the effects of any charge distribution that may be present in the propagation path. In vacuum these terms are zero, and in this situation Eq. (1) reduces to the familiar expression for vacuum propagation of electromagnetic radiation.

B. Nonlinear Polarizations

The problem of propagation in a medium is completely specified when the relations between the polarization and the magnetization and the optical fields, termed constitutive relations, are given. In general these relations can be quite complicated. However, for many situations encountered in the laboratory a set of simplifying assumptions can be made that lead to an expansion of the induced polarization in a power series in the electric field of the light wave of the form

$$P = \epsilon_0 \chi^{(1)} E + \epsilon_0 \chi^{(2)} E^2 + \epsilon_0 \chi^{(3)} E^3 + \dots + \epsilon_0 Q^{(1)} \nabla E + \epsilon_0 Q^{(2)} \nabla E^2 + \dots \quad (2a)$$

For magnetic materials a similar expansion can be made for the magnetization:

$$M = m^{(1)} H + m^{(2)} H^2 + m^{(3)} H^3 + \dots \quad (2b)$$

The types of expression in Eqs. (2a) and (2b) are generally valid when the optical fields are weak compared to the electric field that binds the electrons in the material and when the coefficients of the various terms in Eqs. (2a) and (2b) are constant over the range of frequencies contained in the individual incident and generated fields. In addition, the wavelength of the radiation must be long compared to the dimension of the scattering centers (the atoms and molecules of the nonlinear medium), so that the charge distributions can be accounted for with a multipole expansion.

When the first of these conditions is not met, as for example with extremely intense radiation, the perturbation series will not converge and all powers of the incident fields must be used. When the second of the conditions is not met, as, for example, can happen when certain resonance conditions are satisfied, each term in the response of the medium will involve convolutions of the optical fields instead of simply powers of the fields multiplied by constants. In these situations the polarizations induced in the medium must be solved for as dynamical variables along with the optical fields.

The coefficients in the various terms in Eqs. (2a) and (2b) are termed the n th-order susceptibilities. The first-order susceptibilities describe the linear optical effects, while the remaining terms describe the n th order nonlinear optical effects. The coefficients $\chi^{(n)}$ are the n th-order electric dipole susceptibilities, the coefficients $Q^{(n)}$ are the n th-order quadrupole susceptibilities, and so on. Similar terminology is used for the various magnetic susceptibilities. For most nonlinear optical interactions, the electric dipole susceptibilities are the dominant terms because the wavelength of the radiation is usually much longer than the scattering centers. These will be the ones considered primarily from now on.

In some situations, however, the electric quadrupole susceptibilities can make a significant contribution to the response. For example, symmetry restrictions in certain classes of materials can prevent some nonlinear processes from occurring by dipole interactions, leaving the quadrupole interactions as the dominant response. This occurs, for example, with second-harmonic generation in media with inversion symmetry. Other situations in which quadrupole interactions are significant are those in which they are enhanced by resonances with appropriate energy levels or for interactions involving sufficiently short-wavelength X radiation.

The nonlinear susceptibilities are tensors and, as such, relate components of the nonlinear polarization vector to various components of the optical field vectors. For example, the second-order polarization is given by

$$P_i = \epsilon_0 \chi_{ijk}^{(2)} E_j E_k,$$

where i , j , and k refer to the spatial directions x , y , and z . The susceptibility tensors χ have the symmetry properties of the nonlinear medium. As a result they can restrict the combinations of vector components of the various optical fields that can be used effectively. In some situations, such as those involving nonlinear optical processes that change the polarization vector of the optical wave, the tensor properties play a central role in the nonlinear interaction. In other situations, however, the tensor properties are important only in determining which combinations of vector components can, or must, be used for the optical

fields and the nonlinear polarizations, and beyond that the nonlinear optical susceptibilities can usually be treated as scalars.

The magnitudes of the nonlinear susceptibilities are usually determined by measurement, although in some cases, most notably for third- and higher-order interactions in atomic gases or second-order interactions in some crystals, they can be calculated using various theories. In some cases they can be estimated with varying degrees of accuracy from products of the linear refractive indices at the various wavelengths involved.

C. Calculation of Optical Fields

In order to calculate the optical fields involved in the nonlinear interactions, we assume that they can be expressed in the form

$$E(r, z, t) = \frac{1}{2} [A(r, z, t)e^{i(kz - \omega t)} + \text{complex conjugate}]. \quad (3)$$

This expression describes a wave of amplitude A propagating in the z direction with a frequency ω and a wave vector $k = 2\pi n/\lambda$, where λ is the wavelength of the light in vacuum, and n is the linear refractive index, which takes into account the effects of the linear susceptibilities. Many nonlinear effects involve the interaction of two or more optical fields at different wavelengths. In such situations the optical field is written as

$$E(r, z, t) = \frac{1}{2} \sum [A_i(r, z)e^{i(k_i n_i z - \omega_i t)} + \text{complex conjugate}]. \quad (4)$$

Here the index for the summation extends over all fields in the problem, including those fields that are generated in the nonlinear interaction as well as those that are incident on the medium.

The nonlinear polarization is written as

$$P(r, z, t) = \frac{1}{2} \sum [P_i(r, z)e^{i(k_i^p z - \omega_i t)} + \text{complex conjugate}]. \quad (5)$$

Here P_i is the amplitude of the nonlinear polarization with frequency ω_i . It is determined by an appropriate combination of optical fields and nonlinear susceptibilities according to the expansion in Eq. (2a). It serves as a source term for the optical field with frequency ω_i . The wave-vector of the nonlinear polarization k_i^p is in general different from the wave-vector of the optical field at the same frequency. This difference plays a very important role in determining the effectiveness of many nonlinear optical interactions.

The amplitudes of the various waves can be calculated from the expressions given in Eqs. (1) and (2), using the fields and polarizations with the forms given in Eqs. (4) and (5). In deriving equations for the field amplitudes, a further simplifying assumption, the slowly varying envelope approximation, is also usually made. This approximation assumes that the field envelopes $A_i(r, z, t)$ vary slowly in time compared to the optical frequency and slowly in space compared to the optical wavelength. It is generally valid for all of the interactions that are encountered in the laboratory. It enables us to neglect the second derivatives of the field amplitudes with respect to z and t when Eqs. (4) and (5) are substituted into Eq. (1). When this substitution is made, we obtain expressions of the following form for the various fields involved in a nonlinear interaction:

$$\nabla_{\perp}^2 A_1 + 2ik_1[\partial A_1/\partial z + (n_1/c)\partial A_1/\partial t] = -\mu_0\omega_1^2 P_1 e^{-i\Delta k_1 z} \quad (6a)$$

$$\nabla_{\perp}^2 A_2 + 2ik_2[\partial A_2/\partial z + (n_2/c)\partial A_2/\partial t] = -\mu_0\omega_2^2 P_2 e^{-i\Delta k_2 z} \quad (6b)$$

\vdots

$$\nabla_{\perp}^2 A_n + 2ik_n[\partial A_n/\partial z + (n_n/c)\partial A_n/\partial t] = -\mu_0\omega_n^2 P_n e^{-i\Delta k_n z}. \quad (6c)$$

Here A_i is the amplitude of the i th optical field, which may be either an incident field or a field generated in the interaction, and Δk_i is the wave-vector mismatch between the i th optical field and the polarization that drives it and is given by

$$\Delta k_i = k_i - k_i^p. \quad (7)$$

In nonlinear interactions that involve the generation of fields at new frequencies, the nonlinear polarization for the generated field will involve only incident fields. In many cases the interactions will be weak and the incident fields can be taken as constants. In other situations, the generated fields can grow to be sufficiently intense that their nonlinear interaction on the incident fields must be taken into account. In still other nonlinear interactions, such as those involving self action effects, an incident field amplitude can occur in its own nonlinear polarization. The set of Eq. (6) along with the constitutive relations in Eq. (2) can be used to calculate the optical fields in most nonlinear interactions.

D. Classifications of Nonlinear Interactions

Nonlinear interactions can be classified according to whether they are driven by electric or magnetic fields, whether they involve the generation of fields at new frequencies or change the propagation properties of the incident fields, and whether or not they involve a transfer

of energy to or from the nonlinear medium. To a certain extent these classifications are overlapping, and examples of each type of interaction can be found in the other categories.

A list of the more common nonlinear interactions is given in Table I. Most of the commonly observed nonlinear optical interactions are driven by the electric fields. They can proceed through multipole interactions of any order, but the dipole interactions are most common unless certain symmetry or resonance conditions are satisfied in the medium. The interactions are classified first according to whether they involve frequency conversion (that is, the generation of a wave at a new frequency) or self-action (that is, whether they affect the propagation of the incident wave). The frequency-conversion interactions are further classified according to whether they are parametric processes or inelastic processes. Parametric processes do not involve a transfer of energy to or from the nonlinear material. The material merely serves as a medium in which the optical waves can exchange energy among themselves. On the other hand, inelastic frequency conversion processes, generally termed stimulated scattering, do involve a transfer of energy to or from the nonlinear medium.

Self-action effects involve changes in the propagation characteristics of a light wave that are caused by its own intensity. They can involve both elastic and inelastic processes and can affect almost any of the propagation characteristics of the light wave, including both absorption and focusing properties. They can also change the spatial, temporal, and spectral distributions of the incident wave, as well as its state of polarization.

Coherent effects involve interactions that occur before the wave functions that describe the excitations in the nonlinear medium have time to get out of phase with one another. They can involve changes in propagation characteristics as well as the generation of new optical signals.

Electrooptic and magneto-optic effects involve changes in the refractive index of the medium caused by external electric or magnetic fields resulting in changes in the phase of the optical wave or in its state of polarization.

III. SPECIFIC NONLINEAR PROCESSES

A. Frequency–Conversion Processes

Frequency conversion processes are those that involve the generation of radiation at wavelengths other than the ones that are contained in the incident radiation. A typical frequency conversion geometry is illustrated in Fig. 3. Pump radiation, consisting of one or more optical waves at one or more frequencies, is incident on a nonlinear medium and interacts with it to generate the new wave. Depending on the interaction, the generated wave can be at a lower

or higher frequency than the waves in the pump radiation, or, in some situations, can be an amplified version of one of the incident waves.

Frequency–conversion interactions can either be parametric processes or stimulated scattering interactions. In parametric frequency conversion, the incident wave generates a new wave at a frequency that is a multiple, or harmonic, of the incident frequency, or, if more than one frequency is present in the incident radiation, a sum or difference combination of the frequencies of the incident fields. Parametric frequency conversion can also involve the generation of radiation in two or more longer-wavelength fields whose frequencies add up to the frequency of the incident field. In parametric frequency conversion, energy is conserved among the various optical waves, with the increase in energy of the generated wave being offset by a corresponding decrease in energy of the incident waves.

Parametric frequency-conversion interactions can occur in any order of the perturbation expansion of Eq. (2). The most commonly observed processes have involved interactions of second and third order, although interactions up to order 11 have been reported. Parametric interactions are characterized by a growth rate for the intensity of the generated wave that depends on a power, or a product of powers, of the intensities of the incident waves, and they are strongly dependent on difference in wave vectors between the nonlinear polarization and the optical fields.

Stimulated scattering generally involves in elastic nonlinear frequency-conversion processes. They are nonlinear counterparts of the various linear inelastic scattering processes that were mentioned earlier. In these processes, an incident wave at frequency ω_{inc} is scattered into a wave with a different frequency, ω_{scat} , with the difference in energy between the incident and scattered photons taken up by excitation or deexcitation of an internal mode of the material. When the internal mode is excited in the nonlinear process, the scattered wavelength is longer than the incident wavelength, while it is shorter than the incident wavelength when the internal mode is deexcited.

Various types of stimulated scattering processes can occur, including stimulated Raman, stimulated Brillouin, and stimulated Rayleigh scattering, each of which involves a different type of internal mode of the medium, as will be described in Section III.A.2. The optical wave generated in stimulated scattering processes is characterized by exponential growth similar to the exponential growth experienced by laser radiation in the stimulated emission process.

Frequency–conversion interactions are used primarily to generate coherent radiation at wavelengths other than those that can be obtained directly from lasers, or to

TABLE I Common Nonlinear Optical Interactions

Nonlinear process	Description	Incident radiation	Generated radiation
Frequency conversion			
Parametric processes			
q th-Harmonic generation	Conversion of radiation from ω to $q\omega$	ω	$q\omega$
Sum- or difference-frequency mixing	Conversion of radiation at two or more frequencies to radiation at a sum- or difference-frequency combination	Three-wave (second order) ω_1, ω_2 Four-wave (third order) $\omega_1, \omega_2,$ $\omega_1, \omega_2, \omega_3$ Six-wave (fifth order) ω_1, ω_2	$\omega_g = \omega_1 \pm \omega_2$ $\omega_g = 2\omega_1 \pm \omega_2$ $\omega_g = \omega_1 \pm \omega_2 \pm \omega_3$ $\omega_g = 4\omega_1 \pm \omega_2$
Parametric down-conversion	Conversion of radiation at frequency ω_1 to two or more lower frequencies	Three-wave (second order) ω_1 Four-wave (third order) ω_1	$\omega_{g1} + \omega_{g2} = \omega_1$ $\omega_{g1} + \omega_{g2} = 2\omega_1$
Optical rectification	Conversion of radiation at frequency ω_1 to an electric voltage	ω_1	$\omega_g = 0 = \omega_1 - \omega_1$
Inelastic processes			
Stimulated scattering	Conversion of radiation at frequency ω_1 to radiation at frequency ω_2 with excitation or deexcitation of an internal mode at frequency ω_0	Stokes scattering ω_1, ω_2 Anti-Stokes scattering ω_1, ω_2	$\omega_2 = \omega_1 - \omega_0$ $\omega_2 = \omega_1 + \omega_0$
Self-action effects			
Self-focusing	Focusing of an optical beam due to a change in the refractive index caused by its own intensity		
Self-defocusing	Defocusing of an optical beam due to a change in the refractive index caused by its own intensity		
Self-phase modulation	Modulation of the phase of an optical wave due to a time-dependent change in the refractive index caused by its own intensity		
Optical Kerr effect	Birefringence in a medium caused by an intensity-induced change in the refractive index		
Ellipse rotation	Intensity-dependent rotation of the polarization ellipse of a light wave due to the nonlinear refractive index		
Raman-induced Kerr effect	Change of refractive index or birefringence of a light wave at one wavelength due to the intensity of a light wave at another wavelength		
Multiphoton absorption	Increase in the absorption of a material at high intensities		
Saturable absorption	Decrease in the absorption of a material at high intensities		
Coherent effects			
Self-induced transparency	High transmission level of a medium at an absorbing transition for a short-duration light pulse with the proper shape and intensity		
Photon echo	Appearance of a third pulse radiated from an absorbing medium following application of two pulses of proper duration, intensity, and separation		
Adiabatic rapid passage	Inversion of population of a medium by application of a light pulse whose frequency is swept quickly through that of an absorbing transition		
Electrooptic effects			
Pockels effect	Birefringence in a medium caused by an electric field and proportional to the electric field strength		
Kerr effect	Birefringence in a medium caused by an electric field and proportional to the square of the electric field strength		
Magneto optic effects			
Faraday rotation	Rotation of the direction of linear polarization of a light beam resulting from changes in the relative velocities of oppositely circular polarizations caused by a magnetic field		
Cotton–Mouton effect	Birefringence induced in a material caused by a magnetic field		
Miscellaneous			
Optical breakdown	Rapid ionization of a material by avalanche production of electrons caused by an intense optical field		

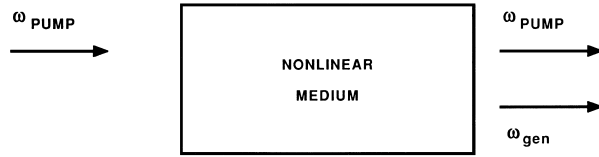


FIGURE 3 Schematic illustration of a typical geometry used in frequency-conversion interactions. Input radiation is supplied at one or more incident frequencies ω_{pump} . The wave at ω_{gen} is generated in the nonlinear interaction.

amplify existing laser radiation. Various frequency-conversion processes have been used to generate coherent radiation at wavelengths ranging from the extreme ultraviolet to the millimeter wave region. The radiation generated through nonlinear frequency-conversion processes has all of the special properties usually associated with laser radiation, such as spatial and temporal coherence, collimation, and narrow bandwidths, and can be tunable if the pump radiation is tunable. The generated radiation can have special properties not possessed by existing laser radiation in a given wavelength range, such as tunability, short pulse duration, or narrow bandwidth, or it can be in a spectral range in which there is no direct laser radiation available. An example of the first is the generation of tunable radiation in spectral regions in the vacuum ultraviolet or the infrared where only fixed-frequency lasers exist. An example of the second is the generation of coherent radiation in the extreme ultraviolet at wavelengths less than 100 nm, a spectral region in which lasers of any kind are just beginning to be developed.

Various frequency-conversion interactions can also be used for spectroscopy, image conversion, control of various properties of laser beams, and certain optical information-processing techniques.

1. Parametric Frequency-Conversion Processes

Parametric frequency-conversion processes include harmonic conversion of various order, various forms of frequency mixing, parametric down-conversion, and optical rectification. They are used primarily for generation of radiation at new wavelengths, although some of the interactions can be used for amplifying existing radiation. Various forms of parametric frequency-conversion interactions have been used to generate coherent radiation at wavelengths ranging from almost the X-ray range at 35.3 nm to the microwave range at 2 mm. Parametric frequency-conversion interactions are usually done in the transparent region of the nonlinear medium. Different interactions, nonlinear materials, and lasers are used to generate radiation in the various wavelength ranges. Some of the more commonly used interactions, types of lasers, and types of nonlinear materials in various wavelength ranges are given in [Table II](#).

a. Second-order effects. Second-order effects are primarily parametric in nature. They include second-harmonic generation, three-wave sum- and difference-frequency mixing, parametric down-conversion, parametric oscillation, and optical rectification. The strongest second-order processes proceed through electric dipole interactions. Because of symmetry restrictions, the even-order electric dipole susceptibilities are zero in materials with inversion symmetry. As a result, the second-order nonlinear interactions are observed most commonly only in certain classes of crystals that lack a center of inversion. However, some second-order processes due to electric quadrupole interactions have been observed in solids (e.g., sodium chloride) and gases that have a center of inversion. In addition, some second-order processes, such as second-harmonic generation, have also been observed in gases in which the inversion symmetry has been lifted with electric or magnetic fields. In these situations the interactions are actually third-order ones in which the frequency of one of the waves is zero.

The nonlinear polarization for the second-order processes can be written as

$$P_{\text{NL}}^{(2)} = 2\epsilon_0 d E^2, \quad (8)$$

where d is a nonlinear optical susceptibility related to $\chi^{(2)}$ in Eq. (2) by $d = \chi^{(2)}/2$ and E is the total electric field as given in Eq. (4). The right-hand side of Eq. (8) has terms that oscillate at twice the frequency of the incident waves and at sum- and difference-frequency combinations and terms that do not oscillate. Each of these terms gives rise to a different nonlinear process. In general they are all present in the response of the medium. The term (or terms) that forms the dominant response of the medium is usually determined by phase matching, as will be discussed later.

The nonlinear polarization amplitude as defined in Eq. (5) can be evaluated for each of the specific nonlinear processes by substituting the form of the electric field in Eq. (4) into Eq. (8) and identifying each term by the frequency combination it contains. The resulting form of the nonlinear polarization for the sum-frequency process $\omega_3 = \omega_1 + \omega_2$ is

$$P_i^{\text{NL}}(\omega_3) = g\epsilon_0 d_{ijk}(-\omega_3, \omega_1, \omega_2) A_j(\omega_1) A_k(\omega_2). \quad (9)$$

Here i, j, k refer to the crystallographic directions of the nonlinear medium and the vector components of the various fields and polarizations, and g is a degeneracy factor that accounts for the number of distinct permutations of the pump fields.

Similar expressions can be obtained for the other second-order processes with minor changes. For second-harmonic generation, $\omega_1 = \omega_2$, while for difference-frequency mixing processes of the form $\omega_3 = \omega_1 - \omega_2$, ω_2 appears with a minus sign and the complex conjugate

TABLE II Wavelength Ranges for Various Nonlinear Parametric Interactions

Nonlinear interaction	Wavelength range	Typical laser	Type of material
Second-harmonic, three-wave sum-frequency mixing	2–5 μm	CO ₂ , CO	Crystal without inversion center
	Visible	Nd:YAG, Nd:glass	Crystal without inversion center
	Ultraviolet (to 185 nm)	Harmonics of Nd, ruby, dye, argon	Crystal without inversion center
Three-wave difference-frequency mixing	Visible to infrared (25 μm)	CO, CO ₂ , Nd, dye, Parametric oscillator, spin-flip Raman	Crystal without inversion center
Parametric oscillation	Visible to infrared	Nd, dye, CO ₂	Crystal without inversion center
Third harmonic, four wave sum-frequency mixing, four-wave difference-frequency mixing	Infrared (3.3 μm) to XUV (57 nm)	CO ₂ , CO, Nd, ruby, dye, excimer	Rare gases, molecular gases, metal vapors, cryogenic liquids
Fifth- and higher-order harmonic and frequency mixing	UV (216 nm)	Nd, harmonics of Nd, excimer	Rare gases, metal vapors
	XUV (106–35.3 nm)		

of A_2 is used. The degeneracy factor g is equal to 1 for second-harmonic generation and 2 for the second-order mixing processes.

Because the susceptibility tensor d has the same symmetry as the nonlinear medium, certain of its components will be zero, and others will be related to one another, depending on the material involved. As a result, only certain combinations of polarization vectors and propagation directions can be used for the incident and generated waves in a given material. For a specific combination of tensor and field components, the nonlinear polarization amplitude can be written as

$$P^{\text{NL}}(\omega_3) = g\varepsilon_0 d_{\text{eff}}(-\omega_3, \omega_1, \omega_2)A(\omega_1)A(\omega_2), \quad (10)$$

where d_{eff} is an effective susceptibility that takes into account the nonzero values of the d tensor and the projections of the optical field components on the crystallographic directions and $A(\omega_i)$ is the total optical field at ω_i . The nonlinear polarization amplitudes for the other second-order processes have similar expressions with the changes noted above. The forms for the polarization amplitude

for the various second-order interactions are given in Table III.

The k vector of the nonlinear polarization is given by

$$k_3^p = k_1 \pm k_2, \quad (11)$$

where the plus sign is used for sum-frequency combinations and the minus sign is used for difference-frequency combinations. The wave-vector mismatch of Eq. (7) is also given in Table III for each nonlinear process. The nonlinear polarization amplitudes and the appropriate wave-vector mismatches as given in Table III can be used in Eq. (6) for the field amplitudes to calculate the intensities of the generated waves in specific interactions.

The second-order frequency-conversion processes are used primarily to produce coherent radiation in wavelength ranges where radiation from direct laser sources is not available or to obtain radiation in a given wavelength range with properties, such as tunability, that are not available with existing direct laser sources. They have been used with various crystals and various types of lasers to generate radiation ranging from 185 nm in the ultraviolet

TABLE III Second-Order Nonlinear Polarizations

Nonlinear process	Nonlinear polarization	Wave-vector mismatch	Plane-wave phase-matching condition
Second-harmonic generation	$P(2\omega) = \varepsilon_0 d_{\text{eff}}(-2\omega, \omega, \omega)A^2(\omega)$	$\Delta k = k(2\omega) - 2k(\omega)^a$	$n(2\omega) = n(\omega)^a$
Three-wave sum-frequency mixing	$P(\omega_3) = 2\varepsilon_0 d_{\text{eff}}(-\omega_3, \omega_1, \omega_2) \times A(\omega_1)A(\omega_2)$	$\Delta k = k(2\omega) - k_1(\omega) - k_2(\omega)^b$	$n(2\omega) = n_1(\omega)/2 + n_2(\omega)/2^b$
		$\Delta k = k(\omega_3) - k(\omega_1) - k(\omega_2)$	$n(\omega_3)/\lambda_3 = n(\omega_1)/\lambda_1 + n(\omega_2)/\lambda_2$
Three-wave difference-frequency mixing	$P(\omega_3) = 2\varepsilon_0 d_{\text{eff}}(-\omega_3, \omega_1, -\omega_2) \times A(\omega_1)A^*(\omega_2)$	$\Delta k = k(\omega_3) - k(\omega_1) + k(\omega_2)$	$n(\omega_3)/\lambda_3 = n(\omega_1)/\lambda_1 - n(\omega_2)/\lambda_2$
Parametric down-conversion	$P(\omega_3) = 2\varepsilon_0 d_{\text{eff}}(-\omega_3, \omega_1, -\omega_2) \times A(\omega_1)A^*(\omega_2)$	$\Delta k = k(\omega_3) + k(\omega_2) - k(\omega_1)$	$n(\omega_1)/\lambda_1 = n(\omega_2)/\lambda_2 + n(\omega_3)/\lambda_3$
	$P(\omega_2) = 2\varepsilon_0 d_{\text{eff}}(-\omega_2, \omega_1, -\omega_3) \times A(\omega_1)A^*(\omega_3)$		
Optical rectification	$P(0) = 2\varepsilon_0 d_{\text{eff}}(0, \omega, -\omega)A(\omega_1)A^*(\omega_1)$	$\Delta k = k(0) - k(\omega) + k(\omega) = 0$	Automatically satisfied

^a Type I phase matching.

^b Type II phase matching.

to 2 mm in the microwave region. The properties of some nonlinear crystals used for second-order interactions are given in Table IV, and some specific examples of second-order interactions are given in Table V.

i. Second-harmonic generation. In second-harmonic generation, radiation at an incident frequency ω_1 is converted to radiation at twice the frequency, and one-half the wavelength, of the incident radiation:

$$\omega_2 = 2\omega_1. \quad (12)$$

For this interaction, the incident frequencies in the nonlinear polarization of Eq. (10) are equal, the generated frequency is ω_2 , and $g = 1$. The intensity of the harmonic wave can be calculated for specific configurations from the appropriate equation in Eq. (6) using the nonlinear polarization amplitude $P(2\omega) = \epsilon_0 d_{\text{eff}}(-2\omega, \omega, \omega)A(\omega)^2$ from Table III. The simplest situation is one in which the incident radiation is in the form of a plane wave and the conversion efficiency to the harmonic is small enough that the incident intensity can be regarded as constant. The harmonic intensity, related to the field amplitude by

$$I(2\omega) = cn_2\epsilon_0|A(2\omega)|^2/2, \quad (13)$$

where n_2 is the refractive index at ω_2 , is then given by

$$I(2\omega) = [8\pi^2 d_{\text{eff}}^2 / n_2 n_1^2 c \epsilon_0 \lambda_1^2] I_0(\omega)^2 \text{sinc}^2(\Delta k L / 2). \quad (14)$$

Here $\text{sinc}(x) = (\sin x)/x$, $I_0(\omega)$ is the incident intensity at ω , L is the crystal length, and

$$\Delta k = k(2\omega) - k_1(\omega) - k_2(\omega), \quad (15)$$

where $k(2\omega)$ is the usual wave vector at 2ω , and $k_1(\omega)$ and $k_2(\omega)$ are the wave vectors of the components of the incident wave at ω . If the incident radiation has only polarization component relative to the principal directions of the crystal (i.e., is either an ordinary or an extraordinary ray), $k_1(\omega) = k_2(\omega) = k(\omega)$ and $\Delta k = k(2\omega) - 2k(\omega)$, whereas if the incident radiation has both ordinary and extraordinary polarization components, $k_1(\omega)$ is in general not equal to $k_2(\omega)$.

In the low conversion regime, the harmonic intensity grows as the square of the incident intensity. This is one of the reasons effective harmonic conversion requires the intense fields that are present in laser radiation. The harmonic intensity is also very sensitive to the value of the wave-vector mismatch Δk . Its dependence on Δk at fixed L and on L for different values of Δk is shown in Fig. 4a,b, respectively. If $\Delta k \neq 0$, the harmonic wave gradually gets out of phase with the polarization that drives it as it propagates through the crystal. As a result the harmonic intensity oscillates with distance, with the harmonic wave first taking energy from the incident wave, and then, after a

phase change of π relative to the nonlinear polarization, returning it to the incident wave. The shortest distance at which the maximum conversion occurs is termed the coherence length and is given by

$$L_{\text{coh}} = \pi / \Delta k. \quad (16)$$

When $\Delta k \neq 0$, maximum conversion is obtained for crystals whose length is an odd multiple of the coherence length, while no conversion is obtained for crystals whose length is an even multiple of the coherence length.

When

$$\Delta k = 0 \quad (17)$$

the harmonic wave stays in phase with its driving polarization. The process is then said to be phase-matched, and the harmonic intensity grows as the square of the crystal length. Under this condition the harmonic intensity can grow so large that the depletion of the incident radiation must be accounted for. In this case the harmonic and incident intensities are given by

$$I(2\omega_1) = I_0(\omega_1) \tanh^2(L/L_{\text{dep}}) \quad (18a)$$

$$I(\omega_1) = I_0(\omega_1) \text{sech}^2(L/L_{\text{dep}}), \quad (18b)$$

where L_{dep} is a depletion length given by

$$L_{\text{dep}} = [n_2 n_1^2 c \epsilon_0 \lambda_1^2 / 8\pi^2 d_{\text{eff}}^2 I_0(\omega_1)]^{1/2}. \quad (19)$$

When $L = L_{\text{dep}}$, 58% of the fundamental is converted to the harmonic. The harmonic and fundamental intensities are shown as a function of $(L/L_{\text{dep}})^2$ in Fig. 5. Although energy is conserved between the optical waves in second-harmonic generation, the number of photons is not, with two photons being lost at the fundamental for every one that is created at the harmonic.

The wave-vector mismatch Δk occurs because of the natural dispersion in the refractive index that is present in all materials. Effective second-harmonic conversion therefore requires that special steps be taken to satisfy the condition of Eq. (17). This process is termed phase matching. For the noncentrosymmetric crystals used for second-order processes, the most common method of phase matching involves use of the birefringence of the crystal to offset the natural dispersion in the refractive indices. In general this means that the harmonic wave will be polarized differently from the incident wave, with the particular combinations being determined by the symmetry properties of the nonlinear crystal. The value of the wave-vector mismatch can be adjusted by varying the direction of propagation of the various waves relative to the optic axis or by varying the temperature of the crystal for a fixed direction of propagation. The angle at which $\Delta k = 0$ is called the phase-matching angle, and the temperature at which $\Delta k = 0$ when $\theta = 90^\circ$ is called the phase-matching temperature.

TABLE IV Properties of Selected Nonlinear Optical Crystals

Nonlinear material	Symmetry point group	Nonlinear susceptibility d (10^{-12} m/V)	Transparency range (μm)	Effective nonlinearity, $ d_{\text{eff}} $	
				Type I phase matching ^a	Type II phase matching
Ag ₃ AsS ₃ (proustite)	3m	$d_{22} = 22; d_{15} = 13$	0.6–13	$d_{15} \sin \theta - d_{22} \cos \theta \sin 3\phi$	$d_{22} \cos^2 \theta \cos 3\phi$
Te (tellurium)	32	$d_{11} = 649$	4–25	$d_{15} \sin \theta + \cos \theta (d_{11} \cos 3\phi - d_{22} \sin 3\phi)$	$d_{11} \cos^2 \theta$
Tl ₃ AsSe ₃ (TAS)		$d_+ = 40$	1.2–18	d_+	
CdGeAs ₂	$\bar{4}2m$	$d_{36} = d_{14} = d_{25} = 236$	2.4–17	$d_{36} \sin \theta$	$d_{36} \sin 2\theta$
AgGaS ₂	$\bar{4}2m$	$d_{36} = d_{14} = d_{25} = 12$	0.6–13	$d_{36} \sin \theta$	$d_{36} \sin 2\theta$
AgGaSe ₂	$\bar{4}2m$	$d_{36} = d_{14} = d_{25} = 40$	0.7–17	$d_{36} \sin \theta$	$d_{36} \sin 2\theta$
GaAs	$\bar{4}3m$	$d_{36} = d_{14} = d_{25} = 90.1$	0.9–17	$d_{36} \sin \theta$	$d_{36} \sin 2\theta$
LiNbO ₃ (lithium niobate)	3m	$d_{15} = 6.25; d_{22} = 3.3$	0.35–4.5	$d_{15} \sin \theta - d_{22} \cos \theta \sin 3\phi$	$d_{22} \cos^2 \theta \cos 3\phi$
LiIO ₃ (lithium iodate)	6	$d_{31} = 7.5$	0.31–5.5	$d_{31} \sin \theta$	
NH ₄ H ₂ (PO ₄) ₂ (ammonium dihydrogen phosphate, ADP)	$\bar{4}2m$	$d_{36} = d_{14} = d_{25} = 0.57$	0.2–1.2	$d_{36} \sin \theta$	$d_{36} \sin 2\theta$
KH ₂ (PO ₄) ₂ (potassium dihydrogen phosphate, KDP)	$\bar{4}2m$	$d_{36} = d_{14} = d_{25} = 0.5$	0.2–1.5	$d_{36} \sin \theta$	$d_{36} \sin 2\theta$
KD ₂ (PO ₄) ₂ (potassium dideuterium phosphate, KD*P)	$\bar{4}2m$	$d_{36} = d_{14} = d_{25} = 0.53$	0.2–1.5	$d_{36} \sin \theta$	$d_{36} \sin 2\theta$
RbH ₂ (AsO ₄) ₂ (rubidium dihydrogen arsenate, RDA)	$\bar{4}2m$	$d_{36} = d_{14} = d_{25} = 0.47$	0.26–1.46	$d_{36} \sin \theta$	$d_{36} \sin 2\theta$
RbH ₂ (PO ₄) ₂ (rubidium dihydrogen phosphate, RDP)	$\bar{4}2m$	$d_{36} = d_{14} = d_{25} = 0.48$	0.22–1.4	$d_{36} \sin \theta$	$d_{36} \sin 2\theta$
NH ₄ H ₂ (AsO ₄) ₂ (ammonium dihydrogen arsenate, ADA)	$\bar{4}2m$			$d_{36} \sin \theta$	$d_{36} \sin 2\theta$
KD ₂ (AsO ₄) ₂ (potassium dideuterium arsenate, KD*A)	$\bar{4}2m$	$d_{36} = d_{14} = d_{25} = 0.4$	0.22–1.4	$d_{36} \sin \theta$	$d_{36} \sin 2\theta$
CsH ₂ (AsO ₄) ₂ (cesium dihydrogen arsenate, CDA)	$\bar{4}2m$	$d_{36} = d_{14} = d_{25} = 0.48$	0.26–1.43	$d_{36} \sin \theta$	$d_{36} \sin 2\theta$
CsD ₂ (AsO ₄) ₂ (cesium dideuterium arsenate, CD*A)	$\bar{4}2m$	$d_{36} = d_{14} = d_{25} = 0.48$	0.27–1.66	$d_{36} \sin \theta$	$d_{36} \sin 2\theta$
KTiOPO ₄ (potassium titanyl phosphate, KTP)	mm2	$d_{31} = 7; d_{32} = 5.4; d_{33} = 15; d_{24} = 8.1; d_{15} = 6.6$	0.35–4.5	$d_{31} \cos 2\theta + d_{32} \sin 2\theta$	
LiCHO ₈ ·H ₂ O (lithium formate monohydrate, LFM)	mm2	$d_{31} = d_{15} = 0.107$ $d_{32} = d_{24} = 1.25$ $d_{33} = 4.36$	0.23–1.2	$d_{31} \cos 2\theta + d_{32} \sin 2\theta$	
KB ₅ O ₂ ·4H ₂ O (potassium pentaborate, KB5)	mm2	$d_{31} = 0.046$ $d_{32} = 0.003$	0.17–>0.76	$d_{31} \cos 2\theta + d_{32} \sin 2\theta$	
Urea	$\bar{4}2m$	$d_{36} = d_{14} = d_{25} = 1.42$	0.2–1.43	$d_{36} \sin \theta$	$d_{36} \sin 2\theta$
Damage threshold (10^6 W/cm ²)			Uses		
12–40	Harmonic generation, frequency mixing, parametric oscillation in mid infrared				
40–60 (at 5 μm)	Harmonic generation, frequency mixing, parametric oscillation in mid infrared				
32	Harmonic generation, frequency mixing, parametric oscillation in mid infrared				
20–40	Harmonic generation, frequency mixing, parametric oscillation in mid infrared				
12–25	Harmonic generation, frequency mixing, parametric oscillation in mid infrared				
>10	Harmonic generation, frequency mixing, parametric oscillation in mid infrared				
	Harmonic generation, frequency mixing, parametric oscillation in mid infrared; difference frequency generation in far infrared				
50–140	Harmonic generation, frequency mixing, parametric oscillation in near and mid infrared; harmonic generation and frequency mixing in visible and near ultraviolet				
125	Harmonic generation, frequency mixing, parametric oscillation near and in mid infrared; harmonic generation and frequency mixing in visible and near ultraviolet				
500 (60 nsec, 1.064 μm)	Harmonic generation, frequency mixing, parametric oscillation in near infrared; harmonic generation and frequency mixing in visible and near ultraviolet				
400 (20 nsec, 694.3 nm)	Harmonic generation, frequency mixing, parametric oscillation in near infrared; harmonic generation and frequency mixing in visible and ultraviolet				
23,000 (200 psec, 1.064 μm)	Harmonic generation, frequency mixing, parametric oscillation in near and mid infrared; harmonic generation and frequency mixing in visible and near ultraviolet				
500 (10 nsec, 1.064 μm)	Harmonic generation, frequency mixing, parametric oscillation in near and mid infrared; harmonic generation and frequency mixing in visible and near ultraviolet				
20,000 (30 psec, 1.064 μm)	Harmonic generation, frequency mixing in near ultraviolet				
350 (10 nsec, 694.3 nm)	Harmonic generation, frequency mixing in near ultraviolet				
200 (10 nsec, 694.3 nm)	Harmonic generation, frequency mixing in near ultraviolet				
	Harmonic generation, frequency mixing in near ultraviolet				
500 (10 nsec, 1.064 μm)	Harmonic generation, frequency mixing and parametric oscillation in visible and near infrared				
>260 (12 nsec, 1.064 μm)	Harmonic generation, frequency mixing and parametric oscillation in visible and near infrared				
160 (20 nsec, 1.064 μm)	Harmonic generation, frequency mixing and parametric oscillation in visible and near infrared				
>1000	Harmonic generation and frequency mixing in the ultraviolet				
	Harmonic generation and frequency mixing in the ultraviolet				
1.4×10^3	Harmonic generation and frequency mixing in the ultraviolet				

^a θ is the direction of propagation with respect to the optic axis, ϕ is the direction of propagation with respect to the x axis.

TABLE V Performance for Selected Second-Order Frequency Conversion Interactions

Laser	Incident wavelength (μm)	Generated wavelength (μm)	Nonlinear material	Conversion efficiency		Conditions, comments
				Power (%)	Energy (%)	
Second harmonic generation						
CO ₂	10.6	5.3	Ag ₃ AsS ₃ (proustite)	~1		
	9.6	4.8	Te	5		Limited by multiphoton absorption
			Tl ₃ AsSe ₃	40	25	1 J/cm ² pump, 100 nsec
			CdGeAs ₂		27	
			AgGaSe ₂	60	14	2.1 cm, 75 nsec, 30 mJ, 12 MW/cm ²
Nd:YAG	1.064	0.532	GaAs	2.7		Phase matching with stacked plates
			ADP	23		
			KDP	83		346 J output (with Nd:glass amplifiers)
			KD*P		75	30 psec pulses, 10 ¹⁰ W/cm ² pump
			CDA	60		
			CD*A	57		
			LiNbO ₃	30		
			KTP	42		
Nd:glass (silicate)	1.059	0.5295	KDP	92		5-psec pulses
Nd:glass (phosphate)	1.052	0.526	KDP		80	
Ruby	0.6943	0.347	ADP	40		
			RDA	40		
Nd:YAG	0.532	0.266	ADP		85	30-psec pulses, 10 ¹⁰ W/cm ² pump intensity
			KDP	70		50 J output (with Nd:glass amplifiers)
			KD*P	75		30 psec, 10 ¹⁰ W/cm ² pump intensity
Dye lasers	0.560–0.620	0.280–0.310	ADP, KDP	8–9		
	0.460–0.600	0.230–0.300	LFM	1–2		
	0.434–0.500	0.217–0.250	KB5	0.2–2		
Ar ⁺	0.5145	0.2572	ADP, KDP	30		300 mW, intracavity, cw
Three-wave sum-frequency mixing ($\omega_3 = \omega_1 + \omega_2$)						
CO ₂	9.6, 4.8	3.53	Tl ₃ AsSe ₃	1.6		1 J/cm ² pump, overall efficiency from 9.6 μm
Nd:YAG	1.064, 0.532	0.3547	KDP	55		41 J output (with Nd:glass amplifiers)
			RDP	21		
Nd:YAG	1.064, 0.266	0.2128	ADP	60		Measured with respect to power at 266 nm
			KDP			1-kW output, 20-nsec pulses
			KB5	0.002		
Nd:YAG,	1.064,	0.208–0.234	ADP	10		
Dye (second harmonic)	0.258–0.300					
Nd:YAG,	0.266	0.1966–0.199	KB5	3		40 kW, 5 nsec at 0.1966
Dye	0.745–0.790					
Dye	0.622–0.652,					
	0.310–0.326,					
	0.740–0.920,					
(Second harmonic)	0.232–0.268;	0.185–0.2174	KB5	8–12		
Difference-frequency mixing ($\omega_3 = \omega_1 - \omega_2$)						
Nd:YAG,	0.532,	2.8–5.65	LiIO ₃	0.92		Measured with respect to total power
Dye	0.575–0.660					

continues

TABLE V (continued)

Laser	Incident wavelength (μm)	Generated wavelength (μm)	Nonlinear material	Conversion efficiency		Conditions, comments
				Power (%)	Energy (%)	
Nd:YAG,	1.064,	1.25–1.6	LiIO ₃	0.7		70 W output
Dye	0.575–0.640					
Ruby	0.6943					
Dye	0.84–0.89	3–4	LiNbO ₃	0.015		6-kW output pulses
	0.80–0.83	4.1–5.2	LiIO ₃	6×10^{-3}		100-W output pulses
	0.78–0.87	3.2–6.5	Ag ₃ AsS ₃	2.7×10^{-3}		110-W pulses at 5 μm
	0.73–0.75	10–13	Ag ₃ AsS ₃	2×10^{-4}		100 mW at 10 μm
	0.74–0.818	4.6–12	AgGaS ₂	2×10^{-4}		300 mW at 11 μm
	0.72–0.75	9.5–17	GaSe	0.075		300 W, 20 nsec at 12 μm
Dye	0.81–0.84,	50–500	LiNbO ₃	10^{-5} – 10^{-6}		10–100 mW output
	0.81–0.84					
Dye	0.440–0.510,	1.5–4.8	LiIO ₃	10^4		400-mW-output power
	0.570–0.620					
Dye	0.586,	8.7–11.6	AgGaS ₂	5×10^{-4}		100 μW output power
	0.570–0.620					
Dye		300–2 mm	ZnTe, ZnSe			
Argon,	0.5145,	2.2–4.2	LiNbO ₃	10^{-3}		1 μW continuous
Dye	0.560–0.620					
CO ₂	Various lines between 9.3 and 10.6 10.6	70–2 mm	GaAs	10^{-3}		10-W-output power
CO ₂		90–111	InSb	4×10^{-8}		2- μW -output power
InSn spin-flip	11.7–12					
Parametric oscillators ($\omega_1 \rightarrow \omega_2 + \omega_3$)						
Nd:YAG	0.266	0.420–0.730	ADP	25		100-kW-output power, 30-psec pulses
Nd:YAG	0.472, 0.532	0.55–3.65	LiNbO ₃	50	30	
	0.579, 0.635					
Nd:YAG	1.064	1.4–4.0	AgGaS ₂		16	20-nsec pulses
	1.064	1.4–4.4	LiNbO ₃	40		
Ruby	0.6943	66–200	LiNbO ₃	10^{-4}		
Ruby	0.6943	0.77–4	LiIO ₃	1–10		2- to 100-kW output
Ruby	0.347	0.415–2.1	LiNbO ₂	8		10-kW output
Hydrogen fluoride	2.87	4.3–4.5	CdSe	10		800-W output, 300 nsec
		8.1–8.3				

In the simplest situation, termed type I phase matching, the incident radiation is polarized as an ordinary or an extraordinary ray, depending on the properties of the nonlinear material, and the harmonic radiation is polarized orthogonally to the incident radiation as shown in Fig. 6a. The wave-vector mismatch is then given by

$$\Delta k = k(2\omega) - 2k(\omega) = 4\pi[n(2\omega) - n(\omega)]/\lambda_1. \quad (20)$$

Phase-matching entails choosing conditions such that the refractive index for the field at the fundamental wavelength equals the refractive index for the field at the harmonic,

that is, $n(2\omega) = n(\omega)$. For example, in a material in which the extraordinary refractive index is less than the ordinary refractive index, the pump wave is polarized as an ordinary ray and the harmonic wave is polarized as an extraordinary ray. The variation of the refractive index of these rays with wavelength is shown in Fig. 6b. The ordinary index is independent of the propagation direction, but the extraordinary index depends on the angle θ between the optic axis and the propagation direction according to the relation

$$[1/n_{2\omega}(\theta)]^2 = (\cos^2 \theta)/(n_{2\omega}^o)^2 + (\sin^2 \theta)/(n_{2\omega}^e)^2. \quad (21)$$

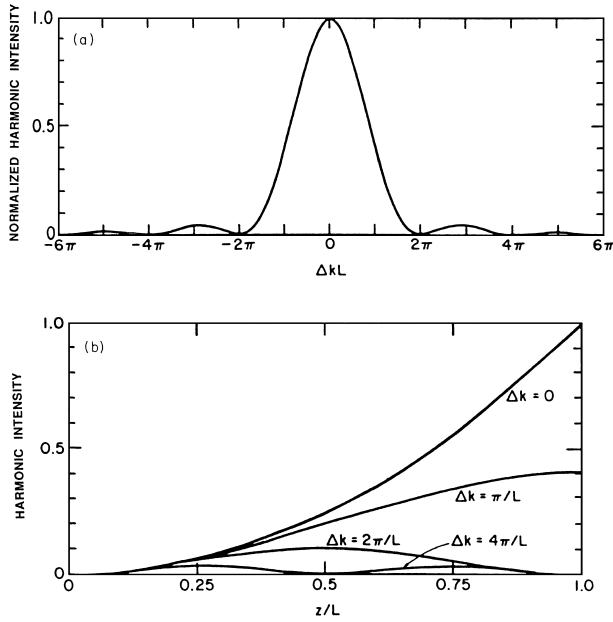


FIGURE 4 (a) Variation of second harmonic intensity with Δk at fixed L . (b) Variation of second harmonic intensity with L for various values of Δk .

As θ varies from 0 to 90° , $n_{2\omega}(\theta)$ varies from $n_{2\omega}^o$ to $n_{2\omega}^e$. In this situation, phase matching consists of choosing the propagation direction such that $n_{2\omega}(\theta) = n_\omega^o$.

In type II phase matching, the incident radiation has both ordinary and extraordinary polarization components, while the harmonic ray is an extraordinary ray (or vice versa, if the nonlinear medium has the property $n^o < n^e$). In type II phase matching, the more general relation among

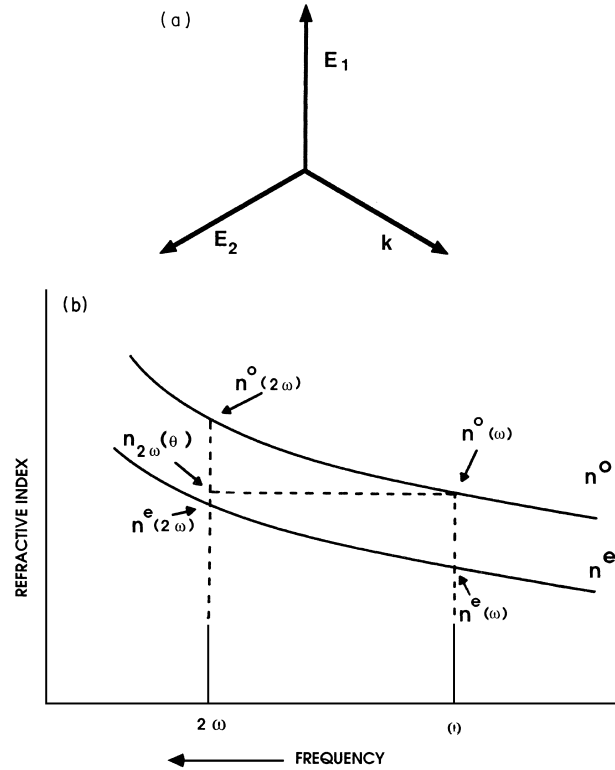


FIGURE 6 (a) Combination of directions of propagation (k), fundamental polarization (E_1), and second harmonic polarization (E_2) for type I second-harmonic generation. (b) Typical variation of the ordinary (n^o) and the extraordinary (n^e) refractive indices with wavelength for a negative uniaxial crystal, showing conditions necessary for type I phase matching of second-harmonic generation. The refractive index of the extraordinary ray at the harmonic varies between n^o and n^e as the angle of propagation relative to the optic axis is varied. Phase matching occurs when $n_{2\omega}(\theta) = n_\omega^o$.

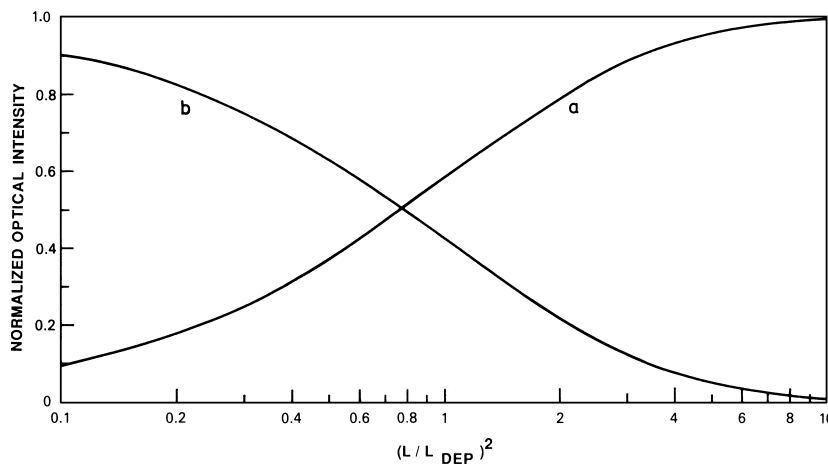


FIGURE 5 Variation of the second harmonic (a) and fundamental (b) intensity with $(L/L_{\text{dep}})^2$, where L_{dep} is the second-harmonic depletion length, for plane waves at perfect phase matching. [Reproduced from Reintjes, J. (1985). Coherent ultraviolet and vacuum ultraviolet sources. In "Laser Handbook," Vol. 5 (M. Bass and M. Sticht, eds.), North-Holland, Amsterdam.]

the k vectors of Eq. (15) must be used, and the relationship between the various refractive indexes required for phase matching is more complicated than that given in Eq. (20) for type I, but the propagation direction is again chosen such that Eq. (17) is satisfied.

Second-harmonic generation has been used to generate light at wavelengths ranging from 217 nm (by doubling radiation from a dye laser at 434 nm) to 5 μm , by doubling radiation from a CO_2 laser at 10.6 μm . Examples of specific results are given in Table V. The usefulness of a particular material for second-harmonic generation is determined by a combination of the magnitude of its nonlinear susceptibility, its ability to phase match the process and its degree of transparency to the harmonic and fundamental wavelengths. The short wave limit of second-harmonic conversion is currently set by an inability to phase match the harmonic process at shorter wavelengths, while the infrared limit of 5 μm is determined by increasing absorption at the fundamental.

Second-harmonic conversion is commonly used with high-power pulsed lasers such as Nd:YAG, Nd:glass, ruby, or CO_2 to generate high-power fixed-frequency radiation at selected wavelengths in the infrared, visible, and ultraviolet. It is also a versatile and convenient source of tunable ultraviolet radiation when the pump radiation is obtained from tunable dye lasers in the visible and near infrared.

Second-harmonic conversion efficiencies range from the order of 10^{-6} for conversion of continuous wave radiation outside the laser cavity to over 90% for high-power pulsed radiation. The laser power required for high efficiency depends on the configuration and nonlinear material. A typical value for 75% conversion of radiation from an Nd:glass laser at 1.06 μm to 530 nm in a 1-cm-long crystal of potassium dihydrogen phosphate (KDP) is 10^8 W/cm^2 . Intensities for conversion in other materials and with other lasers are noted in Table V. Conversion of radiation from high-power pulsed lasers is typically done in collimated-beam geometries with angle-tuned crystals chosen for their high damage thresholds. Conversion of radiation from lower-power pulsed lasers or continuous-wave (cw) lasers is often done in focused geometries in materials with large nonlinear susceptibilities. In many situations the maximum conversion efficiency that can be achieved is not determined by the available laser power but by additional processes that affect the conversion process. Some of these competing processes are linear or nonlinear absorption, imperfect phase matching because of a spread of the laser divergence or bandwidth, incomplete phase matching because of intensity-dependent changes in the refractive index (see Section III.B.1), or damage to the nonlinear crystal caused by the intense laser radiation.

ii. Three-wave sum-frequency mixing. Three-wave sum-frequency mixing is used to generate radiation at higher frequencies, and therefore shorter wavelengths, than those in the pump radiation according to the relation

$$\omega_3 = \omega_1 + \omega_2, \quad (22)$$

where ω_3 is the frequency of the radiation generated in the interaction and ω_1 and ω_2 are the frequencies of the pump radiation. Calculation of the generated intensity is similar to that of second-harmonic generation with the appropriate form of the nonlinear polarization as given in Table III used in Eq. (6). At low conversion efficiencies, the generated intensity has the same $\sin^2(\Delta k L/2)/(\Delta k L/2)^2$ dependence on the wave-vector mismatch as has second-harmonic conversion, but the phase-matching condition is now $\Delta k = k(\omega_3) - k(\omega_1) - k(\omega_2) = 0$, which is equivalent to

$$n(\omega_3)/\lambda_3 = n(\omega_1)/\lambda_1 + n(\omega_2)/\lambda_2. \quad (23)$$

The intensity of the generated radiation grows as the product of the incident pump intensities at low conversion efficiency. In the three-wave sum-frequency mixing interaction, one photon from each of the pump waves is annihilated for each photon created in the generated wave. Complete conversion of the total radiation in both pump waves at perfect phase matching is possible in principle if they start with an equal number of photons. Otherwise the conversion will oscillate with pump intensity or crystal length, even at exact phase matching.

Three-wave sum-frequency mixing is done in the same types of materials as second-harmonic generation. It is used to generate both tunable and fixed-frequency radiation at various wavelengths ranging from the infrared to the ultraviolet. It allows radiation to be generated at shorter wavelengths in the ultraviolet than can be reached with second harmonic conversion if one of the pump wavelengths is also in the ultraviolet.

Examples of specific three-wave sum-frequency mixing interactions are given in Table V. This interaction has been used to produce radiation at wavelengths as short as 185 nm in potassium pentaborate (KB_5), with the cut-off being determined by the limits of phase matching. Three-wave sum-frequency mixing can also be used to improve the efficiency of the generation of tunable radiation, as compared to second-harmonic generation, by allowing radiation from a relatively powerful fixed-frequency laser to be combined with radiation from a relatively weak tunable laser. Another application is in the generation of the third harmonic of certain fixed frequency lasers such as Nd:YAG, Nd:glass, or CO_2 through two second-order processes: second-harmonic conversion of part of the laser fundamental followed by sum-frequency mixing of the unconverted fundamental with the second harmonic. Under certain conditions this process is more efficient than

direct third-harmonic conversion (see Section III.A.1.b). Three-wave sum-frequency mixing has also been used for up-conversion of infrared radiation to the visible, where it can be measured by more sensitive photoelectric detectors or photographic film.

iii. *Three-wave difference-frequency mixing.* Three-wave difference-frequency mixing is used to convert radiation from two incident waves at frequencies ω_1 and ω_2 to a third wave at frequency ω_3 according to the relation

$$\omega_3 = \omega_1 - \omega_2. \quad (24)$$

Just as for sum-frequency mixing, the generated intensity at low conversion efficiency grows as the product of the pump intensities at ω_1 and ω_2 . In this situation a photon is created at both ω_3 and ω_2 for every photon annihilated at ω_1 . The wave-vector mismatch is $\Delta k = k(\omega_3) - [k(\omega_1) - k(\omega_2)]$, and the phase-matching condition is

$$n(\omega_3)/\lambda_3 = n(\omega_1)/\lambda_1 - n(\omega_2)/\lambda_2. \quad (25)$$

The materials used for difference-frequency mixing are of the same type as those used for second-harmonic and sum-frequency mixing. Difference-frequency mixing is generally used to produce coherent radiation at longer wavelengths than either of the pump wavelengths. It is used most often to produce tunable radiation in the infrared from pump radiation in the visible or near infrared, although it can also be used to generate radiation in the visible. It has been used to produce radiation at wavelengths as long as 2 mm in GaAs, with the limit being set by a combination of the increasing mismatch in the diffraction of the pump and generated radiation and the increasing absorption of the generated radiation in the nonlinear medium at long wavelengths.

iv. *Parametric down-conversion.* Parametric down-conversion is used to convert radiation in an optical wave at frequency ω_1 into two optical waves at lower frequencies ω_2 and ω_3 according to the relation

$$\omega_1 = \omega_2 + \omega_3. \quad (26)$$

This process is illustrated schematically in Fig. 7. The wave at ω_1 is termed the pump wave, while one of the waves at ω_2 or ω_3 is termed the signal and the other the idler. When $\omega_2 = \omega_3$ the process is termed degenerate para-



FIGURE 7 Schematic illustration of the waves used in parametric down-conversion. Radiation at ω_1 is supplied, and radiation at ω_2 and ω_3 is generated in the nonlinear interaction.

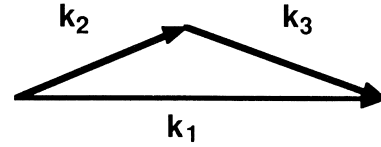


FIGURE 8 Off-axis phase-matching diagram for parametric down-conversion. Direction of arrows indicates direction of propagation of the various waves.

metric down-conversion and is the opposite of second-harmonic generation, whereas if $\omega_2 \neq \omega_3$, the process is called nondegenerate parametric down-conversion and is the opposite of sum-frequency generation. The individual values of ω_2 and ω_3 are determined by the phase-matching condition, which for plane-wave interactions is given by

$$\Delta k = k(\omega_1) - k(\omega_2) - k(\omega_3) = 0 \quad (27)$$

and can be varied by changing an appropriate phase-matching parameter of the crystal such as the angle of propagation or the temperature. Offaxis phase matching, as illustrated in Fig. 8, is also possible. The phase-matching condition of Eq. (27) is the same as that for the sum-frequency process $\omega_2 + \omega_3 = \omega_1$. The relative phase of the waves involved determines whether the sum process or the parametric down-conversion process will occur.

In the absence of pump depletion, the amplitudes of the waves generated in a parametric down-conversion process are given by

$$A(\omega_2) = A_0(\omega_2) \cosh \kappa z + i A_0^*(\omega_3) (\omega_2 n_3 / \omega_3 n_2)^{1/2} \sinh \kappa z \quad (28a)$$

$$A(\omega_3) = A_0(\omega_3) \cosh \kappa z + i A_0^*(\omega_2) (\omega_3 n_2 / \omega_2 n_3)^{1/2} \sinh \kappa z, \quad (28b)$$

where

$$\kappa^2 = [2\omega_2\omega_3 d_{\text{eff}}^2 / n_2 n_3 c^2] I_0(\omega_1) \quad (29)$$

and $A_0(\omega_i)$ and $I_0(\omega_i)$ are the incident field amplitude and intensity, respectively, at ω_i . At low pump intensities the generated field amplitudes grow in proportion to the square root of the pump intensity, while at high pump intensities the growth of the generated waves is exponential.

If there is no incident intensity supplied at ω_2 or ω_3 , the process is termed parametric generation or parametric oscillation, depending on the geometry. For this situation the initial intensity for the generated waves arise from the zero-point radiation field with an energy of $h\nu/2$ per mode for each field. If the interaction involves a single pass through the nonlinear medium, as was shown in Fig. 7, the process is termed parametric generation. This geometry is typically used with picosecond pulses for generation of tunable infrared or visible radiation from pump radiation

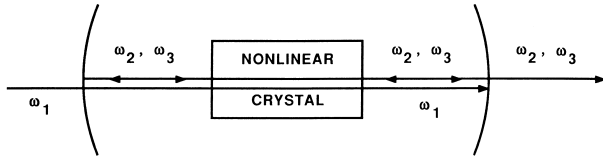


FIGURE 9 Illustration of a doubly resonant cavity for a parametric oscillator.

in the ultraviolet, visible, or near infrared. Amplification factors of the order of 10^{10} are typically required for single-pass parametric generation.

Parametric down-conversion can also be used with a resonant cavity that circulates the radiation at either or both of the generated frequencies, as shown in Fig. 9. In this geometry the process is termed parametric oscillation. If only one of the generated waves is circulated in the cavity, it is termed singly resonant, whereas if both waves are circulated, the cavity is termed doubly resonant. In singly resonant cavities the wave that is circulated is termed the signal, while the other generated wave is termed the idler. Optical parametric oscillators are typically used with pump radiation from Q-switched lasers with pulses lasting several tens of nanoseconds, allowing several passes of the generated radiation through the cavity while the pump light is present.

One of the primary uses of parametric down-conversion is the generation of tunable radiation at wavelengths ranging from the visible to the far infrared. Its wavelength range is generally the same as that covered by difference-frequency mixing, although it has not been extended to as long a wavelength. Tuning is done by varying one of the phase-matching parameters such as angle or temperature. A typical tuning curve for a parametric oscillator is shown in Fig. 10. As the phase-matching condition is changed from the degenerate condition, ω_2 increases and ω_3 decreases in such a way as to maintain the relation in Eq. (26). The extent of the tuning range for a given combination of pump wavelength and nonlinear material is generally set by the limits of phase matching, although absorption can be important if ω_3 is too far in the infrared. Radiation in different wavelength ranges can be produced by using different nonlinear materials and different pump sources. Parametric down-conversion has been used to produce radiation at wavelengths ranging from the visible to $25\ \mu\text{m}$. Some of the combinations of pump sources, nonlinear materials, and tuning ranges are listed in Table V.

Parametric down-conversion can also be used to amplify radiation at ω_2 or ω_3 . In this arrangement radiation is supplied at both the pump wavelength and the lower-frequency wave to be amplified, which is termed the signal. The process is similar to difference-frequency mixing,

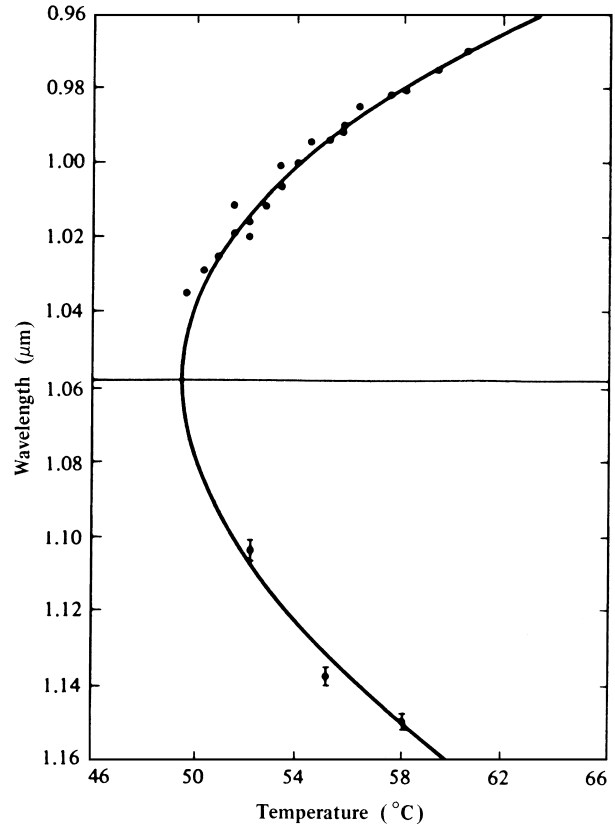


FIGURE 10 Tuning curve for a LiNbO_3 parametric oscillator pumped by radiation from the second harmonic of a Nd laser at 529 nm. [Reproduced from Giordmaine, J. A., and Miller, R. C. (1965). *Phys. Rev. Lett.* **14**, 973.]

and differs from the difference-frequency process only in that the incident intensity in the signal wave is considerably less than the pump intensity for parametric amplification, whereas the two incident intensities are comparable in difference-frequency mixing. In principle, very high gains can be obtained from parametric amplification, with values up to 10^{10} being possible in some cases.

Optical rectification is a form of difference-frequency mixing in which the generated signal has no carrier frequency. It is produced through the interaction

$$\omega_3 = 0 = \omega_1 - \omega_1. \quad (30)$$

It does not produce an optical wave, but rather an electrical voltage signal with a duration corresponding to the pulse duration of the pump radiation. Optical rectification has been used with picosecond laser pulses to produce electrical voltage pulses with durations in the picosecond range, among the shortest voltage pulses yet produced. These pulses do not propagate in the bulk of the nonlinear crystal as do the optical waves, and observation or use of them requires coupling to an appropriate microwave strip line on the nonlinear crystal.

b. Third- and higher-order processes. Third- and higher-order parametric processes are also used for frequency conversion. They have been used to generate radiation ranging from 35.5 nm, almost in the soft X-ray range, to about 25 μm in the infrared. The most commonly used interactions of this type are given in Table VI, along with the form of the nonlinear polarization, the wavevector mismatch, and the plane-wave phase-matching conditions. These interactions include third- and higher-order harmonic conversion, in which an incident wave at frequency ω_1 is converted to a wave at frequency $q\omega_1$, and various forms of four- and six-wave frequency mixing in which radiation in two or more incident waves is converted to radiation in a wave at an appropriate sum- or difference-frequency combination as indicated in Table VI. Four-wave parametric oscillation, in which radiation at frequency ω_1 is converted to radiation at two other frequencies, ω_2 and ω_3 , according to the relation

$$2\omega_1 = \omega_2 + \omega_3 \quad (31)$$

has also been observed.

Third-order and higher odd-order processes can be observed with electric dipole interactions in materials with any symmetry. They are used most commonly in materials that have a center of symmetry, such as gases, liquids, and some solids, since in these materials they are the lowest-order nonzero nonlinearities allowed by electric dipole transitions. Fourth-order and higher even-order processes involving electric dipole interactions are allowed only in crystals with no center of symmetry, and, although they have been observed, they are relatively inefficient and are seldom used for frequency conversion.

The intensity generated in these interactions can be calculated from Eq. (6) using the appropriate nonlinear polarization from Table VI. In the absence of pump depletion, the intensity generated in q th harmonic conversion varies as the q th power of the incident pump intensity, while the intensity generated in the frequency-mixing interactions varies as the product of the incident pump intensities with each raised to a power corresponding to its multiple in the appropriate frequency combination. The generated intensity in the plane wave configuration has the same $[(\sin \Delta k L/2)/(\Delta k L/2)]^2$ dependence on the wave-vector mismatch as do the second-order processes. The plane-wave phase-matching conditions for each of these interactions is the same as for the second-order interactions, namely, $\Delta k = 0$, but the requirements on the individual refractive indexes depend on the particular interaction involved, as indicated in Table VI.

Third- and higher-order frequency conversion are often done with beams that are tightly focused within the nonlinear medium to increase the peak intensity. In this situation, optimal performance can require either a positive or neg-

ative value of Δk , depending on the interaction involved. Phase-matching requirements with focused beams for the various interactions are also noted in Table VI.

The isotropic materials used for third- and higher-order parametric processes are not birefringent, and so alternative phase-matching techniques must be used. In gases, phase matching can be accomplished through use of the negative dispersion that occurs near allowed transitions as shown in Fig. 11. Normal dispersion occurs when the refractive index increases with frequency and is encountered in all materials when the optical frequency falls below, or sufficiently far away from, excited energy levels with allowed transitions to the ground state. Anomalous dispersion occurs in narrow regions about the transition frequencies in which the refractive index decreases with frequency, as shown in Fig. 11a. Negative dispersion occurs in a wavelength range above an allowed transition in which the refractive index, although increasing with frequency, is less than it is below the transition frequency. Regions of negative dispersion occur in restricted wavelength ranges above most, but not all, excited levels in many gases. Examples of the energy-level structures that give positive and negative dispersion for third-harmonic generation are shown in Fig. 11b.

Phase matching can be accomplished by using a mixture of gases with different signs of dispersion. In this situation each component makes a contribution to the wave-vector mismatch in proportion to its concentration in the mixture. The value of the wave-vector mismatch can be controlled by adjusting the relative concentration of the two gases until the appropriate phase-matching condition is met for either collimated or focused beams, as shown in Fig. 12.

Phase matching can also be done in single-component media with focused beams, provided that the dispersion of the medium is of the correct sign for the interaction and wavelengths involved. With this technique the wave-vector mismatch depends on the density of the gas, and the pressure is adjusted until the proper wave-vector mismatch is achieved. Alternatively, phase matching in a single-component medium can be done by choosing the pump and generated frequencies to lie on either side of the transition frequency so that the phase-matching condition is satisfied. This technique is usually used in gases with plane-wave pump beams.

A fourth method for phase matching is the use of non-collinear waves, as shown in Fig. 13. This technique can be used for sum-frequency processes in media with negative dispersion and for difference-frequency processes in media with positive dispersion. It is commonly used, for example, in liquids for the difference frequency process $\omega_4 = 2\omega_1 - \omega_2$.

The conversion efficiency can be increased significantly if resonances are present between certain energy levels of

TABLE VI Nonlinear Polarizations for Third- and Higher-Order Parametric Frequency Conversion Processes

Nonlinear interaction	Process	Nonlinear polarization	Wave-vector mismatch	Plane-wave-matching condition	Dispersion requirement for focused beams in infinitely long media
q th Harmonic generation	$\omega_q = q\omega_1$	$P(q\omega) = \epsilon_0 \chi(-q\omega, \omega, \omega, \dots, \omega) \times A^q(\omega)/2^{q-1}$	$\Delta k = k(q\omega) - qk(\omega)$	$n(q\omega) = n(\omega)$	$\Delta k < 0$
Four-wave sum-frequency mixing	$\omega_4 = 2\omega_1 + \omega_2$	$P(\omega_4) = 3\epsilon_0 \chi(-\omega_4, \omega_1, \omega_1, \omega_2) \times A^2(\omega_1)A(\omega_2)/4$	$\Delta k = k(\omega_4) - 2k(\omega_1) - k(\omega_2)$	$n(\omega_4)/\lambda_4 = 2n(\omega_1)/\lambda_1 + n(\omega_2)/\lambda_2$	$\Delta k < 0$
	$\omega_4 = \omega_1 + \omega_2 + \omega_3$	$P(\omega_4) = 3\epsilon_0 \chi(-\omega_4, \omega_1, \omega_2, \omega_3) \times A(\omega_1)A(\omega_2)A(\omega_3)/2$	$\Delta k = k(\omega_4) - k(\omega_1) - k(\omega_2) - k(\omega_3)$	$n(\omega_4)/\lambda_4 = n(\omega_1)/\lambda_1 + n(\omega_2)/\lambda_2 + n(\omega_3)/\lambda_3$	$\Delta k < 0$
Four-wave difference-frequency mixing	$\omega_4 = 2\omega_1 - \omega_2$	$P(\omega_4) = 3\epsilon_0 \chi(-\omega_4, \omega_1, \omega_1, -\omega_2) \times A^2(\omega_1)A^*(\omega_2)/4$	$\Delta k = k(\omega_4) - 2k(\omega_1) + k(\omega_2)$	$n(\omega_4)/\lambda_4 = 2n(\omega_1)/\lambda_1 - n(\omega_2)/\lambda_2$	$\Delta k = 0^a, \Delta k \leq 0^b$
	$\omega_4 = \omega_1 + \omega_2 - \omega_3$	$P(\omega_4) = 3\epsilon_0 \chi(-\omega_4, \omega_1, \omega_2, -\omega_3) \times A(\omega_1)A(\omega_2)A^*(\omega_3)/2$	$\Delta k = k(\omega_4) - k(\omega_1) - k(\omega_2) + k(\omega_3)$	$n(\omega_4)/\lambda_4 = n(\omega_1)/\lambda_1 + n(\omega_2)/\lambda_2 - n(\omega_3)/\lambda_3$	$\Delta k = 0^a, \Delta k \leq 0^b$
	$\omega_4 = \omega_1 - \omega_2 - \omega_3$	$P(\omega_4) = 3\epsilon_0 \chi(-\omega_4, \omega_1, -\omega_2, -\omega_3) \times A(\omega_1)A^*(\omega_2)A^*(\omega_3)/2$	$\Delta k = k(\omega_4) - k(\omega_1) + k(\omega_2) + k(\omega_3)$	$n(\omega_4)/\lambda_4 = n(\omega_1)/\lambda_1 - n(\omega_2)/\lambda_2 - n(\omega_3)/\lambda_3$	$\Delta k > 0$
Four-wave parametric oscillation	$2\omega_1 \rightarrow \omega_2 + \omega_3$	$P(\omega_2) = 3\epsilon_0 \chi(-\omega_2, \omega_1, \omega_1, -\omega_3) \times A^2(\omega_1)A^*(\omega_3)/4$	$\Delta k = k(\omega_2) + k(\omega_3) - 2k(\omega_1)$	$2n(\omega_1)/\lambda_1 = n(\omega_2)/\lambda_2 + n(\omega_3)/\lambda_3$	$\Delta k = 0^a, \Delta k \leq 0^b$
		$P(\omega_3) = 3\epsilon_0 \chi(-\omega_3, \omega_1, \omega_1, -\omega_2) \times A^2(\omega_1)A^*(\omega_2)/4$	$\Delta k = k(\omega_6) - 4k(\omega_1) - k(\omega_2)$	$n(\omega_6)/\lambda_6 = 4n(\omega_1)/\lambda_1 + n(\omega_2)/\lambda_2$	$\Delta k < 0$
Six-wave sum-frequency mixing	$\omega_6 = 4\omega_1 + \omega_2$	$P(\omega_6) = 5\epsilon_0 \chi(-\omega_6, \omega_1, \omega_1, \omega_1, \omega_1, \omega_2) \times A^4(\omega_1)A(\omega_2)/16$	$\Delta k = k(\omega_6) - 4k(\omega_1) + k(\omega_2)$	$n(\omega_6)/\lambda_6 = 4n(\omega_1)/\lambda_1 - n(\omega_2)/\lambda_2$	$\Delta k < 0^a, \Delta k > 0^d$
Six-wave difference-frequency mixing	$\omega_6 = 4\omega_1 - \omega_2$	$P(\omega_6) = 5\epsilon_0 \chi(-\omega_6, \omega_1, \omega_1, \omega_1, \omega_1, -\omega_2) \times A^4(\omega_1)A^*(\omega_2)/16$			

^a Phase matching in mixtures or by angle.^b Phase optimization in single-component media.^c Requirement for optimized conversion.^d Positive dispersion allowed but not optimal.

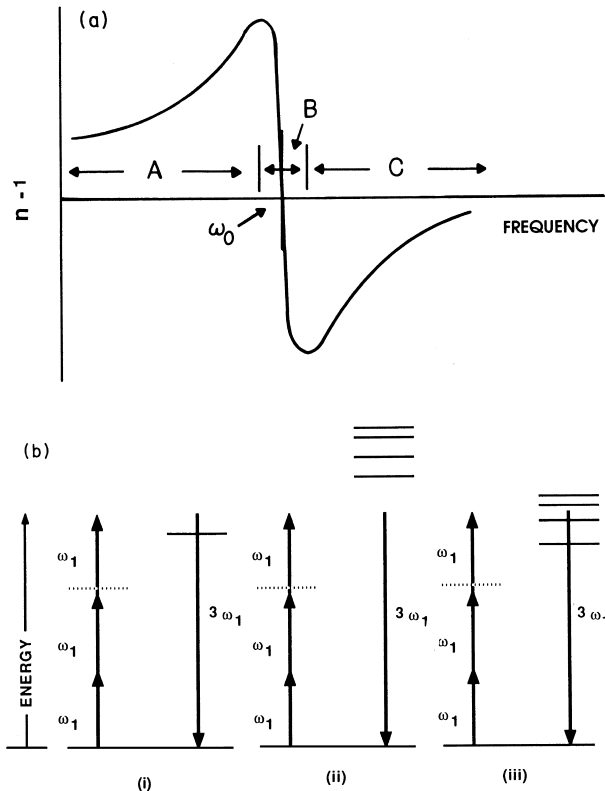


FIGURE 11 (a) Variation of the refractive index near an allowed transition showing regions of normal dispersion (A), anomalous dispersion (B), and negative dispersion (C). (b) Energy level diagrams for third-harmonic generation that provide negative dispersion (i) and positive dispersion (ii and iii).

the medium and the incident and generated frequencies or their sum or difference combinations. This increase in conversion efficiency is similar to the increase in linear absorption or scattering that occurs when the incident wavelength approaches an allowed transition of the medium.

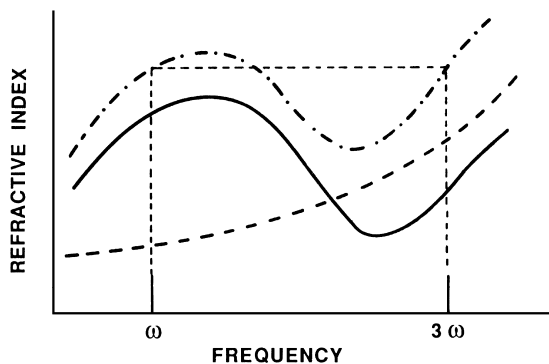


FIGURE 12 Illustration of the use of mixtures for phase matching. The solid curve gives the refractive index variation of a negatively dispersive medium, the dashed curve shows a positively dispersive medium, and the chain curve shows the refractive index of a mixture chosen so that the refractive index at ω is equal to that at 3ω .

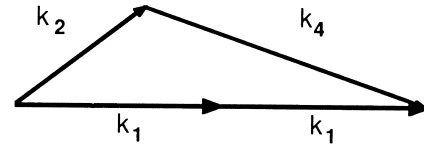


FIGURE 13 Off-axis phase-matching diagram for four-wave mixing of the form $\omega_4 = 2\omega_1 - \omega_2$.

For the nonlinear interactions, however, a much greater variety of resonances is possible. Single-photon resonances occur between the incident or generated frequencies and allowed transitions just as with linear optical effects. The effectiveness of these resonances in enhancing nonlinear processes is limited, however, because of the absorption and dispersion that accompanies them.

Resonances in the nonlinear effects also occur when multiples, or sum or difference combinations, of the incident frequencies match the frequency of certain types of transitions. The most commonly used of these resonances is a two-photon resonance involving two successive dipole transitions between levels whose spacing matches twice the value of an incident frequency or a sum or difference of two input frequencies, as indicated in Fig. 14. In single atoms and in centrosymmetric molecules, the energy levels involved in such two-photon resonances are of the same parity, and transitions between them are not observable in linear spectroscopy, which involves single-photon transitions. Near a two-photon resonance, the nonlinear susceptibility can increase by as much as four to eight orders of magnitude, depending on the relative linewidths of the two-photon transition and the input radiation, resulting in a dramatic increase in the generated power as the input frequency is tuned through the two-photon resonance. An example of the increase in efficiency that is observed as the pump frequency is tuned through a two-photon resonance is shown in Fig. 15. Other higher-order resonances are also possible, but they have not been used as commonly as the two-photon resonances. Resonantly enhanced third-harmonic generation and four-wave mixing have proven very useful in allowing effective conversion of tunable

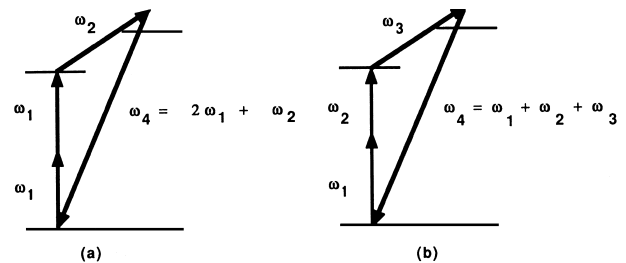


FIGURE 14 Level diagrams showing two-photon resonances at (a) $2\omega_1$ and (b) $\omega_1 + \omega_2$ in four-wave sum-frequency mixing.

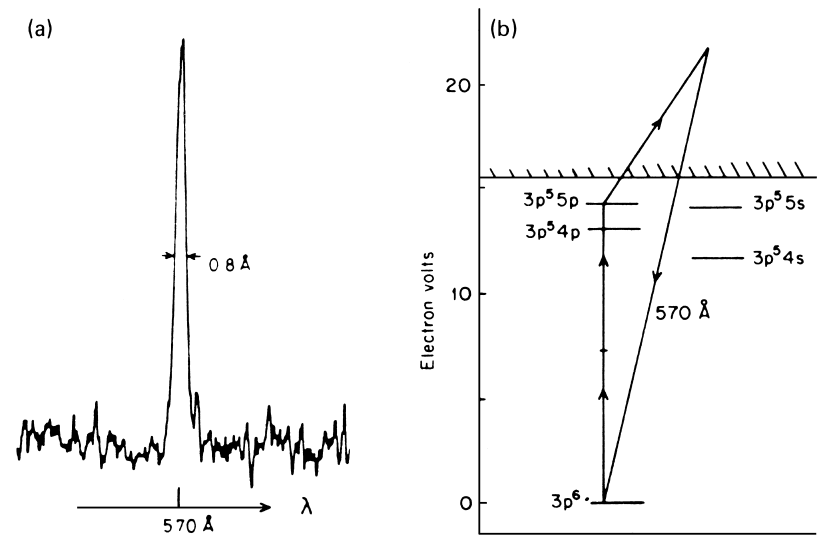


FIGURE 15 (a) Third-harmonic conversion of radiation from an Xe laser showing enhancement as the wavelength of the laser is tuned through a two-photon resonance. (b) The resonant-level structure. [Reproduced from Hutchinson, M. H. R., *et al.* (1976). *Opt. Commun.* **18**, 203. Copyright North-Holland, Amsterdam.]

radiation from dye lasers to the vacuum ultraviolet to providing high-brightness, narrow-band sources of radiation for high-resolution spectroscopy and other applications.

Some of the applications of third- and higher-order frequency conversion are given in Table VII. The q th harmonic generation is used to produce radiation at a frequency that is q times the incident frequency. The most commonly used interaction of this type is third-harmonic conversion. It has been used to produce radiation at wavelengths ranging from the infrared to the extreme ultraviolet. Third-harmonic conversion of radiation from high power pulsed lasers such as CO₂, Nd:glass, Nd:YAG, ruby, and various rare-gas halide and rare-gas excimer lasers has

been used to generate fixed-frequency radiation at various wavelengths ranging from 3.5 μ m to 57 nm, as indicated in Table VII. It has also been used with dye lasers to generate radiation tunable in spectral bands between 110 and 200 nm. The extent of the spectral bands generated in this manner is determined by the extent of the negative dispersion region in the nonlinear materials.

Four-wave sum- and difference-frequency mixing interactions of the form

$$\omega_4 = 2\omega_1 \pm \omega_3 \quad (32a)$$

and

$$\omega_4 = \omega_1 + \omega_2 \pm \omega_3, \quad (32b)$$

TABLE VII Selected Results for Third- and Higher-Order Frequency Conversion Processes

Interaction	Laser	Pump wavelength (nm)	Generated wavelength (nm)	Nonlinear material	Efficiency (%)
Third harmonic	CO ₂	10.6 μ m	3.5 μ m	CO (liquid), CO (gas), BCl ₃ , SF ₆ , NO, DCl	8 (liquid CO)
	Nd:YAG	1.064 μ m	354.7	Rb, Na	10
	Nd:YAG	354.7	118.2	Xe	0.2
	Xe ₂	170	57	Ar	
	Dye	360–600	120–200	Xe, Kr, Sr, Mg, Zn, Hg	Up to 1
Fifth harmonic	Nd:YAG	266	53.2	He, Ne, Ar, Kr	10 ⁻³
	XeCl	308	61.6	He	
	KrF	249	49.8	He	
	ArF	193	38.6	He	
Seventh harmonic	Nd:YAG	266	38	He	10 ⁻⁴
	KrF	249	35.5	He	10 ⁻⁹

where ω_1 , ω_2 , and ω_3 are input frequencies, are also commonly used to produce radiation in wavelength ranges that are inaccessible by other means. These processes can be favored over possible simultaneous third-harmonic generation by the use of opposite circular polarization in the two pump waves, since third-harmonic conversion with circularly polarized pump light is not allowed by symmetry. They have been used to generate radiation over a considerable range of wavelengths in the vacuum ultraviolet, extreme ultraviolet, and the mid infrared. In particular, they have been used to generate tunable radiation over most of the vacuum ultraviolet range from 100 to 200 nm.

These interactions can be used to generate tunable radiation in resonantly enhanced processes, thereby increasing the efficiency of the process. In this situation the pump frequency at ω_1 or the sum combination $\omega_1 + \omega_2$ is adjusted to match a suitable two-photon resonance, while the remaining pump frequency at ω_3 is varied, producing the tunable generated radiation, as illustrated in Fig. 16. The difference-frequency processes $\omega_4 = 2\omega_1 - \omega_3$ and $\omega_4 = \omega_1 + \omega_2 - \omega_3$ can be optimized with focused beams in media with either sign of dispersion. As a result, their usefulness is not restricted to narrow wavelength ranges above dispersive resonances, and they have been used to generate tunable radiation over extensive ranges in the vacuum ultraviolet between 140 and 200 nm in the rare gases Xe and Kr. Tunable radiation generated in this manner in Xe between 160 and 200 nm is illustrated in Fig. 17.

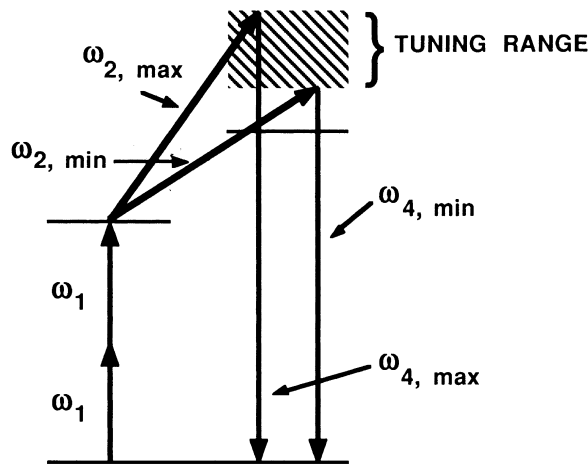


FIGURE 16 Level diagram for producing tunable radiation with two-photon resonantly enhanced four-wave sum-frequency mixing. Pump radiation at ω_1 is tuned to the two-photon resonance, pump radiation at ω_2 is tuned over the range $\omega_{2,\min}$ to $\omega_{2,\max}$, and the frequency of the generated radiation tunes over the range $(2\omega_1 + \omega_{2,\min})$ to $(2\omega_1 + \omega_{2,\max})$.

The difference-frequency processes

$$\omega_4 = 2\omega_1 - \omega_3 \quad (33a)$$

and

$$\omega_4 = \omega_1 \pm \omega_2 - \omega_3 \quad (33b)$$

have also been used to generate tunable radiation in the infrared by using pump radiation from visible and near infrared lasers. In some interactions all the frequencies involved in the mixing processes are supplied externally and in others some of them are generated in a stimulated Raman interaction (see Section III.A.2). Because the gases used for these nonlinear interactions are not absorbing at far-infrared wavelengths, it can be expected that they will ultimately allow more efficient generation of tunable far-infrared radiation using pump radiation in the visible and near infrared than can be achieved in second-order interactions in crystals, although they have not yet been extended to as long wavelengths. Ultimately the limitations on conversion can be expected to arise from difficulties with phase matching and a mismatch between the diffraction of the pump and generated wavelengths. To date, the four-wave difference-frequency mixing interactions have been used to produce coherent radiation at wavelengths out to 25 μm .

Resonances between Raman active molecular vibrations and rotations and the difference frequency combination $\omega_1 - \omega_3$ can also occur. When the four-wave mixing process $2\omega_1 - \omega_3$ or $\omega_1 + \omega_2 - \omega_3$ is used with these resonances it is termed coherent anti-Stokes Raman scattering (CARS). The resonant enhancement that occurs in the generated intensity as the pump frequencies are tuned through the two-photon difference-frequency resonance forms the basis of CARS spectroscopy (see Section IV.A).

Various forms of higher-order interactions are also used for frequency conversion. These consist primarily or harmonic conversion up to order seven- and six-wave mixing interactions of the form $\omega_6 = 4\omega_1 \pm \omega_2$, although harmonic generation up to order 11 has been reported. Generally, the conversion efficiency in the higher-order processes is lower than it is in the lower-order ones, and the required pump intensity is higher. As a result, higher-order processes have been used primarily for the generation of radiation in the extreme ultraviolet at wavelengths too short to be reached with lower-order interaction. The pump sources have for the most part involved mode-locked lasers with pulse durations under 30 psec and peak power levels in excess of several hundred megawatts.

Fifth-harmonic conversion has been used to generate radiation at wavelengths as short as 38.6 nm using radiation from an ArF laser at 193 nm, seventh-harmonic conversion has been used to generate radiation at wavelengths as short as 35.5 nm with radiation from a KrF laser, and ninth-harmonic conversion has been used to generate radiation at

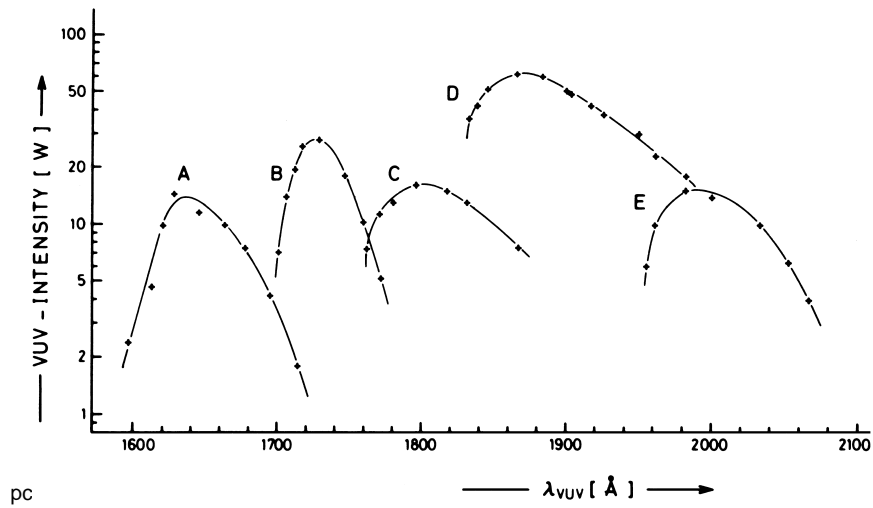


FIGURE 17 Variation of intensity of vacuum ultraviolet radiation (VUV) generated in xenon through the process $2\omega_1 - \omega_2$ in the range 160–200 nm. The different regions A–E correspond to pump radiation obtained from different dye lasers. The radiation is continuously variable within each of the regions. [Reproduced from Hilbig, R., and Wallenstein, R. (1982). *Appl. Opt.* **21**, 913.]

117.7 nm with radiation from a Nd:glass laser at $1.06 \mu\text{m}$. Radiation at various wavelengths in the extreme ultraviolet between 38 and 76 nm using fifth- and seventh-harmonic generation and six-wave mixing of radiation from harmonics of a Nd:YAG laser has also been generated.

Observed conversion efficiencies for many of the third- and higher-order processes are noted in [Table VII](#). They range from about 10% for third-harmonic conversion of Nd:YAG laser radiation in rubidium ($1.064 \mu\text{m} \rightarrow 354.7 \text{ nm}$) or CO_2 laser radiation in liquid CO ($9.6 \mu\text{m} \rightarrow 3.2 \mu\text{m}$) to values of the order of 10^{-11} for some of the higher-order processes. The pump intensities used vary between several hundred kilowatts per square centimeter for resonantly enhanced processes to 10^{15} W/cm^2 for nonresonant processes.

The largest conversion efficiencies that can be achieved with third- and higher-order processes are generally less than those that can be obtained with second-order interactions, because competing processes that limit efficiency are more important for the higher-order interactions. Some of the important limiting processes are listed in [Table VIII](#), along with situations for which they are important. As a result the higher-order processes are most useful in generating radiation in spectral regions, such as the vacuum ultraviolet or far infrared, that are inaccessible by the second-order interactions, or for certain applications such as phase conjugation or spectroscopy.

2. Stimulated Scattering Processes

Simulated scattering processes are nonlinear interactions in which an incident wave at frequency ω_{inc} is converted

to a scattered wave at a different frequency ω_{scat} . The difference in photon energy between the incident and scattered frequencies is taken up or supplied by the nonlinear medium, which undergoes a transition between two of its internal energy levels, as illustrated in [Fig. 18](#).

If the medium is initially in its ground state, the scattered wave is at a lower frequency (longer wavelength) than the incident wave, and the medium is excited to one of its internal energy levels during the interaction. In this situation the frequency shift is termed a Stokes shift, in analogy to the shift to lower frequencies that is observed in fluorescence (which was explained by Sir George Stokes), and the scattered wave is termed a Stokes wave. The incident (laser) and scattered (Stokes) frequencies are related by

$$\omega_S = \omega_L - \omega_0, \quad (34)$$

where ω_L and ω_S are the frequencies of the laser and Stokes waves and ω_0 is the frequency of the internal energy level of the medium.

If the medium is initially in an excited state, the scattered wave is at a higher frequency (shorter wavelength) than the incident wave and the medium is deexcited during the interaction, with its energy being given to the scattered wave. In this situation the scattered wave is termed an anti-Stokes wave (shifts to higher frequencies are not possible in fluorescence, as explained by Stokes), and the frequency shift is termed the anti-Stokes shift. The laser and anti-Stokes frequencies are related by

$$\omega_{AS} = \omega_L + \omega_0 \quad (35)$$

Various types of stimulated scattering processes are possible, each involving a different type of internal

TABLE VIII Competing Processes for Third- and Higher-Order Frequency Conversion^a

Competing process	Effect on conversion efficiency	Conditions under which competing process can be expected to be important
Linear absorption of generated radiation	Loss of generated power Reduction of improvement from phase matching	UV or XUV generation in ionizing continuum of nonlinear medium Generated wavelength close to allowed transition
Nonlinear absorption of pump radiation	Limitation on product NL Loss of pump power Saturation of susceptibility Disturbance of phase-matching conditions Self-focusing or self defocusing	Two-photon resonant interactions
Stark shift	Saturation of susceptibility Self-focusing or self defocusing	Resonant or near-resonant interactions, with pump intensity close to or greater than the appropriate saturation intensity
Kerr effect	Disturbance of phase-matching conditions Self-focusing or self defocusing	Nonresonant interactions Near-resonant interactions when the pump intensity is much less than the saturation intensity
Dielectric breakdown, multiphoton ionization	Disturbance of phase-matching conditions Saturation of susceptibility Loss of power at pump or generated wavelength	Conversion in low-pressure gases at high intensities Tightly focused geometries

^a [Reproduced from Reintjes, J. (1985). Coherent ultraviolet and vacuum ultraviolet sources. In "Laser Handbook," Vol. 5 (M. Bass and M. L. Stitch, eds.), North-Holland, Amsterdam.]

excitation. Some of the more common ones are listed in Table IX, along with the types of internal excitations and the types of materials in which they are commonly observed. Stimulated Brillouin scattering involves interactions with sound waves in solids, liquids, or gases or ion-acoustic waves in plasmas, and stimulated Rayleigh scattering involves interactions with density or orientational fluctuations of molecules. Various forms of stimulated Raman scattering can involve interaction with molecular vibrations, molecular rotations (rotational Raman scattering), electronic levels of atoms or molecules

(electronic Raman scattering), lattice vibrations, polaritons, electron plasma waves, or nondegenerate spin levels in certain semiconductors in magnetic fields. The magnitude of the frequency shifts that occur depends on the combination of nonlinear interaction and the particular material that is used. Orders of magnitude for shifts in different types of materials for the various interactions are also given in Table IX.

Stimulated scattering processes that involve Stokes waves arise from third-order nonlinear interactions with nonlinear polarizations of the form

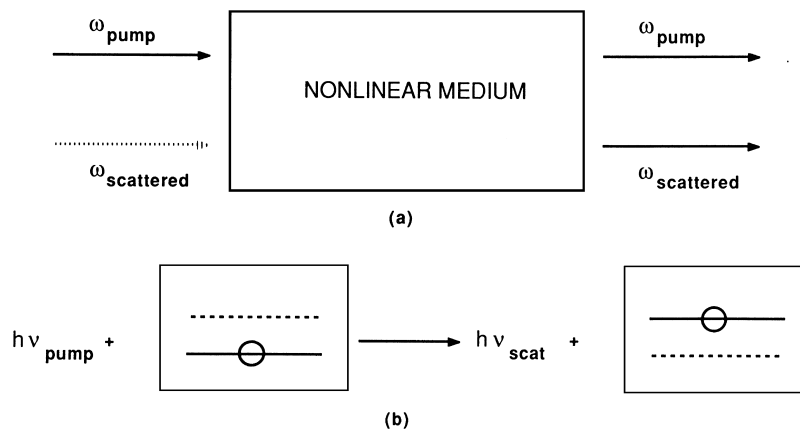


FIGURE 18 (a) Schematic illustration of a typical stimulated scattering interaction. The scattered wave can be supplied along with the pump radiation or it can be generated from noise in the interaction. (b) Representation of a stimulated Stokes scattering interaction, illustrating the transitions in the nonlinear medium.

TABLE IX Stimulated Scattering Interactions

Interaction	Internal mode	Type of material	Typical range of frequency shift $\Delta\nu$ (cm ⁻¹)
Stimulated Raman scattering	Molecular vibrations	Molecular gases	600–4150
		Molecular liquids	600–3500
	Molecular rotations	Molecular gases	10–400
	Electronic levels	Atomic and molecular gases, semiconductors	7000–20,000
	Lattice vibrations	Crystals	10–400
	Polaritons	Crystals	10–400
	Electron plasma waves	Plasmas	100–10,000
Stimulated Brillouin scattering	High-frequency sound waves	Solids, liquids, gases	100 MHz–10 GHz
	Ion-acoustic waves	Plasmas	0.1–10 MHz
Stimulated Rayleigh scattering	Density fluctuations, orientational fluctuations	Gases, liquids	0.1–1

$$P(\omega_S) = \frac{3}{2}\epsilon_0\chi^{(3)}(-\omega_S, \omega_L, -\omega_L, \omega_S)|A_L|^2 A_S. \quad (36)$$

The nonlinear susceptibility involves a two-photon resonance with the difference-frequency combination $\omega_L - \omega_S$. As with other two-photon resonances, the energy levels that are involved are of the same parity as those of the ground state in atoms or molecules with a center of inversion. When the two-photon resonance condition is satisfied, the susceptibility for the stimulated scattering interactions is negative and purely imaginary and can be written as

$$\chi^{(3)}(-\omega_S, \omega_L, -\omega_L, \omega_S) = -i\chi''^{(3)}(-\omega_S, \omega_L, -\omega_L, \omega_S). \quad (37)$$

Stimulated scattering processes can be in the forward direction, involving a scattered wave that propagates in the same direction as the incident laser, as shown in Fig. 18, or in the backward direction, involving a scattered wave that propagates in the opposite direction to the incident laser, as shown in Fig. 19. The field amplitude of the wave generated in the forward direction by a susceptibility of the type in Eq. (37) is described by the equation

$$dA_S/dz = (3\omega_S/4n_S c)\chi''|A_L|^2 A_S. \quad (38)$$

When pump depletion is negligible, the intensity of the Stokes wave is given by



FIGURE 19 Backward stimulated scattering in which the scattered radiation propagates in the direction opposite to the pump radiation. The scattered wave either can be supplied along with the pump radiation, or can be generated from noise in the interaction.

$$I_S = I_S(0)e^{gL}, \quad (39)$$

where the quantity g is the gain coefficient of the process given by

$$g = (3\omega_S/n_S n_L c^2 \epsilon_0)\chi''. \quad (40)$$

The waves generated in stimulated scattering processes have exponential growth, in contrast with the power-law dependences of the waves generated in the parametric interactions. The exponential gain is proportional to the propagation distance and to the intensity of the pump radiation. Phase matching is generally not required for stimulated scattering processes, since the phase of the material excitation adjusts itself for maximum gain automatically.

Photon number is conserved in stimulated scattering interactions, with one photon being lost in the pump wave for every one created in the Stokes wave. The energy created in the Stokes wave is smaller than that lost in the pump wave by the ratio ω_S/ω_L , termed the Manly–Rowe ratio, and the difference in photon energy between the pump and Stokes waves represents the energy given to the medium.

When the energy in the Stokes wave becomes comparable to that in the incident laser pump, depletion occurs and the gain is reduced. In principle, every photon in the incident pump wave can be converted to a Stokes photon, giving a maximum theoretical conversion efficiency of

$$\eta_{\max} = [I_S(L)/I_L(0)] = \omega_S/\omega_L. \quad (41)$$

In practice, the efficiency is never as high as indicated in Eq. (41) because of the lower conversion efficiency that is present in the low-intensity spatial and temporal wings of most laser beams. Photon conversion efficiencies of over 90% and energy-conversion efficiencies over 80% have, however, been observed in certain stimulated Raman and Brillouin interactions.

a. Stimulated Raman scattering. Stimulated Raman scattering can occur in solids, liquids, gases, and plasmas. It involves frequency shifts ranging from several tens of reciprocal centimeters for rotational scattering in molecules to tens of thousands of reciprocal centimeters for electronic scattering in gases. Forward stimulated Raman scattering is commonly used for generation of coherent radiation at the Stokes wavelength, amplification of an incident wave at the Stokes wavelength, reduction of beam aberrations, nonlinear spectroscopy (see Section IV.C), and generation of tunable infrared radiation through polariton scattering. Backward stimulated Raman scattering can be used for wave generation, amplification, pulse compression and phase conjugation (see Section IV.A).

The susceptibility for stimulated Raman scattering in molecules or atoms is given by

$$\chi = -\left[\frac{i}{6\Gamma\hbar^3} \left(1 - \frac{i\Delta}{\Gamma} \right) \right] \times \left[\sum \mu_{0i}\mu_{i2} \left\{ \frac{1}{\omega_{i0} - \omega_L} + \frac{1}{\omega_{i0} + \omega_S} \right\} \right]^2, \quad (42)$$

where $\Delta = \omega_0 - (\omega_L - \omega_S)$ is the detuning from the Raman resonance, Γ the linewidth of the Raman transition, ω_{i0} the frequency of the transition from level i to level 0, and μ_{0i} the dipole moment for the transition between levels 0 and i .

Gain coefficients and frequency shifts for some materials are given in Table X.

Amplification of an incident Stokes wave generally occurs for exponential gains up to about e^8 to e^{10} , corresponding to small signal amplifications of the order of 3000 to 22,000, although under some conditions stable gains up to e^{19} can be obtained. Raman amplifiers of this type are used to increase the power in the Stokes beam. When the pump beam has a poor spatial quality due to phase or amplitude structure, Raman amplification can be used to transfer the energy of the pump beam to the Stokes beam without transferring the aberrations, thereby increasing the effective brightness of the laser system.

Generation of a new wave at the Stokes frequency can be done in a single pass geometry, as shown in Fig. 20a, or in an oscillator configuration at lower gains, as shown in Fig. 20b. The threshold for single-pass generation depends on the degree of focusing that is used but is generally certain for gains above about e^{23} (10^{10}). Once the threshold for single-pass generation is reached, the process generally proceeds to pump depletion very quickly. As a result, the Stokes frequency that is generated in a Raman oscillator usually involves the Raman-active mode with the highest gain.

TABLE X Stimulated Raman Shifts and Gain Coefficients at 694.3 nm

Material	$\Delta\nu_R$ (cm ⁻¹)	$g \times 10^3$ (cm/MW)
Liquids		
Carbon disulfide	656	24
Acetone	2921	0.9
Methanol	2837	0.4
Ethanol	2928	4.0
Toluene	1002	1.3
Benzene	992	3
Nitrobenzene	1345	2.1
N ₂	2326	17
O ₂	1555	16
Carbon tetrachloride	458	1.1
Water	3290	0.14
Gases		
Methane	2916	0.66 (10 atm, 500 nm)
Hydrogen	4155 (vibrational)	1.5 (above 10 atm)
	450 (rotational)	0.5 (above 0.5 atm)
Deuterium	2991 (vibrational)	1.1 (above 10 atm)
N ₂	2326	0.071 (10 atm, 500 nm)
O ₂	1555	0.016 (10 atm, 500 nm)

If sufficient conversion to the Stokes wave takes place, generation of multiple Stokes waves can occur. In this situation the first Stokes wave serves as a pump for a second Stokes wave that is generated in a second stimulated Raman interaction. The second Stokes wave is shifted in frequency from the first Stokes wave by ω_0 . If sufficient intensity is present in the original pump wave, multiple Stokes waves can be generated, each shifted from the preceding one by ω_0 as illustrated in Fig. 21. Stimulated Raman scattering can thus be used as a source of coherent radiation at several wavelengths by utilizing multiple Stokes shifts, different materials, and different pump wavelengths.

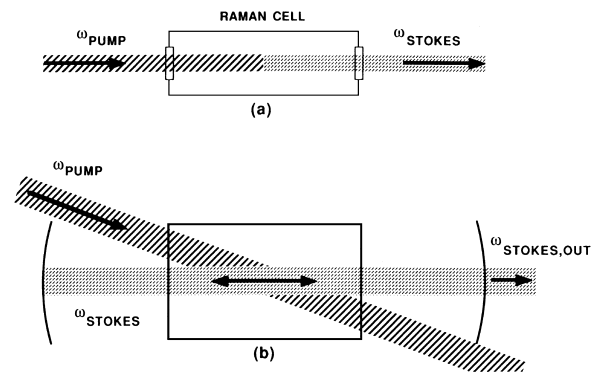


FIGURE 20 (a) Generation of a Stokes wave in single-pass stimulated Raman scattering. (b) Generation of a Raman-Stokes wave in a Raman laser oscillator.

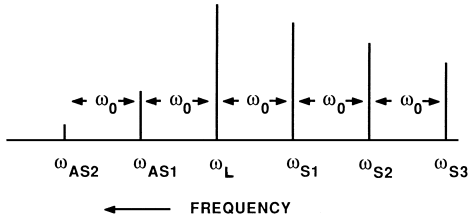


FIGURE 21 Schematic illustration of the spectrum produced by multiple Stokes and anti-Stokes scattering.

Continuously tunable Stokes radiation can be generated by using a tunable laser for the pump radiation or in some situations by using materials with internal modes whose energies can be changed. An example of this type of interaction is the generation of tunable narrow-band infrared radiation from spin-flip Raman lasers in semiconductors, as illustrated in Fig. 22. In these lasers the energy levels of electrons with different spin orientations are split by an external magnetic field. The Raman process involves a transition from one spin state to the other. For a fixed pump wavelength—for example, from a CO₂ laser at 10.6 μm —the wavelength of the Stokes wave can be tuned by varying the magnetic field, which determines the separation of the electron energy levels. Multiple Stokes shifts and anti-Stokes shifts can also be obtained. A list of materials, laser sources, and tuning ranges for various spin-flip lasers is given in Table XI. Radiation that is unstable in bands between 5.2 and 16.2 μm has been generated in this manner with linewidths as narrow as 100 MHz. This is among the narrowest bandwidth infrared radiation available and has been used for high resolution spectroscopy.

Backward stimulated Raman scattering involves the generation or amplification of a Stokes wave that travels in the opposite direction to the pump wave, as was

TABLE XI Spin-Flip Raman Lasers

Pump laser (wavelength, μm)	Material	Tuning range (μm)	Raman order
NH ₃ (12.8)	InSb	13.9–16.8	I Stokes
CO ₂ (10.6)	InSb	9.0–14.6	I, II, III, Stokes I, Anti-Stokes
CO (5.3)	InSb	5.2–6.2	I, II, III, Stokes I, Anti-Stokes
HF	InAs	3–5	

shown in Fig. 19. Backward stimulated Raman scattering requires radiation with a much narrower bandwidth than does forward Raman scattering and is usually observed for laser bandwidths less than about 10 GHz. The backward traveling Stokes wave encounters a fresh undepleted pump wave as it propagates through the nonlinear medium. As a result, the gain does not saturate when the pump wave is depleted, as it does in the forward direction. The peak Stokes intensity can grow to be many times of the incident laser intensity, while the duration of the Stokes pulse decreases relative to that of the incident pump pulse, allowing conservation of energy to be maintained. This process, illustrated in Fig. 23, is termed pulse compression. Compression ratios of the order of 50:1, producing pulses as short as a few nanoseconds, have been observed.

Anti-Stokes Raman scattering involves the generation of radiation at shorter wavelengths than those of the pump wave. Anti-Stokes scattering can occur in one of two ways. The more common method involves a four-wave difference frequency mixing process of the form

$$\omega_{\text{AS}} = 2\omega_{\text{L}} - \omega_{\text{S}} \quad (43a)$$

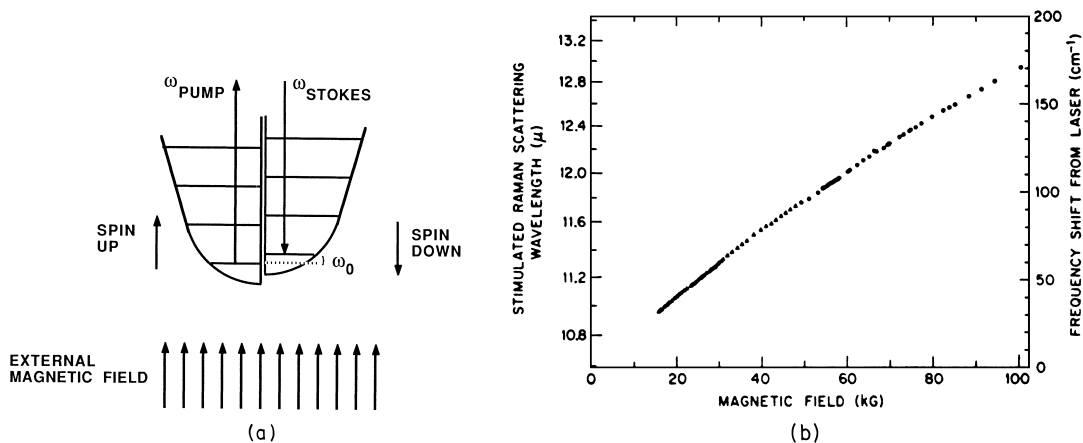


FIGURE 22 (a) Level diagram for a spin-flip Raman laser. (b) Typical tuning curve for an InSb spin-flip Raman laser pumped by a CO₂ laser at 10.6 μm . [Part (b) reproduced from Patel, C. K. N., and Shaw, E. D. (1974). *Phys. Rev. B* 3, 1279.]

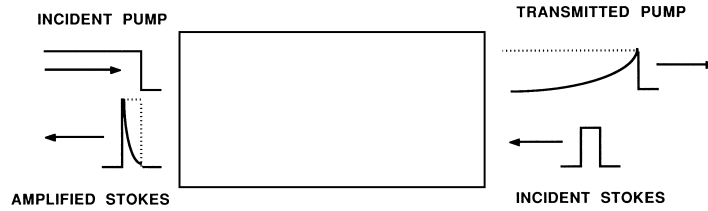


FIGURE 23 Illustration of pulse compression with backward stimulated Raman scattering. Backward-traveling Stokes pulse sweeps out the energy of the pump in a short pulse.

(see Section III.A.I.b) in media without a population inversion. In this interaction the Stokes radiation is generated with exponential gain through the stimulated Raman interaction as described above. The anti-Stokes radiation is generated by the four-wave mixing process, using the Stokes radiation as one of the pump waves. The anti-Stokes radiation that is generated through this interaction grows as part of a mixed mode along with the Stokes radiation. It has the same exponential gain as does the Stokes radiation, but the amplitude of the anti-Stokes radiation depends on the phase mismatch just as for other four-wave mixing interactions. The anti-Stokes generation is strongest for interactions that are nearly phase matched, although neither wave has exponential gain at exact phase matching. For common liquids and gases that have normal dispersion, the anti-Stokes process is not phase matched in the forward direction but is phase matched when the Stokes and anti-Stokes waves propagate at angles to the pump wave, as shown schematically in Fig. 24a. The anti-Stokes radiation is thus generated in cones about the pump radiation in most of these materials. The opening angle of the cone depends on the dispersion in the medium and on the frequency shift. It is of the order of several tens of milliradians for molecular vibrational shifts in liquids and can be considerably smaller for rotational shifts in molecular gases. An example of anti-Stokes emission in H_2 is shown in Fig. 24b. Here the pump radiation was at 532 nm. The anti-Stokes radiation at 435.7 nm was concentrated near the phase-matching direction of about 7 mrad, with a dark band appearing at the exact phase-matching direction and bright emission appearing in narrow bands on either side of phase matching.

Multiple anti-Stokes generation can occur through interactions of the form

$$\omega_{AS,n} = \omega_1 + \omega_2 - \omega_3, \quad (43b)$$

where $\omega_1, \omega_2, \omega_3$ are any of the Stokes, anti-Stokes, or laser frequencies involved in the interaction that satisfy the relations

$$\omega_1 - \omega_3 = \omega_0 \quad (44a)$$

$$\omega_{AS,n} - \omega_2 = \omega_0. \quad (44b)$$

Just as with multiple Stokes generation, the successive anti-Stokes lines are shifted from the preceding one by ω_0 as shown in Fig. 21. Multiple Stokes and anti-Stokes Raman scattering in molecular gases have been used to generate radiation ranging from 138 nm in the ultraviolet to wavelengths in the infrared. Some of the combinations of lasers and materials are listed in Table XII.

In media with a population inversion between the ground state and an excited Raman-active level, radiation at the anti-Stokes wavelength can be produced through a process similar to the one just described for Stokes generation in media that start from the ground state. This combination, illustrated in Fig. 25, is termed an anti-Stokes Raman laser. It has been used to generate

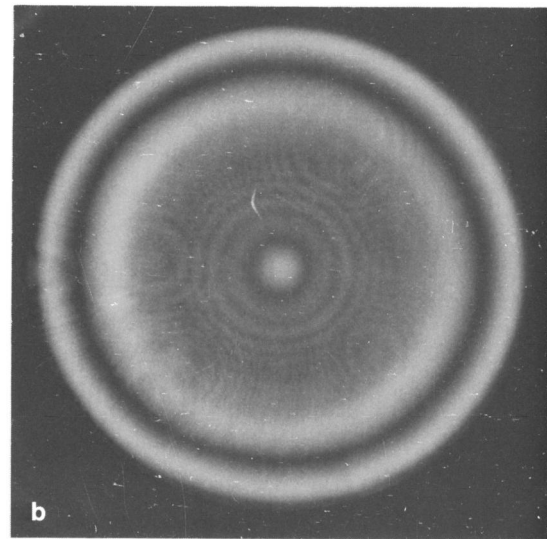
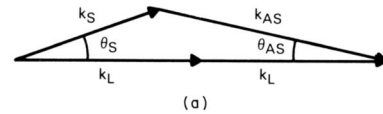


FIGURE 24 (a) Off-axis phase-matching diagram for anti-Stokes generation. (b) Anti-Stokes rings produced in stimulated Raman scattering from hydrogen. The dark band in the anti-Stokes ring is caused by suppression of the exponential gain due to the interaction of the Stokes and anti-Stokes waves. [Reproduced from Duncan, M. O., *et al.* (1986). *Opt. lett.* **11**, 803. Copyright © 1986 by the Optical society of America.]

TABLE XII Wavelength Range (nm) of Tunable UV and VUV Radiation Generated by Stimulated Raman Scattering in H_2 ($\Delta\nu = 4155 \text{ cm}^{-1}$) with Dye Lasers^a

Process (order)	Pump wavelength (nm)					
	600–625 (Rh 101) ^b	570–600 (Rh B) ^c	550–580 (Rh 6G) ^b	300–312.5 (Rh 101, SH) ^b	275–290 (Rh 6G, SH) ^b	548 (Fluorescein 27) ^c
AS (13)						138
AS (12)						146
AS (11)						156
AS (10)						167
AS (9)	185.0–187.3					179
AS (8)	200.4–203.1	196.9–200.4				194
AS (7)	218.6–221.8	214.4–218.6	194.5–198.1			211
AS (6)	240.4–244.3	235.4–240.4	211.6–215.9			231
AS (5)	267.1–271.9	261.0–267.1	256.7–263.0	184.8–189.5		256
AS (4)	300.4–306.6	292.7–300.4	287.3–295.3	200.2–205.7	188.7–189.5	286
AS (3)	343.3–351.3	332.2–343.3	326.3–336.6	218.3–224.9	204.8–213.0	325
AS (2)			377.5–391.4	240.1–248.1	223.8–233.7	376
AS (1)				266.7–276.6	246.8–258.8	
S (1)				342.7–359.1	310.5–329.7	
S (2)					356.5–382.1	

^a From Reintjes, J. (1985). Coherent ultraviolet and vacuum ultraviolet sources, in "Laser Handbook," Vol. 5 (M. Bass and M.L. Stitch, eds.), p. 1, North-Holland, Amsterdam.]

^b [Data from Wilke and Schmidt (1979). *Appl. Phys.* **18**, 177.]

^c [Data from Schomburg *et al.* (1983). *Appl. Phys. B* **30**, 131.]

radiation at wavelengths ranging from 149 to 378 nm using transitions in atomic gases such as Tl, I, or Br.

b. Stimulated Brillouin scattering. Stimulated Brillouin scattering (SBS) involves scattering from high-frequency sound waves. The gain for SBS is usually

greatest in the backward direction and is observed most commonly in this geometry, as shown in Fig. 26a. The equations describing backward SBS are

$$dI_B/dz = -g_B I_L I_B \quad (45a)$$

$$dI_L/dz = -g_B (\omega_L/\omega_B) I_L I_B, \quad (45b)$$

where the intensity gain coefficient g_B is given by

$$g_B = \frac{\omega_B^2 \rho_0 (\partial \epsilon / \partial \rho)^2}{4\pi c^3 n \nu \Gamma_B \epsilon_0}, \quad (45c)$$

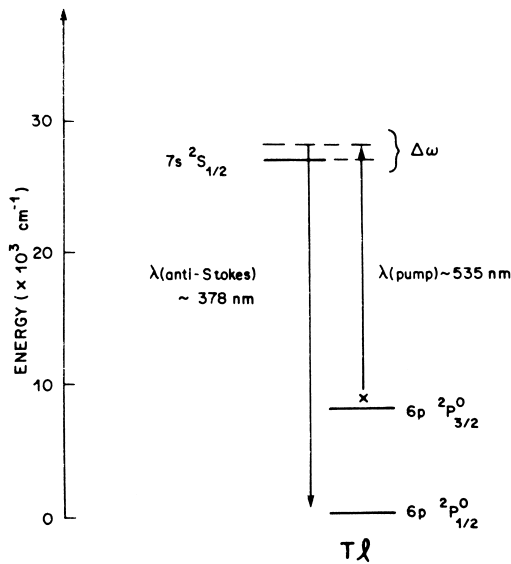


FIGURE 25 Level diagram of anti-Stokes Raman laser in Tl. [Reprinted from White, J. C., and Henderson, D. (1982). *Opt. Lett.* **7**, 517. Copyright ©1982 by the Optical society of America.]

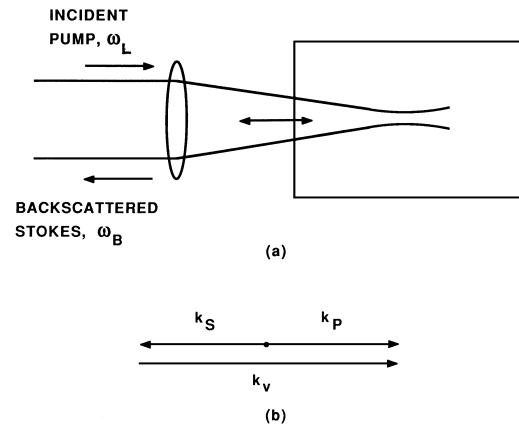


FIGURE 26 (a) Schematic diagram of a typical configuration used for stimulated Brillouin scattering. (b) The k -vector diagram for stimulated Brillouin scattering.

where ρ is the density, v the velocity of sound, $\partial\epsilon/\partial\rho$ the change of the dielectric constant with density, and Γ_B is the linewidth of the Brillouin transition. The wave-vector diagram is shown in Fig. 26b. The incident and generated optical waves are in the opposite directions, and the wave vector of the acoustic wave is determined by the momentum matching condition

$$\mathbf{k}_v = \mathbf{k}_L - \mathbf{k}_S \approx 2\mathbf{k}_L, \quad (46)$$

where \mathbf{k}_v is the \mathbf{k} vector of the sound wave and \mathbf{k}_L and \mathbf{k}_S are the wave vectors of the laser and scattered waves. Because the speed of light is so much greater than the speed of sound, the magnitude of the \mathbf{k} vector of the incident wave is almost equal to that of the scattered wave. The corresponding frequency shift of the scattered wave, termed the Brillouin shift, is equal to the frequency of the sound wave generated in the interaction and is given by

$$\Delta\omega_B = 2\omega_L v/c, \quad (47)$$

where v is the velocity of sound and we have used the approximation that $\mathbf{k}_L = -\mathbf{k}_S$.

Stimulated Brillouin scattering can be observed in liquids, gases, and solids, and also in plasmas, in which the scattering is from ion-acoustic waves. The SBS shift generally ranges from hundreds of megahertz in gases to several tens of gigahertz in liquids, depending on the particular material and laser wavelength involved. The SBS shifts for some materials are given in Table XIII. The threshold intensity for SBS ranges from 10^8 to 10^{10} W/cm², depending on the nonlinear material and the focusing conditions, and maximum reflectivities (the ratio of the intensity of the scattered wave to that of the incident pump wave) commonly exceed 50%. The acoustic waves generated in these interactions are among the most intense high-frequency sound waves generated.

Because the response time of the acoustic phonons is relatively long (of the order of 1 nsec in liquids and up to several hundred nanoseconds in some gases), SBS is ob-

served most commonly with laser radiation with a bandwidth less than 0.1 cm^{-1} . Generally, SBS has a higher gain than stimulated Raman scattering in liquids and usually dominates the interaction for wave generation when the laser radiation consists of narrow-band pulses that are longer than the response time of the acoustic phonon. Stimulated Raman scattering is commonly observed in these materials only for relatively broad-band radiation (for which the SBS interaction is suppressed), for short-duration pulses (for which SBS does not have time to grow), or at the beginning of longer-duration, narrow-band pulses.

In liquids and gases, SBS is used most commonly for phase conjugation (see Section IV.C) and for pulse compression in a manner similar to that described above for stimulated Raman scattering. SBS in solids can also be used for these purposes but is less common because the materials can be easily damaged by the acoustic wave that is generated in the medium.

B. Self-Action Effects

Self-action effects are those that affect the propagation characteristics of the incident light beam. They are due to nonlinear polarizations that are at the same frequency as that of the incident light wave. Depending on the particular effect, they can change the direction of propagation, the degree of focusing, the state of polarization or the bandwidth of the incident radiation, as was indicated in Table I. Self-action effects can also change the amount of absorption of the incident radiation. Sometimes one of these effects can occur alone, but more commonly two or more of them occur simultaneously.

The most common self-action effects arise from third-order interactions. The nonlinear polarization has the form

$$P(\omega) = \frac{3}{4}\epsilon_0\chi^{(3)}(-\omega, \omega, -\omega, \omega)|A|^2 A. \quad (48)$$

The various types of self-action effects depend on whether the susceptibility is real or imaginary and on the temporal and spatial distribution of the incident light. Interactions that change the polarization vector of the radiation depend on the components of the polarization vector present in the incident radiation, as well as on the tensor components of the susceptibility.

The real part of the nonlinear susceptibility in Eq. (48) gives rise to the spatial effects of self-focusing and self-defocusing, spectral broadening, and changes in the polarization vector. The imaginary part of the susceptibility causes nonlinear absorption.

1. Spatial Effects

The real part of the third-order susceptibility in Eq. (48) causes a change in the index of refraction of the material

TABLE XIII Stimulated Brillouin Shifts and Gain Coefficients at 1.064 μm

Material	$\Delta\nu_B$ (GHz)	g (cm/MW)
Carbon disulfide	3.84	0.13–0.16
Methanol	2.8	0.014
Ethanol	3	0.012
Toluene	3.9	0.013
Benzene	4.26	0.021
Acetone	3	0.019
<i>n</i> -Hexane	2.8	0.023
Cyclohexane	3.66	0.007
Carbon tetrachloride	1.9	0.007
Water	3.75	0.006

according to the relation

$$n = n^L + n_2 \langle E^2 \rangle = n^L + \frac{1}{2} n_2 |A|^2, \quad (49)$$

where

$$n_2 = (3/4n^L) \chi'. \quad (50)$$

In these equations, $\langle E^2 \rangle$ is the time average of the square of the total electric field of Eq. (3), which is proportional to the intensity, n^L is the linear refractive index, χ' is the real part of χ , and n_2 is termed the nonlinear refractive index.

Self-focusing occurs as a result of a combination of a positive value of n_2 and an incident beam that is more intense in the center than at the edge, a common situation that occurs, for example, as a result of the spatial-mode structure of a laser. In this situation the refractive index at the center of the beam is greater than that at its edge and the optical path length for rays at the center is greater than that for rays at the edge. This is the same condition that occurs for propagation through a focusing lens, and as a result the light beam creates its own positive lens in the nonlinear medium. As the beam focuses, the strength of the nonlinear lens increases, causing stronger focusing and increasing the strength of the lens still further. This behavior results in catastrophic focusing, in which the beam collapses to a very intense, small spot, in contrast to the relatively gentle focusing that occurs for normal lenses, as illustrated in Fig. 27.

Self-focusing can occur in any transparent material at sufficiently high intensities and has been observed in a wide range of materials, including glasses, crystals, liquids, gases, and plasmas. The mechanism causing self-focusing varies from one material to another. In solids and some gases the nonlinear index is due to interaction with the electronic energy levels which causes a distortion of the electron cloud, which results in an increase in the refractive index. In materials such as molecular liquids with anisotropic molecules, the nonlinear index arises from orientation of the molecules so that their axis of easy

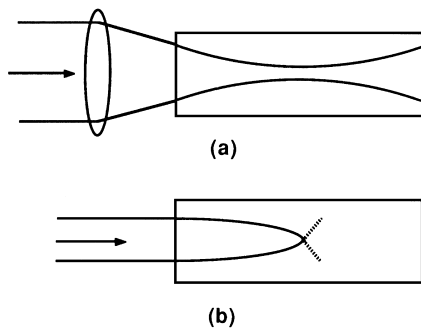


FIGURE 27 Schematic of focal region produced by (a) a normal and (b) a self-focusing lens.

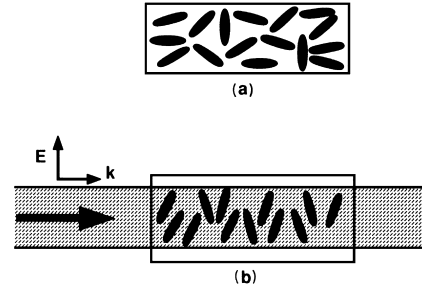


FIGURE 28 (a) Random molecular orientation in a liquid with anisotropic molecules produces an isotropic refractive index. (b) Partial alignment with a laser beam produces a nonlinear index and optically induced birefringence.

polarization is aligned more closely along the polarization vector of the incident field, as shown in Fig. 28. In such materials the molecules are normally arranged randomly, resulting in an isotropic refractive index. When the molecules line up along the optical field, the polarizability increases in that direction, resulting in both an increase in the refractive index for light polarized in the same direction as the incident field and, because the change in the refractive index is less for light polarized perpendicular to the incident radiation, birefringence. This effect is termed the optical Kerr effect, and the materials in which it occurs are termed Kerr-active. Self-focusing is observed most commonly in these materials.

A nonlinear index can also arise from electrostriction, in which the molecules of the medium move into the most intense regions of the electric field. The resulting increase in density causes an increase in the refractive index near the regions of intense fields. Because of the relatively slow response time of moving molecules, electrostriction has a longer time constant than molecular orientation and is typically important only for pulses that last for several tens to several hundreds of nanoseconds or longer.

Self-focusing in plasmas occurs because of a form of electrostriction in which the electrons move away from the most intense regions of the beam. Because the electrons make a negative contribution of the refractive index, the change in their density distribution results in a positive lens.

In order for a beam to self-focus, the self-focusing force must overcome the tendency of the beam to increase in size due to diffraction. This requirement leads to the existence of a critical power defined by

$$P_c = 0.04 \epsilon_0 \lambda^2 c / n_2. \quad (51)$$

For incident powers above the critical power, the self-focusing overcomes diffraction and a beam of radius a focuses at a distance given by

$$z_f = \frac{0.369ka^2}{\sqrt{P/P_c} - 0.858}. \quad (52)$$

TABLE XIV Self-Focusing Parameters for Selected Materials

Material	$n_2 \times 10^{-22}$ (MKS units)	Critical power at 1.064 μm (kW)
Carbon disulfide	122	35.9
Methanol	41	106
Ethanol	32	136
Toluene	29	152
Benzene	21	207
Acetone	3.4	1274
<i>n</i> -Hexane	2.6	1645
Cyclohexane	2.3	1880
Carbon tetrachloride	1.8	2194
Water	1.4	3038
Cesium vapor	-2.9×10^{-16} N	

For incident powers below the critical power, the self-focusing cannot overcome the spreading due to diffraction and the beam does not focus, although it spreads more slowly than it would in the absence of the nonlinear index. Values of the nonlinear index and the critical powers for some materials are given in Table XIV.

When the power in the incident beam is just above the critical power, the entire beam focuses as described above in a process termed whole-beam self-focusing. When the incident power exceeds the critical power significantly, the beam usually breaks up into a series of focal regions, each one of which contains one or a small number of critical powers. This behavior is termed beam break-up, and the resulting focusing behavior is termed small-scale self-focusing. If the incident beam contains a regular intensity pattern, the distribution of focal spots can be regular. More commonly, however, the pattern is random and is determined by whatever minor intensity variations are one the incident beam due to mode structure, interference patterns, or from imperfections or dust in or on optical materials through which it has propagated. An example of a regular pattern of self-focused spots developed on a beam with diffraction structure is shown in Fig. 29.

Once the self-focusing process has started, it will continue until catastrophic focus is reached at the distance given in Eq. (52). The minimum size of the focal point is not determined by the third-order nonlinear index but can be determined by higher-order terms in the nonlinear index, which saturate the self-focusing effect. Scuh saturation has been observed in atomic gases. For self-focusing in liquids, it is thought that other mechanisms, such as nonlinear absorption, stimulated Raman scattering, or multiphoton ionization, place a lower limit on the size of the focal spots. Minimum diameters of self-focal

spots are of the order of a few micrometers to a few tens of micrometers, depending on the material involved.

If the end of the nonlinear medium is reached before the catastrophic self-focal distance of Eq. (52), the material forms an intensity-dependent variable-focal-length lens. It can be used in conjunction with a pinhole or aperture to form an optical limiter or power stabilizer.

When the incident beam has a constant intensity in time, the focal spot occurs in one place in the medium. When the incident wave is a pulse that varies in time, the beam focuses to a succession of focal points, each corresponding to a different self-focal distance according to Eq. (52). This gives rise to a moving self-focal point, which, when observed from the side and integrated over the pulse duration, can have the appearance of a continuous track. In special cases the beam can be confined in a region of small diameter for many diffraction lengths in an effect termed self-trapping. This happens, for example, when the nonlinear index is heavily saturated, as, for example, in an atomic transition. When the pulse duration is short compared to the response time of the nonlinear index, which can vary from several picoseconds in common liquids to several hundred picoseconds in liquid crystals, the back end of the pulse can be effectively trapped in the index distribution set up by the front of the pulse. This behavior is termed dynamic self-trapping.

Self-focusing in solids is generally accompanied by damage in the catastrophic focal regions. This results in an upper limit on the intensity that can be used in many optical components and also results in a major limitation on the intensity that can be obtained from some solid-state pulsed lasers.

Because of the tensor nature of the nonlinear susceptibility, the intensity-dependent change in the refractive index is different for light polarized parallel and perpendicular to the laser radiation, resulting in optically induced birefringence in the medium. The birefringence can be used to change linearly polarized light into elliptically polarized light. This effect forms the basis of ultrafast light gates with picosecond time resolution that are used in time-resolved spectroscopy and measurements of the duration of short pulses. The birefringence also results in the rotation of the principal axis of elliptically polarized light, an effect used in nonlinear spectroscopy and the measurement of nonlinear optical susceptibilities.

2. Self-Defocusing

Self-defocusing results from a combination of a negative value of n_2 and a beam profile that is more intense at the center than at the edge. In this situation the refractive index is smaller at the center of the beam than at the edges, resulting a shorter optical path for rays at the center than

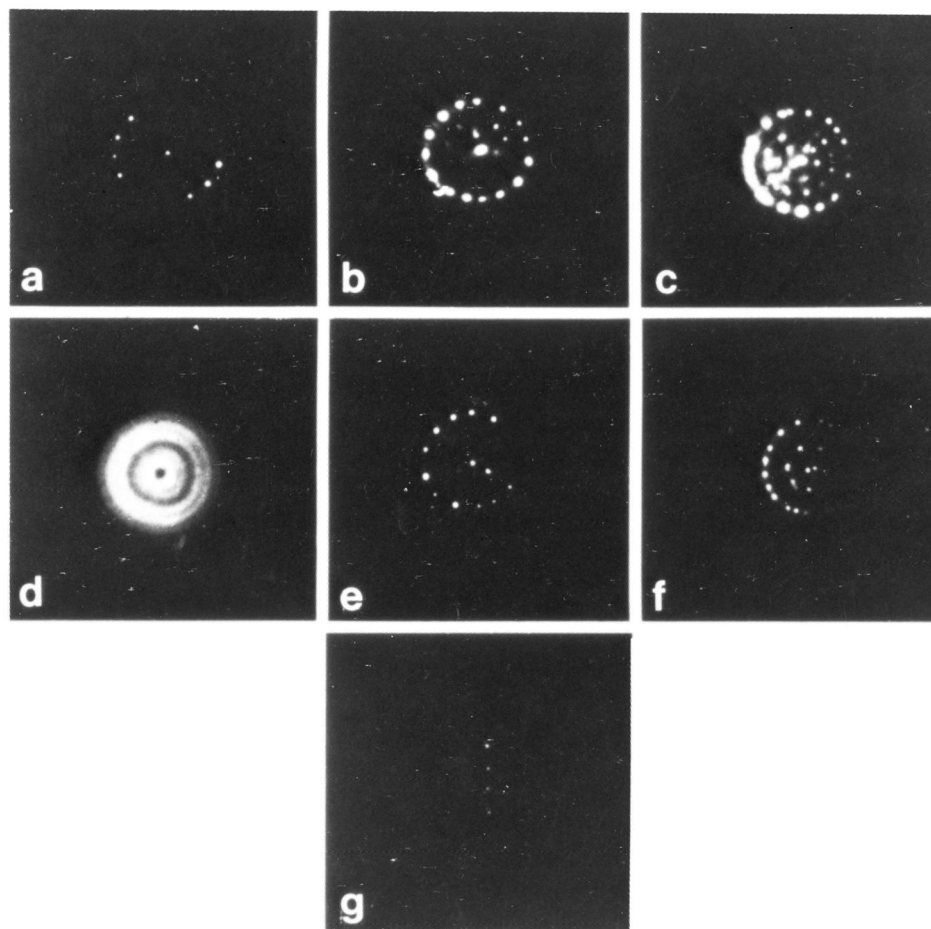


FIGURE 29 Pattern of self-focal spots (a–c, e–f) obtained from a beam with the ring pattern in (d). (g) Pattern produced by a beam with a horizontal straight edge. [Reproduced from Campillo, A. J., *et al.* (1973). *Appl. Phys. Lett.* **23**, 628.]

for those at the edge. This is the same condition that exists for propagation through a negative-focal-length lens, and the beam defocuses.

Negative values of the nonlinear refractive index can occur because of interaction with electronic energy levels of the medium when the laser frequency is just below a single-photon resonance or just above a two-photon resonance. Generally, self-defocusing due to electronic interactions is observed only for resonant interactions in gases and has been observed in gases for both single- and two-photon resonant conditions. A more common source of self-defocusing is thermal self-defocusing or, as it is commonly called, thermal blooming, which occurs in materials that are weakly absorbing. The energy that is absorbed from the light wave heats the medium, reducing its density, and hence its refractive index, in the most intense regions of the beam. When the beam profile is more intense at the center than at the edge, the medium becomes a negative lens and the beam spreads. Thermal blooming

can occur in liquids, solids, and gases. It is commonly observed in the propagation of high-power infrared laser beams through the atmosphere and is one of the major limitations on atmospheric propagation of such beams.

3. Self-Phase Modulation

Self-phase modulation results from a combination of a material with a nonlinear refractive index and an incident field amplitude that varies in time. Because the index of refraction depends on the optical intensity, the optical phase, which is given by

$$\phi = kz - \omega t = \frac{2\pi}{\lambda} \left[n^L + \frac{1}{2} n_2 |A(t)|^2 \right] z - \omega t, \quad (53)$$

develops a time dependence that follows the temporal variation of the optical intensity. Just as with other situations involving phase modulation, the laser pulse develops spectral side bands. Typical phase and frequency variations are

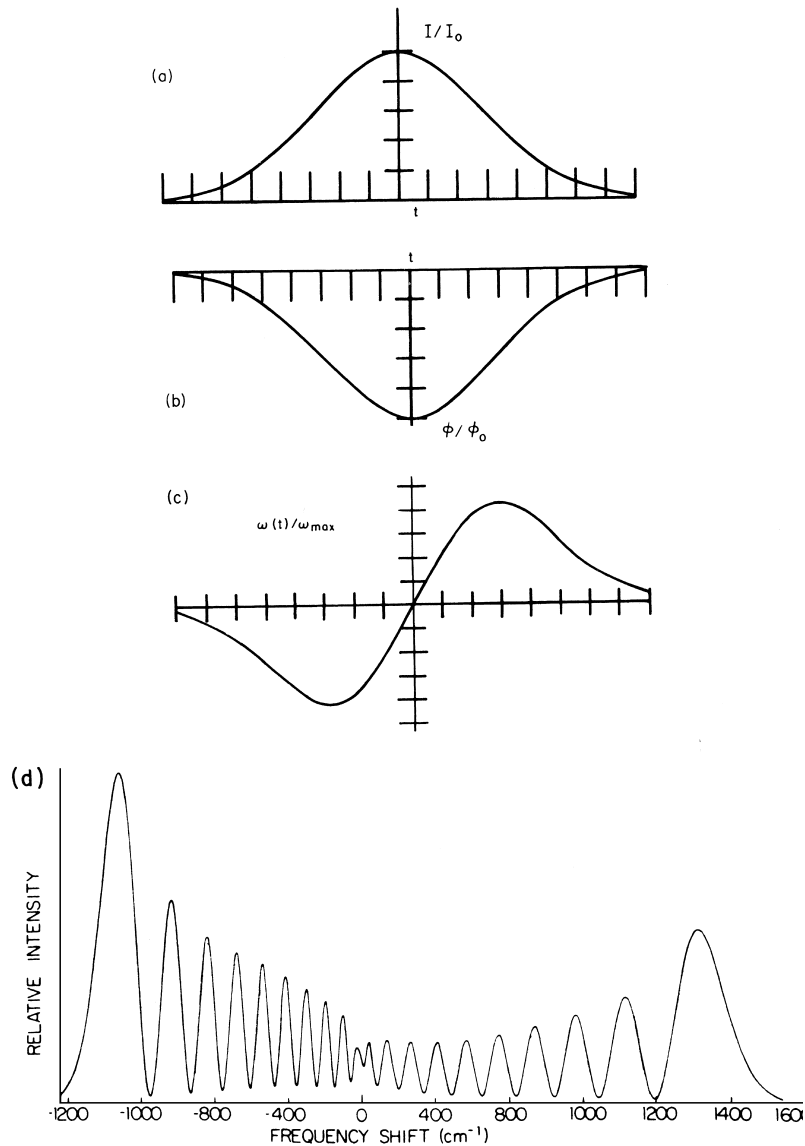


FIGURE 30 Temporal variation of (a) intensity, (b) phase, and (c) frequency produced by a medium with a positive nonlinear refractive index. (d) Spectrum of a self-phase-modulated pulse. The asymmetry results from a rise time that was 1.4 times faster than the full time. [Part (d) reproduced from Reintjes, J. (1984). "Nonlinear Optical Parametric Processes in Liquids and Gases," Academic Press, Orlando, FL.]

shown in Fig. 30 for a pulse that has a bell-shaped profile. The phase develops a bell-shaped temporal dependence, and the frequency, which is the time derivative of the phase, undergoes an oscillatory behavior as shown. For a medium with a positive n_2 , the down-shifted part of the spectrum is controlled by the leading part of the pulse and the up-shifted part of the spectrum is controlled by the trailing part of the pulse. The number of oscillations in the spectrum is determined by the peak phase modulation, while the extent of the spectrum is determined by the pulse duration and the peak phase excursion. For pulses generated in Q-switched lasers, which last for tens of nanosec-

onds, the spectral extent is relatively small, usually less than a wave number. For picosecond-duration pulses generated in mode-locked lasers, the spectrum can extend for hundreds of wave numbers and, in some materials such as water, for thousands of wave numbers.

In many instances, self-phase modulation represents a detrimental effect, such as in applications to high-resolution spectroscopy or in the propagation of a pulse through a dispersive medium when the self-phase modulation can cause spreading of the pulse envelope. In some instances, however, self-phase modulation can be desirable. For example, the wide spectra generated from

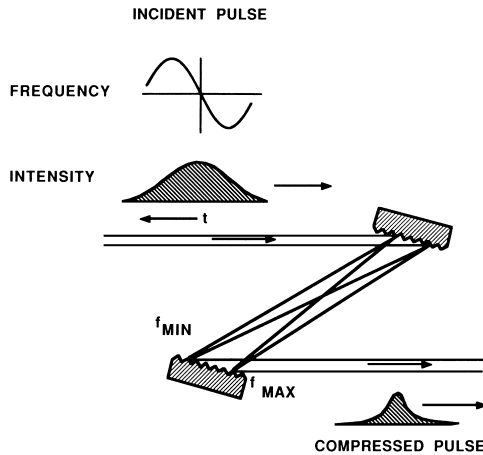


FIGURE 31 Pulse compression produced by a self-phase-modulated pulse and a pair of diffraction gratings. The dispersion of the gratings causes the lower-frequency components of the pulse to travel a longer path than that of the higher-frequency components, allowing the back of the pulse to catch up to the front.

picosecond pulses in materials such as water have been used for time-resolved spectroscopic studies. In other situations, the variation of the frequency in the center of a pulse, as illustrated in Fig. 30c, can be used in conjunction with a dispersive delay line formed by a grating pair to compress phase-modulated pulses in an arrangement as shown in Fig. 31, in a manner similar to pulse compression in chirped radar.

Pulse compressions of factors of 10 or more have been achieved with self-phase modulation of pulses propagated freely in nonlinear media resulting in the generation of subpicosecond pulses from pulses in the picosecond time range. This technique has also been very successful in the compression of pulses that have been phase-modulated in glass fibers. Here the mode structure of the fiber prevents the beam break-up due to self-focusing that can occur for large phase modulation in nonguided propagation. Pulse compressions of over 100 have been achieved in this manner and have resulted in the generation of pulses that, at 8 fsec (8×10^{-15} sec), are the shortest-duration optical pulses yet produced. Self-phase modulation in fibers, coupled with anomalous dispersion that occurs for infrared wavelengths longer than about $1.4 \mu\text{m}$, has also been used to produce solitons, which are pulses that can propagate long distances without spreading because the negative dispersion in the fiber actually causes the phase-modulated pulse to narrow in time. Soliton pulses are useful in the generation of picosecond-duration pulses and in long-distance propagation for optical-fiber communication.

4. Nonlinear Absorption

Self-action effects can also change the transmission of light through a material. Nonlinear effects can cause ma-

terials that are strongly absorbing at low intensities to become transparent at high intensities in an effect termed saturable absorption or, conversely, they can cause materials that are transparent at low intensities to become absorbing at high intensities in an effect termed multiphoton absorption.

Multiphoton absorption can occur through absorption of two, three, or more photons. The photons can be of the same or different frequencies. When the frequencies are different, the effect is termed sum-frequency absorption. Multiphoton absorption can occur in liquids, gases, or solids. In gases the transitions can occur between the ground state and excited bound states or between the ground state and the continuum. When the transition is to the continuum, the effect is termed multiphoton ionization. Multiphoton absorption in gases with atoms or symmetric molecules follow selection rules for multiple dipole transitions. Thus two-photon absorption occurs in these materials between levels that have the same parity. These are the same transitions that are allowed in stimulated Raman scattering but are not allowed in a single-photon transitions of linear optics. Multiphoton absorption in solids involves transitions to the conduction band or to discrete states in the band gap. In semiconductors, two- or three-photon absorption can be observed for near-infrared radiation, while for transparent dielectric materials multiphoton absorption is generally observed for visible and ultraviolet radiation. Multiphoton absorption increases with increasing laser intensity and can become quite strong at intensities that can be achieved in pulsed laser beams, often resulting in damage of solids or breakdown in gases. This can be one of the major limitations on the intensity of radiation that can be passed through semiconductors in the infrared or other dielectric materials in the ultraviolet.

The simplest form of multiphoton absorption is two-photon absorption. It is described by an equation of the form

$$dI/dz = -\beta I^2, \quad (54)$$

where β is the two-photon absorption coefficient. This equation has a solution of the form

$$I(L) = I_0/(1 + \beta I_0 L) \quad (55)$$

for the intensity transmitted through a material of length L , where I_0 is the intensity incident on the material at $z = 0$. The form of the solution is indicated graphically in Fig. 32. Note that the transmission as a function of distance is quite different from that encountered for linear absorption, for which the transmitted intensity decreases with distance as $e^{-\alpha L}$. In the limit of large values of $\beta I_0 L$, the transmitted intensity approaches the constant value $1/\beta L$, independent of the incident intensity. Two-photon absorption can thus be used for optical limiting. It can also be

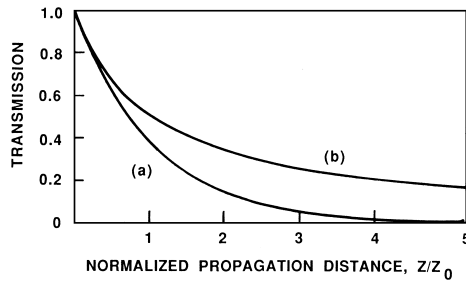


FIGURE 32 Dependence of transmitted intensity with distance for (a) linear absorption and (b) two-photon absorption. For (a) $z_0 = \alpha^{-1}$, while for (b) $z_0 = (\beta I_0)^{-1}$.

used for spectroscopy of atomic and molecular levels that have the same parity as the ground state and are therefore not accessible in linear spectroscopy. Finally, two-photon absorption can be used in Doppler-free spectroscopy in gases to provide spectral resolutions that are less than the Doppler width.

5. Saturable Absorption

Saturable absorption involves a decrease in absorption at high optical intensities. It can occur in atomic or molecular gases and in various types of liquids and solids, and it is usually observed in materials that are strongly absorbing at low light intensities. Saturable absorption occurs when the upper state of the absorbing transition gains enough population to become filled, preventing the transfer of any more population into it. It generally occurs in materials that have a restricted density of states for the upper level, such as atomic or molecular gases, direct bandgap semiconductors, and certain liquids, such as organic dyes. Saturable absorption in dyes that have relaxation times of the order of several picoseconds have been used to mode-lock solid-state and pulsed-dye lasers, producing optical pulses on the order of several tens of picoseconds or less. These dyes have also been used outside of laser cavities to shorten the duration of laser pulses. Saturable absorbers with longer relaxation times have been used to mode-lock cw dye lasers, producing subpicosecond pulses. Saturable absorption can also be used with four-wave mixing interactions to produce optical phase conjugation. Saturation of the gain in laser amplifiers is similar to saturable absorption, but with a change in sign. It is described by the same equations and determines the amount of energy that can be produced by a particular laser. In a linear laser cavity, the gain can be reduced in a narrow spectral region at the center of a Doppler-broadened line, which forms the basis of a spectroscopic technique known as Lamb dip spectroscopy that has a resolution less than the Doppler width.

C. Coherent Optical Effects

Coherent nonlinear effects involve interactions that occur before the wave functions that describe the excitations of the medium have time to relax or dephase. They occur primarily when the nonlinear interaction involves one- or two-photon resonances, and the duration of the laser pulse is shorter than the dephasing time of the excited state wave functions, a time that is equivalent to the inverse of the linewidth of the appropriate transition. Coherent nonlinear optical interactions generally involve significant population transfer between the states of the medium involved in the resonance. As a result, the nonlinear polarization cannot be described by the simple perturbation expansion given in Eq. (2), which assumed that the population was in the ground state. Rather, it must be solved for as a dynamic variable along with the optical fields.

Virtually any of the nonlinear effects that were described earlier can occur in the regime of coherent interactions. In this regime the saturation of the medium associated with the population transfer generally weakens the nonlinear process involved relative to the strength it would have in the absence of the coherent process.

Coherent interactions can also give rise to new nonlinear optical effects. These are listed in Table XV, along with some of their characteristics and the conditions under which they are likely to occur.

Self-induced transparency is a coherent effect in which a material that is otherwise absorbing becomes transparent to a properly shaped laser pulse. It occurs when laser pulses that are shorter than the dephasing time of the excited-state wave functions propagate through materials with absorbing transitions that are inhomogeneously broadened. In self-induced transparency, the energy at the beginning of the laser pulse is absorbed and is subsequently reemitted at the end of the laser pulse, reproducing the original pulse with a time delay. In order for all the

TABLE XV Coherent Nonlinear Interactions

Interaction	Conditions for observation
Self-induced transparency	Resonant interaction with inhomogeneously broadened transition; pulse duration less than dephasing time
Photon echoes	Resonant interaction with inhomogeneously broadened transition; pulse duration less than dephasing time; two pulses spaced by echo time τ , with pulse areas of $\pi/2$ and π , respectively
Adiabatic following	Near-resonant interaction; pulse duration less than dephasing time
Adiabatic rapid passage	Near-resonant interaction; pulse duration less than dephasing time; frequency of pulse swept through resonance

energy that is absorbed from the beginning of the pulse to be reemitted at the end, the pulse field amplitude must have the special temporal profile of a hyperbolic secant. In addition, the pulse must have the correct “area,” which is proportional to the product of the transition dipole moment and the time integral of the field amplitude over the pulse duration. The pulse area is represented as an angle that is equivalent to the rotation angle from its initial position of a vector that describes the state of the atom or molecule. In self-induced transparency, the required pulse area is 2π , corresponding to one full rotation of the state vector, indicating that the medium starts and ends in the ground state. If the incident pulse has an area greater than 2π it will break up into multiple pulses, each with area 2π , and any excess energy will eventually be dissipated in the medium. If the initial pulse has an area less than 2π it will eventually be absorbed in the medium.

Self-induced transparency is different from ordinary saturated absorption. In saturated absorption, the energy that is taken from the pulse to maintain the medium in a partial state of excitation is permanently lost to the radiation field. In self-induced transparency, the energy given to the medium is lost from the radiation field only temporarily and is eventually returned to it.

Photon echoes also occur in materials with inhomogeneously broadened transitions. In producing a photon echo, two pulses are used with a spacing of τ , as shown in Fig. 33. The first pulse has an area of $\pi/2$ and produces an excitation in the medium. The second pulse has an area of π and reverses the phase of the excited-state wave functions after they have had time to dephase. Instead of continuing to get further out of phase as time progresses, the wave functions come back into phase. At a time τ after the second pulse, the wave functions are again all in phase and a third pulse, termed the echo, is produced. Photon echoes are observed most easily when the pulsed spacing τ is larger than the inhomogeneous dephasing time caused, for example, by Doppler broadening, but smaller than the homogeneous dephasing time caused by collisions or population decay.

Other coherent interactions include optical nutation and free induction decay, in which the population oscillates be-

tween two levels of the medium, producing oscillations in the optical fields that are radiated, and adiabatic rapid passage, in which a population inversion can be produced between two levels in a medium by sweeping the frequency of a pulse through the absorption frequency in a time short compared to the dephasing time of the upper level.

D. Electrooptic and Magneto-optic Effects

Electrooptic and magneto-optic effects involve changes in the refractive index of a medium caused by an external electric or magnetic field. These are not normally thought of as nonlinear optical effects but are technically nonlinear optical processes in which the frequency of one of the fields is equal to zero. Various electrooptic and magneto-optic effects can occur depending on the situation. Some of these were listed in Table I.

In the Pockels effect, the change in the refractive index is proportional to the external electric field, whereas in the quadratic Kerr effect the change in refractive index is proportional to the square of the electric field. The Pockels effect occurs in solids without inversion centers, the same types that allow second-order nonlinear effects. The quadratic Kerr effect occurs in materials with any symmetry and is commonly used in liquids with anisotropic molecules such as nitrobenzene. The Faraday and Cotton–Mouton effects produce changes in the refractive index that are proportional to the magnetic field strength. Electrooptic and magneto-optic effects generally cause different changes in the refractive indexes in different directions relative to the applied field or the crystal axes, resulting in field-induced birefringence.

Electrooptic and magneto-optic effects are commonly used in light modulators and shutters. The field-dependent changes in the refractive index can be used directly for phase or frequency modulation. This is most commonly done with the Pockels effect, and units with modulation frequencies of up to the order of 1 GHz are commercially available. Field-induced birefringence is used to change the polarization state of a light beam and can be used for both modulation and shuttering. A light shutter can be constructed with either a Pockels cell or a Kerr cell

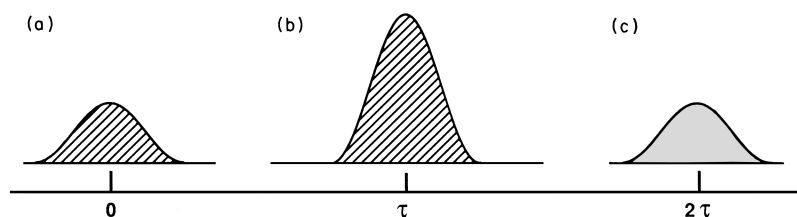


FIGURE 33 Arrangement of pulses for generating photon echoes: (a) with an area of $\pi/2$, and (b) with an area of π , are supplied spaced by time τ ; and (c) the echo, is generated in the interaction at a time τ after the second pulse.

by adjusting the field so that the polarization of the light wave changes by 90° . These are commonly used for producing laser pulses with controlled durations or shapes and for Q-switching of pulsed lasers. The Faraday effect produces a birefringence for circular polarization, resulting in the rotation of the direction of polarization of linearly polarized light. When adjusted for 45° and combined with linear polarizers, it will pass light in only one direction. It is commonly used for isolation of lasers from reflections from optical elements in the beam.

IV. APPLICATIONS

In the previous sections the basic nonlinear optical interactions have been described, along with some of their properties. In the following sections we shall describe applications of the various nonlinear interactions. Some applications have already been noted in the description given for certain effects. Here we shall describe applications that can be made with a wide variety of interactions.

A. Nonlinear Spectroscopy

Nonlinear spectroscopy involves the use of a nonlinear optical interaction for spectroscopic studies. It makes use of the frequency variation of the nonlinear susceptibility to obtain information about a material in much the same way that variation of the linear susceptibility with frequency provides information in linear spectroscopy. In nonlinear spectroscopy the spectroscopic information is obtained directly from the nonlinear interaction as the frequency of the pump radiation is varied. This can be contrasted with linear spectroscopy that is done with radiation that is generated in a nonlinear interaction and used separately for spectroscopic studies.

Nonlinear spectroscopy can provide information different from that available in linear spectroscopy. For example, it can be used to probe transitions that are forbidden in single-photon interactions and to measure the kinetics of excited states. Nonlinear spectroscopy can allow measurements to be made in a spectral region in which radiation needed for linear spectroscopy would be absorbed or perhaps would not be available. It can also provide increased signal levels or spectral resolution.

Many different types of nonlinear effects can be used for nonlinear spectroscopy, including various forms of parametric frequency conversion (harmonic generation, four-wave sum- and difference-frequency mixing), degenerate four-wave mixing, multiphoton absorption, multiphoton ionization, and stimulated scattering. Some of the effects that have been used for nonlinear spectroscopy are given in Table XVI, along with the information that is provided and the quantities that are varied and detected.

An example of improved spectral resolution obtained through nonlinear spectroscopy is the Doppler-free spectroscopy that can be done with two-photon absorption, as illustrated in Figs. 34 and 35. In this interaction, two light waves with the same frequency are incident on the sample from opposite directions. As the frequency of the incident light is swept through one-half of the transition frequency of a two-photon transition, the light is absorbed and the presence of the absorption is detected through fluorescence from the upper state to a lower one through an allowed transition. Normally, the spectral resolution of absorption measurements is limited by the Doppler broadening caused by the random motion of the atoms. Atoms that move in different directions with different speeds absorb light at slightly different frequencies, and the net result is an absorption profile that is wider than the natural width of the transition. In the nonlinear measurement, however, the atom absorbs one photon from each beam coming from opposite directions. A moving atom sees the frequency of one beam shifted in one direction by the same amount as the frequency of the other beam is shifted in the opposite direction. As a result, each atom sees the same sum of the frequencies of the two beams regardless of its speed or direction of motion, and the absorption profile is not broadened by the Doppler effect. This type of spectroscopy can be used to measure the natural width of absorption lines underneath a much wider Doppler profile. An example of a spectrum obtained with two-photon Doppler-free spectroscopy is shown in Fig. 35.

Nonlinear spectroscopy can also be used to measure the frequency of states that have the same parity as the ground state and therefore are not accessible through linear spectroscopy. For example, S and D levels in atoms can be probed as two-photon transitions in multiphoton absorption or four-wave mixing spectroscopy. Nonlinear spectroscopy can also be used for spectroscopic studies of energy levels that are in regions of the spectrum in which radiation for linear spectroscopy either is absorbed or is not available. Examples of such applications are spectroscopy of highly excited levels in the vacuum ultraviolet or extreme ultraviolet or of levels in the far infrared. Nonlinear effects can also be used for spectroscopic studies of excited levels, providing information on their coupling strength to other excited states or to the continuum. Time-resolved measurements can also be made to determine the kinetics of excited states such as lifetimes, dephasing times, and the energy-decay paths.

One of the most extensively used nonlinear processes for spectroscopy is coherent anti-Stokes Raman scattering (CARS). This is a four-wave mixing process of the form $\omega_{AS} = 2\omega_L - \omega_S$, where ω_L and ω_S are the laser and Stokes frequencies that are provided in the incident radiation and ω_{AS} is the anti-Stokes frequency that is generated

TABLE XVI Application of Nonlinear Optics to Spectroscopy^a

Nonlinear interaction	Quantity varied	Quantity measured	Information obtained
Multiphoton absorption Atom + $2\omega_{\text{laser}} \rightarrow \text{atom}^{*b}$	ω_{laser}	Fluorescence from excited level	Energy levels of states with same parity as ground state; sub-Doppler spectral resolution
Multiphoton ionization Atom + $2\omega_{\text{laser}} + \mathcal{E}^a \rightarrow \text{atom}^+$	ω_{laser}	Ionization current	Rydberg energy levels
Atom (molecule) + $(\omega_1 + \omega_2) \rightarrow \text{atom}^+ (\text{molecule}^+)^c$	ω_2	Ionization current	Even-parity autoionizing levels
Four-wave mixing Sum-frequency mixing $\omega_4 = 2\omega_1 + \omega_2$	ω_2	Optical power at ω_4	Energy structure of levels near ω_4 , e.g., autoionizing levels, matrix elements, line shapes
Third-harmonic generation $\omega_3 = 3\omega_1$	ω_1	Optical power at ω_3	Energy levels near $2\omega_1$ with same parity as ground state
Difference-frequency mixing $\omega_4 = 2\omega_1 - \omega_2$ (CARS)	ω_2	Optical power at ω_4	Raman energy levels Solids Polaritons Lattice vibrations Gases Measurement of nonlinear susceptibilities Solids Liquids Concentrations of liquids in solutions Concentrations of gases in mixtures Temperature measurements in gases Measurements in flames, combustion diagnostics Interference of nonlinear susceptibilities Time-resolved measurements Lifetimes, dephasing times CARS background suppression
Four-frequency CARS $\omega_4 = \omega_1 + \omega_2 - \omega_3$	$\omega_1 - \omega_3$ $\omega_2 - \omega_3$	Optical power at ω_4	Same as CARS, CARS background suppression
Raman-induced Kerr effect $\omega_{2,y} = \omega_1 - \omega_1 + \omega_{2,x}$	ω_2	Polarization changes at ω_2	Same as CARS, CARS background suppression
Higher-order Raman processes $\omega_4 = 3\omega_1 - 2\omega_2$	ω_2	Optical power at ω_4	Sames as CARS
Coherent Stokes-Raman spectroscopy $\omega_4 = 2\omega_1 - \omega_2 (\omega_1 < \omega_2)$	ω_2	Optical power at ω_4	Same as CARS
Coherent Stokes scattering $\omega_S = \omega_1(t) - \omega_1(t) + \omega_S(t) + \omega_1(t + \Delta t) - \omega_1(t + \Delta t)$	Δt	Optical power at $\omega_S(t + \Delta t)$	Lifetimes, dephasing times, resonant contributions to $\chi^{(3)}$
Raman gain spectroscopy $\omega_S = \omega_L - \omega_L + \omega_S$	ω_S	Gain or loss at ω_S	Ramann energy levels
Saturation spectroscopy	ω_{laser}	Induced gain or loss	High-resolution spectra

^a (From Reintjes, J. F. (1984). "Nonlinear Parametric Processes in Liquids and Gases," pp. 422–423, Academic Press, New York.)^b An atom or molecule in an excited state is designated by *.^c An ionized atom or molecule is designated by + and \mathcal{E}^a designates ionizing energy supplied by an external electric field.

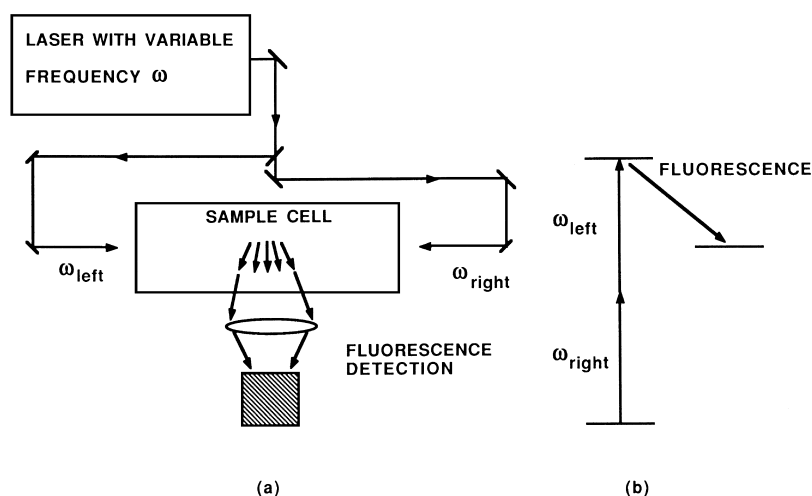


FIGURE 34 The use of two-photon absorption with counterpropagating beams for Doppler-free spectroscopy. (a) Experimental arrangement. (b) Level diagram of transitions used.

in the interaction. The process is resonantly enhanced when the difference frequency $\omega_L - \omega_S$ is tuned near a Raman-active mode with frequency $\omega_0 = \omega_L - \omega_S$, and the resulting structure in the spectrum provides the desired spectroscopic information. The Raman-active modes can involve transitions between vibrational or rotational levels in molecules or lattice vibrations in solids. Effective use of the CARS technique requires that phase-matching conditions be met. This is usually done by angle phase matching, as was shown in Fig. 24a. Coherent anti-Stokes Raman scattering offers advantages over spontaneous Raman scattering of increased signal levels and signal-to-noise ratios in the spectroscopy of pure materials, because the intensity of the generated radiation depends on the product $(NL)^2$, where N is the density and L is the length of the interaction region as compared with the (NL) dependence of the radiation generated in spontaneous Raman scattering. As a result, spectra can be obtained more quickly and with higher resolution with CARS. Coherent anti-Stokes Raman scattering has been used for measurements of Raman spectra in molecular gases, Raman spectra in flames for combustion diagnostics and temperature measurements, and in spatially resolved measurements for molecular selective microscopy, an example of which is shown in Fig. 36.

Although CARS offers the advantage of increased signal levels in pure materials for strong Raman resonances, it suffers from the presence of a nonresonant background generated by all the other levels in the medium. In order to overcome these limitations, several alternative techniques have been developed involving interactions, such as the Raman-induced Kerr effect, in which the sample develops birefringence, causing changes in the polarization of the incident waves as their frequency difference is tuned

through a resonance, and four-frequency CARS, in which the variation of the CARS signal near a resonance with one set of frequencies is used to offset the background in the other set.

B. Infrared Up-Conversion

Three-wave and four-wave sum-frequency mixing has been used for conversion of radiation from the infrared to the visible, where photographic film can be used and photoelectric detectors offer the advantages of improved sensitivity and increased temporal resolution. Infrared up-conversion has been used for up-conversion of infrared images to the visible and for time-resolved studies of the rotational spectra of molecules formed in explosions, giving information as to the time dependence of the temperature.

C. Optical Phase Conjugation

Optical phase conjugation, also referred to as time reversal or wavefront reversal, is a technique involving the creation of an optical beam that has the variations in its wavefront, or phase, reversed relative to a reference beam. If the optical field is represented as the product of an amplitude and complex exponential phase,

$$E = Ae^{i\phi}, \quad (56)$$

then the process of reversing the sign of the phase is equivalent to forming the complex conjugate of the original field, an identification that gives rise to the name phase conjugation. When optical phase conjugation is combined with a reversal of the propagation direction, it allows for compensation of distortions on an optical beam, which develop as a result of propagation through distorting media,

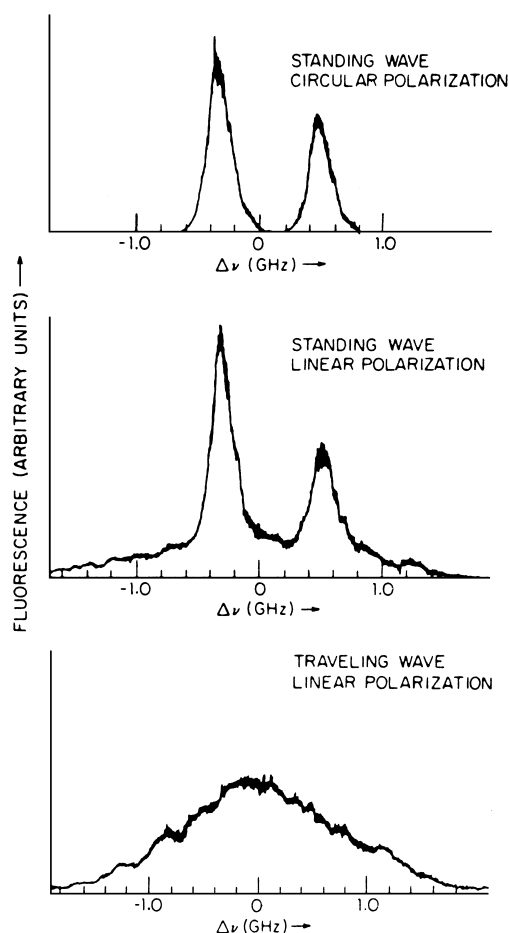


FIGURE 35 Example of high-resolution spectrum in sodium vapor obtained with two-photon Doppler-free spectroscopy. [Reproduced from Bloembergen, N., and Levenson, M. D. (1976). Doppler-free two-photon absorption spectroscopy. In "High Resolution Laser Spectroscopy" (K. Shimoda, ed.), p. 355, Springer, New York.]

for example, the atmosphere, or imperfect or low-quality optical components such as mirrors, lenses, or windows. Such distortions are familiar to people who have looked through old window glass, through the air above a hot radiator, or in the apparent reflections present on the highway on a hot day. Optical phase conjugation can also be used for holographic imaging and can allow images to be transmitted through multimode fibers without degradation due to the difference in phase velocity among the various modes. It can also be used in various forms of optical signal processing such as correlators and for spectroscopy.

The concept of correction of distortions by optical phase conjugation is as follows. When a wave propagates through a distorting medium, its phase contour acquires structure that will eventually diffract, leading to increased spreading, reduced propagation length, reduced focal-spot intensity, and image distortion. The basic idea of compen-

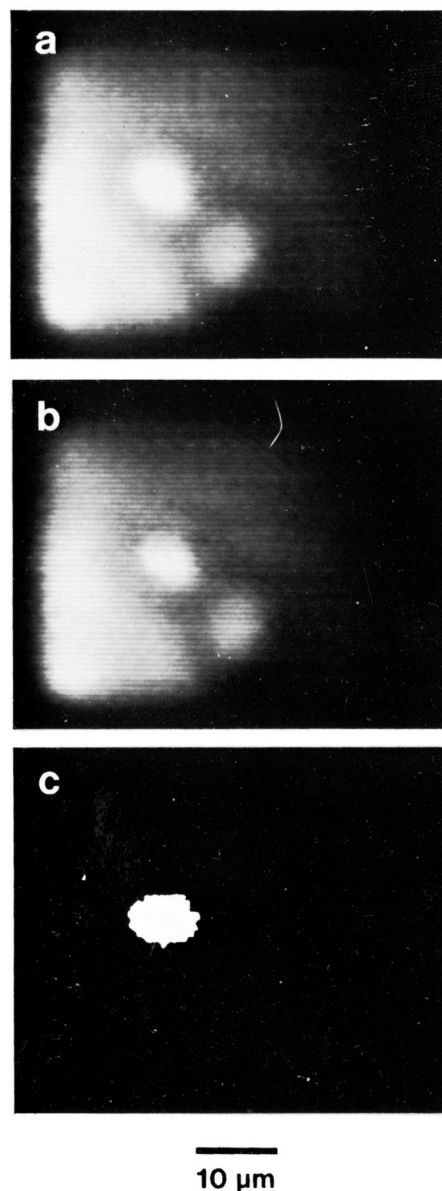


FIGURE 36 Molecular selective microscopy using coherent anti-Stokes Raman scattering. (a) On-resonant picture of deuterated and nondeuterated liposomes with two liposomes visible are shown. Only the deuterated liposomes have a Raman resonance for the radiation used. (b) The same picture as that in (a), but the two pump waves have been detuned from the Raman resonance. (c) The nonresonant signal has been subtracted from the resonant one, leaving only the deuterated liposome visible. [Reproduced from Duncan, M. D. (1984). *Opt. Comm.* **50**, 307. Copyright © North-Holland, Amsterdam.]

sation of distortions is to prepare at the entrance of the medium a beam whose wave front is distorted in such a way that the distortion introduced by the medium cancels the one that would develop in the beam, resulting in an undistorted wave front at the exit of the medium.

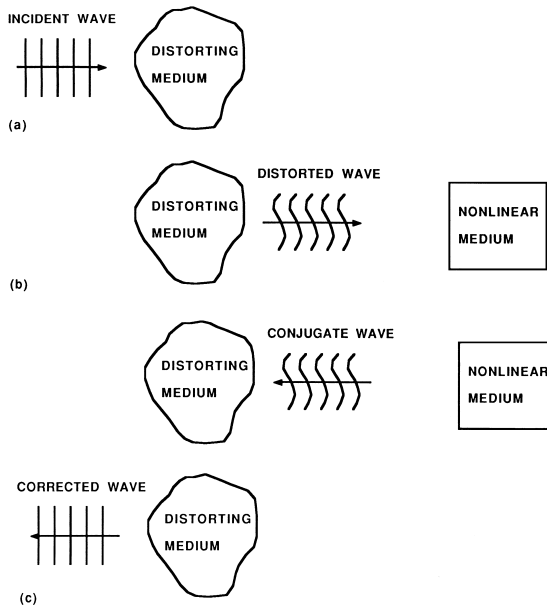


FIGURE 37 The use of optical phase conjugation for compensation of distortions. (a) An initially smooth beam propagates from right to left through a distorting medium. (b) After emerging from the distorting medium, the acquired phase variations are conjugated and the direction of propagation is reversed in a suitable nonlinear interaction. (c) After the second pass through the nonlinear medium, the beam emerges with a smooth wavefront.

The wave front required for this compensation to occur is the conjugate of the wave front obtained by propagation of an initially plane wave through the medium. It is usually obtained by propagating a beam through the distorting medium and then into a nonlinear medium that produces a phase-conjugate beam that propagates in the reverse direction, as illustrated in Fig. 37b. Various nonlinear interactions can be used for optical phase conjugation, including degenerate four-wave mixing in transparent, absorbing, and amplifying media and various forms of stimulated scattering such as Raman, Brillouin, and Rayleigh and stimulated emission. The two most widely used techniques are degenerate four-wave mixing and stimulated Brillouin scattering.

Degenerate four-wave mixing configured for phase conjugation is illustrated in Fig. 38. Here two strong waves, A_1 and A_2 , are incident on the nonlinear medium from opposite directions. The wave front carrying the information to be conjugated is the probe wave A_3 and the conjugate wave A_4 that is generated in the interaction $\omega_4 = \omega_1 + \omega_2 - \omega_3$ propagates in the direction opposite to A_3 . Because the waves are present in counterpropagating pairs and all the waves have the same frequency, the process is automatically phase matched regardless of the angle between the probe wave and the pump waves. Three-wave difference-frequency mixing of the form $\omega_3 = \omega_1 - \omega_2$ can also be

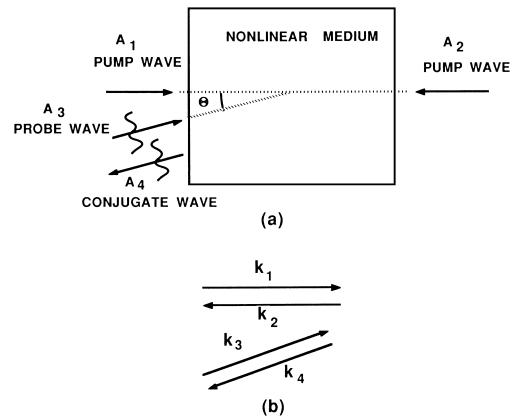


FIGURE 38 (a) Configuration for the use of degenerate four-wave mixing for phase conjugation. The frequencies of all waves are equal; the two pump waves, A_1 and A_2 , propagate in opposite directions; and the angle θ between the probe wave and the forward pump wave is arbitrary. (b) The k vector diagram is shown.

used in crystals for phase conjugation, but the phase-matching requirements in nonlinear crystals restrict the angles that can be used, and hence the magnitude of the distortions that can be corrected.

The generation of conjugate waves in nonlinear interactions involves the creation of volume holograms, or diffraction gratings, in the medium through interference of the incident waves. The interference that occurs, for example, between the waves A_1 and A_3 in degenerate four-wave mixing is illustrated in Fig. 39a. The backward wave is created by the scattering of the pump wave A_2 off the

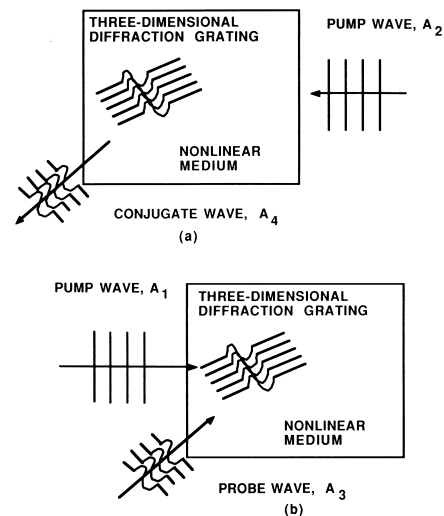


FIGURE 39 Illustration of the formation of phase conjugate waves by scattering from the interference patterns formed by pairs of the incident waves. (a) Interference between the probe wave and the forward pump wave forming a contoured grating in the nonlinear medium. (b) Scattering of the backward pump wave from the contoured grating to form the phase-conjugate wave.

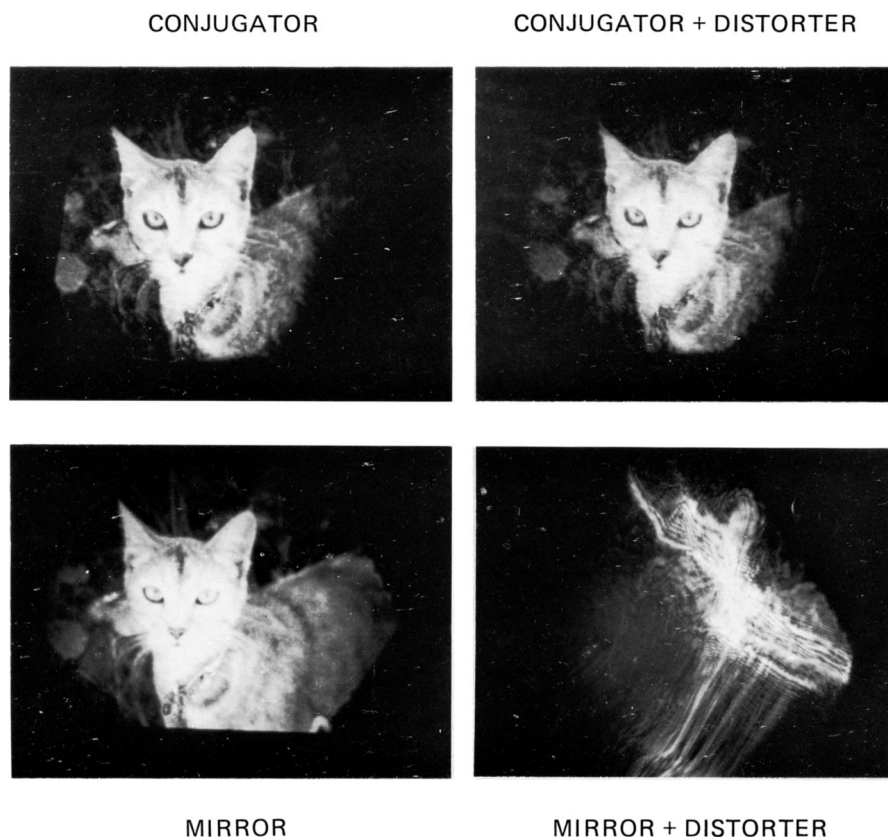


FIGURE 40 Example of image reconstruction with phase conjugation. The unaberrated image is shown at the lower left using a plane mirror and at the upper left using a conjugate mirror. The image at the lower right shows the effect of an aberrator (a distorting piece of glass) on the image obtained with a normal mirror, while the image at the upper right shows the corrected image obtained with the aberrator and the conjugate mirror. [From Feinberg, J. (1982). *Opt. Lett.* 7, 488. Copyright ©1982 by the Optical Society of America.]

diffraction grating created by the interference of waves A_1 and A_3 . The information as to the distortions on the incoming wave appears as bending of the contours of the grating. It is transferred to the phase contour of the backward wave as it is created but with its sense reversed. A similar interference occurs between all pairs of the incident waves, and each interference pattern creates a contribution to the backward-propagating phase-conjugate wave.

In stimulated Brillouin scattering, the incident wave serves as both the pump wave for the nonlinear process and the distorted wave to be conjugated. When the incident wave is highly aberrated, different components of the pump wave interfere in the focus, producing regions of relatively high and low gain. The wave that is ultimately generated in the Brillouin process is the one with the highest average gain. This wave in turn is one that is the phase conjugate of the incident wave, because it has the most constructive interference in the focus. The interference of the various components of the pump beam in the focus can also be viewed as the creation of diffraction gratings

or holograms in a manner similar to that described for degenerate four-wave mixing. The phase-conjugate wave is then produced by scattering of a backward-traveling wave from these diffraction gratings.

Optical phase conjugation has been used for reconstruction of images when viewed through distorting glass windows or other low-quality optical components, for removal

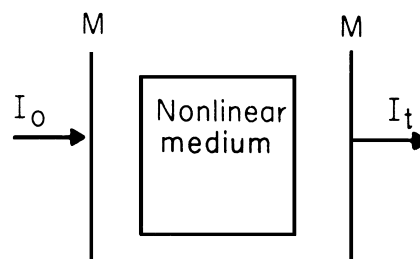


FIGURE 41 Illustration of nonlinear Fabry-Perot cavity for optical bistability. The end mirrors are plane and parallel, and the medium in the center exhibits either a nonlinear refractive index or saturable absorption.

of distortions from laser beams, for holographic image reconstruction, and for image up-conversion. An example of image reconstruction using optical phase conjugation is shown in Fig. 40.

D. Optical Bistability

Nonlinear optical effects can also be used to produce bistable devices that are similar in operation to bistable electric circuits. They have potential use in optical memories and the control of optical beams and can perform

many of the functions of transistors such as differential gain, limiting, and switching. Bistable optical elements all have a combination of a nonlinear optical component and some form of feedback. A typical all-optical device consisting of a Fabry–Perot (FP) cavity that is filled with a nonlinear medium is illustrated in Fig. 41. The FP cavity has the property that when the medium between its plates is transparent, its transmissivity is high at those optical wavelengths for which an integral number of wavelengths can be contained between the plates, and its reflection is high for other wavelengths.

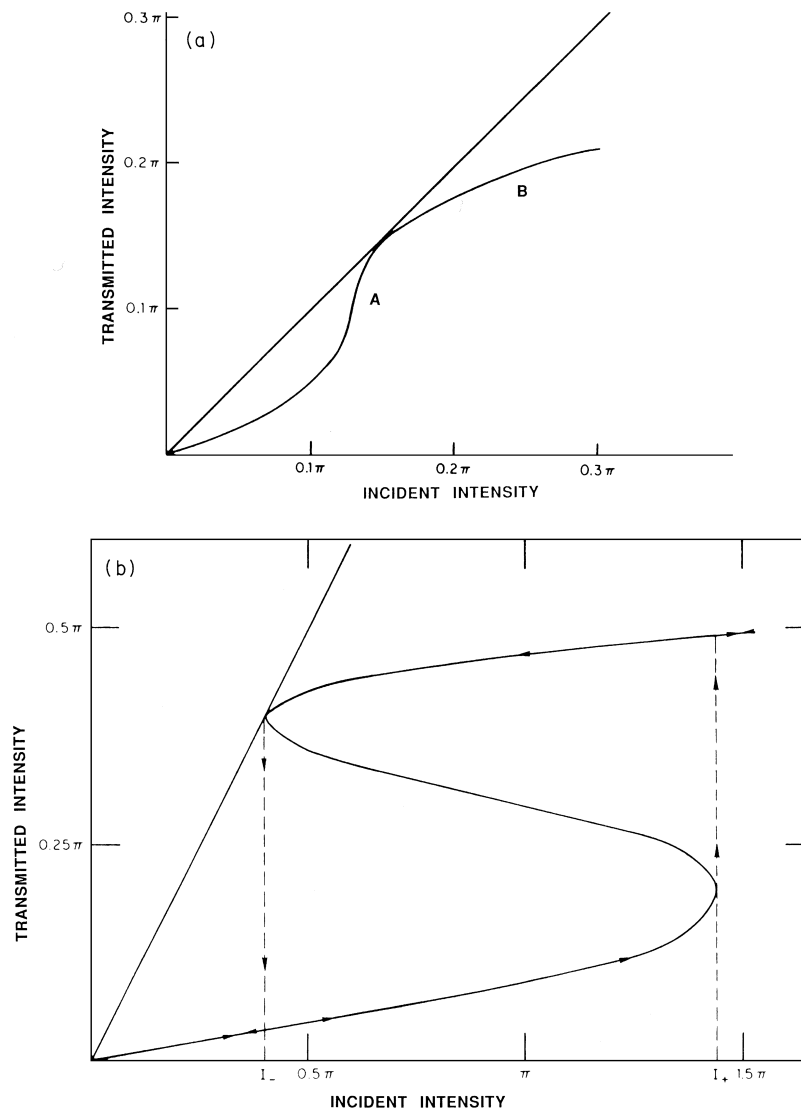


FIGURE 42 Illustration of nonlinear transmission with a nonlinear Fabry–Perot cavity using a medium with a nonlinear refractive index. (a) The transmission curve does not show hysteresis and can be used for differential gain (in region A) or optical limiting (in region B). (b) The transmission shows hysteresis and can be used for optical bistability. The difference between the differential gain and bistable operation is determined by the initial detuning from the resonance. [Reproduced from Reintjes, J. (1984). "Nonlinear Optical Parametric Processes in Liquids and Gases," Academic Press, Orlando, Florida.]

The operation of a nonlinear FP cavity can be illustrated by assuming that the nonlinear medium has a refractive index that is a function of the intensity inside the FP cavity. The wavelength of the incident light is chosen to be off resonance so that at low incident intensities the reflectivity is high, the transmission is low, and not much light gets into the FP cavity. As the incident intensity is increased, more light gets into the FP cavity, changing the refractive index of the medium, and bringing the cavity closer to resonance. As a result, the fractional transmission increases and the intensity of the light transmitted through the cavity increases faster than does the intensity of the incident light. When the incident intensity is sufficiently high the cavity is brought into exact resonance, allowing full transmission. For further increases in the incident intensity, the cavity moves beyond the resonance and the fractional transmission decreases.

For conditions in which the transmitted intensity increases faster than the incident intensity, the nonlinear FP cavity can be used for differential gain in a manner similar to transistor amplifiers. An example of the transmission curve under these conditions is shown in Fig. 42a. In this situation, a small modulation on the incident light wave can be converted to a larger modulation on the transmitted wave. Alternatively, the modulation can be introduced on a different wave and can be transferred to the forward wave with a larger value. If the incident intensity is held just below the value at which the transmission increases and the cavity is illuminated by another, weaker wave that changes the transmission level, the nonlinear FP cavity can act as a switch. In the range of intensities for which an increase in intensity moves the FP cavity away from resonance, the nonlinear FP acts as an optical limiter. For certain conditions, the level of light inside the cavity is sufficient to maintain the resonance condition once it has been achieved, resulting in hysteresis of the transmission curve with the incident intensity, and providing the conditions for bistability. An example of the transmission curve under this condition is shown in Fig. 42b.

A medium with saturable absorption can also be used for the nonlinear medium. In this case, the resonance condition is always met as far as the wavelength is concerned, but the feedback from the second mirror, which is required for resonant transmission, is not obtained until sufficient light is present inside the cavity to bleach the nonlinear medium. Devices of this type can show differential gain and optical bistability but not optical limiting.

Bistable optical devices can also be constructed in such a way that the transmitted light signal is converted to an electrical voltage that is used to control the refractive index of the material inside the cavity. These hybrid optical-electrical devices can operate with very low levels of light intensity and as such are especially compatible with integrated optical components.

Optical bistability has also been observed in a wide range of optical configurations, such as laser cavities, degenerate four-wave mixing, and second-harmonic generation. Under some conditions, bistable devices can exhibit chaotic fluctuations in which the light output oscillates rapidly between two extremes.

SEE ALSO THE FOLLOWING ARTICLES

ELECTROMAGNETICS • IMAGING OPTICS • LASERS, DYE • LASERS, ULTRAFAST PULSE TECHNOLOGY • MICROWAVE MOLECULAR SPECTROSCOPY • OPTICAL DIFFRACTION • POLARIZATION AND POLARIMETRY • RAMAN SPECTROSCOPY

BIBLIOGRAPHY

- Bowden, C. M., and Haus, J. (feature eds.) (1989). Nonlinear optical properties of materials. *J. Optical Society of America B, Optical Physics* 6(4), 562–853.
- Fisher, R. (ed.) (1983). "Optical Phase Conjugation," Academic Press, New York.
- Institution of Electrical Engineers (1990). "Properties of Gallium Arsenide," 2nd Ed., *emis Data Review Series No. 2*, Institution of Electrical Engineers, U.K.
- Kobayashi, T. (1989). "Nonlinear Optics of Organics and Semiconductors," Springer Proceedings in Physics Series, Vol. 36, Springer-Verlag, Berlin and New York.
- Levenson, M. D. (1982). "Introduction to Nonlinear Spectroscopy," Academic Press, New York.
- Marder, S. R. (1991). "Materials for Nonlinear Optics," ACS Symposium Series No. 455. Am. Chem. Soc., Washington, D.C.
- Pepper, D. M. (1986). Applications of optical phase conjugation. *Sci. Am.* 254 (January)(6), 74.
- Reintjes, J. (1984). "Nonlinear Optical Parametric Processes in Liquids and Gases," Academic, New York.
- Reintjes, J. (1985). Coherent ultraviolet sources. In "Laser Handbook" (M. Bass and M. L. Stitch, eds.), Vol. 5, pp. 1–202, North-Holland Publ., Amsterdam.
- Shen, Y. R. (1984). "The Principles of Nonlinear Optics," Wiley, New York.
- Shkunov, V. V., and Zeldovich, B. Y. (1985). Optical phase conjugation. *Sci. Am.* 253, (December)(6), 54.



Optical Amplifiers (Semiconductor)

Mark Summerfield

The University of Melbourne

- I. Structure of the Semiconductor Optical Amplifier
- II. Gain
- III. Noise
- IV. Dynamics
- V. Coherent Interactions
- VI. Conclusions

GLOSSARY

Amplified spontaneous emission (ASE) The result of spontaneously emitted photons being amplified as they travel along an optical amplifier. ASE is often the main source of noise in optically amplified systems.

Beat noise Electrical noise generated at a receiver due to the mixing of optical signal power and optical ASE noise power at the detector.

Carrier density The number of electron–hole pairs per unit volume within the semiconductor material.

Carrier lifetime The average time for which an electron–hole pair exists before recombining via either a spontaneous or stimulated process.

Carrier rate equation A differential equation describing the physical processes leading to the formation and recombination of carriers within the semiconductor material. The carrier rate equation is the most impor-

tant equation governing the behavior of semiconductor optical amplifiers.

Fabry–Perot amplifier An amplifier in which there is some feedback due to reflections at both ends of the gain medium so that signals are amplified by multiple passes through the device. Fabry–Perot amplifiers have the benefit of providing enhanced gain at specific resonant wavelengths.

Gain compression The process by which high power input signals can reduce the gain of an optical amplifier by depleting the population inversion.

Gain spectrum The gain of an optical amplifier as a function of optical frequency or wavelength.

Hole In the valence band of a semiconductor material, which is almost fully populated by electrons, a hole is the absence of an electron at a particular energy state.

Mode A permitted electromagnetic field pattern within an optical waveguide, such as the active region of a semiconductor optical amplifier.

Modulation The variation in time of a property of an electromagnetic signal, such as its power, amplitude, frequency or phase. Also used to refer to the variation in time of a property of an optical amplifier, such as its gain or carrier density.

Population inversion In an optical amplifier, the condition in which the number of carriers in an excited state exceeds the number of carriers in an unexcited, or ground, state, such that a net gain exists in the amplifier.

Spontaneous emission The random generation of a photon due to the spontaneous recombination of an electron-hole pair. Spontaneous emission is the fundamental noise-generating process in an optical amplifier.

Stimulated emission The generation of a new photon due to recombination of an electron-hole pair in the presence of an existing photon. The new photon is a coherent replica of the existing photon. Stimulated emission is the fundamental gain process in an optical amplifier.

Transparency In an optical amplifier, the condition in which the number of carriers in an excited state exactly balances the number of carriers in an unexcited, or ground, state, such that there is no net gain or loss in the amplifier.

Traveling wave amplifier An amplifier in which there is no feedback due to reflections at the ends of the gain medium, so that signals are amplified during a single pass through the device. Traveling wave amplifiers have the advantage that their gain is less wavelength-dependent than Fabry-Perot amplifiers, and they have optimum noise performance.

Wavelength division multiplexing (WDM) In optical communications, a technique whereby the information-carrying capacity of optical fibers is greatly increased by sending multiple modulated signals on optical carriers of differing wavelengths. Optical amplifiers are a key enabling technology for WDM systems, since they enable multiple signals to be amplified simultaneously without the need for conversion of all individual signals back into electronic form.

SEMICONDUCTOR OPTICAL AMPLIFIERS have received less attention in the literature than amplifiers using optical fibers as the gain medium. Fiber amplifiers, including Erbium-doped fiber amplifiers and Raman amplifiers are currently transforming the world of optical fiber telecommunications as they become incorporated as line amplifiers in high-capacity wavelength-division multiplexed (WDM) transmission systems. However, despite the dominance of fiber amplifiers in the telecommunications marketplace, semiconductor optical amplifiers continue to show potential for application in a range of

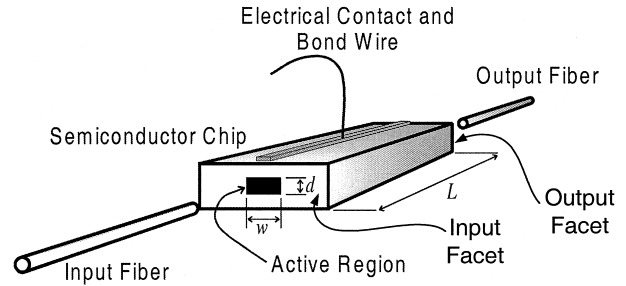


FIGURE 1 Structure of a semiconductor optical amplifier.

areas such as high-speed optical switching, optical integrated circuits, and optical wavelength converters. This article provides an overview of the technology and operation of semiconductor optical amplifiers, and highlights a number of applications.

I. STRUCTURE OF THE SEMICONDUCTOR OPTICAL AMPLIFIER

The basic structure of a Semiconductor Optical Amplifier (SOA) is shown schematically in Fig. 1. The SOA is essentially a semiconductor laser in which lasing has been suppressed. Light usually enters the device from one of the two end facets. It is convenient to think of the SOA as having an input facet and an output facet, as shown. Amplification of light occurs in the *active region*, which is a rectangular waveguide of width w and depth d within the semiconductor chip. For most applications it is desirable that this waveguide supports only a single transverse mode. Typical values of w and d are between 1 and 10 μm . When an electric current is injected into the electrical contact, a population inversion is achieved within the active. The mechanism by which gain is produced from this electrical injection is described in greater detail in Section II.A.

Choice of the length L of the active region depends upon a number of factors. In general, a longer device can produce higher gain. However, fabricating longer devices may reduce the yield and hence increase the device cost due to the higher probability of a defect occurring somewhere within the active region. Furthermore, high-gain SOAs may suffer from poor performance due to residual facet reflectivities that can cause large ripples in the gain spectrum, instability and even lasing within the SOA. The effects of facet reflectivities, and common techniques used to minimize reflections are discussed in Section II.C. SOAs are typically fabricated with active region lengths ranging from less than 400 μm up to more than 2 mm.

Devices packaged for use in telecommunications applications are usually coupled to single-mode fibers at the input and output facets. In order to obtain efficient

coupling of light between the fiber pigtails and the active region of the SOA, it is necessary to match the mode fields of the fiber and SOA waveguides as closely as possible, as well as achieve accurate alignment of the fiber to the active region. The use of lensed fibers or separate spherical lenses, as well as active alignment techniques, can help to improve the coupling; however, obtaining very high coupling efficiency remains a challenging and labor-intensive task. The coupling loss per facet is typically 3–4 dB.

II. GAIN

A. Gain Mechanism

Gain is achieved in SOAs, as in other laser media, by stimulated emission of photons. The lasing material in an SOA is a direct band-gap semiconductor. A direct band-gap semiconductor crystal can be considered as a quasi two-level laser system, as illustrated in Fig. 2. The active region of the SOA is typically an intrinsic (i.e., undoped) semiconductor. In the absence of current injection, the population of the conduction band is small, being generated only by the thermal excitation of electrons from the valence band. In this case, as shown in Fig. 2 (a), a photon with energy $h\nu$ equal to the band-gap energy E_g will most likely be absorbed, resulting in the promotion of an electron from the heavily populated valence band, to the lightly populated conduction band. However, if electrical current is injected into the SOA, it acts as a pump source supplying electrons to the conduction band, and removing them from the valence band leaving it populated with holes. If sufficient current is supplied, a population inversion results as shown in Fig. 2 (b). In this case, the incident photon will most likely induce a stimulated emission event, in which a conduction-band electron recombines with a valence-band hole, generating in the process an additional photon, which is an exact duplicate of the original. This process,

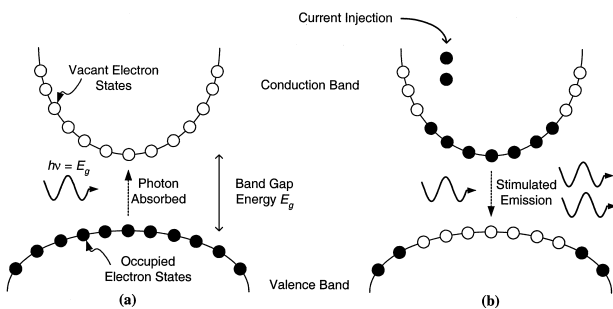


FIGURE 2 The conduction and valence bands of a direct band-gap semiconductor crystal act as the excited and ground state energy levels in a quasi two-level laser system. (a) Thermal equilibrium. (b) Pumping via injected current.

repeated many times as the photons propagate along the active waveguide region of the SOA, produces an exponential increase in the intensity of light which is the origin of the very high gain available from SOAs.

The exponential growth in the number of photons resulting from stimulated emission produces a gain G_0 of the amplifier, given by the ratio of the output power P_{out} of the SOA to the input power P_{in} , which is also exponential with length. The gain is given by

$$G_0 = \frac{P_{out}}{P_{in}} = e^{(g_0 - \gamma_{sc})L}. \quad (1)$$

In Eq. (1), g_0 is the differential gain representing the increase in photon number per unit length. The parameter γ_{sc} is an empirical term describing optical losses, i.e., reduction in photon number per unit length, within the SOA. These losses are mainly due to scattering of photons out of the guided mode due to imperfections in the semiconductor lattice, and roughness in the edges of the active-region waveguide. The expression in Eq. (1) in fact only gives the *small-signal gain*, which is the gain observed when the input power P_{in} is very small. The effect of large-signal input is considered in Section D.

The differential gain g_0 must, in general, be calculated from a detailed model of the semiconductor band structure. However, it turns out that g_0 is a very nearly linear function of the density of electrons N_c injected into the conduction band by the pump source. An empirical expression often used to represent the differential gain is

$$g_0 = \Gamma a (N_c - N_0). \quad (2)$$

Equation (2) represents the gain due to stimulated emission in the SOA. The parameter a has units of area, and is commonly known as the *gain cross section*. It is a property of the particular semiconductor alloy from which the SOA is fabricated, and can be determined from the slope of the relationship between the differential gain and the carrier density. The parameter Γ is the *optical confinement factor*, which represents the extent to which the propagating optical field is actually confined within the active region. It is a property of the waveguide structure of the SOA, and depends upon the material properties and the geometry of the device. Since most SOAs support only a single guided mode, and this mode has a significant fraction of its power in the evanescent field outside of the active region, Γ is a constant somewhat less than unity (typically between 0.3 and 0.6). Finally, the parameter N_0 is the transparency carrier density. It represents the carrier density at which the populations of the valence and conduction bands are balanced such that the rate of absorption exactly balances the rate of stimulated emission. At this carrier density, there is no net gain or loss of photons via stimulated emission or absorption—the SOA is “transparent.”

The second term in Eq. (2), γ_{sc} , is an empirical term describing other sources of optical loss within the SOA. These losses are mainly due to scattering of photons out of the guided mode due to imperfections in the semiconductor lattice, and roughness in the edges of the active-region waveguide.

The density of electrons N_c in the conduction band is dependent upon the rate at which the SOA is pumped via the injected current—which increases the carrier density—and upon the rate at which carriers (electrons and holes) recombine, which decreases the carrier density. There are generally four significant mechanisms whereby carriers may recombine. First, there are nonradiative processes that occur when a single carrier interacts with a feature of the semiconductor crystal lattice, such as a defect. Since only a single carrier is initially involved, the probability of such an event occurring is simply proportional to the carrier density. Second, a pair of carriers (an electron and a hole) may spontaneously recombine. This process is proportional to the square of the carrier density, and results in the emission of a photon. Third, three carriers may interact to exchange energy and momentum, such as in a collision, culminating in the recombination of an electron–hole pair. Processes of this type are known as Auger recombination, and are rare; however the probability of a three-carrier interaction is proportional to the cube of the carrier density, and thus the rate of Auger recombination may be relatively high in a strongly pumped semiconductor in which N_c may be large. The final recombination process is stimulated emission, whereby a photon interacts with an electron–hole pair that recombine, emitting in the process an additional photon coherent with the initiating photon. This is the origin of the gain process described by Eq. (2), and the rate at which it occurs is directly proportional to the carrier density, and to the number of photons available, i.e., the intensity of the light.

Overall, the carrier density can thus be described by the following differential equation, known as the *carrier rate equation*:

$$\frac{dN_c}{dt} = \frac{I}{qV} - (A_{nr}N_c + BN_c^2 + CN_c^3) - \frac{\Gamma a(N_c - N_0)}{h\nu} P. \quad (3)$$

Equation (3) is the key equation that describes the most important characteristics of the SOA, including the gain and carrier dynamics. In this equation, I is the injected current, q is the charge on the electron (i.e., 1.6×10^{-19} C), $V = wdL$ is the volume of the active region, A_{nr} is a coefficient describing the rate of non-radiative recombination, B is the radiative recombination coefficient, C is the Auger recombination coefficient, h is Planck's constant (i.e., 6.63×10^{-34} Js), ν is the optical frequency, and P is the optical irradiance (i.e., the optical power divided by the cross-sectional area of the active region).

Equation (3) can be explained in simple terms as follows. The rate of change of the carrier density (i.e., dN_c/dt) is equal to the rate at which new carriers are injected, *minus* the rate at which carriers are removed via nonradiative, radiative, and Auger recombination processes, *minus* the rate at which carriers are removed via stimulated recombination. Note that the rate of stimulated recombination is proportional to the differential gain given by Eq. (2). This is because the creation of a new photon by stimulated emission (i.e., the gain process) corresponds precisely to the loss of one carrier through stimulated recombination.

An important point to appreciate regarding Eq. (3) is that it is position-dependent. The optical power is, in general, dependent upon position within the amplifier—in fact, it increases along the length of the SOA. Consequently, the carrier density N_c must also be position-dependent. For small values of input optical power (the so-called *small signal* regime), the effect on the carrier density may be negligible, in which case the gain is effectively described by Eqs. (1) and (2). However, for larger input optical power the consequent reduction in N_c along the length of the amplifier may be significant, resulting in a measurable reduction in the gain of the SOA. This point will be addressed further in Section D, in which the subject of *gain compression* is discussed.

In some applications, the carrier density may not vary a great deal with position or time. In such cases a simplification may be made to Eq. (3) by linearizing the recombination term around the operating carrier density. The resulting variation on the carrier rate equation is normally written as

$$\frac{dN_c}{dt} = \frac{I}{qV} = \frac{N_c}{\tau_s} - \frac{\Gamma a(N_c - N_0)}{h\nu} P. \quad (4)$$

The constant τ_s is known as the *spontaneous carrier lifetime* and may be interpreted as the average time for which a carrier “lives” in the conduction band before spontaneously recombining via any one non-radiative, radiative or Auger recombination processes. The spontaneous carrier lifetime is typically on the order of 1 ns.

Table I lists some typical ranges for the values of SOA parameters introduced in this section, valid for long-wavelength devices (i.e., 1300 to 1550 nm) typically designed for telecommunications applications.

B. Spectral Properties

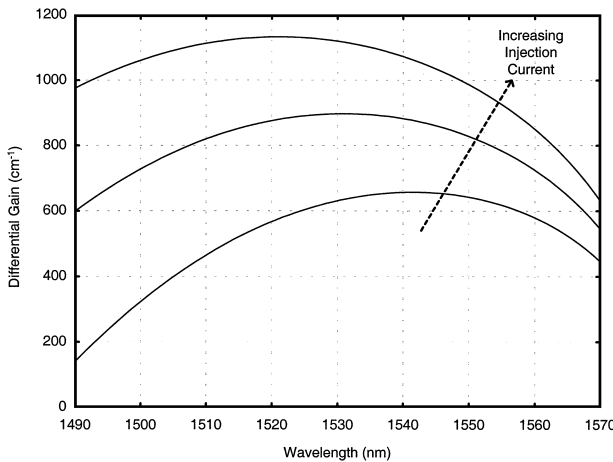
In the preceding section, the most important properties of SOAs were introduced without regard to any dependence upon the wavelength of the light being amplified. However, it is apparent from Fig. 2 that an electron in the conduction band may occupy any one of a wide range of energy levels. Consequently, a range of transitions

TABLE I Typical Values of Some SOA Material Parameters

Parameter	Description	Typical values
L	Amplifier length	400–2000 μm
$A_x = wd$	Active region cross-sectional area	$\sim 10^{-9} \text{ cm}^2$
a	Material gain cross section	$(2.7\text{--}3.0) \times 10^{-16} \text{ cm}^2$
N_0	Transparency carrier density	$(1.0\text{--}2.0) \times 10^{18} \text{ cm}^{-3}$
γ_{sc}	Material scattering loss	10–40 cm^{-1}
Γ	Confinement factor	0.3–0.6
A_{nr}	Nonradiative recombination coefficient	$\sim 10^8 \text{ s}^{-1}$
B	Radiative recombination coefficient	$\sim 10^{-11} \text{ cm}^3 \text{ s}^{-1}$
C	Auger recombination coefficient	$\sim 10^{-28} \text{ cm}^6 \text{ s}^{-1}$

from the conduction band to the valence band is available, corresponding to the emission of photons with a range of different energies, and hence different optical frequencies. The probability of emission of a photon with a given energy—either through spontaneous or stimulated emission—depends in general upon the *density of states*, i.e., the number of available energy levels close to the photon energy, and the carrier density N_c , i.e., the number of carriers available to populate those states. A description of the calculation of the density of states is complex, and beyond the scope of this article, however the general result is that the differential gain g_0 introduced in Eq. (2) is a function of optical wavelength as well as carrier density.

Figure 3 shows the differential gain of a nominally 1550 nm SOA as a function of wavelength for three different values of injected current, calculated using the density-of-states method. The device has a band-gap energy corresponding to 1580 nm; however, the gain peak

**FIGURE 3** Calculated gain spectra of a 1550-nm SOA for three different values of injected current.

always occurs at a higher energy (i.e., shorter wavelength). As the injection current is increased, the carrier density increases and more states are filled, resulting in a shift in the gain peak to shorter wavelengths. The gain decreases on the low energy (long wavelength) side of the gain peak because the density of states is lower. It decreases on the high-energy (short-wavelength) side because although the density of states is greater, the number of carriers is not sufficient to fully occupy those states.

As already mentioned, calculation of the gain based on the density of states is complex. Furthermore, it is difficult to obtain the required material parameters to perform such a calculation for a specific device. Consequently, it is convenient to have an empirical model that may be employed when the simple linear gain model described by Eq. (2) is inadequate. It is common in such cases to replace the linear gain model with a parabolic model, which may be expressed in terms of wavelength as

$$g(N_c, \lambda) = a(N_c - N_0) - b_2[\lambda - (\lambda_0 - b_3(N_c - N_0))]^2. \quad (5)$$

In Eq. 5 the parameter λ_0 is the wavelength of the gain peak at transparency. The parameter b_2 describes the width of the gain spectrum, while the parameter b_3 represents the shift in the gain peak to shorter wavelengths that occurs as the carrier density is increased. To within an order of magnitude, typical values of these parameters for bulk semiconductor materials are $b_2 \sim 0.1 \text{ cm}^{-1} \text{ nm}^{-2}$ and $b_3 \sim -1 \times 10^{-17} \text{ nm cm}^3$.

Close inspection of Fig. 3 reveals that the gain spectrum is not symmetrical around the gain peak if a wavelength range of many tens of nanometers is considered. Consequently, more sophisticated empirical models such as cubic approximations have occasionally been employed.

C. Facet Reflectivity

In an ideal traveling-wave SOA amplification of the optical signal is achieved in a single pass of the amplifier. Indeed, an SOA can be viewed simply as a semiconductor laser that has had its end-facet mirrors removed to prevent optical feedback from occurring, and hence inhibit the onset of lasing. However, in a Fabry–Perot semiconductor laser the mirrors typically comprise only the cleaved end-facets of the device, which provide a reflectivity of around 30 to 40% simply as a result of the refractive index difference between the semiconductor material and the surrounding air. Consequently, additional design and/or manufacturing steps are required in order to produce a device *without* reflections at the facets.

The most common technique employed to reduce facet reflectivity is to apply an *anti-reflection (AR)* coating to each end facet. As illustrated in Fig. 4, an AR coating is a

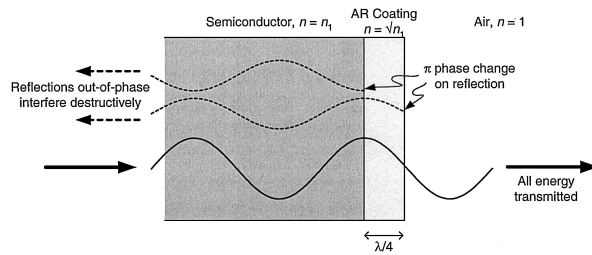


FIGURE 4 Operation of anti-reflection coating.

thin film of material a quarter wavelength thick deposited on the end facets of the SOA that has a refractive index ideally equal to the square root of the refractive index of the semiconductor. This results in two points of reflection: at the interface between the semiconductor and the AR coating; and at the interface between the AR coating and air. These two reflected waves interfere destructively, so that there is no net reflection from the coating. Conversely, the transmitted waves interfere constructively, so that ideally all light is transmitted.

In practice, however, AR coatings are imperfect. The main drawback of a single-layer AR coating is that it is wavelength-specific—it works best only for light whose quarter-wavelength corresponds to the thickness of the film. Furthermore, attenuation of light in the coating material, surface roughness, diffraction of light exiting the SOA waveguide, and polarization dependence of the modal index within the waveguide all result in further reduction of the effectiveness of the AR coating. Figure 5 shows the measured fiber-to-fiber gain spectrum of a typical AR-coated SOA. The device exhibits up to 4 dB of Fabry–Perot gain ripple due to residual facet reflectivity. The ripple is substantially smaller around 1550 nm—the center of

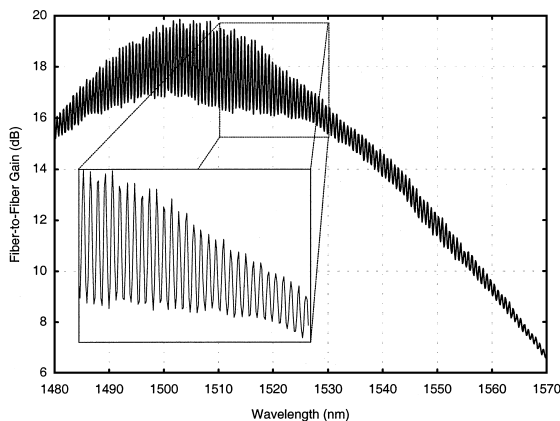


FIGURE 5 Gain spectrum of an SOA with imperfect AR coatings. Inset shows detail of the wavelength range from 1510 to 1530 nm, illustrating the approximately sinusoidal ripple resulting from residual Fabry–Perot cavity reflections.

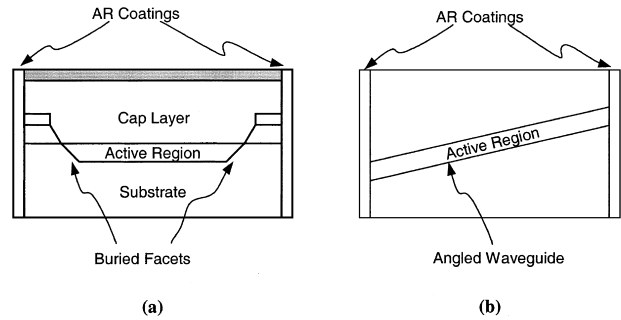


FIGURE 6 (a) SOA structure with buried facets (side view, cross section); (b) SOA structure with angled facets (top view, plan).

the conventional low-loss telecommunications window—and it is evident that the AR coatings have been optimized for this wavelength. Under the operating conditions shown, the gain peak is at about 1505 nm; however, at a lower bias current, or under more strongly saturated operation (see Section D) the gain peak would coincide with the minimum reflection wavelength.

A more broadband AR coating can be obtained by using multiple layers; however, for SOAs two additional techniques are more commonly employed to reduce the impact of reflections. These are the use of *buried facets* and *angled facets*, as illustrated in Fig. 6. In the case of the buried facet device, shown in Fig. 6(a), the index difference between the active region and the passive region is relatively small, reducing the magnitude of any reflections. Additionally, the interface may be angled as shown to direct residual reflections upwards, away from the waveguide. The end-facets of the device are AR coated to reduce reflections at that point; however, due to diffraction of light leaving the active region, the impact of residual reflections is minimal because the majority of reflected light is not directed back into the active region. In an angled facet device, shown in Fig. 6(b), the device is cut from the original wafer in such a way that the active region waveguide is angled with respect to the end facets. Consequently, residual reflected light is directed substantially away from the waveguide preventing repeated amplification and reflection that leads to the Fabry–Perot ripple in the gain characteristic. Typically the angle is around 10 to 15°. The disadvantage of these techniques for reducing the reflections is that generally alignment of optical fibers to the device is made more difficult, resulting in increased coupling losses and a higher cost of packaging.

SOAs with significant residual facet reflectivity have been intentionally fabricated in the past. The perceived benefit of these so-called *Fabry–Perot (FP) Amplifiers* (as distinct from traveling-wave amplifiers) is that there is an enhancement of the gain at the peaks of the FP spectral response due to the signal effectively experiencing

multiple passes of the active region. The disadvantages are that the optical signals to be amplified must be aligned to these peaks, the optical bandwidth of each signal is limited by the width of the peaks, and the noise performance of the amplifier is degraded. With improved fabrication techniques, it is now possible to obtain high yield for amplifiers up to 2 mm or more in length, and enhanced gain is achieved simply by using a longer active region. At the same time, amplifiers with very high gain require greatly improved suppression of reflections to prevent large residual gain ripple, and the onset of lasing. For example, to achieve a gain ripple of less than 1 dB in a device with 30 dB of facet-to-facet gain requires that the reflectivity of each facet be below 6×10^{-5} .

D. Gain Compression

It was mentioned in Section A that, in general, the optical power and hence the carrier density N_c are functions of position along the SOA. This is illustrated in Fig. 7 for the simplest case—which is also the one of greatest interest for amplification of optical signals—in which a weak signal of power P_{in} is injected into one end of the amplifier, and then amplified along the length of the active region to produce a strong output signal of power P_{out} . In order to determine the net gain G of the amplifier in this case, it is necessary to integrate the optical power along the

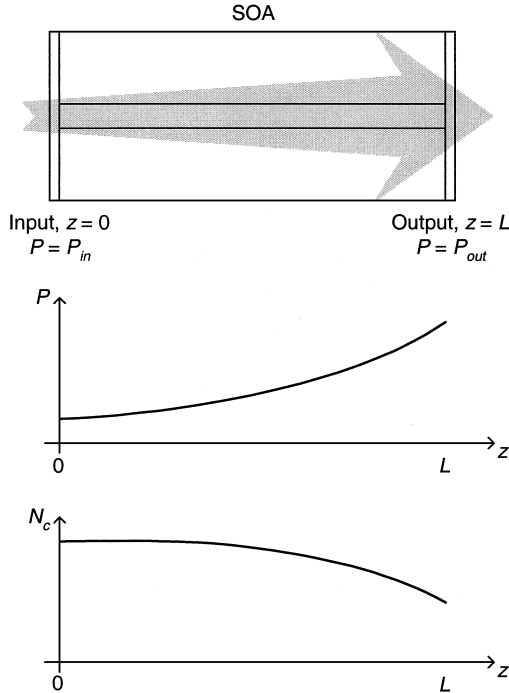


FIGURE 7 Growth of the optical field, and consequent reduction of the carrier density, along the length of an SOA.

length of the amplifier according to the following differential equation, noting that both the optical irradiance P , and the carrier density N_c are implicitly functions of the position z :

$$\frac{dP}{dz} = [g(N_c) - \gamma_{sc}]P. \quad (6)$$

The gain $g(N_c)$ can be found as a function of the optical irradiance by solving Eq. 4 in the steady-state (i.e., setting $dN_c/dt = 0$). The result is

$$g(P) = \frac{g_0}{1 + P/P_{sat}}, \quad (7)$$

where g_0 is the unsaturated gain, given by Eq. (2) in the absence of an optical signal, and P_{sat} is the *saturation irradiance*, given by $P_{sat} = h\nu / \Gamma a \tau_s$. The value of $A_x P_{sat}$ (i.e., the corresponding total power, as opposed to the irradiance) is typically in the range of 1 to 10 mW. Note that if $P \ll P_{sat}$, then $g(P) \approx g_0$, and the solution to Eq. (6), integrated from $z = 0$ to $z = L$ is simply given by Eq. (1), as would be expected.

In many cases of practical interest the magnitude of the differential gain g may be much larger than the losses represented by γ_{sc} . Consequently, the loss term may be neglected in Eq. (6), and integration along the length of the SOA yields the following transcendental equation:

$$\overline{P_{out}} e^{\overline{P_{out}}} = G_0 \overline{P_{in}} e^{\overline{P_{in}}}, \quad (8)$$

where G_0 is the small-signal gain (given by Equation (1) with $\gamma_{sc} = 0$), $\overline{P_{in}} = P_{in} / (A_x P_{sat})$ is the input power to the SOA normalized to the saturation power, and $\overline{P_{out}} = P_{out} / (A_x P_{sat})$ is similarly the output power normalized to the saturation power. Assuming the small-signal gain is known (e.g., from an experimental measurement), Eq. (8) can be solved numerically for any given value of P_{in} to determine P_{out} , and hence the net gain $G = P_{out} / P_{in}$.

Similarly, a more complex transcendental solution to Eq. (6) can be found for the case in which γ_{sc} is not negligible:

$$\ln \frac{1 - \xi(1 + \overline{P_{in}})}{1 - \xi(1 + \overline{P_{out}})} = \xi \ln \frac{G_0 \overline{P_{in}}}{\overline{P_{out}}}, \quad (9)$$

where the parameter ξ is the normalized internal loss, given by $\xi = \gamma_{sc} / g_0$, and all other parameters have the same definitions as in Eq. (8). Note, however, that in order to solve Eq. (9) it is also implicitly necessary to know the length L of the SOA in order to compute g_0 given the signal gain G_0 (or vice-versa).

As an example, consider the case of a 500- μm -long amplifier for which the facet-to-facet small-signal gain has been measured to be 20 dB, with a saturation power of 10 mW. If the internal losses are assumed to be negligible,

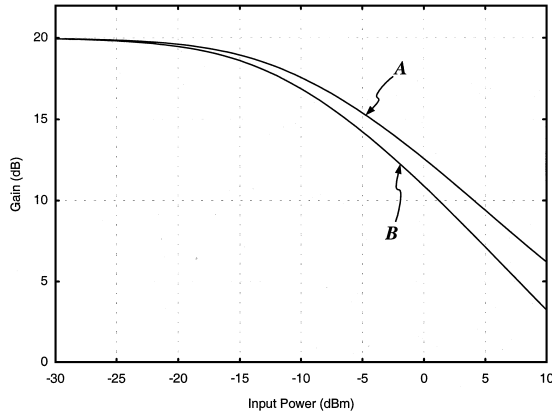


FIGURE 8 Calculated gain compression curves for an SOA with a 20-dB facet-to-facet gain and a saturation power of 10 mW, (A) assuming negligible internal losses and (B) assuming internal losses of $\gamma_{sc} = 34 \text{ cm}^{-1}$.

i.e., $\gamma_{sc} = 0$, then the gain as a function of input power may be calculated using Eq. (8). The result is the curve labeled A in Fig. 8. If, however, the internal losses in fact have the more realistic value of $\gamma_{sc} = 34 \text{ cm}^{-1}$, then using Eq. (9) the curve B is instead obtained. Both curves show the same general behavior—as the input power is increased, the gain of the amplifier is compressed due to depletion of carriers. Thus this process is generally referred to as *gain compression* or *gain saturation*. Curves such as those shown in Fig. 8 are known as *gain compression characteristics*. It is apparent that at small input power there is little difference between the two models, and it is generally simpler to use Eq. (8), which requires fewer amplifier parameters. However, if the amplifier is to be used in a regime of stronger gain compression, then the influence of the internal loss becomes more significant and cannot be neglected. In general, as the gain is further compressed, the losses (which are independent of the input power) become relatively more important, and the simple model represented by Eq. (8) thus tends to overestimate the gain.

For applications involving the amplification of one optical signal at a single, predetermined, wavelength, the simple analysis presented earlier is often sufficient to characterize the behavior of an SOA. However, it is apparent from Fig. 3 that the gain compression characteristics are wavelength-dependent. In particular, it is clear that a given reduction in the carrier density results in a greater compression of the gain at shorter wavelengths than at longer wavelengths. Accordingly, in the models represented by Eqs. (8) and (9) the saturation power must be considered to be wavelength-dependent. Specifically, P_{sat} is smaller at shorter wavelengths, and larger at longer wavelengths. If the gain compression behavior of the amplifier in re-

sponse to multiple input signals at different wavelengths is required, it is no longer possible to use these simple models. Instead it is generally necessary to employ numerical integration techniques to solve a system of equations such as Eqs. (4), (5), and (6).

III. NOISE

A. Amplified Spontaneous Emission

Whether in the electrical or optical domain, amplification is never obtained for free—invariably the gain is associated with some level of added noise. In the case of an SOA, the fundamental origin of noise is *spontaneous emission*. As has already been discussed in Section II.A, with reference to recombination processes, electrons and holes spontaneously recombine resulting in the emission of a photon over timescales on the order of a nanosecond. Since the process is random it results in optical noise emitted over the complete gain spectrum of the SOA. Indeed, it was Einstein who first showed, on the basis of fundamental thermodynamic considerations, that in any inverted medium, spontaneous emission must occur at a rate that is proportional to the differential gain provided by the population inversion. In other words, spontaneous emission noise is an unavoidable fact of life.

The most important parameter of an SOA that determines its noise performance is the *spontaneous emission factor* (also sometimes known as the *population inversion factor*) n_{sp} . This parameter is a measure of the amount of spontaneous emission noise produced, relative to the gain. A lower value of n_{sp} corresponds to a less noisy amplifier, and the minimum possible value in theory—never achieved in practice in an SOA—is $n_{sp} = 1$. As with the differential gain, a detailed calculation of n_{sp} requires knowledge of the density of states. Figure 9 shows the spontaneous emission factor calculated in this way, as a function of wavelength for a 1550-nm SOA corresponding to the middle differential gain curve of Fig. 3. Note that n_{sp} is smaller at longer wavelengths, closer to the band gap energy. This is because the available conduction band energy states are more densely populated, and hence the effective population inversion is more complete. A common empirical expression for n_{sp} , based on the parameters introduced in Section II, is given by

$$n_{sp} = \frac{N_c}{N_c - N_0}. \quad (10)$$

At first glance, this simple expression seems at odds with the result shown in Fig. 9 since it is independent of the optical wavelength. However, considering the gain curves

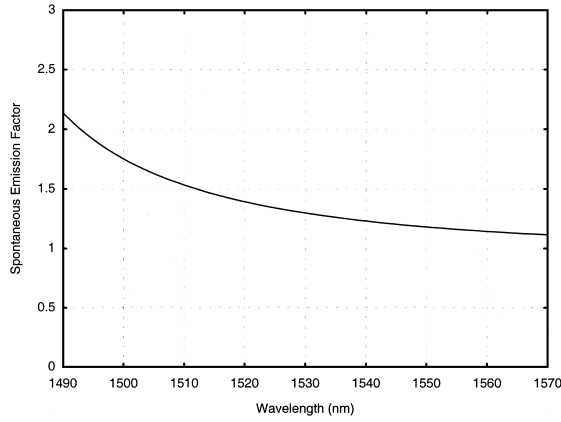


FIGURE 9 Spontaneous emission factor n_{sp} of a 1550 nm SOA.

shown in Figure 3 it is apparent that the carrier density at which the amplifier becomes transparent is higher at shorter wavelengths, due to the shift in the gain curves towards longer wavelengths as the carrier density decreases. Thus in effect N_0 is wavelength dependent, being larger at shorter wavelengths, and so Eq. (10) is in agreement with the more detailed calculation presented in Fig. 9. With this modification, Eq. (10) can be employed along with the gain models presented in Section II to produce useful models of the combined gain and noise properties of SOAs.

Spontaneous emission events may occur anywhere within the SOA, and those spontaneously emitted photons that are captured in the guided mode of the active region propagate along the amplifier, and are amplified along with the signal. Thus the optical noise present at the output in fact comprises *amplified spontaneous emission* (ASE) noise. This is illustrated in Fig. 10, in which Fig. 10 (a)

shows a schematic arrangement of an SOA with a small input signal of power P_{in} followed by an optical bandpass filter with optical bandwidth $\Delta\nu$. A schematic picture of the input optical spectrum is shown in Fig. 10 (b). At the output of the SOA the spectrum shown in Fig. 10 (c) comprises the amplified signal, of power $P_{out} = GP_{in}$, accompanied by a broad spectrum of ASE noise. After optical filtering, the noise is limited to a bandwidth $\Delta\nu$ surrounding the signal. For an ideal optical amplifier, the total ASE noise power in this band, in a single polarization state, is given by

$$P_{ASE} = n_{sp}(G - 1)h\nu\Delta\nu. \quad (11)$$

This expression is only strictly valid when the population inversion, and hence n_{sp} , is constant along the length of the amplifier. However, since the greatest contribution to the ASE noise is made by spontaneous emission occurring at the input end of the amplifier, which experiences the largest gain, in practice Eq. (11) proves highly accurate if a value of n_{sp} valid close to the amplifier input is used.

B. Signal-to-Noise Ratio and Noise Figure

In many applications of practical interest, the total optical ASE noise power given by Eq. (11) is not the most appropriate measure of the noise introduced by the SOA. In optical communications systems, for example, it is more useful to know the *electrical noise power* that accompanies the signal after detection. Modern communications systems generally employ *intensity modulation with direct detection* (IMDD), in which the information is represented directly by the power of the optical signal. Upon reception with a photodetector, the signal power is converted into photocurrent. However, since direct detection is a square-law process with respect to the optical field amplitude, conversion of the optical noise into the electrical domain is somewhat more complicated.

The process of direct detection of an optical carrier accompanied by noise is shown schematically in the power spectral diagrams of Fig. 11. The signal power and average ASE power are converted through the process of photodetection into a DC photocurrent. As a result of the square-law nature of the photodetector, the optical carrier “beats” with components of the ASE noise to generate a broadband, random current component known as *signal-spontaneous beat noise*. If the optical ASE power spectrum is flat, then the signal-spontaneous beat noise power spectral density is flat, and limited to a bandwidth of $\Delta\nu/2$. Additionally, components of the ASE noise “beat” with other components of the ASE noise to generate a further random current component known as *spontaneous-spontaneous beat noise*. For a flat optical ASE power spectrum, the spontaneous-spontaneous beat noise power

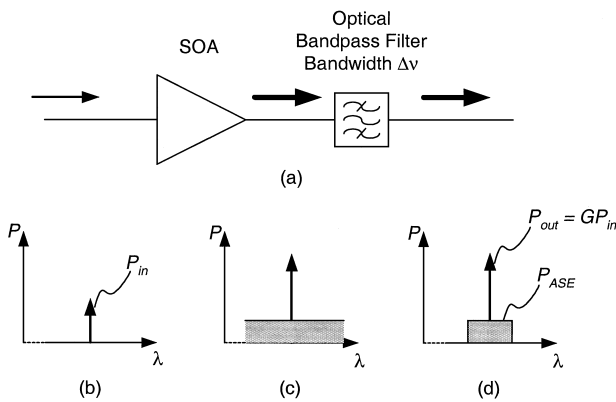


FIGURE 10 Generation of Amplified Spontaneous Emission noise in an SOA. (a) schematic diagram of amplifier arrangement, (b) schematic diagram of input optical power spectrum, (c) schematic diagram of optical spectrum at SOA output, (d) schematic diagram of optical spectrum at bandpass filter output.

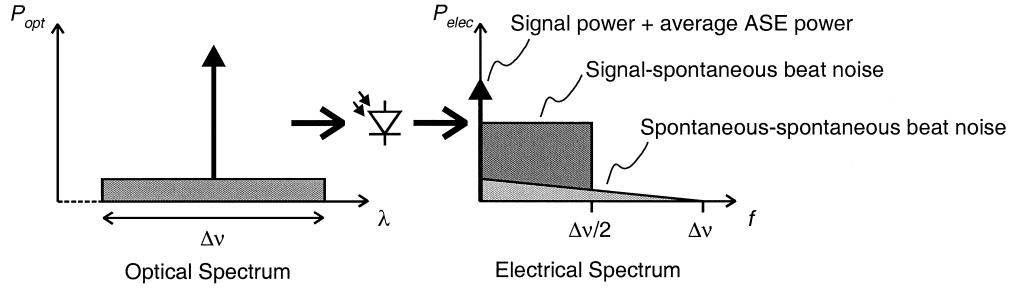


FIGURE 11 Generation of beat noise via direct detection of an optical signal accompanied by ASE noise.

spectral density is triangular, being maximum near DC and tailing off to zero at a frequency of $\Delta\nu$.

Consider the system comprising an amplifier and filter shown in Fig. 10. If this arrangement is followed by an ideal photodetector, with a quantum efficiency of $\eta = 1$, then the complete combination may be analyzed as an ideal optically preamplified receiver. In the following, all optical powers are represented as their photocurrent equivalents, e.g., $i_{in} = P_{in}q/h\nu$, and it is assumed that the electrical detection system has a bandwidth B_e . The photocurrent equivalent of the spontaneous emission power i_{sp} is

$$i_{sp} = n_{sp}(G - 1)q\Delta\nu. \quad (12)$$

The signal power S dissipated in a reference load of $1\ \Omega$ is

$$S = (Gi_{in})^2. \quad (13)$$

The noise power generated in the process of optical detection consists of three components, being the shot noise N_{shot} , the signal-spontaneous beat noise N_{s-sp} , and the spontaneous-spontaneous beat noise N_{sp-sp} , given by

$$N_{shot} = 2Beq(Gi_{in} + i_{sp}) \quad (14)$$

$$N_{s-sp} = 4Gi_{in}i_{sp}\frac{B_e}{\Delta\nu} \quad (15)$$

$$N_{sp-sp} = \frac{i_{sp}^2 B_e(2\Delta\nu - B_e)}{\Delta\nu^2}. \quad (16)$$

Note that in a practical system there is also an additional constant noise component N_{th} due to Johnson noise, and other noise sources in the receiver electronics; however this is not pertinent to the present analysis of the noise introduced by the amplifier.

It has become common practice to define the signal-to-noise ratio SNR_{out} at the output of the amplifier to be the SNR that would be observed at the electrical terminals of an ideal photoreceiver, i.e.,

$$SNR_{out} = \frac{S}{N_{shot} + N_{s-sp} + N_{sp-sp}}. \quad (17)$$

In many cases of practical interest, the signal power is much greater than the ASE noise power. Under these conditions, the terms not including Gi_{in} in Eqs. (14), (15), and (16) may be neglected, and the output SNR may be approximated as

$$SNR_{out} \approx \frac{Gi_{in}}{4i_{sp}B_e/\Delta\nu + 2qB_e}. \quad (18)$$

The *noise figure* of an electrical or optical component is a measure of the degree of degradation that is experienced by a signal passing through the component due to added noise. Typically, it is the ratio of the SNR at the input to the (reduced) SNR at the output, under a specified set of conditions and assumptions. There are a number of definitions of the noise figure of an optical amplifier in use; however, by far the most widespread is based on an input reference signal that is noise free, other than the fundamental shot-noise generated in the process of photodetection. In this case, the reference input SNR is given by

$$SNR_{in} = \frac{i_{in}}{2qB_e}. \quad (19)$$

Using Equations (19), (17), and (12), the noise figure F can be written as

$$F = \frac{SNR_{in}}{SNR_{out}} = \frac{2n_{sp}(G - 1) + 1}{G}. \quad (20)$$

A commonly used figure of merit is the noise figure of the amplifier when the gain is large. When $G \gg 1$, Eq. (20) may be approximated as

$$F = 2n_{sp}. \quad (21)$$

Thus, for even the very best possible high-gain amplifier with complete inversion ($n_{sp} = 1$), the noise figure F is never less than 2 (3.01 dB). This is a fundamental result for an IMDD system that can be arrived at more generally using Heisenberg's uncertainty principle as the starting point. While this so-called *quantum limit* has been approached in other types of optical amplifier (most notably erbium-doped fiber amplifiers), it is never achieved in SOAs because, as Fig. 9 shows, the spontaneous emission factor never reaches the theoretical limit of $n_{sp} = 1$.

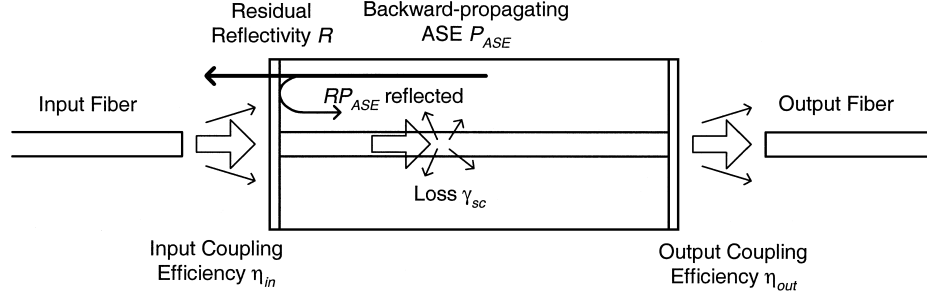


FIGURE 12 Schematic diagram illustrating the impact of coupling losses, residual facet reflectivity and internal losses on noise performance.

Additionally, as discussed in the following section, there are a number of other factors that, in practical devices, result in further increases in the noise figure.

C. “Excess” Noise

In the preceding section, the most fundamental origins of optical noise were addressed. In this section further nonideal properties of practical SOAs that further impact upon their noise performance are discussed.

Figure 12 shows schematically an SOA coupled to optical fibers at the input and output. A number of properties of the amplifier are highlighted that degrade the noise performance in practice. These are the input coupling efficiency η_{in} , the residual facet reflectivity R and the internal losses γ_{sc} . Together these contribute to an additional factor χ contributing to the noise figure that is commonly known as the *excess noise factor*. Each contribution will be discussed in turn in the following paragraphs.

Due to mismatches in the mode field profiles, and difficulties of precise alignment, power is invariably lost when optical fibers are coupled to SOAs. This loss of power is represented in Fig. 12 by the input and output coupling efficiencies η_{in} and η_{out} . The coupling efficiency is always less than one, and in practice for SOAs is typically less than 0.5. While coupling terms contribute to a reduction in the net gain of the packaged device—the fiber-to-fiber gain is lower than the SOA chip gain by a factor of $\eta_{in}\eta_{out}$ —only the input coupling efficiency impacts on the noise performance. This is because at the output facet, the signal, and the accompanying noise are attenuated equally by the coupling loss. At the input, however, if there is optical power P_{in} present in the fiber, only $\eta_{in}P_{in}$ will actually enter the SOA. Consequently, according to Eq. (17), the output SNR is reduced by the same factor.

Residual facet reflectivity also contributes to degradation in noise performance. Assuming the SOA has been designed to be a traveling-wave amplifier, the residual reflectivity is low, e.g., $R < 10^{-4}$. However, the facet-to-facet

gain of the SOA may be very high. Backward-traveling ASE noise, which is comparable in power to the forward-traveling noise, undergoes a small reflection from the front facet, as shown in Fig. 12. This reflected noise power experiences the full gain G of the amplifier, resulting in an additional noise component of power $RG P_{ASE}$ at the output facet. Of course, further reflections also contribute additional noise in the same manner; however, if $RG \ll 1$, as is the case in a traveling-wave amplifier, this first-order contribution is by far the most important.

The final factor illustrated in Fig. 12 that contributes to excess noise is the internal loss, represented by γ_{sc} . In a medium where the losses are negligible compared to the gain, Eq. (11) accurately describes the total output ASE noise power. However, if internal loss is significant, then the gain G is affected by these losses, whereas the initial generation of spontaneously emitted photons is not. The overall effect of this is that the noise is slightly higher, relative to the gain, than predicted by Eq. (11), by a factor of $g/(g - \gamma_{sc})$.

Taking into consideration all of these non-ideal factors, the practical high-gain noise figure of an SOA is given by

$$F_{SOA} \approx 2\chi n_{sp}, \quad (22)$$

where the excess noise factor χ is:

$$\chi = \frac{(1 + RG)}{\eta_{in}} \frac{g}{g - \gamma_{sc}}. \quad (23)$$

Consequently, in practice the noise figure of an SOA in the high-gain region is typically in the range 6 dB to 10 dB—much higher than the theoretical minimum of 3 dB predicted by Eq. (20).

IV. DYNAMICS

A. Carrier Lifetime

The main factor that determines the maximum response speed of an SOA is the carrier lifetime. The spontaneous

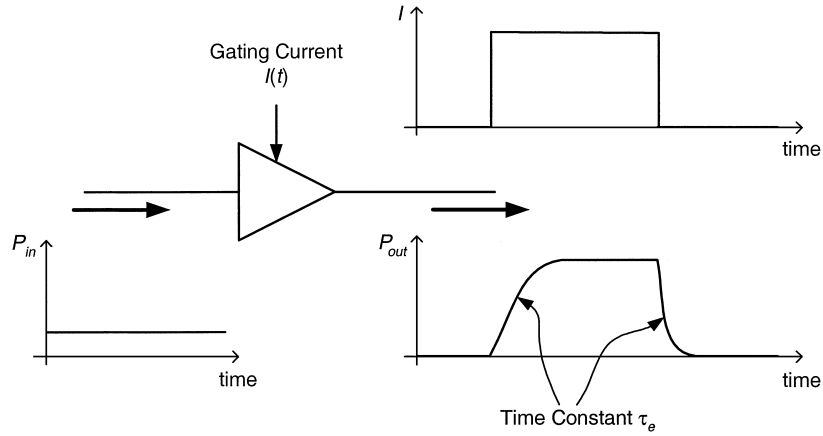


FIGURE 13 Operation of an SOA as an optical gate.

carrier lifetime τ_s was introduced in Eq. (4), where it was stated that this parameter is typically on the order of 1 ns. Consequently, it is reasonable to suppose that the maximum modulation bandwidth of an SOA would be on the order of 1 GHz. However, SOAs have been demonstrated as all-optical switching devices, wavelength converters, and regenerators among other applications, operating at speeds of tens of gigahertz. The explanation for this lies in the process of stimulated emission. In the presence of an optical field, carriers can be transported between the valence and conduction bands not only through the processes of current injection and spontaneous emission, but also through stimulated emission and absorption. Stimulated processes themselves are almost instantaneous, and the rate at which such events occur is proportional to the number of photons present to initiate them.

Looking again at Eq. (4), and ignoring the carrier-independent terms corresponding to electrical pumping by the injection current I and optical pumping due to the incomplete inversion represented by the term proportional to $N_0 P$, the carrier rate equation has the basic form:

$$\frac{dN_c}{dt} \sim -\left(\frac{1}{\tau_s} + \frac{\Gamma a P}{h\nu}\right)N_c. \quad (24)$$

A first-order differential equation of this form has an *effective time constant* τ_e given by

$$\tau_e^{-1} = \tau_s^{-1} + \frac{\Gamma a P}{h\nu}. \quad (25)$$

The second term in this expression increases in proportion to the optical power, resulting in a consequent reduction in the effective carrier lifetime. In practical devices, injected optical power on the order of 10 mW may result in reduction of the carrier lifetime to an effective value of 100 ps or less.

B. Current Modulation

An important application of semiconductor optical amplifiers is in the area of optical gating and switching. The process of optical switching by current modulation is shown schematically in Fig. 13. A CW signal of power P_{in} is injected into an SOA, to which the bias current I is modulated by the square gating signal shown. The gain of the amplifier responds to these changes in injected current according to Eq. (3). This response is approximately exponential, with an effective time constant as given in Eq. (25). Thus the SOA acts as an optical gate that is “off” when the bias current is low, and “on” when the bias current is high. Note that according to Eq. (25) the effective lifetime is shorter when the optical intensity in the SOA is higher, thus switching from “on” to “off” is typically faster than switching from “off” to “on.” Switching from the “on” to “off” states can be achieved in sub-nanosecond timescales, whereas switching from the “off” to “on” states may take a few nanoseconds.

C. Optical Modulation

Section II.D shows how gain compression in a semiconductor optical amplifier manifests itself as a strong dependence of gain on input power level, as shown in Fig. 8. In some applications, this dependence of gain on power level can lead to problems. However, there are some applications where gain compression can be used to advantage. One such application is in so-called cross-gain modulation wavelength converters. Cross-gain modulation wavelength converters (sometimes called XGM wavelength converters) take advantage of the gain compression effect. The operation of XGM wavelength converters is based on optical modulation, or more precisely, optical cross-modulation. Figure 14 outlines the operating principle of an XGM wavelength converter.

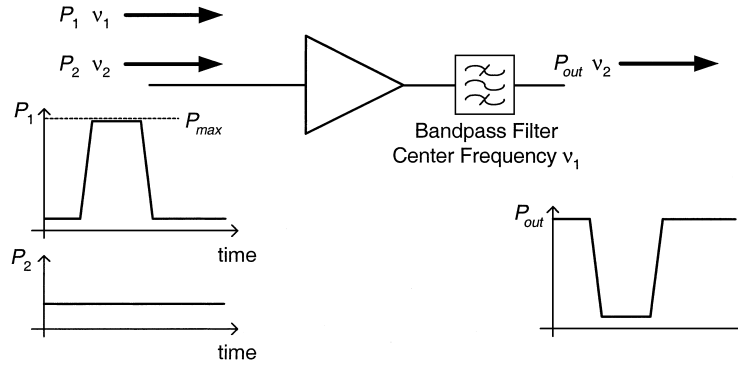


FIGURE 14 Schematic of XGM wavelength converter.

The wavelength converter has two input optical signals. One of these signals is a CW signal optical of power P_2 and frequency ν_2 , and the other is an intensity-modulated signal of power P_1 at optical frequency ν_1 . The modulated input signal has a maximum power level P_{max} that drives the SOA well into the gain compression regime. For the device with gain compression characteristics shown in Fig. 8, P_{max} typically lies in the region of 0 to 10 dBm. At this power level, the gain of the SOA is 20 dB when the input signal is “off” and around 10 to 5 dB when the input signal is “on.” Thus the gain of the amplifier is modulated by the optical input signal, and undergoes excursion of around 10 dB or more.

In the arrangement shown in Fig. 14, the optical input signal at frequency ν_1 optically modulates the gain of the SOA. The output power P_{out} of the CW input signal at frequency ν_2 is modulated by the changing gain of the SOA. Thus the data on the input signal at ν_1 are transferred to the signal at ν_2 . As shown in Fig. 14, when the input signal at ν_1 goes “high,” the output signal at ν_2 goes “low” and vice versa. Thus the wavelength converter inverts the data. An optical band stop filter at the output of the SOA removes the signal at ν_1 . Note that since the amplifier bias current is always “on,” and the CW input at frequency ν_2 is always present even when the signal at frequency ν_1 is “off,” the effective lifetime remains short according to Eq. (25). Thus all transitions between the high and low output levels may be very fast, unlike the case of optical gating discussed in Section B, in which the “off” to “on” transitions are slower than “on” to “off” transitions. Wavelength conversion of optical signals using XGM in SOAs has been reported at bit-rates in excess of 40 Gb/s.

V. COHERENT INTERACTIONS

A. Refractive Index and Optical Phase

In the previous sections of this article, the optical intensity and the relationship between the carrier density N_c and the

amplifier gain have been considered in detail. The refractive index of an SOA is also dependent upon the carrier density, and this fact leads to device properties, such as chirp and phase modulation, that are important in applications such as optically amplified transmission, wavelength conversion, and digital signal regeneration.

The group refractive index for light traveling in an InGaAsP SOA operating in the 1300- to 1600-nm wavelength range has been measured to be in the range 3.5 to 4. The change in refractive index relative to the change in carrier density, i.e., dn/dN_c , has been measured at around $-1 \times 10^{-20} \text{ cm}^3$ to $-2 \times 10^{-20} \text{ cm}^3$. Note that the refractive index decreases as the carrier density increases. These values may seem small, however they are sufficient to enable a total phase shift along the SOA corresponding to a full cycle of the optical field for relatively modest variations in carrier density (e.g., $0.2 \times 10^{-18} \text{ cm}^{-3}$) even in short devices (e.g., 500 μm).

The dependence of the gain and refractive index on the carrier density are related to the well-known *linewidth enhancement factor* α for semiconductor lasers according to

$$\alpha = -\frac{4\pi}{\lambda} \frac{dn/dN_c}{dg/dN_c}. \quad (26)$$

The quantity dn/dN_c is not strongly wavelength dependent; however, the quantity dg/dN_c —which is essentially the gain cross-section a defined previously—is strongly wavelength dependent, as shown in Fig. 3. Specifically, dg/dN_c is larger at shorter wavelengths. Consequently α is greater at longer wavelengths. Over the wavelength range shown in Fig. 3 the linewidth enhancement factor would typically vary between $\alpha = 4$ at shorter wavelengths, and $\alpha = 12$ or more at longer wavelengths.

B. Optical Phase Modulation

When the carrier-density in an SOA is modulated either via a change in the drive current, or a change in the optical power level, the refractive index of the device changes and

this results in a change in phase shift of the field passing through the amplifier. This effect is optical phase modulation. In some situations it has undesirable and detrimental outcomes and in some situations it can be used to advantage. An example of a detrimental effect is in intensity-modulated direct detection communications systems, where phase modulation by an intensity-modulated signal causes a chirp or optical frequency shift of the optical carrier. This chirp can result in increased pulse spreading in dispersive media such as optical fibers. On the other hand cross-phase modulation, where an intensity-modulated signal modulates the phase of another signal can be used to advantage in so-called cross-gain modulation wavelength converters.

In a cross-phase modulation wavelength converter (XPM wavelength converter), an SOA is placed in one arm of an optical interferometer. Intensity modulation on an optical signal injected into the SOA causes the phase of a second optical signal to change. At the output of the interferometer, the phase-shifted second optical signal is added coherently with a non-phase-shifted signal at the same frequency. The net result is that the output intensity of the second optical signal is dependent on the phase shift in the optical modulator and is therefore dependent on the optical power of the first signal. Consequently, data on the first optical frequency is transferred to the second optical frequency.

C. Four-Wave Mixing

Consider the situation depicted in Fig. 15 (a), in which two optical carriers with different optical frequency are simultaneously injected into an SOA. Each optical carrier

may be represented by its time-varying complex electric field, as illustrated by the schematic spectrum shown in Fig. 15 (b):

$$E_1(t) = A_1 e^{j2\pi\nu_1 t} \quad (27)$$

$$E_2(t) = A_2 e^{j2\pi\nu_2 t}, \quad (28)$$

where A_1 and A_2 are the complex amplitudes of the two carriers, and ν_1 and ν_2 are their optical frequencies. Assume that the field amplitudes are normalized so that the optical irradiance within the SOA is given by the squared magnitude of the field, e.g., $P_1 = |E_1(t)|^2 = |A_1|^2$. Upon detection by a square law device, such as a photodetector, the total time-varying power that would be observed in this optical field is

$$P(t) = |A_1|^2 + |A_2|^2 + A_1 A_2^* e^{-j\Omega t} + A_1^* A_2 e^{j\Omega t}, \quad (29)$$

where

$$\Omega = 2\pi(\nu_2 - \nu_1) \quad (30)$$

is the angular difference frequency between the two optical fields.

This process of “beating” between the two optical fields occurs within the SOA, because the carrier density indeed responds to the square of the electric field, just like the photoelectrons in a photodetector. If Eq. (29) is substituted into the rate equation, Eq. (4), then the resulting differential equation has an approximate small-signal solution of the form

$$N_c(t) = \overline{N}_c + \Delta N_c e^{j\Omega t} + \Delta N_c^* e^{-j\Omega t}. \quad (31)$$

In other words, the carrier density oscillates around a mean value (determined by the total average power) at

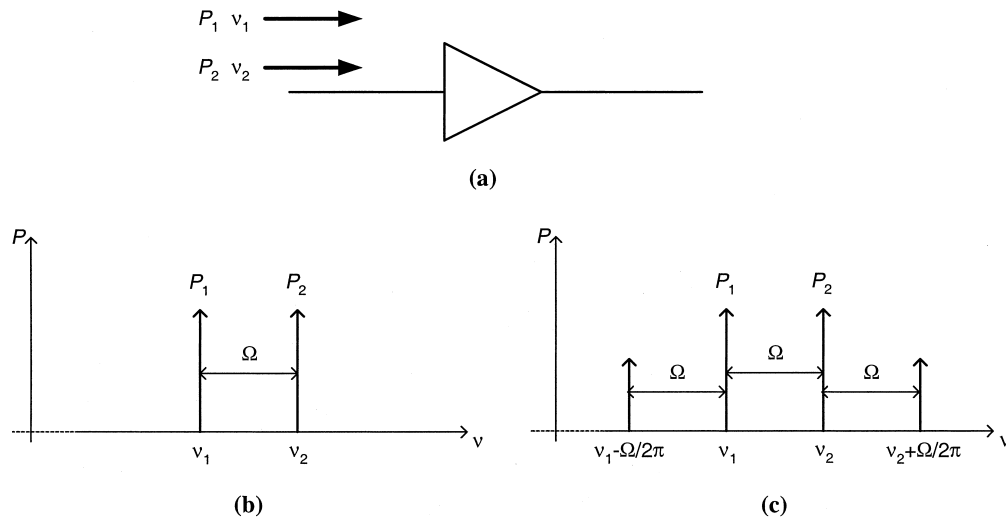


FIGURE 15 Illustration of four-wave mixing in an SOA. (a) Two optical carriers at slightly different wavelengths injected into an SOA. (b) Schematic diagram of input optical spectrum. (c) Schematic diagram of output optical spectrum.

a frequency equal to the difference in frequency of the two optical carriers. The magnitude of these carrier oscillations is given by the expression

$$\Delta N_c = -\frac{(\bar{N}_c - N_0)A_1^*A_2/P_e}{1 + j\Omega\tau_e}, \quad (32)$$

where P_e is an effective saturation irradiance (given by $P_e = h\nu/\Gamma a\tau_e$), and τ_e is the effective carrier lifetime defined in Eq. (25). Thus the magnitude of the carrier oscillations is proportional to the gain, to the magnitude of the electric fields, and inversely proportional to an effective saturation power. Furthermore, Eq. (32) has the form of a single-pole response, with a 3-dB bandwidth determined by the inverse of the effective carrier lifetime. As discussed in Section IV.A, in the presence of a strong optical

field the effective carrier lifetime can be less than 100 ps, and so the bandwidth of this process may be on the order of tens of gigahertz.

The two optical carriers therefore propagate through the SOA in the presence of a time-varying gain and refractive index. This results in amplitude and phase modulation of both carriers, and the generation of new optical signals in the form of intermodulation products at optical frequencies of $\nu_1 - \Omega/2\pi$ and $\nu_2 + \Omega/2\pi$, as shown schematically in Fig. 15 (c). This nonlinear process is known as *four-wave mixing*. A potential application of this process is wavelength conversion. Data encoded on the probe signal are transferred to the newly generated intermodulation products, and these signals can be separated from the pump and probe signals using optical filtering.

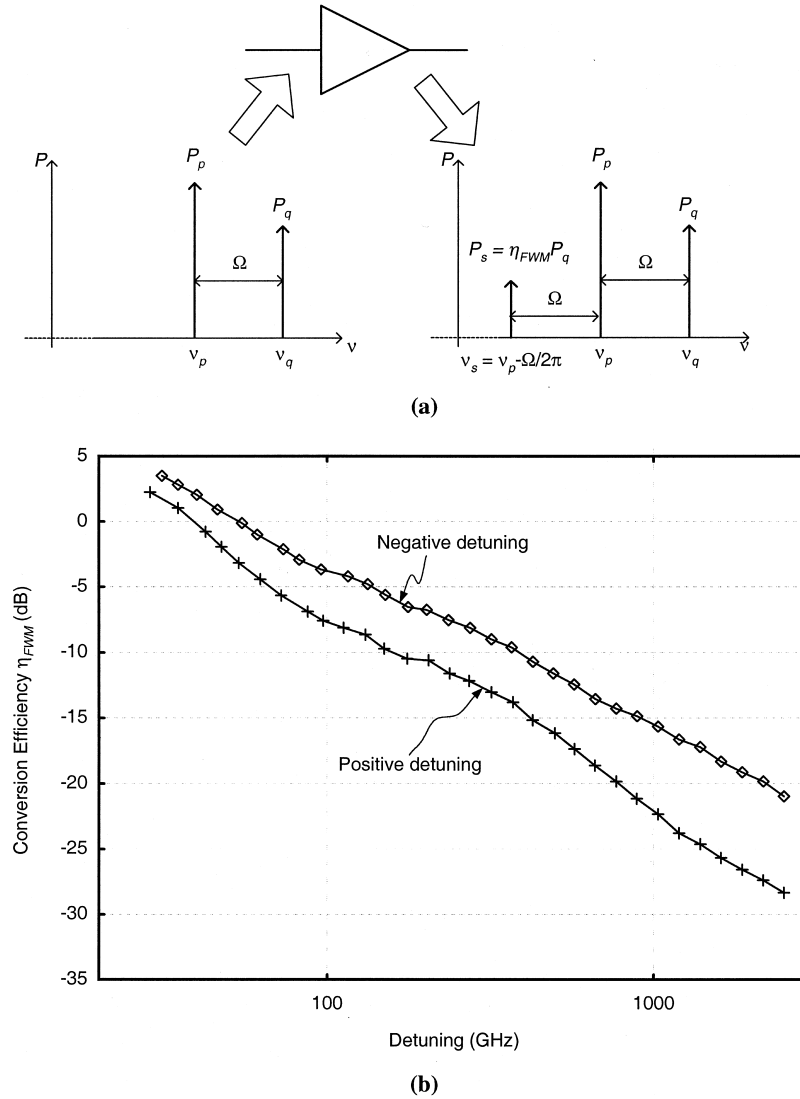


FIGURE 16 Frequency response of four-wave mixing in an SOA. (a) Schematic arrangement of the frequency domain pump-probe measurement. (b) Four-wave mixing conversion efficiency as a function of pump-probe detuning.

The mechanism of carrier oscillations outlined above is the most significant one of a number of processes that lead to four-wave mixing in SOAs. This mechanism is sometimes known as an *interband process*, because it results from the oscillation of carriers between the valence band and conduction band. Additionally, there are a number of *intraband processes* that result in similar oscillations of the gain and refractive index experienced by the propagating field. Such processes do not correspond to a change in the overall carrier population, but rather to a redistribution of the carriers amongst the energy states of the conduction band. Such processes may be very fast—typical time constants are sub-picosecond—and so in principle may be extremely broadband, e.g., into the terahertz range. The most significant intraband processes are *spectral hole burning* and *carrier heating*. However, even these mechanisms are some orders of magnitude weaker than the interband carrier oscillations, and so despite their very broad bandwidth, in practice they result in little more than a perturbation of the dominant interband response.

Figure 16 shows the results of an experimental measurement of the frequency response of the four-wave mixing process in an SOA. The experimental technique used is a frequency domain pump-probe measurement, illustrated schematically in Fig. 16(a). Two optical carriers are injected into the SOA, a strong *pump signal* P_p and a much weaker *probe signal* P_q . The frequency detuning Ω is defined to be the probe frequency minus the pump frequency, i.e., $\Omega = \nu_q - \nu_p$. The power P_s of the four-wave mixing product at the frequency $\nu_s = \nu_p - \Omega$ is measured, and the conversion efficiency, defined as $\eta_{FWM} = P_s/P_q$ is determined. The results of this measurement, for a range of positive and negative detuning, are shown in Fig. 16(b). The response is dominated by the roll-off determined by the single-pole characteristic of the interband oscillations, as defined by Eq. (32). The difference between the response curves for positive and negative detuning is due to the differing contributions made by the gain modulation and index modulation for the two different conversion directions, as well as to interactions between the dominant interband processes and the intraband processes.

VI. CONCLUSIONS

Semiconductor optical amplifier technology is now relatively mature in single device form, but more developments are expected as applications drive the need for integration of multiple SOAs on a single semiconductor chip. Semiconductor optical amplifiers generally cannot compete with fiber amplifiers for power amplification, loss compensation and low noise preamplification in long-haul

fiber optic transmission. However, they have a number of advantages, such as their compact size and low power consumption, that may make them attractive in more cost-sensitive applications such as metropolitan area and access networks, where tradeoffs in gain, noise performance, and output power may be made.

The relatively fast dynamic response of SOAs opens up a rich range of applications other than simply providing gain. As an optical gating device, the SOA's switching time of less than 1 ns means that it is one of the fastest gating technologies available. Arrays of fast SOA optical gates may find application in future optical packet-switched telecommunications systems, where conventional optical switches cannot provide the required switching speed. Future optical packet switches may use large numbers of integrated SOAs and other devices.

Unless deliberate and complex steps are taken in the device design, or special controlling circuitry is employed, SOAs exhibit an extremely strong nonlinear response. In conventional signal gain applications this can be a disadvantage, potentially leading to signal distortion and crosstalk. However, the strong nonlinear response of SOAs can be used to advantage in certain signal-processing applications. One important example is optical wavelength conversion. A wavelength converter shifts modulated optical signals from one optical carrier frequency to another. There are a number of possible structures for wavelength converters based on SOAs, including, cross-gain modulation, cross-phase modulation, and four-wave mixing. Like optical packet switching, optical wavelength conversion in dense wavelength-division multiplexed systems will require multiple devices, and a key challenge to device and sub-system designers will be to devise techniques to integrate large numbers of SOAs along with other functional components.

ACKNOWLEDGMENT

The author would like to thank Rodney S. Tucker, of the University of Melbourne, for his helpful comments on, and contributions to, this article.

SEE ALSO THE FOLLOWING ARTICLES

LASERS, SEMICONDUCTOR • OPTICAL FIBER COMMUNICATIONS • PHOTONIC BANDGAP MATERIALS

BIBLIOGRAPHY

- Agrawal, G. P., and Dutta, N. K. (1993). "Semiconductor Lasers," Kluwer Academic Publishers, Dordrecht.
- Agrawal, G. P. (1995). Semiconductor laser amplifiers. In "Semiconductor Lasers: Past, Present and Future" (G. P. Agrawal, eds.), pp. 243–283, AIP Press, Woodbury, New York.

- Henry, C. H. (1986). "Theory of spontaneous emission noise in open resonators and its application to lasers and optical amplifiers," *J. Lightwave Technol.* **LT-4**, 288–297.
- Olsson, N. A. (1989). "Lightwave systems with optical amplifiers," *J. Lightwave Technol.* **7**, 1071–1082.
- Olsson, N. A. (1992). "Semiconductor optical amplifiers," *Proceedings of the IEEE* **80**, 375–382.
- O'Mahony, M. J. (1988). "Semiconductor laser optical amplifiers for use in future fiber systems," *J. Lightwave Technol.* **6**, 531–544.
- Westbrook, L. D. (1986). "Measurements of dg/dN and dn/dN and their dependence on photon energy in $\lambda = 1.5 \mu\text{m}$ in GaAsP laser diodes," *IEE Proc. J.* **133**, 135–142.
- Yamamoto, Y., and Mukai, T. (1989). "Fundamentals of optical amplifiers," *Optical Quantum Electr.* **21**, S1–S14.
- Yariv, A. (1997). "Optical Electronics in Modern Communications," Oxford University Press, New York.



Optical Detectors

Robert H. Kingston

Massachusetts Institute of Technology

- I. Optical Detection: Thermal and Electronic
- II. Fundamental Limits to Electronic Detector Performance
- III. Vacuum Photodetectors
- IV. Semiconductor Photodetectors
- V. Graph of Detector Characteristics
- VI. Enhanced Detection Techniques
- VII. Optical Imaging Detectors

GLOSSARY

Detectivity, D The inverse of NEP . The higher the detectivity, the higher the sensitivity.

Electron charge, q 1.6×10^{-19} C.

Electron volt A convenient unit of energy for detector calculations, an electron volt, eV, is the energy gained by an electron in a one volt change in potential.

$$1 \text{ eV} = 1.6 \times 10^{-19} \text{ J.}$$

Frequency, ν , f Optical or infrared wave frequency, ν .
Electrical current or voltage frequency, f .

Noise equivalent power, NEP That optical power which produces a detector output current equal to the total noise current. A detector with input power, NEP , produces an output signal-to-noise ratio of unity.

Photon energy, $h\nu$ The quantized energy absorbed or emitted from an electromagnetic wave, given by $h\nu = 1.24/\lambda$ (in μm) eV.

Quantum efficiency, η The ratio of the number of elec-

trons produced in the detector to the number of photons or quanta of energy, $h\nu$, incident.

Responsivity, \mathfrak{R} The detector output current per unit optical power, given by $\mathfrak{R} = \eta G \lambda$ (in μm)/1.24 in A/W , with G the internal detector electron gain.

Shot noise The current fluctuation or noise produced by the random excitation of photoelectrons.

Specific detectivity, D^* The detectivity, D , normalized to a detector area of 1 cm^2 and a bandwidth of 1 Hz.

Thermal energy, kT At room temperature (300 K) $kT = 0.026 \text{ eV}$.

OPTICAL DETECTORS sense and measure electromagnetic radiation in the broad spectral region from the longest infrared wavelength atmospheric “window” at $\lambda = 10 \mu\text{m}$ to the ultraviolet end of the visible spectrum at $\lambda = 0.3 \mu\text{m}$. The detection process may depend upon the heat generated by the incident radiation (thermal detection) or by the photoexcitation of electrons or

holes (electronic detection). The two most important performance parameters are the sensitivity and temporal response. In addition to the use of specific detector types, system performance may be enhanced by using heterodyne detection or optical amplification. Detectors may be grouped or integrated on chip into two-dimensional imaging arrays, the most useful being the charge-coupled device (CCD) array.

I. OPTICAL DETECTION: THERMAL AND ELECTRONIC

Optical and infrared detectors fall into two distinct categories, either thermal or electronic. In thermal detectors, optical power is converted to heat, and the associated temperature increase is measured electrically, for example, by a thermocouple. In contrast, electronic detectors depend upon the photoexcitation of an individual electron and its consequent emission into a vacuum (the *external* photoelectric effect) or its flow in a semiconductor material (the *internal* photoelectric effect). Thermal detectors, with appropriate absorbing coatings, are sensitive to all wavelengths, but have moderate sensitivity and relatively slow response times, typically greater than a millisecond. Electronic detectors approach single “photon” detection capability, have response times as short as a nanosecond, but must be tailored to the wavelength region of interest, because of the threshold wavelength of the photoelectric effect.

The performance of any optical detector is measured by two critical parameters, the sensitivity and the temporal response. The sensitivity can be expressed in several ways. The usual measure is the *noise equivalent power* or *NEP*, which is defined as that optical power which produces an electrical signal voltage (or current) equal to the root mean square (rms) noise voltage (or current). The ratio of the signal voltage to the noise voltage, designated $(S/N)_V$, is the critical performance parameter in a measurement or in an operating system. In the simplest case, the fractional precision of a measurement is the inverse of $(S/N)_V$. In terms of the *NEP*, the signal-to-noise ratio (except for the special case of signal-induced noise discussed in Section II.A) may be written

$$\left(\frac{S}{N}\right)_V = \frac{v_s}{\sqrt{v_n^2}} = \frac{i_s}{\sqrt{i_n^2}} = \frac{P}{NEP} \quad (1)$$

where the subscripts, *s* and *n*, refer to the signal and noise, respectively. It should be emphasized that the output electrical *voltage* from a detector or detection circuit is proportional to the optical *power*. Thus the output electrical power is proportional to the *square* of the optical power, in

contrast to most radio and microwave frequency systems where the electrical output *power* is proportional to the input radio frequency *power*. The significance of this will be apparent below.

Two other useful concepts are the detectivity, *D*, and the specific detectivity, *D**. The detectivity, in W^{-1} , is the inverse of the *NEP* and as such, the larger the number the better the sensitivity or ability to detect a weak signal. Finally, *D** is the value of *D* normalized to a detector area of 1 cm² and an output electrical bandwidth of 1 Hz. For a given detector of area, *A*, operating at output bandwidth, *B*, the specific detectivity is given by

$$D^* = D\sqrt{AB} = \sqrt{AB}/NEP: A \text{ in cm}^2, B \text{ in Hz}, \quad (2)$$

and has the unusual dimensions of cm $\sqrt{\text{Hz}}/\text{W}$. *D** is a valid measure of detector performance when, as is often the case, the *mean-square* noise current is proportional to detector *area* as well as to the electrical output *bandwidth*. Since $D = 1/NEP$ is proportional to the square root of the noise current, it varies inversely as the square root of the area-bandwidth product. *D** is most appropriate to infrared detectors, especially in the Long-wavelength infrared (LWIR) or 10- μm wavelength region.

A. Wavelength Region

The wavelength region considered here extends from approximately 10 μm in the LWIR to 3 μm at the ultraviolet (UV) end of the visible spectrum. Although somewhat arbitrary, these wavelengths include all the major transmission bands of the atmosphere, from the 8- to 12- μm “window” in the LWIR region to the edge of the UV absorption bands at the violet end of the visible region. These limits also have significance in terms of the photon energy. In the 10- μm infrared wavelength region, the photon energy, $h\nu$ (in electron volts, given by $1.24/\lambda(\lambda \text{ in } \mu\text{m})$), is about 0.1 eV, which is the order of the room-temperature thermal energy, $kT = 0.026 \text{ eV}$. In the UV, photon energies are greater than 4 eV, and the radiation is energetic enough to ionize individual atoms or molecules. Detection of such “ionizing” radiation is not considered in this article.

B. Thermal Detectors

Historically, the first and also the most straightforward way of detecting optical radiation is by measuring the temperature rise of an object which is absorbing the incident optical energy. In fact, infrared radiation, at wavelengths longer than about 0.7 μm , was discovered by Herschel in the 19th century using a prism to refract the light and a blackened thermometer bulb as a detector. Modern thermal detectors convert the measured temperature rise to an electrical signal using several different devices. These include

the thermocouple or thermopile whose voltage output is a function of temperature, the bolometer with a temperature-sensitive resistance, and the pyroelectric detector which has a dielectric polarization which is temperature dependent. In the last case, the current is proportional to dP/dt , the rate of change of polarization, and as a result, the pyroelectric is especially useful for short pulses or high-frequency modulation. Despite this small advantage for the pyroelectric device, thermal detectors are generally inferior to electronic detectors in sensitivity as well as in frequency response. First, fundamental temperature fluctuations in the device produce significant noise compared to electronic sensors. Second, the temporal response is determined by the ratio of the heat capacity of the device to its thermal conductance to the surroundings, the latter by radiation or heat flow along the metallic leads. Such time constants are significantly longer than electron transit times in typical photoelectric devices.

A thermal detector can be modeled as shown in Fig. 1. Optical or infrared power, P , is focused onto a small detector of area, A , in a chamber at temperature, T . The resultant increase in temperature, ΔT , produces a signal voltage output, v_s , at the output of an electrical sensing network, N . The sensing network might be a simple amplifier in the case of a thermocouple or a resistance bridge in the case of a bolometer. The fundamental limit in the NEP is governed by the temperature fluctuation of the detector element, which in turn is determined by the ambient temperature, T , the heat capacity, C_T , and the heat conductance, G_T . The signal voltage may be written

$$v_s \propto \Delta T = K \frac{P}{G_T} \frac{1}{\sqrt{1 + 4\pi^2 f^2 \tau^2}} \quad (3)$$

$$\tau = \frac{C_T}{G_T} \text{ with } C_T \text{ in joules/K, } G_T \text{ in watts/K,}$$

where f is the electrical frequency and the response is identical to that of a low-pass RC circuit. There is also

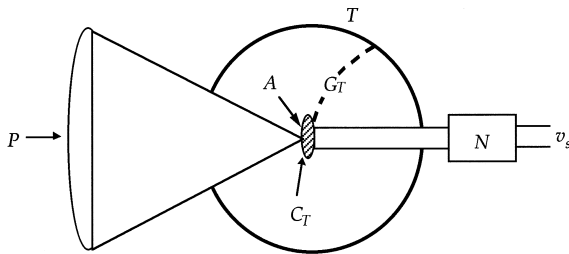


FIGURE 1 A thermal detector of optical radiation. The detector of area, A , is illuminated by optical power, P , and is enclosed in a chamber at temperature, T . G_T and C_T are the thermal conductance and heat capacity of the element, and the network, N , converts the temperature change to a signal voltage, v_s .

a random voltage or noise produced by the inherent temperature fluctuation of the detector element. These fluctuations arise from the statistical variation of the emitted and absorbed equilibrium radiation. From thermal statistics, the total mean square fluctuation of a body with heat capacity, C_T , is given by

$$\overline{\delta T^2} = kT^2/C_T, \quad (4)$$

which, for a bandwidth, B , of the output electrical circuit, becomes

$$\overline{\delta T^2} = 4kT^2 B/G_T. \quad (5)$$

The limiting NEP for a thermal detector is obtained by setting the signal-induced temperature change equal to the rms temperature fluctuation,

$$\Delta T = NEP/G_T = \sqrt{\overline{\delta T^2}} = \sqrt{4kT^2 B/G_T}$$

$$\therefore NEP = \sqrt{4kT^2 G_T B}. \quad (6)$$

Obviously, the smallest NEP or the highest detectivity occurs at the lowest thermal conductance, G_T , but at the longest time response, given by C_T/G_T . This limiting sensitivity occurs when the conductance is solely caused by the radiative flow between the detector surface and the surrounding chamber. With unit emissivity over the sensitive area, corresponding to complete absorption of the incident power, the thermal conductance becomes, from the Stefan–Boltzmann law,

$$G_T = \frac{d}{dT}(\sigma T^4 A) = 4\sigma T^3 A, \quad (7)$$

from which the noise equivalent power is

$$NEP = 4\sqrt{\sigma T^5 AB} = 4\sqrt{(\sigma T^4)(kT)AB} \quad (8)$$

a form consistent with the definition of D^* in Eq. (2). As a benchmark, let the detector area be 1 mm^2 , and the bandwidth be 1 kHz . Then,

$$\sigma T^4 \text{ at } 300 \text{ K} = 0.046 \text{ W/cm}^2$$

$$kT = (kT)_{\text{eV}}(q_{\text{electron}})$$

$$= (0.026)(1.6 \times 10^{-19}) \text{ J.}$$

$$\begin{aligned} NEP &= 4\sqrt{(\sigma T^4)(kT)AB} \\ &= 4\sqrt{(0.046)(4.2 \times 10^{-21})(10^{-2})(10^3)} \\ &= 1.8 \times 10^{-10} \text{ W.} \end{aligned} \quad (9)$$

Electric noise in the temperature sensing circuitry can increase this value by as much as an order of magnitude, depending upon the type of detector. In addition, if faster response is required, the thermal conductance must be increased by a solid conducting path to the chamber walls,

resulting in a reduced sensitivity. In contrast a perfect photon counter, as discussed below, would have an NEP of

$$\begin{aligned} NEP &= 2h\nu B = 2(h\nu)_{\text{eV}} q_{\text{electron}} B \\ &= 2[1.24/\lambda(\mu\text{m})](1.6 \times 10^{-19})(10^3) \\ &= [4.0 \times 10^{-16}/\lambda(\mu\text{m})] \text{ W} \end{aligned} \quad (10)$$

which is many orders of magnitude smaller than that of the thermal detector. Since electronic detectors approach photon-counting performance in the visible, they are far superior to thermal detectors in this spectral region. Albeit, for moderate sensitivity and response speed, the broad spectral response and simplicity make thermal detectors an attractive choice.

At the longer wavelengths, especially in the LWIR at $10 \mu\text{m}$, even electronic detectors are limited by the thermal or “background” radiation, and thermal detectors can have comparable sensitivity, while limited in frequency response. For either detector type, cooling to the order of 100 K or less yields optimum operation. Fluctuation of the incident infrared power then produces the limiting noise in either case. In specialized applications, usually in laboratory measurements where frequency response is not critical, cooling to liquid helium temperatures, 4 K, provides thermal detector sensitivity comparable to that of electronic detectors, and the broad spectral response obviates the necessity of specialized semiconductor fabrication to match the desired wavelength threshold. Of particular interest is the superconducting bolometer, whose rapid resistance change with induced temperature yields a strong electrical signal compared to typical electric circuit noise.

C. Electronic Detectors

The internal photoelectric effect, just as infrared radiation, was also first observed in the 19th century, when certain minerals such as selenium or lead sulfide were found to increase their electrical conductivity in the presence of light. These photoconductors depend upon the photoexcitation of bound electrons and/or holes into the conduction and/or valence bands of the material. Then, at the turn of the century, external photoemission was discovered in vacuum diodes. As first explained by Einstein, the photoelectric effect was found to have a threshold wavelength determined by the relation $h\nu = hc/\lambda \geq E$, where E is the energy required for the electron to exit the material. In the case of a semiconductor, the excitation energy, E , is that of the gap between the valence and conduction bands or the ionization energy of an impurity in the material. The electronic detector family has two main branches, the first being the vacuum photodiode and its more useful

adaptation, the vacuum photomultiplier. Following these are the semiconductor devices, the photoconductor and the photodiode. The former behaves as a light-controlled variable resistance, while the semiconductor photodiode, almost identical to the vacuum photodiode in characterization, is a high impedance current source with current magnitude proportional to the incident radiation.

Electronic detectors offer the ultimate in frequency response, as high as tens of gigahertz, and especially in the visible, approach photon-counting or quantum-limited performance. As such, they offer magnitudes of improvement in sensitivity over thermal devices. In the limit of photon-counting performance, the signal measurement fluctuation or noise is produced by the random production of photo electrons. In many cases, electrical noise in the postdetection amplifier, rather than “photon” noise, limits the sensitivity.

II. FUNDAMENTAL LIMITS TO ELECTRONIC DETECTOR PERFORMANCE

Figure 2(a) depicts a typical photodiode detector circuit, in this case ac coupled to an amplifier. The equivalent circuit, in Figure 2(b), shows a current source proportional to incident optical power, a parallel “dark” current, i_D , a noise current source which represents the “shot” or photon noise associated with these currents, and a resistance and capacitance attributable to both the detector and the electrical amplifier. For the vacuum or semiconductor diode, the responsivity, \mathfrak{R} , determined for a given wavelength, λ , and quantum efficiency, η , is

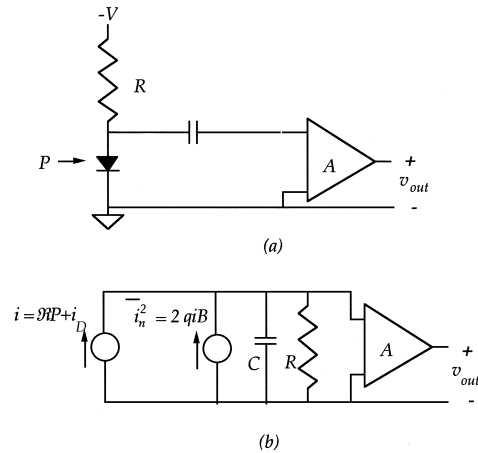


FIGURE 2 (a) Typical circuit diagram for an electronic detector, in this case, an ac coupled reverse-biased semiconductor photodiode. (b) Equivalent circuit for electronic detector. The noise current shown is that of a simple photodiode. Detectors with internal gain have a modified noise term.

$$\mathfrak{R} = \eta q / h\nu = \eta \lambda (\mu\text{m}) / 1.24 \text{ A/W} \quad (11)$$

and its magnitude is seen to be of the order of unity in the near infrared region for unit quantum efficiency. Quantum efficiencies are generally in the range from 0.2 to 1.0 for semiconductor diodes, when operated at wavelengths shorter than the threshold determined by the energy gap. Typical vacuum photodiodes are generally limited to the visible region, since η falls below 0.001 beyond a wavelength of about $1 \mu\text{m}$. The dark current, arising from the thermal generation of electrons in both the vacuum and semiconductor devices, is usually negligible or can be made so by cooling the detector.

The ultimate sensitivity of an electronic detector is then determined by its responsivity and the various noise processes in the detector and the following electric circuitry. First to be considered is the fundamental “photon” or signal-noise contribution.

A. Quantum or Photon Limits

The photoexcitation process in an electronic detector occurs at an average rate proportional to the instantaneous incident optical power as determined by the responsivity, \mathfrak{R} . The electron production process is random, however, and obeys Poisson statistics, which state that in any interval of time in which the expected or average value of the electron count is n , then the actual probability of k events during this time interval is given by

$$p(k, n) = n^k e^{-n} / k! \quad \text{and} \quad \sum_{k=0}^{\infty} p(k, n) = 1; \quad \bar{k} = \sum_{k=0}^{\infty} k p(k, n) = n. \quad (12)$$

This distribution is illustrated in Fig. 3, for $n=2$ and $n=20$, which shows the evolution to a Gaussian distribution for large samples. For so-called photon-counting systems, where individual photoelectron events can be observed, such as in a photomultiplier, this random distribution may be observed directly in the presence of a constant optical power. In the more usual case, where the photoelectron production rate is much greater than the system

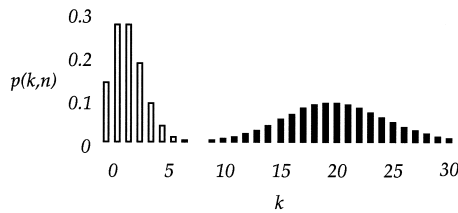


FIGURE 3 The Poisson distribution for $n=2$ (open bars) and $n=20$ (solid bars).

bandwidth, producing a large number of events per measurement interval, the fluctuation of the output current is Gaussian about the mean and has a mean square value and probability distribution given by

$$\overline{(i - \bar{i})^2} = \bar{i}_n^2 = 2q\bar{i}B; \quad p(i) = \frac{e^{-(i - \bar{i})^2 / 2\bar{i}_n^2}}{\sqrt{2\pi\bar{i}_n^2}}. \quad (13)$$

Called “shot noise,” this current fluctuation applies to both thermal and laser or coherent radiation. (Lead “shot” for a shotgun cartridge is produced by dropping a stream of molten droplets in a “shot tower.” The droplets solidify before striking the base of the tower, producing a random set of sound impulses due to the random variation in size and thus arrival time.) Although thermal radiation power fluctuates with an exponential probability distribution (Rayleigh power), the rate is so high, at frequencies corresponding to the spectral width, that the detector and its circuitry average this fluctuation to a negligible value. A very special form of radiation, so-called “squeezed state” or “sub-Poisson,” can be produced in nonlinear optical processes such as optical parametric oscillation. The fluctuation in the detector output can be less than that expected from a Poisson process. This special case is not treated here.

Equation (13) is applicable to “signal-noise-limited” detector performance where the noise is the quantum-limited fluctuation of the signal itself. Using Eqs. (1), (11), and (13), the signal-to-noise ratio becomes

$$\left(\frac{S}{N}\right)_V = \frac{i_s}{\sqrt{\bar{i}_n^2}} = \frac{i_s}{\sqrt{2q\bar{i}_s B}} = \frac{\sqrt{\bar{i}_s}}{\sqrt{2qB}} = \frac{\sqrt{\eta q P / h\nu}}{\sqrt{2qB}} = \sqrt{\frac{P}{NEP_{SL}}}; \quad NEP_{SL} = \frac{2 h\nu B}{\eta} \quad (14)$$

and increases as the square root of the signal power because of the increase of noise with signal power. The signal-limited NEP_{SL} is useful as a benchmark, as in Eq. (10) above, but when a detector is really signal or photon-noise limited, pulse-counting techniques may be used to monitor individual photo-events, and $(S/N)_V$ loses its significance.

Photon-noise also applies to “background-noise-limited” detection, where the total optical flux from the surroundings establishes the output noise level, and the signal is now a small fraction of this “background” power. An example is the detection of a small temperature change in a room-temperature object using its thermal radiation in the 8- to 12- μm wavelength region. The small change in power produced by the temperature change produces the signal current, while there is a much larger total detector current in the presence of the thermal background power.

Another example would be a weak pulsed laser signal observed coming from a bright sunlit cloud background. The signal-to-noise expression for the background-limited case is

$$\begin{aligned} \left(\frac{S}{N}\right)_V &= \frac{i_S}{\sqrt{2q i_B B}} = \frac{\eta q P_S / h\nu}{\sqrt{2q(\eta q P_B / h\nu)B}} \\ &= \frac{P_S}{NEP_{BL}}; \quad NEP_{BL} = \sqrt{\frac{2 h\nu B P_B}{\eta}} \end{aligned} \quad (15)$$

with P_B the background power. Since the effective integration time, τ_B , for an electrical network with bandwidth, B , is $\tau_B = 1/2 B$, the NEP_{BL} may also be written

$$\begin{aligned} NEP_{BL} &= \sqrt{\frac{2 h\nu B P_B}{\eta}} = \sqrt{\frac{h\nu P_B}{\eta \tau_B}} \\ &= \sqrt{\frac{(\eta P_B / h\nu) T (h\nu / \eta)^2}{\tau_B^2}} = \frac{h\nu}{\eta \tau_B} \sqrt{N_B} \\ \therefore \eta(NEP_{BL}) &= \frac{h\nu}{\tau_B} \sqrt{N_B} \end{aligned} \quad (16a)$$

which states that the rms fluctuation in the measured optical power is the average power per photon times the square root of the number of photoevents.

In the case of thermal radiation in the 8- to 12- μm atmospheric “window,” consider the NEP_{BL} for a 1-mm² detector with unit quantum efficiency operating at a 1 kHz bandwidth. By cooling the detector, internal thermal electron generation is made negligible, and the background power from the object establishes the noise. The radiance of a unit emissivity object at 300 K is 460 W/m² for the full spectrum and becomes 130 W/m² after passing through an 8- to 12- μm filter. The irradiance at the detector is then the object radiance reduced by the factor $(2f/\#)^2$, where $f/\#$ is the f -number of the optical system. For the most sensitive or “fastest” system, $f/\# = 0.5$, and the irradiance at the image is maximum and equal to the radiance of the object. The noise equivalent power then becomes, using an average wavelength of 10 μm ,

$$\begin{aligned} NEP_{BL} &= \sqrt{\frac{2 h\nu B P_B}{\eta}} \\ &= \sqrt{\frac{2(1.24/10)(1.6 \times 10^{-19})10^3[(130)(10^{-6})]}{1}} \\ &= 3.6 \times 10^{-11} \text{ W} \end{aligned} \quad (16b)$$

which is within an order of magnitude of the result for the thermal detector, Eq. (9). This is not too surprising, since the temperature fluctuations of the thermal detector are determined by the independent random fluctuations in the emitted and absorbed radiation. It is thus the fluctuation

of the blackbody radiation which limits the sensitivity in either case. Increasing the $f/\#$ of the optical receiver lowers P_B quadratically and the NEP linearly, but the signal power striking the detector drops quadratically, resulting in a lower value of $(S/N)_V$. Since background noise is additive, that is, independent of the true signal, the signal-to-noise ratio is proportional to signal power and satisfies Eq. (1). Unless the signal has a broad spectrum, as in thermal imaging, a background-limited system should have the narrowest possible spectral filter for maximum sensitivity. If the signal is at the longer infrared wavelengths, it may be necessary to cool the filter to prevent its thermal emission from reaching the detector.

B. Postdetection Electrical Noise

1. Analog or Continuous Signal Mode

With the exception of photon counters in the visible and near infrared and background-limited detectors in the LWIR, additive electric circuit noise usually dominates the noise output from a detection system. The noise arises in the load resistance, R , of Fig. 2(b) as well as in the elements of the following amplifier. Sometimes called Johnson noise, after its discoverer in the late 1920s, the current fluctuation arises from the random thermal motion of electrons in the resistors and transistors in the circuit. The resistor noise can be represented by a parallel current source whose mean square value and Gaussian probability distribution are

$$\overline{i_n^2} = 4kTB/R; \quad p(i) = \frac{e^{-i^2/2\overline{i_n^2}}}{\sqrt{2\pi\overline{i_n^2}}}. \quad (17)$$

The capacitance in Fig. 2(b) is a critical parameter, since for a given bandwidth, B , the higher the capacitance, the lower the value of R , and thus the higher the noise current. Consider a detector operating at a bandwidth of 1 MHz, with a typical combined device and circuit shunt capacitance of 1 pf. The required load resistance, $R = 1/2 \cdot BC = 157,000$ ohms and the rms noise current is

$$\begin{aligned} \sqrt{\overline{i_n^2}} &= \sqrt{\frac{4kTB}{R}} = \sqrt{\frac{4(0.026)(1.6 \times 10^{-19})10^6}{1.57 \times 10^5}} \\ &= 3.3 \times 10^{-10} \text{ A} \end{aligned} \quad (18)$$

Setting the signal current, $i = \Re(NEP)$, equal to the noise current yields

$$\begin{aligned} NEP &= \frac{\sqrt{\overline{i_n^2}}}{\Re} = \frac{\sqrt{\overline{i_n^2}}}{\eta q / h\nu} = \frac{3.3 \times 10^{-10}}{\eta \lambda (\mu\text{m}) / 1.24} \\ &= \frac{4.1 \times 10^{-10}}{\eta \lambda (\mu\text{m})} \text{ W} \end{aligned} \quad (19)$$

using Eq. (11). This may be compared with the signal-noise-limited value,

$$\begin{aligned} NE_{PSL} &= \frac{2h\nu B}{\eta} = \frac{2(1.24/\lambda(\mu\text{m}))(1.6 \times 10^{-19})10^6}{\eta} \\ &= \frac{4 \times 10^{-13}}{\eta\lambda(\mu\text{m})} \text{ W} \end{aligned} \quad (20)$$

The decreased sensitivity in Eq. (19) applies for all electronic detectors, unless they have internal electron multiplication, such as in the photomultiplier. As shown in the graph of Fig. 10, this degradation becomes less pronounced at the higher bandwidths because of a smaller increase of the circuit noise with bandwidth.

Fortunately, the bandwidth-dictated load resistance, R , can be supplied by the amplifier and the net noise decreased substantially because of the gain in the amplifier first stage. One of the preferred amplifier types is the transimpedance amplifier, which presents an appropriate input resistance using negative feedback, but has an equivalent input noise current smaller than that of a real resistor. Commercial low noise amplifiers for optical detection are usually characterized by an effective input noise current, sometimes called the “spot” noise current, $(i_n)_{\sqrt{f}}$, usually quoted in $\text{pA}/\sqrt{\text{Hz}}$, with $1 \text{ pA} = 10^{-12} \text{ A}$. Since the amplifier now includes the load resistance, R , of Fig. 2(b), its contribution is included in the effective input noise current. The amplifier limited noise equivalent power, NE_{AL} , is then

$$NE_{AL} = (i_n)_{\sqrt{f}} \sqrt{B}/\Re \quad (21)$$

Typical values of $(i_n)_{\sqrt{f}}$ are in the range of 0.1 to 10 $\text{pA}/\sqrt{\text{Hz}}$, increasing with the full available bandwidth of the amplifier and with the detector capacitance. For the ultimate in sensitivity, the detector and its amplifier should be integrated on a single chip or in a custom package, to minimize the capacitance of the coupling circuit.

2. Sampled or Integrated Signal Mode

Sampled data systems, such as imaging CCDs, use an integrating amplifier rather than a broadband analog amplifier. As shown in Fig. 4, there is no load resistor for the detector. Instead, switch S_1 is closed for an integration time, T , and at the end of this sampling time, S_1 is opened

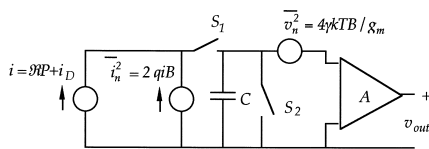


FIGURE 4 Equivalent circuit for an integrating or sampling mode amplifier.

and the change in the output voltage measured. Closing S_2 then resets the capacitor voltage to zero, and after another period, T , the next sample charge is measured. Added to the photon-noise fluctuation of the detector current are two sources of electrical noise. First, although resistor noise has been eliminated, there is still a thermal fluctuation in the reset “zero” voltage of the capacitor. Specifically, when S_2 opens, thermal energy in the short circuit loop is converted to potential energy, and the reset voltage has a sample-to-sample mean square fluctuation given by $\overline{v_n^2} = kT/C$. Second, the amplifier has an equivalent input noise voltage, shown here as that of a field effect transistor (FET), the most common input device for a low-noise amplifier.

Rather than the noise equivalent power as a figure of merit, the common descriptor for a sampling detector is the noise equivalent electron count, NEE . The output voltage is a measure of the apparent charge transferred to the capacitor, and the three sources of the charge fluctuation may each be categorized as a fluctuation in the electron count, N . The photon noise term is simply $(NEE)_{\text{shot}} = \sqrt{N} = \sqrt{(\eta P/h\nu)T}$, from the Poisson statistics of the photoexcitation process. The NEE count associated with the capacitor thermal noise is $(NEE)_C = C\sqrt{\overline{v_n^2}}/q = \sqrt{kTC}/q$, and that of the amplifier FET noise, $(NEE)_{FET} = (C/q)\sqrt{4\gamma kTB/g_m}$, where γ depends upon the particular transistor and is of order unity and g_m is the FET transconductance. Ideally, the values of $(NEE)_C$ and $(NEE)_{FET}$ would be comparable to or smaller than \sqrt{N} , yielding photon-counting performance. For a capacitance of 1 pf at room temperature, the $(NEE)_C$ is 400, while the contribution from the amplifier, $(NEE)_{FET}$, for $\gamma = 1$, $g_m = 10 \text{ mS}$, and $B = 1 \text{ MHz}$, is only 8. The thermal noise of the capacitor is thus limiting, but fortunately a technique known as “correlated double sampling” may be used to eliminate this term.

Specifically, the output voltage from the amplifier is sampled immediately after switch, S_2 , is opened and then again at the end of the sampling time when S_1 is opened and before S_2 closes again. The net voltage change is then independent of the thermal fluctuations of the initial capacitor voltage. Again, in this type of circuit, it is essential to combine the amplifier input transistor on the same chip as the detector or detector array.

C. Temporal Response

In addition to the obvious limitations of the detector capacitance to speed of response, both vacuum and semiconductor devices have specific constraints associated with the dynamics of the photoexcited electron (or hole). In both the vacuum and the semiconductor photodiode, the

frequency response is roughly the inverse of the transit time across either the vacuum or the semiconductor space-charge region. Fortunately, this transit time can be well less than a nanosecond so that both types of diodes can be operated at frequencies as high as 10 GHz. In the case of the semiconductor diode, the requisite thin space-charge region results in high capacitance with its consequent degradation of noise performance. In contrast, the photoconductor response time is determined by the recombination or decay time of the photoexcited carriers in the uniform photoconductor medium. As such, it is a property of the particular material rather than a designable characteristic of the device structure.

In all of these cases, it should be emphasized that the photoexcitation process itself is instantaneous. Once the electromagnetic field of the incident radiation arrives at the photocathode or penetrates the semiconductor, electrons are immediately excited at an average rate proportional to energy density but with a random distribution in time. Thus, even a picosecond (10^{-12} seconds) pulse of light of energy, $Nh\nu$, would produce ηN electrons within the picosecond, but the observed number would have an rms fluctuation of $\sqrt{\eta N}$. This, of course, assumes that the active photoexcitation region is smaller in extent than the spatial extent of the optical pulse, in this case, of order $300\text{ }\mu\text{m}$.

III. VACUUM PHOTODETECTORS

The external photoelectric effect utilized in the vacuum photodiode and photomultiplier is generally effective in the visible part of the spectrum. Quantum efficiencies as high as 25% occur in the blue, but fall to the order of 1% or less in the near infrared at a wavelength of about $1\text{ }\mu\text{m}$.

A. Photodiode

The simple vacuum photodiode is mainly of historical interest since it operates at reasonably high voltages, 100 V or more, is fragile and bulky, and has inferior performance to that of a semiconductor photodiode, especially because of its limited wavelength response. One possible advantage is the large photosensitive area available without commensurate large capacitance. The relatively high voltage operation offers fast response time with electrode spacings of many millimeters and thus low capacitance. It is helpful to calculate the frequency response since it is applicable to the more complicated but much more usable photomultiplier tube, which offers electron multiplication that effectively overcomes external electric circuit noise. In the simple representation of a photodiode in Fig. 5, an emitted photoelectron travels from the photocathode to

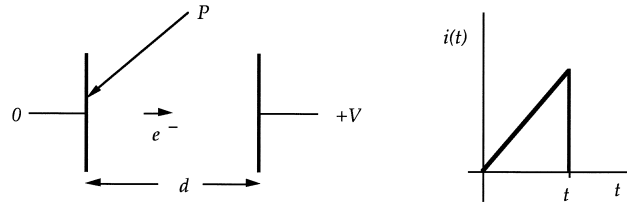


FIGURE 5 Vacuum photodiode and impulse response.

the anode in a time, $\tau = d\sqrt{2m/qV}$, and the current pulse as shown is triangular since the current flow in the external circuit is proportional to the electron velocity which increases linearly at constant acceleration. This is the impulse response of the detector, and the squared magnitude of the Fourier transform yields the frequency response. In this case, the half-power point in the response occurs at a frequency $f_{co} = 1/2\tau$. For a voltage of 300 V and an electrode spacing of 5 mm, the transit time is 1 nanosecond and f_{co} is then 500 MHz. The capacitance for 1 cm^2 electrodes would be only 0.2 pf.

B. Photomultiplier

The photomultiplier, as shown in Fig. 6, is almost universally used as a “photon counter,” that is, the internal electron multiplication produces an output electrical pulse whose voltage is large compared with the output electric circuit noise. Each pulse in turn is the result of an individual photoexcited electron. The numbered electrodes, 1–8, called dynodes, are each successively biased about 100 V positive with respect to the preceding electrode, and an accelerated electron typically produces about 5 secondary electrons as it impacts the dynode. The final current pulse collected at the output electrode, the anode, would in this case contain $5^8 \approx 400,000$ electrons. The secondary emission multiplication process is random, the value of the dynode multiplication factor is close to Poisson distributed from electron to electron. The output pulse amplitude thus fluctuates. For a secondary emission ratio of $\delta = 5$, the rms fractional pulse height fluctuation is $1/\sqrt{\delta - 1} = 0.5$. Since the mean pulse height can be well above the output circuit noise, the threshold for a pulse count may be

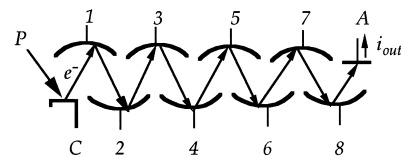


FIGURE 6 Photomultiplier structure. Photocathode, C, dynodes, 1–8, collecting anode, A.

set low enough to detect essentially all photoelectrons. The first dynode is the critical step in the multiplication process. From Eq. (12), the probability of zero secondary emission for $\delta = 5$ is $p(0,5) = 5^0 e^{-5} / 0! = e^{-5} \approx 10^{-2}$. Thus, only about 1% of the photoelectrons is missed, which is inconsequential where the quantum efficiency or photon to electron conversion is at most 25%. If the photomultiplier is used in the analog mode, that is, at count rates much higher than the bandwidth, then the mean square shot noise at the output is increased by a noise factor, $\Gamma = \delta / (\delta - 1)$ or 1.25 for our example. The expression for the output noise current is $i_n^2 = 2q\Gamma G^2 iB$ with $i = \eta q P / h\nu$, the photocathode current.

The impulse response and consequent operating bandwidth arise from two phenomena, the temporal spreading of the electron packet reaching the anode and the shape of the induced current pulse as the charge flows from the last dynode to the anode. The latter effect is the same as that in the simple diode and, by careful design, the transit time spread can be minimized, resulting in bandwidths as high as many gigahertz. Available photomultiplier tubes usually use a circular electrode arrangement rather than the schematic linear description of Fig. 6. This results in a more compact package and a consequent reduction in transit time spread.

Dark current can be a problem for multiplier tubes operating in the near infrared, since the lower work function or barrier for the photoemission process also allows more thermally excited electrons to escape. Cooling of the photocathode can eliminate this shortcoming, although semiconductor devices discussed below are usually preferable at near infrared wavelengths. For the detection of very short pulses, such as in an optical radar system, dark current is not as critical a parameter, and hybrid cathode designs such as the “transferred electron” photocathode can yield quantum efficiencies of the order of 20% out to wavelengths beyond 1 μm . These semiconductor-based cathodes are two-element structures and require a bias voltage of the order of a volt to excite electrons into a higher energy “valley” where they require less photon energy for photoemission. The process results in significantly higher dark currents and generally precludes these detectors for photon counting applications.

In addition to multielectrode photomultiplier tubes, electron multiplication can be obtained using a reverse-biased silicon diode as the anode. At photocathode-to-anode voltages of several thousand volts, the impinging photoelectron can produce about 1000 hole-electron pairs. This gain of 1000 will again have fluctuations of the order of the square root of the gain, about 30 or 3%. This type of multiplication as well as modifications of the photomultiplier electrode structure are used extensively in imaging devices discussed in Section VI.

IV. SEMICONDUCTOR PHOTODETECTORS

Semiconductor devices, the most common photodetectors, are of two main types. The photoconductor is effectively a radiation-controlled resistance, while the photodiode is similar to the vacuum version. The photoconductor is limited to specialized applications, while the photodiode and its internal multiplication form, the avalanche photodiode, have universal application.

A. Photoconductor

The photoconductor, as shown in Fig. 7, depends upon the creation of holes or electrons in a uniform bulk semiconductor material, and the responsivity, temporal response, and wavelength cutoff are unique to the individual semiconductor. An intrinsic photoconductor utilizes “across-the-gap” photoionization or hole-electron pair creation. An extrinsic photoconductor depends upon the ionization of impurities in the material and in this case only one carrier, either hole or electron, is active. The same is true for a quantum-well photoconductor, in which electrons or holes can be photoexcited from a small potential well in the narrower band-gap regions of the semiconductor. The quantum efficiency for the structure in the figure is determined by the absorption coefficient, α , and may be written as $\eta = (1 - R)[1 - e^{-\alpha d}]$, where R is the reflection coefficient at the top surface. Carriers produced by the radiation, P , flow in the electric field and contribute to this current flow for a time, τ_r , the recombination time. The value of the current is

$$i = (qV/L) \frac{\eta P \tau_r}{h\nu} = (\tau_r / \tau_t) \frac{\eta q P}{h\nu} = \frac{\eta q G_p P}{h\nu} = \Re P;$$

$$\Re = \frac{\eta q G_p}{h\nu}, \quad G_p = (\tau_r / \tau_t) \quad (22)$$

where τ_t is the carrier transit time through the device and G_p is the photoconductive gain. In the first term of the equation, the fraction in parentheses is the current produced by a carrier of charge, q , moving at a velocity, v , between electrodes spaced at a distance, L . The second term is the average number of carriers, the rate of production times the lifetime. The transit time from one end of the

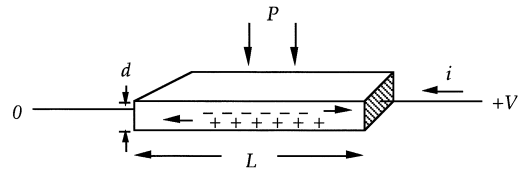


FIGURE 7 Photoconductor. Intrinsic operation shown. In extrinsic mode, mobile electrons are photoexcited from fixed positive donor impurities, or holes from negative acceptors.

structure to the other is $\tau_t = L/v = L/\mu E = L^2/\mu V$, with μ the carrier mobility. The photoconductive gain will be greater than unity if the lifetime exceeds the transit time. In this case, an unrecombined excess carrier leaving the material is immediately replaced by a new carrier at the opposite end, maintaining charge neutrality. A 1-mm-long detector might have a transit time as short as a microsecond, so that for a 1-ms recombination time, the gain could be as high as 1000. The current noise, produced by the signal-induced random (g)eneration of carriers as well as the random (r)ecombination, is called *g-r* noise and is given by

$$\overline{i_n^2} = 4qG_p^2 \frac{\eta q P}{h\nu} B \quad (23)$$

not too different from the noise current of the photomultiplier above. The circuit of Fig. 2(b) is applicable to the photoconductor, but the dark current associated with residual thermally excited carriers can produce enough noise so that the detector is still not signal or photon noise limited, even with gains of 1000 or greater. The dark current noise is given by $\overline{i_n^2} = 4qG_p i_D$. In addition, the dark resistance shunts the external load, R , in Fig. 2(b), reducing the signal voltage.

Photoconductors have two general uses. First, extrinsic, or impurity-doped, materials such as germanium and silicon can be operated at extremely long wavelengths using an appropriate low ionization energy impurity, although they usually require cryogenic cooling. Second, the ease of fabrication (no *p-n* junction required) makes them economically attractive for the middle and long wavelength infrared regions where the ultimate in sensitivity is not required. Intrinsic photoconductor materials include lead sulfide, lead selenide, cadmium sulfide, and mercury cadmium telluride, while germanium and silicon are the usual hosts for extrinsic photoconductors with impurities such as arsenic, copper, gold, and indium.

B. Photodiodes

Semiconductor photodiodes are conveniently divided into three operating categories. First is the reverse-biased, or sometimes confusingly called “photoconductive,” mode. Second is the nonbiased or photovoltaic mode, where the open-circuit voltage of the diode is the signal source. Last is the avalanche mode, where internal gain occurs due to carrier impact ionization in the space-charge region.

1. Reverse-Biased Diode

The most common detector in the most used mode, the reverse-biased diode, Fig. 8, utilizes hole-electron pair

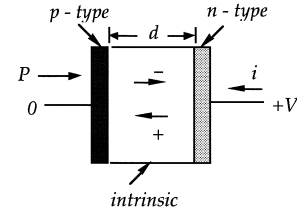


FIGURE 8 A reverse-biased semiconductor photodiode, in this case, a *p-i-n* structure.

creation in the space-charge region, usually fabricated of intrinsic material. Although not feasible in all types of semiconductors, the *p-i-n* configuration allows the thickest space-charge region, yielding, in turn, low capacitance and high quantum efficiency because of the extended photon interaction distance. The optimum diode also has thin and heavily doped *p* and *n* regions to minimize absorption in the outer layers and reduce series resistance. The quantum efficiency, η , is identical in form to that for the photoconductor above, with d now the thickness of the intrinsic layer. As shown in the figure, the preferred incidence of the radiation is on the *p* side of the structure, since this produces the most free carrier pairs at the left of the layer, and thus most of the current is carried by electrons moving to the right, the holes traveling the shorter distance to the left. With rare exceptions, the electron velocity is much greater than that of the hole, and the frequency response is correspondingly higher. The current impulse for an electron produced at the left of the region is a square pulse, since the electron moves at constant velocity in the constant field of the junction. In fact, most high-speed diodes are designed to operate at voltages where the electron velocity has reached a saturation value. Higher voltages can cause breakdown or, as discussed below, avalanche multiplication. The frequency cutoff is again very close to the value $f_{co} = 1/2\tau$, as in the vacuum diode. As an example, a silicon *p-i-n* photodiode might have the following constants and operating performance:

$$\lambda = 0.8 \mu\text{m}; \quad A = 0.1 \times 0.1 \text{ mm};$$

$$\alpha = 10^3 \text{ cm}^{-1}; \quad d = 10 \mu\text{m}; \quad \eta = 0.63.$$

$$v_{\text{electron}} = 3 \times 10^6 \text{ cm/s at 5 V bias}$$

$$f_{co} = 1.5 \text{ GHz}; \quad C = 0.1 \text{ pf.}$$

The many variations of the diode structure of Fig. 8 include the so-called Schottky barrier diode, a simple metal-semiconductor structure, usually gold on gallium arsenide. Here the metal is biased negative with respect to an *n*-type or intrinsic semiconductor with a heavily doped *n*-type contact layer. Gallium arsenide has a 10-fold larger

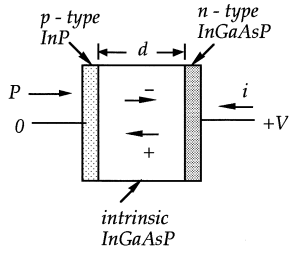


FIGURE 9 A heterostructure photodiode. The *InP* *p*-type layer is completely transparent to the incoming radiation, thus improving the quantum efficiency.

absorption coefficient at $0.8 \mu\text{m}$, for example, and a 3-fold higher saturated electron velocity. For the same quantum efficiency as the silicon diode above, the active layer can be $1 \mu\text{m}$ thick, the f_{co} becomes 50 GHz, but the capacitance increases to 1.0 pf.

A second variation in photodiode design is the heterostructure. Figure 9 shows such a structure, an *InP-InGaAsP* diode. There are two special features of this design. First, the quaternary compound, *InGaAsP*, by varying the composition, can be “tuned” to operate at wavelength cutoffs from 0.9 to $1.7 \mu\text{m}$, while still maintaining a perfect lattice match to the underlying *InP* upon which it is grown epitaxially. Second, the *InP* has a wavelength cutoff at about $1 \mu\text{m}$, so that it is completely transparent to any wavelengths between 1 and $1.7 \mu\text{m}$, completely encompassing the low attenuation bands in optical fibers. In addition to this unique example, there is a large family of binary, ternary, and quaternary semiconductors which may be fabricated as either homojunctions (single material) or heterojunctions. Their wavelength cutoffs extend from the visible to wavelengths beyond the atmospheric window at $12 \mu\text{m}$. In addition to the III–V valence (*Al, Ga, In*)-(As, *P, Sb*) family, there is a II–VI family of ternary compounds, (*Cd, Hg, Pb, Sn*)-(S, *Se, Te*), which is used mainly in the infrared and LWIR region. Mercury cadmium telluride, $\text{Hg}_x\text{Cd}_{(1-x)}\text{Te}$, is the most popular material in this family, tuneable from $0.8 \mu\text{m}$ at $x = 0$ to effectively infinite wavelength, since at $x = 1$, *HgTe* is a semimetal. Figures 8 and 9 are of course schematic. Most devices are formed by diffusion or epitaxial growth and use overlaid metal electrode connections.

Overall, the *NEP* of a reverse-biased photodiode is determined by amplifier noise at visible and near infrared wavelengths. At longer wavelengths, dark current must be reduced by cooling and in many applications, the detector becomes background limited. Particularly in the infrared, the diode is best operated in the photovoltaic rather than reverse-biased mode, since tunneling and avalanche breakdown occur at ever decreasing voltages as the energy gap decreases.

2. Photovoltaic Mode

To avoid the noise produced by excess leakage currents, many diodes in the infrared region are operated as open-circuited voltage sources. Light shining on the active junction region produces a current flow in the *reverse* direction, and under open-circuit conditions, the charge buildup creates a *forward* voltage which produces an equal and opposite current flow, resulting in a net of zero. (This is the principle of the solar cell.) Using the standard semiconductor diode voltage-current equation, the open circuit diode voltage becomes

$$v = \frac{kT}{q} \ln \left(\frac{i + I_s}{I_s} \right) \xrightarrow{i \ll I_s} \frac{kT}{q I_s} i = R_J i \quad (24)$$

where I_s is the saturation current of the diode and i is the photoinduced current, $\Re P$. The junction resistance, $R_J = kT/qI_s$, is the zero-bias resistance of a *p-n* junction and the parallel noise current source is given by the Johnson noise expression, $\overline{i_n^2} = 4kTB/R_J$, since the junction is in thermal equilibrium. This noise may be shown to be the same as the uncorrelated shot noise of two equal and opposite currents, I_s , flowing across the junction. The total mean square noise current then becomes this junction Johnson noise term plus the effective mean square input noise current of the amplifier, as discussed in Section II.B above. Diodes which normally operate in the reverse-biased mode, such as those of silicon and gallium arsenide, are characterized by their saturation or dark current. Diodes which have high reverse currents due to tunneling or avalanche are characterized by their zero-bias resistance, and a standard figure of merit is the *RA* product, the resistance multiplied by the diode area. Assuming surface leakage is negligible, the resistance should be inversely proportional to area, and for a given material, the *RA* product as a function of temperature is a convenient characterization. This behavior allows the use of Eq. (2) and

$$\begin{aligned} D^* &= D\sqrt{AB} = \frac{\sqrt{AB}}{NEP} = \frac{\Re\sqrt{AB}}{\overline{i_n^2}} = \frac{\Re\sqrt{AB}}{\sqrt{4kTB/R}} \\ &= \Re\sqrt{\frac{RA}{4kT}}; \quad A \text{ in cm}^2, \quad B \text{ in Hz} \end{aligned} \quad (25)$$

using Eq. (19) for the value of *NEP*. This of course assumes that the dominant noise source is the diode rather than the following amplifier. Actually, cooling a photovoltaic diode can increase the junction resistance to the extent that the operation is identical to that of a reverse-biased diode. Conversely, at the longer infrared wavelengths, the ultimate performance is determined by the thermal background, except when observing a narrow spectral source such as laser radiation.

3. Avalanche Photodiode

Akin to the vacuum photomultiplier, semiconductor diodes can utilize internal electron (and hole) multiplication to obtain current gain prior to external amplification. In Fig. 8, if an electron produced at the left side of the junction gains enough energy from the electric field to produce an additional hole-electron pair, this process can be repeated as the electrons and their progeny move across the space-charge region. This “avalanche” process is complicated by the production of extra hole-electron pairs by the leftward moving holes. In fact, avalanche “breakdown” occurs when the net electron and hole gains across the space-charge region add to unity, since the avalanche becomes self-sustaining. In the ideal multiplication process, the ionization probability for holes is zero and the electron current grows exponentially in space across the junction. The gain becomes $M = e^{\alpha d}$ where α is the ionization coefficient for the electrons, a sharply increasing function of electric field in the junction above a threshold of about 10^5 V/cm. When α for the holes is zero, the mean-square noise current for high gain is twice the simple shot-noise value, because of the randomness in the multiplication process, and is given by $\overline{i_n^2} = 4qM^2iB = 4qM^2\mathfrak{N}PB$. Silicon avalanche photodiodes have very small ionization coefficients for holes and this equation is valid up to gains of about 200, where excess noise due to hole multiplication becomes appreciable. This is about the maximum usable gain, in any event, since even slight variations in material parameters over the junction area can produce breakdown in selected regions because of the rapid variation of the ionization coefficient with field. Silicon APDs can thus decrease the NEP_{AL} by about a factor of 100, the gain of 200 mitigated by the net increased noise. Unfortunately, most of the other diode materials of interest, such as the III–V and II–VI semiconductors, generally have roughly equal hole and electron ionization coefficients. For exact equality, the noise current expression becomes $\overline{i_n^2} = 2qM^3iB = 2qM^3\mathfrak{N}PB$ and with a typical low-noise amplifier, the NEP_{AL} may be reduced by about a factor of 10, at a gain of about $M = 20$. At higher gains, the net $(S/N)_V$ then decreases again because of the increasing multiplication noise. In special heterostructures, with separate absorption and multiplication regions, SAM devices, the light may be absorbed in one type of material and multiplication occur in an adjacent layer. This allows the absorption to be optimized for the wavelength of interest while using the lowest noise avalanche material.

Avalanche photodiodes may also be used as photon counters in the so-called Geiger mode, named after the Geiger–Müller gas-filled tube, used in the detection of high-energy particles. As in the gas tube, the diode is biased at or above the normal avalanche breakdown voltage,

and a breakdown (or discharge) is triggered by a photoionization in the active region. Application of this mode is limited since the breakdown is followed by a “dead” time when no photoevents may be observed. This dead time lasts until the active region is swept clear of carriers, and standby voltage is restored to the electrodes after the initial current surge. The detector thus can not only miss closely spaced events but also can be triggered by “dark” current events. As an analog signal detector, it would have severe saturation at a photoevent rate approximately the inverse of the dead time.

V. GRAPH OF DETECTOR CHARACTERISTICS

Most of the above material is summarized in the graph shown in Fig. 10. The quantum efficiency, η , is taken as unity. The NEP for a particular detector should be increased by $1/\eta$ from the graphical value. Typical semiconductor quantum efficiencies are of the order of 0.5, vacuum photocathodes from 0.25 to as low as 0.001.

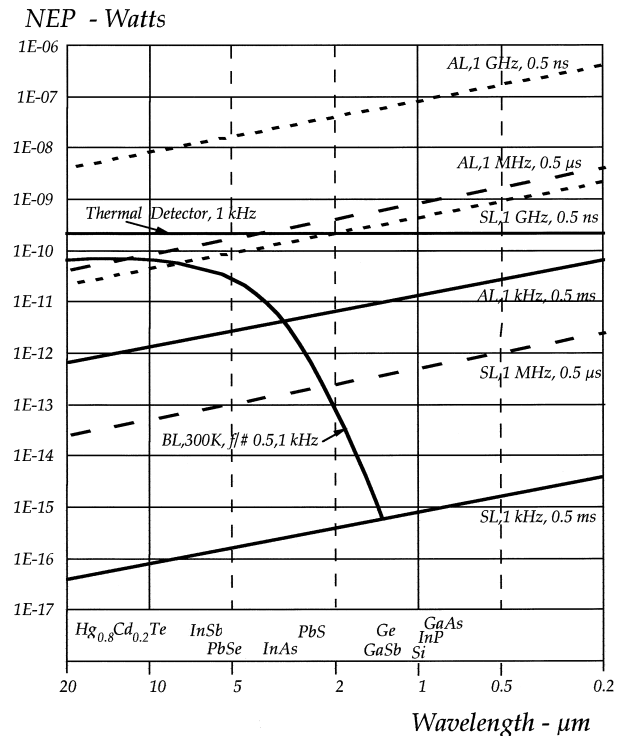


FIGURE 10 Graph of detector performance indicating mode of operation, bandwidth, and pulse response. SL, signal or photon noise limited. BL, background limited. AL, amplifier noise limited. Various semiconductors are centered at their approximate wavelength cutoff.

A. Thermal Detection

The thermal detector wavelength-independent behavior is the theoretical limit for a 1-kHz bandwidth with an area of 1 mm² using Eq. (7), and as such is probably one or two orders of magnitude better than actual available performance. In particular, the associated millisecond response time would require extra heat conductance beyond the radiation term with consequent increase in the *NEP* from Eq. (6). The electric sensing circuitry adds additional noise. Cooling of the detector can offer significant improvement, as discussed above. In the LWIR region, particularly in imaging arrays, thermal detector performance is close enough to that of semiconductor electronic detectors to offer an attractive alternative.

B. Electronic Detection

1. Signal or Photon Noise Limited (SL)

The theoretical detection limit is plotted for three bandwidths, 1 kHz, 1 MHz, and 1 GHz. Also shown are the corresponding pulsewidths, 0.5 ms, 0.5 μ s, and 0.5 ns. This ultimate performance is essentially only obtainable with the vacuum photomultiplier and then only within a factor of 4 (25% quantum efficiency) at wavelengths of 0.5 μ m and shorter. Avalanche photodiodes can also approach the theoretical limit at high bandwidths, as discussed in Section V.B.3. Although the avalanche diode in the Geiger mode can count individual photoevents, it is not operable in a continuous or analog signal mode.

2. Background Noise Limited (BL)

Background-limited operation is usually applicable for broad spectral band signals in the infrared, in particular thermal sensing. The curve shown for a 1-kHz bandwidth, 1 mm² detector area, is calculated for a signal at the specified wavelength, λ , with 300 K background radiation at all wavelengths shorter than λ . The optical system has an $f/\# = 0.5$, or a numerical aperture of unity. The *NEP* is proportional to $1/(f/\#)$, but the signal strength at the detector falls off as $1/(f/\#)^2$, so that the lowest $f/\#$ yields the best overall performance. It should also be noted from Eq. (15) that the *NEP* varies as $1/\sqrt{\eta}$ in the special case of background-limited operation. The *NEP* also decreases as square root of the area, as in the thermal detector case. For narrow spectral band signals such as laser radiation, an appropriate optical filter should be used to reduce the background so that the system becomes signal or amplifier noise limited. Specific calculations are necessary to establish the ultimate sensitivity in this case.

3. Amplifier Noise Limited (AL)

The three curves for amplifier-limited operation have been calculated assuming typical amplifier normalized noise currents of 0.3, 1, and 3 pA/ $\sqrt{\text{Hz}}$ for the respective bandwidths 1 kHz, 1 MHz, and 1 GHz. The applicable value depends upon the specifications of the chosen amplifier. Generally, the best performance is obtained by integrating the amplifier and detector on the chip. Avalanche multiplication can supply pre-amplification current gain and reduce the expected *NEP* by a factor from 10 to 100. At 1 GHz bandwidth, this can bring performance to within an order of magnitude of photon noise-limited behavior.

VI. ENHANCED DETECTION TECHNIQUES

Two techniques may be used to enhance the performance of the above detectors, yielding improved sensitivity but restricted to special applications. These are heterodyne detection and optical amplification. Both techniques approach photon-noise-limited performances but only when receiving a single diffraction-limited mode of the receiver aperture.

A. Heterodyne Detection

Akin to the standard detection process in radio, television, and radar systems, optical heterodyne detection uses a “local” single-frequency oscillator and results in an electrical signal output at an intermediate frequency (i.f.) which is an amplitude and phase replica of the original signal. This process is also known as coherent detection. As shown in Fig. 11, a beamsplitter is used to allow simultaneous illumination of the detector by the incoming signal and the local oscillator. Using, for example, a 90% transmission beamsplitter, 10% of the local oscillator power strikes the detector, 90% of the signal. A limiting attribute of this system is the angular sensitivity of the receiver which is the diffraction-limited field of view associated

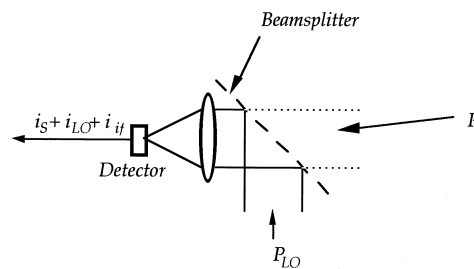


FIGURE 11 Heterodyne detection. The incoming signal power, P , mixes with the single-frequency local oscillator power, P_{LO} . The receiver beamwidth is the diffraction-limited angle determined by aperture diameter.

with the local oscillator distribution. In the figure, if the local oscillator power is circular and uniform over the lens, then the far-field receiver pattern is the classic Airy pattern, with half-power full-width angle of approximately, λ/D , with D the beam diameter. This diffraction-limited or single-spatial-mode behavior produces extremely narrow receiver beams for typical apertures, 10 microradians for a 10-cm aperture at a 1- μm wavelength, for example. As compensation for this shortcoming, the sensitivity of a heterodyne system is close to the photon-noise limit. The three components of detector current are the signal, local oscillator, and i.f. terms given by

$$\begin{aligned} i_S &= \Re P; \quad i_{LO} = \Re P_{LO}; \\ i_{if} &= 2\sqrt{i_{LO} i_S} \cos(2\pi f_{if} t + \phi), \end{aligned} \quad (26)$$

with the i.f. $f_{if} = \nu_{LO} - \nu_S$. With sufficient local oscillator power, the shot noise produced by the local oscillator current dominates any amplifier noise and the (S/N) ratio becomes

$$\begin{aligned} \left(\frac{S}{N}\right)_P &= \frac{\overline{i_{if}^2}}{\overline{i_n^2}} = \frac{(2\sqrt{i_{LO} i_S})^2 \cos^2 2\pi f_{if} t}{2q i_{LO} B} \\ &= \frac{2i_{LO} i_S}{2q i_{LO} B} = \frac{i_S}{qB} = \frac{\eta q P}{h\nu q B} = \frac{\eta P}{h\nu B} \end{aligned} \quad (27)$$

and the NEP is then $h\nu B/\eta$. This is one-half the value for direct detection as given in Eq. (14), completely consistent since the i.f. bandwidth of a signal is twice its baseband or envelope-detected value. A pulse of width, T , has a direct-detected bandwidth of $1/2T$, but an i.f. bandwidth of $1/T$. Note also that the signal-to-noise ratio is defined in terms of the output electrical powers rather than voltages. This maintains the proportionality of (S/N) to optical signal power. In the coherent or heterodyne case, the electrical and optical amplitudes as well as powers are directly proportional. This ideal performance requires a frequency and amplitude stable optical local oscillator, obviously a laser. A balanced mixer circuit using two detectors may be used to reduce the effects of local oscillator power fluctuations. Assuming the availability of a suitable stable laser source and the feasibility of operation with extremely narrow angular beamwidths, heterodyne offers the ultimate in sensitivity as well as the ability to measure phase and frequency. As such, it is most applicable in Doppler radar and high-resolution spectroscopy.

Although a heterodyne system is photon-noise limited, as indicated by Eq. (27), the noise produced by the local oscillator is additive. Thus, unlike the photon counter of Section II.A, there is always a noise floor equivalent to an uncertainty of one photon per sampling time, $T = 1/B$, even with zero signal. In a sense, a heterodyne or coherent system overmeasures the incoming signal by extracting

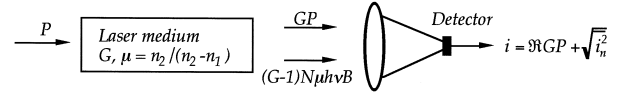


FIGURE 12 Schematic representation of an optical amplifier.

both amplitude and phase simultaneously. In the perfect photon counter, measurement of only the energy yields an uncertainty of \sqrt{N} in a sampling time $T = 1/2B$, while a perfect coherent detector has a measurement uncertainty of $\sqrt{2N + 1}$, where N is the expected number of photons in the measurement interval. This attribute makes heterodyne detection unattractive for radiometry at (S/N) near unity or lower.

B. Optical Amplification

As the a in *Laser* indicates, amplification rather than oscillation is perfectly feasible in a laser system. Rather than amplifying the photoexcited electron stream, why not amplify the signal light wave? Figure 12 illustrates a schematic optical amplifier system and indicates the complications of such a simple concept. First, the laser gain medium not only amplifies the optical wave with gain, G , but adds spontaneous emission noise, $(G-1)N\mu h\nu B$. N is the number of spatial or diffraction-limited modes, counting both polarizations, intercepted by the detector, and μ is the inversion factor of the laser, which has a minimum of unity for complete inversion, that is, an empty ground state. At the output of the laser, the effective ratio of optical signal power to optical noise power is simply $P/h\nu B$ for a diffraction-limited field of view at high gain and in the ideal case of complete inversion with $\mu = 1$. Unfortunately, the optical power has to be converted to an electrical output, and there are two significant noise contributions. First is the expected zero signal contribution, mainly due to the fluctuation of the spontaneous emission optical noise power. Second is a mixing or “beating” between the optical signal and the optical noise, similar to the mixing process in a heterodyne system. The result is a rather complicated signal-to-noise expression,

$$\left(\frac{S}{N}\right)_v = \frac{P_S}{2\mu h\nu B \sqrt{\frac{N\Delta\nu}{2B} + \frac{P_S}{\mu h\nu B}}}; \quad B \ll \Delta\nu. \quad (28)$$

Here $\Delta\nu$ is the lesser of the laser amplifier bandwidths or that of a narrow spectral band optical filter between the laser and the detector. The first term under the square root sign is associated with the zero signal additive noise and the second arises from the beating between the signal and noise fields. Generally, the narrow field of view required to obtain a low N , and bandwidths, B , much less than the spectral bandwidth, $\Delta\nu$, make optical amplification

unattractive because of the large value of the first term under the square root sign. In addition, the inversion factor, μ , is typically much greater than unity for semiconductor lasers, although it can approach unity for optically pumped four-level laser systems. Serendipitously, fiber-optic communication systems utilize single-mode transmission, $N = 1$, and have data bandwidths of the order of 10 GHz, which is comparable to available optical fiber spectral filters at 50 GHz bandwidth. An optically pumped erbium-doped glass fiber laser, operating near $\lambda = 1.5 \mu\text{m}$, has almost complete inversion giving $\mu = 1$. Bit error rates of 10^{-6} and smaller require a signal-to-noise voltage of the order of 10, and as a result, the second term under the square root sign in Eq. (28) must be of order 100, while the first term is the order of unity. The result in this limit is

$$\left(\frac{S}{N}\right)_v = \frac{1}{2} \sqrt{\frac{P_s}{h\nu B}} \quad (29)$$

which is comparable to the photon-noise-limited result of Eq. (14). In this very special application, optical amplification performance approaches ideal detection. In this and any photon-noise-limited application, it is essential to note that the output noise is a function of the signal input. This not only means that the signal-to-noise eventually increases as the square root of the signal power, but that in the absence of signal, the noise is quite small. Thus a digital communications receiver has negligible noise output in the presence of a “zero.” The threshold setting for the digital comparator circuit may thus be lowered commensurately.

VII. OPTICAL IMAGING DETECTORS

A. Vacuum Tube

Historically the first practical optical imaging device, the vacuum envelope image tube has been almost completely replaced by the silicon CCD array. The principal photoemissive form, the image orthicon, has a thin large-area photocathode, and light impinging from the outside of the layer produces an emitted photoelectron pattern which is accelerated to a plane semi-insulating sheet “retina.” Secondary emission then produces a positive charge pattern on the thin sheet, which is “read-out” by a raster scanned electron beam. Electrons in the beam are deposited on the retina until the positive charge is neutralized, and then the leftover electrons are collected and the resultant video current used to produce an image. Typical television rates are 30 frames per second, and the electrical bandwidth is several megahertz for the standard TV format. Orthicon tubes, depending upon external photoemission, suffer from the rapid fall-off in quantum efficiency from 25% in the green

to as low as 1% in the red. The vidicon also uses an electron beam raster scan, but in this case the retina is a photoconductor layer. The back surface of the layer is charged uniformly by a scanned electron beam. During the ensuing frame time, an optical image on the front surface causes charge to leak through the material to a front conducting base layer. When the beam returns to an individual pixel area, the current required to recharge the area flows through the conducting base layer and becomes the video signal. Vidicons, since they utilize the internal photoeffect, can operate throughout the spectrum even to the LWIR region, if the retina is cryogenically cooled.

Both types of tubes have the advantage of frame-to-frame integration; the charge on the retina builds up (or decays in the vidicon) over the full frame time. With special low-noise readout techniques and/or impact ionization gain in the orthicon retina, photon-counting performance can be approached, subject to the available quantum efficiency. Electron gain may also be obtained using two-dimensional electron multiplier “channel” plates, an array of bonded hollow tubes with secondary emission gain along the length of the tube. The disadvantages of vacuum image tubes are bulk and fragility, as well as poor metric capability. The image shape, dimension, and position are determined by the electron beam scanning precision and are thus sensitive to circuit drive voltages as well as stray ambient magnetic fields. In contrast, a solid-state detector array has fixed dimension and pixel location, allowing nondistorting image conversion and precise metric location.

B. Solid-State Arrays and the CCD

As shown in Fig. 13, there are several ways to obtain an optical image using a semiconductor detector: scanning a single element, scanning a line array, or “staring” with a two-dimensional array. Scanning is usually accomplished by mechanical motion of the whole optical system or by a moving mirror in front of a fixed camera. Operation and readout of the single-element scanner is straightforward. Depending upon the number of elements, the linear

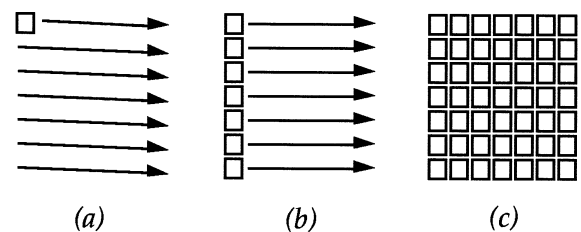


FIGURE 13 Two-dimensional imaging. (a) Single-element scan. (b) Line array scan. (c) Staring array.

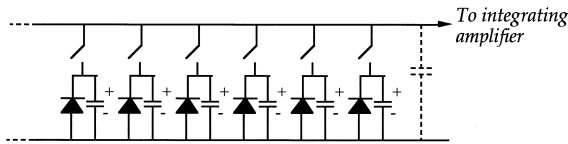


FIGURE 14 Linear array readout. Sequential sampling is usually carried out on chip.

array can be hard-wired to a set of amplifiers and the data parallel processed. Alternatively, especially in long arrays, sequential readout of successive elements can be realized using a circuit such as that shown in Fig. 14, where the switching is performed by an integrated transistor array. This type of readout assumes integration and storage of the optically induced charge over the full readout cycle and thus limits the scanning rate. Similar circuit techniques can be used for a two-dimensional array, but in either case, the *NEE* is strongly dependent upon the combined capacitance of the individual diode storage elements and the line to the amplifier. For example, if the net value of capacitance is 1 pf, and the amplifier FET transconductance is 5 mS with $\gamma = 1$, then at room temperature, using double correlated sampling, the limiting value of *NEE* is 25. Without this special sampling technique, *kTC* noise results in an *NEE* of 400. For operation in the LWIR region such as in thermal imaging, background-limited detection governs and the *NEE* then becomes the square root of the total electron count during a frame time. As a result, the capacitance should be sufficient to store the expected background-induced charge. The pixel-to-pixel uniformity must be better than the noise or external image processing may be required to obtain satisfactory images. LWIR thermal imagers using room temperature-integrated bolometric (thermal) detector/transistor readout arrays also offer competitive performance. Systems using either type of array are characterized by their noise equivalent temperature change, $NE\Delta T$, that temperature change which produces a (S/N) of unity.

The ultimate in sensitivity, especially in the visible and near infrared, is obtained with the CCD. Figure 15 shows

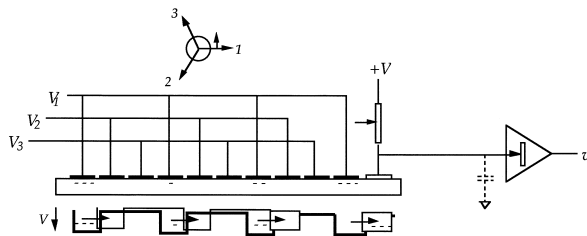


FIGURE 15 Charge-coupled device (CCD) operation. Photoexcited packets of electrons are three-phase clocked to the integrating amplifier at the right.

an array of electrodes on a semiconductor substrate. The electrode structure could be metal-oxide-semiconductor (MOS) on silicon or a metal Schottky barrier on a material such as gallium arsenide. A three-phase voltage is applied to the array such that a “wave” of positive voltage moves to the right on the electrode structure. Electrons stored under the electrodes move with the wave as indicated by the motion of the potential wells in the voltage diagram at the bottom of the figure. For imaging, the CCD array, with fixed voltage, is exposed to light for a frame time, and photo-excited electrons collect under the nearest positive electrode. At the end of the frame time, each charge packet is then “clocked” toward the sensing circuit on the right, where the charge is transferred to the input of an FET integrating amplifier. The clocking and readout technique is shown in Fig. 16. After exposure to the optical image, each vertical column is transferred horizontally to the readout array. The readout array then transfers the charge vertically down to the output sensor. The next column then moves into the readout array and the process is repeated until the whole frame has been sensed. The only capacitance in the sensing circuit is that of the output electrode and associated FET. Since these are both integrated on the semiconductor chip, values of 0.1 pf or less are quite feasible. As a result, using double correlated sampling, *NEEs* less than unity have been measured at kilohertz clock rates while at the standard TV rate of 4 MHz, an *NEE* of 5 to 10 is obtainable. The CCD array thus may be operated as a photon counter, but with the observation or sampling time determined by the frame rate.

The CCD depictions of Figs. 15 and 16 are quite rudimentary, and actual devices require either transparent electrodes or a transparent back plane. In addition, tailoring of the impurity doping in the semiconductor

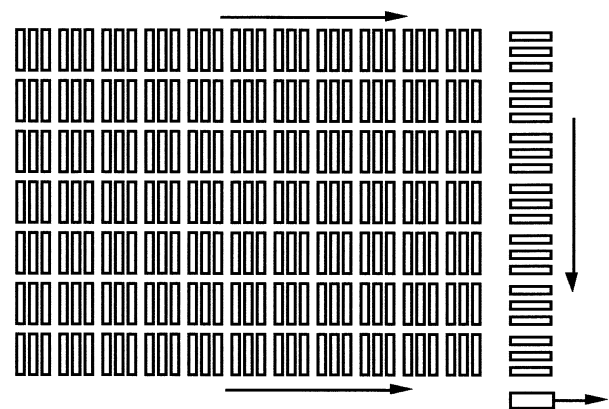


FIGURE 16 Readout technique for a two-dimensional CCD array. Vertical columns of charge packets are sequentially transferred to the output column, then clocked down to the charge-sensing amplifier electrode.

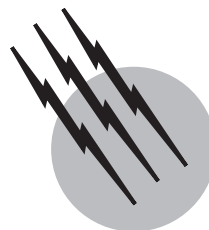
can isolate the electrons from the electrode interface, preventing electron trapping and consequent poor transfer efficiency. This technique distinguishes a “buried channel” from a “surface channel” CCD. Impurity doping gradients along the direction of the electron flow can also implement two-phase clocking. Silicon CCDs are the imaging detector of choice for TV cameras. Three-color operation is obtainable by overlaying color filters on adjacent columns or using three separate arrays with appropriate external filters. The dimensional uniformity and stability of the solid array makes color registration simple in this case, and of course the complete lack of image distortion affords superb metric capability. Silicon and gallium arsenide CCDs have infrared wavelength cutoff at about $1\text{ }\mu\text{m}$ and a limited number of other semiconductors are feasible for operation in this mode. A hybrid form of CCD combines the standard electron storage and clocking with adjacent metal/semiconductor, PtSi/Si, detectors on chip, which supply electrons photoexcited at wavelengths out to several micrometers. Although the quantum efficiency is only a few percent, the low *NEE* makes the platinum-silicide (PtSi) CCD an attractive alternative to discrete device arrays.

SEE ALSO THE FOLLOWING ARTICLES

IMAGING OPTICS • LASERS, SEMICONDUCTOR • OPTICAL AMPLIFIERS • PHOTONIC BANDGAP MATERIALS • RADIOMETRY AND PHOTOMETRY • VACUUM TECHNOLOGY

BIBLIOGRAPHY

- Adesida, I., and Coleman, J. (1997). Optoelectronic devices; Photodetectors, *In* “Handbook of Photonics” (M. Gupta, ed.), pp. 301–322. CRC Press, New York.
- Caniou, J. (1999). “Passive Infrared Detection,” Kluwer Academic, Boston.
- Crowe, *et al.* (1993). Detectors. *In* “The Infrared and Optical Systems Handbook” (J. Accetta and D. Shumaker, eds.), Vol. 3, pp. 175–283. Environmental Institute of Michigan, Ann Arbor, MI.
- Kingston, R. (1995). “Optical Sources, Detectors, and Systems,” Academic Press, New York.
- Kruse, P., and Skatrud, D., eds. (1997). “Uncooled Infrared Imaging Arrays and Systems,” Academic Press, New York.
- Saleh, B., and Teich, M. (1991). “Fundamentals of Photonics,” Wiley, New York.
- Sloan, S. (1994). Photodetectors. *In* “Photonic Devices and Systems” (R. Hunsperger, ed.), pp. 171–246. Dekker, New York.



Optical Diffraction

Salvatore Solimeno

University of Naples

- I. History
- II. Mathematical Techniques
- III. Helmholtz–Kirchhoff Integral Theorem
- IV. Diffraction Integrals
- V. Fock–Leontovich Parabolic Wave Equation
- VI. Ring-Shaped and Diffractionless Beams
- VII. Geometric Theory of Diffraction
- VIII. Watson Transform
- IX. Diffraction by Cylindrical Structures
- X. Diffraction Gratings
- XI. Coherent and Incoherent Diffraction Optics
- XII. Recent Developments

GLOSSARY

Airy Pattern Far field diffracted by a circular aperture illuminated by a plane wave or field on the focal plane of an axisymmetric imaging system.

Babinet's principle The sum of the fields diffracted by two complementary screens (corresponding to a situation in which the second screen is obtained from the first one by interchanging apertures and opaque portions) coincides with the field in the absence of any screen.

Boundary diffraction wave Representation introduced by Rubinowicz of the deviation of the diffracted field from the geometrical optics field as a wave originating from the screen boundary.

Diffraction-limited system Optical system in which

aberrations are negligible with respect to diffraction effects.

Far-field pattern Distribution of the field diffracted at an infinite distance from an obstruction.

Fourier methods Analysis of diffraction and aberration effects introduced by optical imaging systems based on the decomposition of the field into a superposition of sinusoids.

Fraunhofer diffraction Diffraction at a very large distance from an aperture.

Fresnel diffraction Diffraction in proximity to an aperture.

Fresnel number of an aperture Parameter related to the deviation of the field behind an aperture from the geometrical optics field.

Geometric theory of diffraction Method introduced by

J. B. Keller for calculating the field diffracted from an obstacle based on a suitable combination of geometrical optics and diffraction formulas for particular obstacles (wedges, half-planes, cylinders, spheres, etc.).

Green's theorem (Helmholtz–Kirchhoff integral theorem) Representation of the field by means of an integral extending over a closed surface containing the field point.

Huygens–Fresnel principle Every point on a primary wavefront serves as the source of secondary wavelets such that the wavefront at some later time is the envelope of these wavelets.

Kirchhoff method Expression of the field amplitude at any point behind a screen by means of an integral extending only over the open areas in the screen and containing the unobstructed values of the field amplitude and its normal derivative.

Modulation transfer function Function measuring the reduction in contrast from object to image, that is, the ratio of image-to-object modulation for sinusoids of varying spatial frequencies.

Optical transfer function Function measuring the complex amplitude of the image transmitted by an optical system illuminated by a unit-amplitude sinusoidal pattern, versus the spatial frequency.

Spatial frequency Representation of an object or an image as a superposition of sinusoids (Fourier components).

Spread function Image of a point source object; its deviation from an impulse function is a measure of the aberration diffraction effects.

ANY DEVIATION of light rays from rectilinear paths which cannot be interpreted as reflection or refraction is called diffraction. Diffraction optics (DO) deals mainly with the description of fields in proximity to caustics, foci, and shadow boundaries of wavefronts delimited by apertures and with the far fields. Imaging systems, diffraction gratings, optical resonators, and holographic and image processing systems are examples of devices that depend for their ultimate performance on DO. Rudimentary solutions to DO problems can be obtained by using the Huygens–Fresnel principle. More accurate solutions are obtained by solving the wave equation with the Helmholtz–Kirchhoff integral formula. The modern geometric theory of diffraction combines ray-optical techniques with the rigorous description of the field diffracted by typical obstacles (wedges, cylinders, spheres) to represent with great accuracy the light field diffracted by a generic obstruction.

I. HISTORY

A. Early Ideas

The phenomenon of diffraction was discovered by the Jesuit Father Francesco Maria Grimaldi and described in the treatise “*Physico Mathesis de Lumine, Coloribus et Iride*,” published in 1665, 2 years after his death. As it stands today, diffraction optics (DO) is the result of a long and interesting evolution originating from the ideas illustrated in 1690 by Christian Huygens in his “*Traité de la Lumière*.” Thomas Young had the great merit of introducing the ideas of wave propagation to explain optical phenomena that had previously been treated, following Newton, by the corpuscular theory. He introduced the principle of interference, illustrated in his three Bakerian lectures read at the Royal Society in 1801–1803. By virtue of this principle, Young was able to compute for the first time the wavelengths of different colors. In 1875, Augustin Jean Fresnel presented to the French Academy the treatise “*La diffraction de la Lumière*,” in which he presented a systematic description of the fringes observed on the dark side of an obstacle and was able to show agreement between the measured spacings of the fringes observed and those calculated by means of the wave theory.

B. Successes of the Wave Theory

A remarkable success of the wave theory of light was recorded in 1835 with the publication in the *Transactions of the Cambridge Philosophical Society* of a fundamental paper by Sir George Biddell Airy, director of the Cambridge Observatory, in which he derived an expression for the image of a star seen through a well-corrected telescope. The image consists of a bright nucleus, known since then as *Airy's disk*, surrounded by a number of fainter rings, of which only the first is usually bright enough to be visible to the eye. Successive developments took advantage of Fresnel's ideas to solve a host of diffraction problems by using propagation through an elastic medium as a physical model. In particular, in 1861 Glebsch obtained the analytic solution for the diffraction of a plane wave by a spherical object.

C. Fourier Optics

Since the initial success of the Airy formula, the theory of diffraction has enjoyed increasing popularity, providing the fundamental tools for quantitatively assessing the quality of images and measuring the ability of optical systems to provide well-resolved images. To deal with this complex situation, Duffieux proposed in 1946 that the imaging

of sinusoidal intensity patterns be examined as a function of their period. Then the optical system becomes known through an optical transfer function (OTF) that gives the system response versus the number of lines of the object per unit length. More recently, optical systems have been characterized by means of the numerable set of object fields that are faithfully reproduced. This approach, based on the solution of Fredholm's integral equations derived from the standard diffraction integrals, has allowed the ideas of information theory to be applied to optical instruments.

II. MATHEMATICAL TECHNIQUES

Diffraction optics makes use of a number of analytical tools:

1. Spectral representations of the fields (plane, cylindrical, spherical wave expansions; Hermite–Gaussian and Laguerre–Gaussian beams; prolate spheroidal harmonics)
2. Diffraction integrals
3. Integral equations
4. Integral transforms (Lebedev–Kontorovich transform, Watson transform)
5. Separation of variables
6. Wiener–Hopf–Fock functional method
7. Wentzel–Kramers–Brillouin (WKB) asymptotic solutions of the wave equations
8. Variational methods

In many cases the solutions are expressed by complex integrals and series, which can be evaluated either asymptotically or numerically by resorting to a number of techniques:

1. Stationary-phase and saddle-point methods
2. Boundary-layer theory
3. Two-dimensional fast Fourier transform (FFT) algorithm

III. HELMHOLTZ–KIRCHHOFF INTEGRAL THEOREM

A. Green's Theorem

Let $u(\mathbf{r})e^{i\omega t}$ be a scalar function representing a time-harmonic ($e^{i\omega t}$) electromagnetic field which satisfies the Helmholtz wave equation

$$\nabla^2 u + \omega^2 \mu \varepsilon u = 0, \quad (1)$$

where $\nabla^2 = \partial^2/\partial x^2 + \partial^2/\partial y^2 + \partial^2/\partial z^2$ is the Laplacian operator and μ and ε are the magnetic permeability (expressed in henrys per meter in mksa units) and the dielectric constant (faradays per meter), respectively, of the medium. It can be shown that the field $u(\mathbf{r}_0)$ at the point \mathbf{r}_0 depends linearly on the values taken by $u(\mathbf{r})$ and the gradient ∇u on a generic closed surface S containing \mathbf{r}_0 (*Green's theorem*):

$$u(\mathbf{r}) = \oint_S \left[G(\mathbf{r}, \mathbf{r}') \frac{\partial u(\mathbf{r}')}{\partial n_0} - u(\mathbf{r}') \frac{\partial G(\mathbf{r}, \mathbf{r}')}{\partial n_0} \right] dS', \quad (2)$$

where $G(\mathbf{r}_0, \mathbf{r}) = \exp(-ikR)/4\pi R$, $R = |\mathbf{R}| = |\mathbf{r} - \mathbf{r}_0|$, and $\partial/\partial n_0$ represents the derivative along the outward normal $\hat{\mathbf{n}}_0$ to S .

For a field admitting the ray optical representation $u(\mathbf{r}) = A(\mathbf{r})e^{-ikS(\mathbf{r})}$ the integral in Eq. (2) reduces to

$$u(\mathbf{r}) = \frac{i}{2\lambda} \oint_S \frac{A(\mathbf{r}')}{R} e^{-ik(R+S)} \times (\cos \theta_i + \cos \theta_d) dS', \quad (3)$$

where θ_i and θ_d are, respectively, the angle between the ray passing through \mathbf{r}' and the inward normal $\hat{\mathbf{n}}_0$ to S , and the angle between the direction $-\hat{\mathbf{R}}$ along which the diffracted field is calculated and $\hat{\mathbf{n}}_0$. In particular, when the surface S coincides with a wavefront, $\theta_i = 0$ and $\cos \theta_i + \cos \theta_d$ reduces to the *obliquity factor* $1 + \cos \theta_d$.

B. Huygens–Fresnel Principle

The integral of Eq. (3) allows us to consider the field as the superposition of many elementary wavelets of the form

$$du(\mathbf{r}) = \frac{e^{-ikR}}{4\pi R} \times \left[\frac{\partial u(\mathbf{r}')}{\partial n_0} + u(\mathbf{r}') \hat{\mathbf{n}}_0 \cdot \mathbf{R} \left(ik + \frac{1}{R} \right) \right] dS'. \quad (4)$$

According to the above equations every point on a wavefront serves as the source of spherical wavelets. The field amplitude at any point is the superposition of the complex amplitudes of all these wavelets. The representation of a field as a superposition of many elementary wavelets is known as the *Huygens principle*, since it was formulated in his “*Traité de la Lumière*” in 1690. Fresnel completed the description of this principle with the addition of the concept of interference.

C. Kirchhoff Formulation of the Diffraction by a Plane Screen

Let us consider the situation in which a plane Π , of equation $z' = \text{const}$, separates the region *I* containing the sources from the homogeneous region *II*, where the field

must be calculated. In this case, replacing the function G in Eq. (2) by

$$G_{\pm}(\mathbf{r}, \mathbf{r}') = \frac{\exp(-ik|\mathbf{r} - \mathbf{r}'|)}{4\pi|\mathbf{r} - \mathbf{r}'|} \pm \frac{\exp(-ik|\mathbf{r}_s - \mathbf{r}'|)}{4\pi|\mathbf{r}_s - \mathbf{r}'|}, \quad (5)$$

with \mathbf{r}_s referring to the specular image of \mathbf{r} with respect to Π , it is easy to verify that Eq. (2) reduces to

$$\begin{aligned} u(\mathbf{r}) &= -2 \int_{\Pi} \int_{\Pi} \frac{\partial u(x', y', z')}{\partial z'} \frac{e^{-ikR}}{4\pi R} dx' dy' \\ &= 2 \int_{\Pi} \int_{\Pi} u(x', y', z') \left(ik + \frac{1}{R} \right) \frac{e^{-ikR}}{4\pi R} \cos \theta_d dx' dy', \end{aligned} \quad (6)$$

use having been made of the relation $\partial G_{+}(\mathbf{r}', \mathbf{t})/\partial n_0 = G_{-}(\mathbf{r}', \mathbf{r}) = 0$ for \mathbf{r}' on Π . In particular, for an aperture Σ on a plane screen illuminated by the incident field $u_i(\mathbf{r})$, Eq. (6) can be approximated by

$$\begin{aligned} u(\mathbf{r}) &= 2 \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} P(x', y') u_i(x', y', z') \left(ik + \frac{1}{R} \right) \\ &\quad \times \frac{e^{-ikR}}{4\pi R} \cos \theta_d dx' dy', \end{aligned} \quad (7)$$

where $P(x', y')$ (the *pupil function*) takes on the value 1 if (x', y') belongs to the aperture and vanishes otherwise. In writing Eq. (7) we have tacitly assumed that the field on the exit aperture coincides with that existing in the absence of the aperture; this approximation, known as the *Kirchhoff principle*, is equivalent to the assumption that a finite exit pupil does not perturb the field on the pupil plane. Since presumably the actual perturbation is significant only near the pupil edge, we expect the error related to the application of Kirchhoff's principle to be neglected provided the aperture is sufficiently large. Exact analysis of the effects produced by some simple apertures (e.g., half-plane) confirms the validity of Kirchhoff's hypothesis for calculating the field near the *shadow boundaries* separating the lit region from the shady side; the error becomes relevant only for field points lying deep in either the lit or the dark regions.

IV. DIFFRACTION INTEGRALS

A. Fresnel and Fraunhofer Diffraction Formulas

Let us consider a field different from zero only on a finite plane aperture Σ . If we indicate by a the radius of the smallest circumference encircling the aperture and assume that $|z - z'| \gg a$, we can approximate the distance

R with $|z - z'| + \frac{1}{2}[(x - x')^2 + (y - y')^2]/|z - z'|$, so that the diffraction integral of Eq. (7) reduces to the *Fresnel formula*

$$\begin{aligned} u(x, y, z) &= \frac{i}{\lambda d} e^{-ikd} \iint_{\Sigma} u(x', y', z') \\ &\quad \times \exp \left\{ -i \frac{k}{2d} [(x - x')^2 + (y - y')^2] \right\} dx' dy', \end{aligned} \quad (8)$$

where $d = |z - z'|$.

When $d \gg a^2\pi/\lambda$, where a is the radius of the aperture, we can neglect the terms of the exponential in the integrand of Eq. (8) proportional to $x'^2 + y'^2$, so that

$$\begin{aligned} u(x, y, z) &= \frac{i}{\lambda d} e^{-ikR_0} \iint_{\Sigma} u(x', y', z') \\ &\quad \times \exp \left(ik \frac{xx' + yy'}{d} \right) dx' dy', \end{aligned} \quad (9)$$

where $R_0 = d + (x^2 + y^2)/2d$. The above equation, referred to as the *Fraunhofer diffraction formula*, allows us to express the far field in terms of the two-dimensional Fourier transform of u on the aperture plane, evaluated for $k_x = kx/d$ and $k_y = ky/d$. The Fraunhofer fields of some typical apertures illuminated by plane waves are plotted in Fig. 1.

B. Rotationally Invariant Fields and Airy Pattern

The far field radiated by a circular aperture can be obtained by assuming A [see Eq. (3)] independent of the angle ϕ and integrating the integral on the right-hand side of Eq. (9) with respect to ϕ , thus obtaining

$$\begin{aligned} u(\rho, z) &= -i \cdot \text{NA} \cdot ka \exp(-ikd - ikd\theta^2/2) \\ &\quad \times \int_0^1 A(ax) J_0(ka\theta x) \exp[-ikS(ax)] x dx, \end{aligned} \quad (10)$$

where $\rho = \sqrt{x^2 + y^2}$, $\theta = \rho/d$, NA is the numerical aperture of the lens, and $u = A \exp(-ikS)$. In particular, for $A = 1$, $S = 0$, Eq. (10) gives

$$u(\rho, z) \propto 2J_1(ka\theta)/ka\theta, \quad (11)$$

where J_1 is the Bessel function of first order. The intensity distribution associated with the above field is known as the *Airy pattern* (see Fig. 1d). Because of the axial symmetry, the central maximum corresponds to a high-irradiance central spot known as the *Airy disk*. Since $J_1(v) = 0$ for $v = 3.83$, the angular radius of the Airy disk is equal to $\theta = 3.83/ka$.

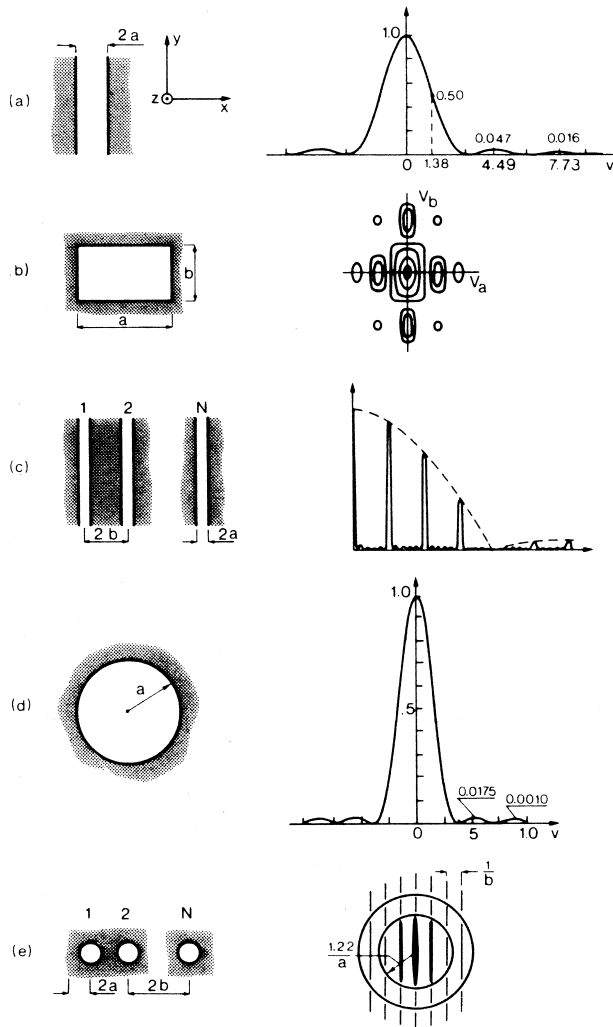


FIGURE 1 Diffraction patterns of typical apertures. The pattern functions $|G(\theta, \phi)|^2$ are respectively proportional to (a) $[\sin(v)/v]^2$, $v = ka \sin \theta \cos \phi$; (b) $\{[\sin(v_a)/v_a][\sin(v_b)/v_b]\}^2$, $v_a = ka \sin \theta \cos \phi$, $v_b = kb \sin \theta \sin \phi$; (c) $\{[\sin(v_a)/v_a][\sin(Nv_b)/\sin v_b]\}^2$, $v_a = ka \sin \theta \cos \phi$, $v_b = kb \sin \theta \sin \phi$; (d) $(2J_1(v)/v)^2$, $v = ka \sin \theta$; (e) $\{[2J_1(v_a)/v_a][\sin(Nv_b)/\sin v_b]\}^2$, $v_a = ka \sin \theta$, $v_b = kb \sin \theta \cos \phi$.

The central spot is surrounded by a series of rings corresponding to the secondary maxima of the function $J_1(v)/v$, which occur when v equals 5.14, 8.42, 11.6, etc. On integrating the irradiance over a pattern region, one finds that 84% of the light arrives within the Airy disk and 91% within the bounds of the second dark ring.

The Airy pattern can be observed at a finite distance by focusing a uniform spherical wave with a lens delimited by a circular pupil. In this case the quantity v is replaced by $k \cdot \text{NA} \cdot \rho$, where $\rho = \sqrt{x^2 + y^2}$.

C. Resolving Power

If we consider the diffraction images of two plane waves, it is customary to assume as a resolution limit the angular

separation at which the center of one Airy disk falls on the first dark ring of the other (*Rayleigh's criterion of resolution*). This gives for the angular resolution

$$\theta_{\min} \simeq 1.22\lambda/D, \quad (12)$$

where, D represents the diameter of the exit pupil.

D. Fields in the Focal Region of a Lens

Imaging systems are designed with the aim of converging a finite conical ray congruence, radiated by a point source on the object plane, toward a focal point on the image plane. In most cases, the field relative to the region between the source and the exit pupil can be calculated by geometrical optics methods, that is, by evaluating the trajectories of the rays propagating through the sequence of refracting surfaces. However, downstream from the exit pupil, we are faced with the unphysical result of a field vanishing abruptly across the shadow boundary surface formed by the envelope of the rays passing through the edge of the *exit pupil*. To eliminate this discontinuity it is necessary to resort to the diffraction integral representation. In particular, the field on the exit aperture can be assumed to coincide with that existing in the absence of the aperture and can be calculated by geometrical optics methods.

When the numerical aperture of the beam entering or leaving the lens is quite large, it is necessary to account for the vector character of the electric field \mathbf{E} . This occurs, for example, in microscope imaging, where the aperture of the beam entering the objective can be very large. As a consequence, for rotationally symmetric lenses, the focal spot for linearly polarized light is not radially symmetric, a fact that affects the resolving power of the instrument.

If we choose a Cartesian system with the z axis parallel to the optic axis and the plane $z=0$ coinciding with the Gaussian image of the object plane $z=z_0 (<0)$ (see Fig. 2), we can show for a field point very close to the Gaussian image that Eq. (3) generalizes to the

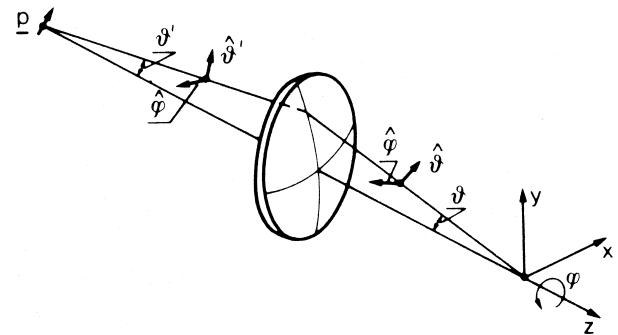


FIGURE 2 Mutual orientation of the spherical coordinate systems relative to the source and the image formed by a lens.

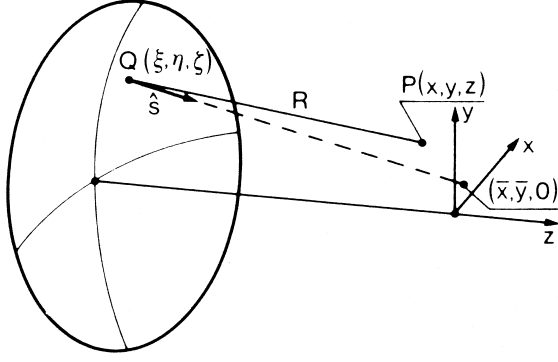


FIGURE 3 Schematic representation of field point $P(x, y, z)$ wavefront $Q(\xi, \eta, \zeta)$, and Gaussian image $(\bar{x}, \bar{y}, 0)$.

Luneburg–Debye integral

$$\mathbf{E}(\mathbf{r}) = i \frac{e^{-ikV_0}}{\lambda} \int \int_{\bar{A}} \mathbf{E}'(\hat{\mathbf{n}}_0) \exp[-ik(p(x - \bar{x}) + q(y - \bar{y}) + rz)] d\Omega, \quad (13)$$

where $\mathbf{E}' = \mathbf{R}e^{ikR}\mathbf{E}$, $p = -n_{0x}$, $q = -n_{0y}$, and $r = -n_{0z}$ are the *direction cosines* of the ray passing through (ξ, η, ζ) (see Fig. 3), while $d\Omega = d\bar{A}/R^2$ is the solid angle that the surface element $d\bar{A}$ of the exit pupil subtends at the focus.

It is now useful to introduce *optical coordinates* $v \bar{u}$ defined by

$$v = k\rho \cdot \text{NA}, \quad \bar{u} = kz \cdot \text{NA}^2, \quad (14)$$

where $\text{NA} = \sin \theta_{\max}$ represents the *numerical aperture*, θ_{\max} indicating the half-aperture in the *image space*. In this way, the Luneburg–Debye integral, for NA sufficiently small and a field linearly polarized along the x axis, reduces to

$$\mathbf{E}(\mathbf{r}) \propto \hat{\mathbf{x}} \exp(-i\bar{u}/\theta_{\max}^2) \times \int_0^{2\pi} d\phi \int_0^1 d\Theta \Theta \mathbf{E}_i(\Theta, \phi) \times \exp[i\nu\Theta \cos(\phi - \psi) + i\bar{u}\Theta^2/2]. \quad (15)$$

A three-dimensional plot of the field amplitude on a focal plane is shown in Fig. 4, and the streamlines of the Poynting vector in a focal region are represented in Fig. 5.

E. Field Near a Caustic

When the field point approaches a *caustic*, then two or more rays having almost equal directions pass through P (see Fig. 6). In this case the diffraction integral I for a two-dimensional field takes the form

$$I \propto \frac{1}{\lambda^{1/2}} \int_{-\infty}^{+\infty} \exp[-ik(as + bs^3)] ds \propto \frac{R^{1/2}}{\lambda^{1/2}} \left(\frac{2}{k\rho_c}\right)^{1/3} \text{Ai}\left[-\left(\frac{k^2\rho_c^6}{4\rho_c^4}\right)^{1/3}\right]. \quad (16)$$

$\text{Ai}(\cdot)$ is the Airy function. A plot of the field in proximity to the caustic is shown in Fig. 7.

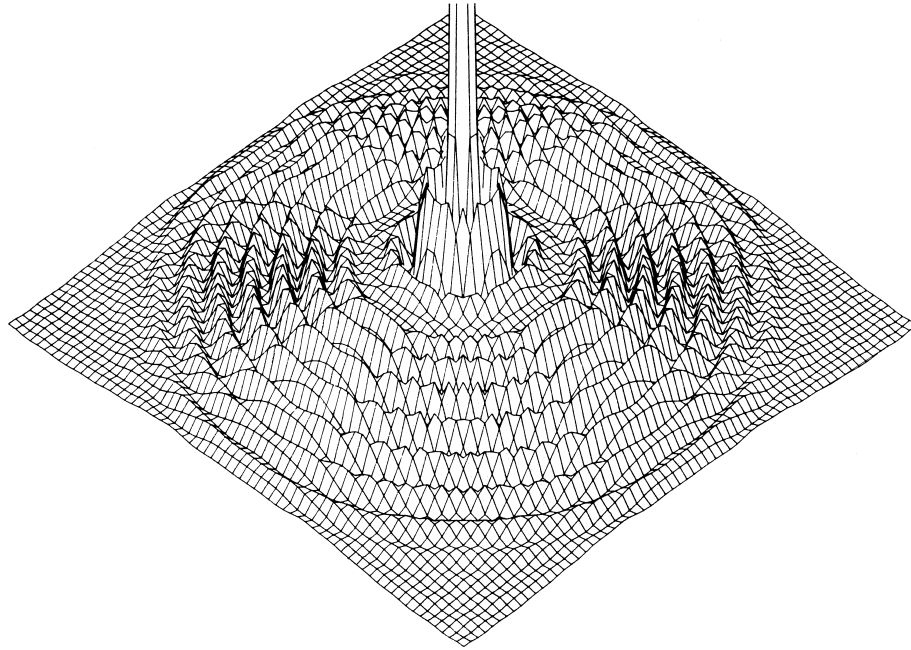


FIGURE 4 Diffraction pattern of a ring-shaped aperture with internal diameter 0.8 times the external one, illuminated by a plane wave at a distance corresponding to a Fresnel number of 15. The three-dimensional plot was obtained by using an improved fast Fourier transform algorithm. [From Luchini, P. (1984). *Comp. Phys. Commun.* **31**, 303.]

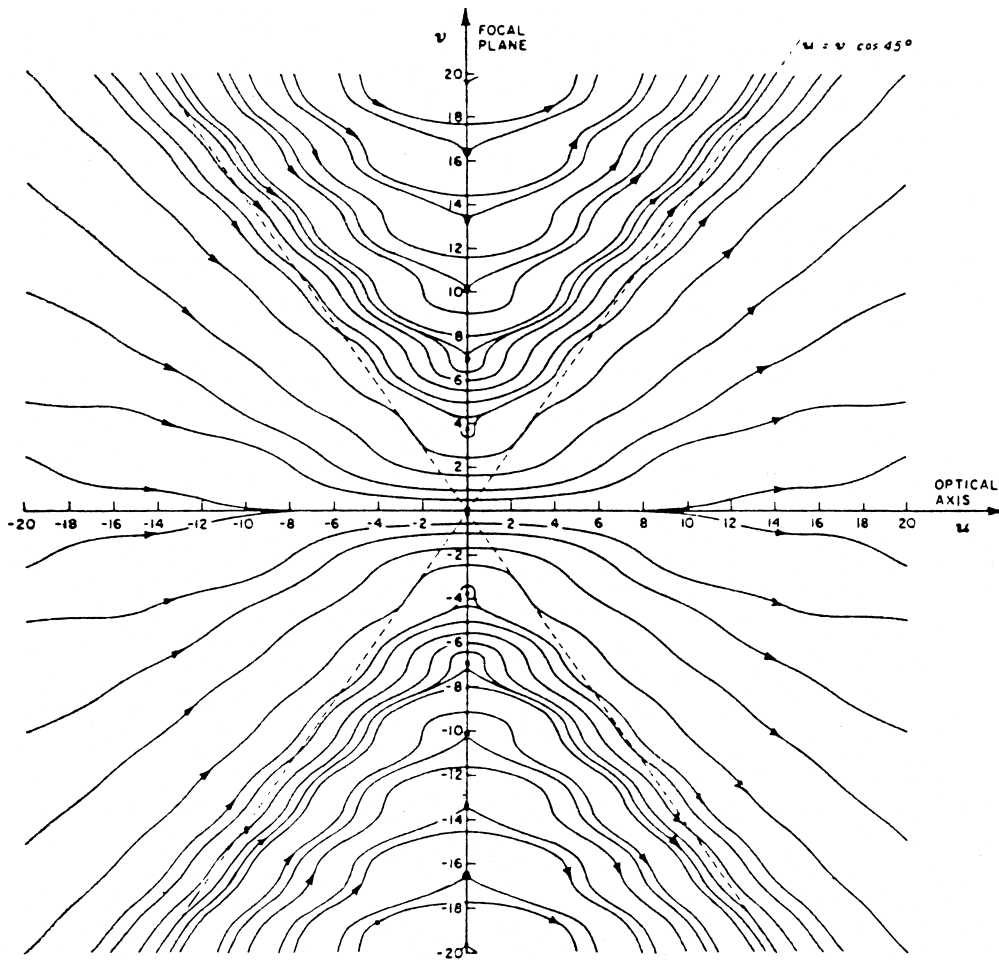


FIGURE 5 Flow lines of the Poynting vector near the focus of an aplanatic system with angular semiaperture of 45. [From Boivin, A., Dow, J., and Wolf, E. (1977). *J. Opt. Soc. Am.* **57**, 1171.]

V. FOCK-LEONTOVICH PARABOLIC WAVE EQUATION

When the wavelength is small compared with all characteristic dimensions, approximate solutions can be obtained by transforming the wave equation to a parabolic one, as shown for the first time by Leontovich in 1944. To this end let us represent a wave $u(x, y, z)$ propagating approximately along the z axis in the form

$$u(x, y, z) = W(x, y, z) e^{-ikz}.$$

Then, plugging the above expression into the wave equation (1) and neglecting the second-order derivative $\partial^2 W / \partial z^2$, we find that the problem reduces to solving the parabolic differential equation

$$\left(\frac{\partial^2}{\partial x^2} + \frac{\partial^2}{\partial y^2} \right) W = -2ik \frac{\partial}{\partial z} W.$$

Accordingly, the propagation consists of the lateral diffusion of the generally complex amplitude W , whereas the longitudinal diffusion is neglected. This is justified by the weak dependence of W on z compared with the dependence on the transverse coordinates.

The above equation is associated with the Green's function

$$G(x - x', y - y', z - z') = \sqrt{\frac{k}{2\pi|z - z'|}} \times \exp \left[ik \frac{(x - x')^2 + (y - y')^2}{2|z - z'|} \right].$$

With the help of this expression for G it can be easily shown that $W e^{-ikz}$ coincides with the field given by the Fresnel formula [see Eq. (8)].

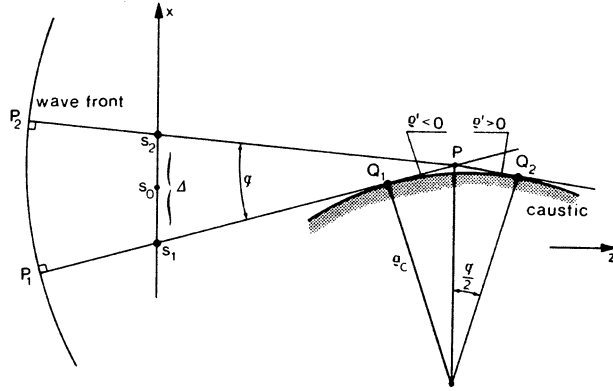


FIGURE 6 Geometry for calculation of the field in proximity to a caustic.

VI. RING-SHAPED AND DIFFRACTIONLESS BEAMS

Intense optical fields whose form shows little change in size over long paths can be obtained by means of either conical lenses or spherical lenses showing spherical aberration together with a single projecting lens (Durnin *et al.*, 1987; Herman and Wiggins, 1991). A conical lens produces fields having a transverse distribution proportional to the zeroth-order Bessel function J_0 .

A. Cylindrical Waves

The Green's function can be represented as a superposition of cylindrical waves

$$u(\rho, z) \propto \exp\left[-i\sqrt{k^2 - \chi^2}|z - z'|\right] J_m(\chi\rho)\chi d\chi.$$

In fact,

$$\begin{aligned} \frac{\exp(-ik|\mathbf{r} - \mathbf{r}'|)}{4\pi|\mathbf{r} - \mathbf{r}'|} &= -\frac{i}{4\pi} \sum_{m=-\infty}^{\infty} \exp[im(\phi - \phi')] \\ &\times \int_0^{\infty} \frac{J_m(\chi\rho)J_m(\chi\rho')}{\sqrt{k^2 - \chi^2}} \\ &\times \exp\left[-i\sqrt{k^2 - \chi^2}|z - z'|\right] \chi d\chi. \end{aligned}$$

Then, representing the field by a Fourier series with respect to ϕ'

$$u(\rho', \phi', z') = \sum_{m=-\infty}^{\infty} C_m(\rho', z') e^{im\phi'},$$

it can be easily shown that

$$\begin{aligned} C_m(\rho, z) &= \int_0^{\infty} H_m[C_m(\rho', z'); \chi] \\ &\times \exp\left[-i\sqrt{k^2 - \chi^2}|z - z'|\right] J_m(\chi\rho)\chi d\chi, \end{aligned}$$

where $H_m[f(\rho'); \chi]$ is the Hankel transform of order m ,

$$H_m[f(\rho); \chi] = \int_0^{\infty} f(\rho) J_m(\chi\rho) \rho d\rho.$$

B. Bessel–Gauss Fields

In the paraxial limit ($\chi \ll k$) the above integral representation for circularly symmetric field ($m = 0$) simplifies as

$$\begin{aligned} u(\rho, z) &= \frac{ik}{z} \exp\left(-ikz - i\frac{k\rho^2}{2z}\right) \int_0^{\infty} u(\rho', 0) \\ &\times \exp\left(-i\frac{k\rho'^2}{2z}\right) J_0\left(\frac{k\rho\rho'}{z}\right) \rho' d\rho'. \end{aligned}$$

In particular, for an incident beam of the form

$$u(\rho, 0) = J_0(k \sin \theta \rho) \exp\left(-\frac{\rho^2}{w_0^2}\right)$$

we have

$$\begin{aligned} u(\rho, z; \omega) &= \frac{w_0}{w(z)} \exp\left[-ik \cos \theta z + i\psi(z)\right. \\ &\left. + \frac{ik}{2q(z)}(\rho^2 + \sin^2 \theta z^2)\right] J_0\left(i\frac{kb \sin \theta}{q(z)} \rho\right), \end{aligned}$$

with $q(z) = z - ib$ and $\psi(z) = \arctan(z/b)$. For $z \rightarrow \infty$ the field becomes peaked along the surface of a cone of aperture θ . This property proves particularly useful for generating conical beams, such as those used for implementing optical traps used in laser cooling apparatus.

VII. GEOMETRIC THEORY OF DIFFRACTION

As we know from experience, the edges of an illuminated aperture shine when observed from the shadow region. This fact was analyzed by Newton, who explained it in terms of repulsion of light corpuscles by the edges. In 1896 Arnold Sommerfeld obtained the rigorous electromagnetic solution of the half-plane diffraction problem. Using this result, it can be shown that the total field splits into a geometrical optics wave and a diffracted wave originating from the edge. In 1917, Rubinowicz recast the (scalar) diffraction integral for a generic aperture illuminated by a spherical wave in the form of a line integral plus a geometrical optics field. Parallel to this development, J. B. Keller (1957) successfully generalized the concept of ray by including those diffracted by the edges of an aperture.

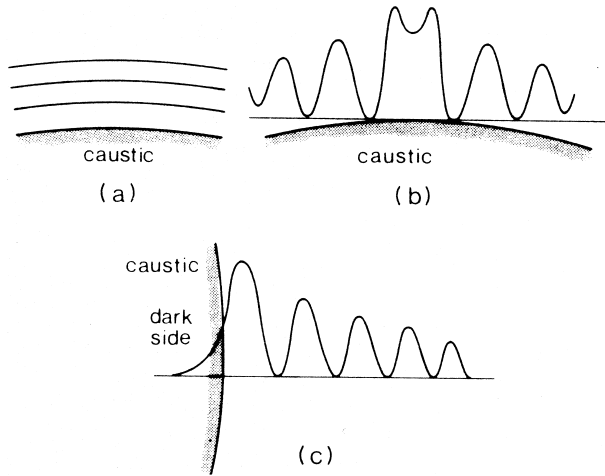


FIGURE 7 Fringes in proximity to a caustic (a) and field amplitude along a ray (b) and a normal to a caustic (c).

In order to emphasize the geometric character of his approach, Keller called it the *geometric theory of diffraction* (GTD).

A. Diffraction Matrix

The asymptotic construction of diffracted fields will be illustrated by using as example a metallic wedge (see Fig. 8) illuminated by a plane wave having the magnetic field parallel to the edge (*s-wave* or *TM wave*). It can be shown that for $\rho \rightarrow \infty$ the diffracted field is given by

$$u(\rho, \phi) \sim D_s(\phi, \phi') \frac{e^{-ik\rho}}{\rho^{1/2}} + \sum_q \exp[-ik\rho \cos(\phi - \phi_q')] \times [U(\phi - \beta_q + 2\pi) - U(\phi - \beta_q)], \quad (18)$$

where $U(x)$ is the unit step function and

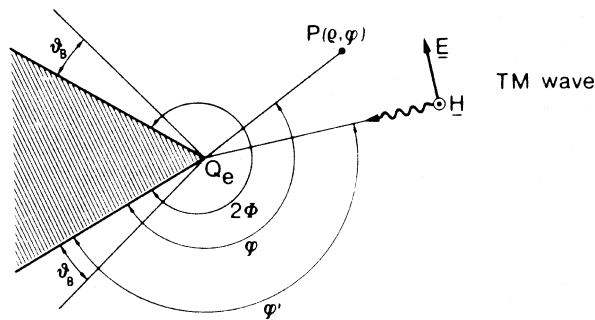


FIGURE 8 Geometry of a wedge-shaped region.

$$D_s(\phi, \phi') = -\frac{e^{-i\pi/4}\lambda^{1/2}}{N} \times \left\{ \frac{1}{\cos(\pi/N) - \cos[(\phi - \phi')/N]} + \frac{1}{\cos(\pi/N) - \cos[(\phi + \phi')/N]} \right\}, \quad (19)$$

where $N = 2\pi/\Phi$ and D_s is the *diffraction coefficient* of the wedge illuminated by an *s-wave*. For a wave having the electric field parallel to the edge (*p-wave*) the coefficient D_p can be obtained from D_s by changing to minus the plus sign before the second term on the right side of Eq. (19).

The above results have been extended to a metallic wedge illuminated by rays forming an angle β with the tangent \hat{e} to the edge at the diffraction point Q_e (see Figs. 9 and 10). It can be shown that the diffraction rays form a half-cone with axis parallel to \hat{e} and aperture equal to the angle β formed by the incident ray with the edge. If we consider the projection of the incident and diffracted rays on a plane perpendicular to the edge at Q_e , the position of the diffracted rays, forming a conical surface, is given by the angle ϕ_e , while the direction of the incident ray is defined by ϕ_e' . The electric component of the edge-diffracted ray can be expressed in the form

$$\mathbf{E}_d(\mathbf{r}) = \{\rho_1/[r(\rho_1 + r)]\}^{1/2} e^{-ikr} \mathbf{D}(\phi, \phi'; \beta) \cdot \mathbf{E}_i(Q_e), \quad (20)$$

where \mathbf{r} represents the distance of the field point P from Q_e , and ρ_1 stands for the distance of P from the focal point along the diffracted ray. The *diffraction matrix* \mathbf{D}

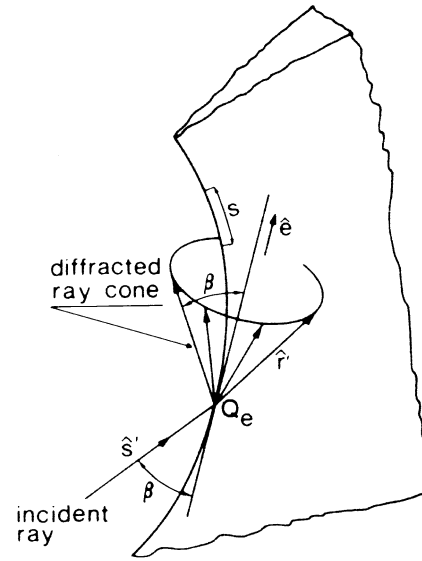


FIGURE 9 Cone of rays diffracted by an edge point Q_e .

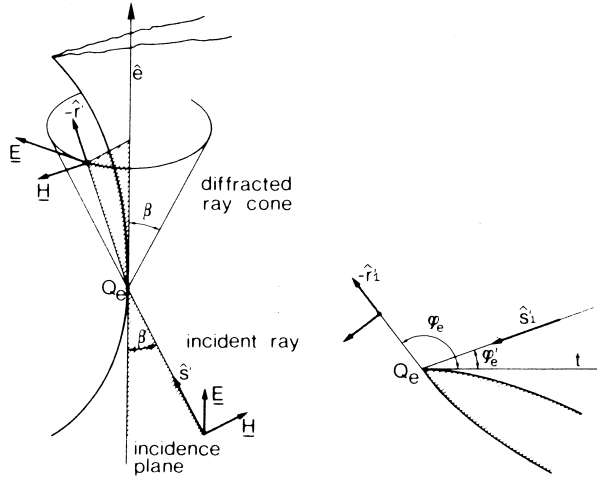


FIGURE 10 Geometry for an edge-diffracted field.

has been derived by Kouyoumjian and Pathak and put in the form

$$\mathbf{D}(\phi, \phi'; \beta) = \hat{\beta}_d \hat{\beta}_i D_p + \hat{\phi} \hat{\phi}' D_s, \quad (21)$$

where $\hat{\beta}_d = \hat{\phi} \times \hat{r}$, $\hat{\beta}_i = \hat{\phi}' \times \hat{s}'$, and D_p and D_s , which are generalizations of the diffraction coefficient (19) for $\beta \neq \pi/2$, are given by

$$D_p = \frac{e^{-in/4} \sin(\pi/N) \lambda^{1/2}}{N \sin \beta} \times \left[-\frac{1}{\cos(\pi/N) - \cos[(\phi - \phi')/N]} \pm \frac{1}{\cos(\pi/N) - \cos[(\phi + \phi')/N]} \right]. \quad (22)$$

B. Diffraction from a Slit

The GTD formalism can be applied conveniently to the calculation of the field diffracted by a slit of width $2a$ and infinity length. For simplicity, we assume a plane incident wave normal to the edges. As a first approximation we take the field on the aperture coincident with the incident field (Kirchhoff approximation). Then, following J. B. Keller, we can say that the field point P at a finite distance is reached by two different rays departing from the two edges and by a geometrical optics ray, if any (see Fig. 11a). The contribution of the diffracted rays can be expressed in the form

$$\mathbf{E}_d(\theta) = \frac{e^{-ik\rho}}{\lambda^{1/2}} \left[\mathbf{D}\left(\frac{3}{2}\pi + \theta, \frac{\pi}{2}; \frac{\pi}{2}\right) + \mathbf{D}\left(\frac{3}{2}\pi - \theta; \frac{\pi}{2}; \frac{\pi}{2}\right) \right] \cdot \mathbf{E}_i. \quad (23)$$

A point that deserves comment is the use of the Kirchhoff approximation, which can be considered valid when the slit width is much larger than the field wavelength. This approximation can be improved, as shown by Keller, by taking into account the multiple diffraction undergone by the rays departing from each and diffracted from the opposite one (see dashed line in Fig. 11b).

VIII. WATSON TRANSFORM

When the field is known on the surface of a cylinder of radius a , i.e., $u = u(a, \phi) = \langle u(a, \phi) e^{-im\phi} \rangle e^{im\phi}$, it can be represented at a generic distance ρ by a series of cylindrical waves,

$$u(\rho, \phi) = \sum_{m=-\infty}^{\infty} \langle u(a, \phi) e^{-im\phi} \rangle \times \frac{H_m^{(2)}(k\rho)}{H_m^{(2)}(\beta)} e^{im\phi}, \quad (24)$$

where $H_m^{(2)}(k\rho) = J_m(k\rho) - iN_m(k\rho)$ represents the Hankel function of second kind, and $J_m(k\rho)$ and $N_m(k\rho)$ are the Bessel and Neuman functions of order m . Treating $\langle u(a, \phi') e^{-im\phi'} \rangle_{\phi'}$ as analytic functions of the index m , we

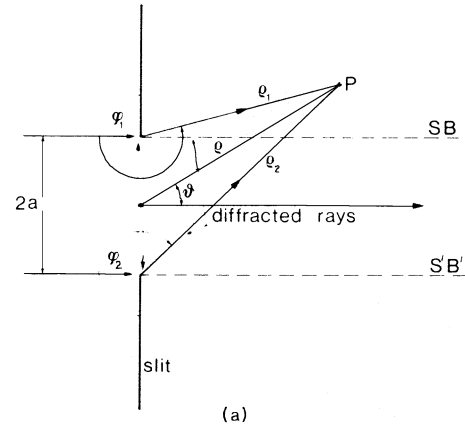


FIGURE 11 Diffraction by a slit of width $2a$ (a) and relative far-field pattern (b) for a p-wave and $ka=8$. The solid curve results from single diffraction; the dashed curve includes the effects of multiple diffraction. The dots represent the exact solution. [From Keller, J. B. (1957). *J. Appl. Phys.* **28**, 426.]

can recast the above series as a contour integral

$$u(\rho, \phi) = -\frac{i}{2} \oint_C \frac{\langle u(a, \phi') e^{-i\nu\phi'} \rangle_{\phi'} e^{i\nu(\phi-\pi)} H_v^{(2)}(k\rho)}{\sin(\nu\pi)} \frac{H_v^{(2)}(k\rho)}{H_v^{(2)}(\beta)} d\nu,$$

where C is a counterclockwise-oriented contour enclosing all the zeros of $\sin(\nu\pi)$ and leaving outside those of $H_v^{(2)}(k\rho)/H_v^{(2)}(\beta)$. The above integral can be immediately justified by noting that its integrand has inside C a series of poles coinciding with the zeros of $\sin(\nu\pi)$, that is, at $\nu = m$. Therefore, by calculating the integral by means of the residue theorem we again obtain the series (24). Taking now into account the recurrence relation $H_{-v}^{(2)}(x) = e^{-i\pi\nu} H_v^{(2)}(x)$ and the fact that $H_v^{(2)}$ has neither zeros nor poles with $\text{Im } \nu = 0$, we can recast the contour integral as

$$u(\rho, \phi) = -i \oint_{C'} \frac{\langle u(a, \phi') \cos[\nu(\pi - \phi + \phi')] \rangle_{\phi'}}{\sin(\nu\pi)} \times \frac{H_v^{(2)}(k\rho)}{H_v^{(2)}(\beta)} d\nu,$$

where C' is a contour running from $-\infty + i\varepsilon$ to $\infty + i\varepsilon$ and completed by a semicircle in the upper half-plane. The integrand of the last integral presents only polar singularities corresponding to the zeros ν_n of $H_v^{(2)}(\beta)$, which are given with good accuracy by

$$\nu_n = -\beta - e^{-i\pi/3} \left(\frac{\beta}{2}\right)^{1/3} x_n,$$

with x_n the n th zero of the Airy function $\text{Ai}(-x)$, $x_1 = 2.338$, $x_2 = 4.088$, $x_3 = 5.521$, $x_4 = 6.787$, $x_5 = 7.944$. Consequently, the field is represented by the rapidly converging series

$$\lim_{\rho \rightarrow \infty} u(\rho, \phi) = -e^{-i5\pi/12} \left(\frac{\beta}{2}\right)^{2/3} \sqrt{\frac{2\pi}{k\rho}} e^{-ik\rho} \times \sum_{n=1}^{\infty} \frac{e^{-i\pi\nu_n/2} \langle u(a, \phi') \cos[\nu_n(\pi - \phi + \phi')] \rangle_{\phi'}}{\sin(\nu_n\pi) \text{Ai}'(-x_n)}.$$

IX. DIFFRACTION BY CYLINDRICAL STRUCTURES

The radiation incident on a dielectric cylinder (e.g., an optical fiber), is scattered in a way that is dependent on the parameters of the cylinder (refractive index n , outer diameter a). The scattering problem is generally treated with the aid of an expansion in a series of Bessel and Hankel functions (cylindrical waves).

For an incident unit-amplitude monochromatic plane wave of wave vector k , directed perpendicularly to the

cylinder axis and linearly polarized along the cylinder axis, the outgoing electric field seen at an angle ϕ from the incident light propagation direction is written as

$$u_d(\rho, \phi) = \sum_{m=-\infty}^{\infty} c_m e^{im\phi}, \quad (25)$$

where the complex coefficients c_n are given by

$$c_m = \frac{n J_m(ka) J'_m(nka) - J'_m(ka) J_m(nka)}{J_m(nka) H_m^{(1)'}(ka) - n J'_m(nka) H_m^{(1)}(ka)}. \quad (26)$$

The above series can be directly used in the calculation provided that $ka \ll 1$, in which case the summation of few terms gives a good approximation of the effective result. When $ka \gg 1$ the series in (25) is not rapidly convergent, and alternative methods of calculation are necessary.

A. Debye Series

The scattered field $u_d(\rho, \phi)$ can also be written as

$$u_d(\rho, \phi) = \frac{1}{2} \sum_{m=-\infty}^{\infty} (-i)^m e^{im\phi} S_m H_m^{(2)}(k\rho), \quad (27)$$

where the scattering matrix S can be expanded in the Debye series

$$S_m(ka) = \frac{H_m^{(1)}(ka)}{H_m^{(2)}(ka)} \left[R_m + T_m^2 \frac{H_m^{(2)}(nka)}{H_m^{(1)}(nka)} \sum_{p=1}^{\infty} r_m^{p-1} \right] = \sum_{p=0}^{\infty} S_m^{(p)}(ka), \quad (28)$$

with

$$R_m = -\frac{\ln' H_m^{(2)}(ka) - n \ln' H_m^{(2)}(nka)}{\ln' H_m^{(2)}(ka) - n \ln' H_m^{(1)}(nka)}$$

$$T_m^2 =$$

$$n \frac{[\ln' H_m^{(2)}(nka) - \ln' H_m^{(1)}(nka)] [\ln' H_m^{(2)}(ka) - \ln' H_m^{(1)}(ka)]}{[\ln' H_m^{(2)}(ka) - n \ln' H_m^{(1)}(nka)]^2}$$

$$r_m = -\frac{H_m^{(2)}(nka) \ln' H_m^{(2)}(ka) - n \ln' H_m^{(2)}(nka)}{H_m^{(1)}(nka) \ln' H_m^{(2)}(ka) - n \ln' H_m^{(1)}(nka)}.$$

Here $\ln' H_m^{(2)}(ka) = H_m^{(2)'}(ka)/H_m^{(2)}(ka)$. Here R_m represents the reflection coefficient for the m th partial wave associated with direct reflection from the surface, the T_m term corresponds to transmission into the sphere and transmission to the outside, and r_m^{p-1} accounts for $p-1$ internal reflections at the surface. The scattered field can thus be written as the *Debye expansion*

$$u_d = u_d^{(0)} + \cdots u_d^{(p)} + \cdots, \quad (29)$$

where

$$u_d^{(0)} = \frac{1}{2} \sum_{m=-\infty}^{\infty} (-i)^m \frac{H_m^{(1)}(ka)}{H_m^{(2)}(ka)} R_m e^{im} H_m^{(2)}(k\rho) \quad (30)$$

and

$$u_d^{(p)} = \frac{1}{2} \sum_{m=-\infty}^{\infty} (-i)^m \frac{T_m^2 H_m^{(1)}(ka) H_m^{(2)}(mka)}{H_m^{(2)}(ka) H_m^{(1)}(mka)} \times r_m^{p-1} e^{im\phi} H_m^{(2)}(k\rho). \quad (31)$$

The p th component of the m th partial wave corresponds to transmission into the cylinder (T_m) followed by a bouncing back and forth between $\rho = a$ and $\rho = 0$ a total of p times with p internal reflections at the surface (r_m^{p-1}) and a final transmission to the outside (T_m). For ka large and $n \simeq 1$ the contributions $u_d^{(p)}$ can be neglected and u_d reduces to $u_d^{(0)}$. The rate of convergence of the Debye series is determined by the damping produced at each internal reflection.

B. Watson–Regge Representation

The p th component $u_d^{(p)}$ of the scattered field is represented by the series of partial waves indicated in (31). The summation in (31) can be considered also as an integration over the complex plane taken around the closed contour C enclosing all the poles due to the zeros of $\sin(\nu\pi)$,

$$u_d^{(p)}(\rho, \theta) = -\frac{i}{2} \int_{-\infty+i\varepsilon}^{\infty+i\varepsilon} \frac{\exp i\left(\frac{1}{2} - p\right)\pi\nu \cos[\nu(\theta + p\pi)]}{\sin \nu\pi} \times S_\nu^{(p)}(ka) H_\nu^{(2)}(k\rho) d\nu$$

To calculate the integral, it is necessary to know the poles of $S_\nu^{(p)}$, which can be shown to be close either to the zeros of $H_\nu^{(2)}(ka)$, ν_n , or to those of $H_\nu^{(1)}(nka)$, ν'_n :

$$\nu_n \simeq -ka - e^{-i\pi/3} \left(\frac{ka}{2}\right)^{1/3} x_n - \frac{i}{\sqrt{n^2 - 1}}$$

$$\nu'_n \simeq -mka + e^{-i\pi/3} \left(\frac{mka}{2}\right)^{1/3} x_n + \frac{n}{(n^2 - 1)^{1/2}},$$

where x_n is the zero of $\text{Ai}(-x)$.

C. First Term of the Debye Expansion

The term $u_d^{(0)}(\rho, \theta)$ is given by

$$u_d^{(0)}(\rho, \theta) = \frac{i}{2} \int_{-\infty+i\varepsilon}^{\infty+i\varepsilon} e^{i\pi\nu/2} \frac{\cos \nu\theta}{\sin \nu\pi} \frac{g_\nu(n, ka, \rho)}{H_\nu^{(2)}(ka)} d\nu, \quad (32)$$

with

$$g_\nu(n, ka, \rho) = H_\nu^{(1)}(k\rho) H_\nu^{(2)}(ka) + R_\nu(ka, nka) H_\nu^{(2)}(k\rho) H_\nu^{(1)}(ka). \quad (33)$$

$R_\nu \rightarrow -1$ for $|\nu| \rightarrow \infty$ in all regions except for the one between $\eta_2 = \pi/2$ and $\eta_2 = -\pi/2$, where $R_\nu \rightarrow 0$. Consequently, in view of the asymptotic representations of $H_\nu^{(1,2)}$, the integrand tends to zero in the first quadrant, so that the path of integration ($i\varepsilon, \infty + i\varepsilon$) may be shifted to the positive imaginary axis ($i\varepsilon, i\infty$) by sweeping across the poles ν'_n , thus giving

$$u^{(0)}(\rho, \theta) = \frac{i}{2} \left(\int_{-\infty+i\varepsilon}^{i\varepsilon} + \int_{i\varepsilon}^{i\infty} \right) \times e^{i\pi\nu/2} \frac{\cos \nu\theta}{\sin \nu\pi} \frac{g_\nu}{H_\nu^{(2)}(ka)} d\nu + \pi \sum_n \exp \left(i \frac{1}{2} \pi \nu'_n \right) \frac{\cos \nu'_n \theta}{\sin \nu'_n \pi} r'_n, \quad (34)$$

where r'_n is the residue of $g_\nu/H_\nu^{(2)}$ relative to the pole located at $\nu = \nu'_n$.

Using now the representation $\sin^{-1} \nu\pi = -2ie^{i\pi\nu} \sum_{m=0}^{\infty} e^{i2\pi m\nu}$ and modifying the integration path $u^{(0)}(\rho, \theta)$ splits the sum into the geometric optic contribution $u_g^{(0)}(\rho, \theta)$ plus a residue series

$$u^{(0)}(\rho, \theta) = u_g^{(0)}(\rho, \theta) + \pi \sum_n \sum_{m=1}^{\infty} \times e^{i[3\pi/2 + 2\pi m]\nu_n} \cos(\nu_n \theta) r_n + \pi \sum_n \sum_{m=0}^{\infty} e^{i[3\pi/2 + 2\pi m]\nu'_n} \cos(\nu'_n \theta) r'_n,$$

with ν_n and ν'_n zeros of $H_\nu^{(2)}(ka)$ and $H_\nu^{(2)}(nka)$, respectively. For refractive index $n > 1$ the residues are negligibly small while $u_g^{(0)}$ is represented by the integral

$$u_g^{(0)}(\rho, \phi) = \int_C e^{3i\pi\nu/2} \cos(\nu_n \theta) \times \frac{R_\nu(nka, ka) H_\nu^{(2)}(k\rho) H_\nu^{(1)}(nka)}{H_\nu^{(1)}(ka)} d\nu, \quad (35)$$

with C a path going around the poles of R_ν and crossing the real axis twice, first between 0 and $-nka$ and then between $-nka$ and $-k\rho$. The integrand of Eq. (35) has a saddle point in each of these intervals

$$\nu_s^{(1)} = -k\rho, \quad \nu_s^{(2)} = -k\rho \sin \phi,$$

corresponding to rays hitting the cylinder with impact parameters p and $\rho \sin \phi$, respectively. Deforming C by letting it pass through these saddle points, we can evaluate the integral asymptotically, giving

$$u_g^{(0)}(\rho, \phi) = \frac{i}{2} \int_C \frac{\exp[i\nu(\phi - \pi)]}{\sin \nu\pi} c_\nu d\nu = c_{\bar{\nu}} e^{i\bar{\nu}\phi} + c_{\bar{\nu}}^* e^{-i\bar{\nu}\phi}. \quad (36)$$

This means that the diffraction pattern in the illuminated zone is equivalent to the pattern obtained from the interference of the light coming from two sources.

X. DIFFRACTION GRATINGS

In the history of physics the diffraction grating stands out as one of the most important instruments. It seems to have been invented by the American astronomer David Rittenhouse in about 1785 and rediscovered some years later by Joseph von Fraunhofer. Great progress in the manufacture of ruled gratings was due to Harry A. Rowland, who in 1882 built a machine for ruling gratings with an uncorrected periodic error much less than $1/300,000$.

When illuminated by a plane wave of assigned wavelength, a reflection (or transmission) grating produces several beams which can be labeled with an integer m ($0, \pm 1, \pm 2, \dots$) that represents the *order of diffraction*. If we indicate with θ the incidence angle and d the period of the grating grooves, the m th diffracted beam forms a reflection (transmission) angle θ_m which satisfies the so-called grating equation

$$\sin \theta_m - \sin \theta = m \frac{\lambda}{d}.$$

In particular, the zeroth, order corresponds to specular reflection found in ordinary mirrors. For m sufficiently large the above relation can be satisfied only by a complex angle $\theta_m = i\theta_m'' + \text{sgn}(m)(\pi/2)$. This means that the grating produces evanescent waves.

Let $u_{(\text{in})}(\mathbf{r}) = \exp(-i\mathbf{k}_{(\text{in})} \cdot \mathbf{r})$ be a generic scalar component of unit amplitude of the incident beam and $u_{(\text{d})}(\mathbf{r})$ be the corresponding diffracted wave by a plane grating lying on the plane $x-y$ with grooves parallel to the y axis. In view of the periodicity of the reflection boundary it can be shown that $u_{(\text{d})}(\mathbf{r}) \exp(ik_{(\text{in})x}x)$ is a periodic function of x having period d . The grating surface divides the space in two regions 1 and 2. Assuming that the incident field $u_{(\text{in})}(\mathbf{r})$ is contained in region 1, we have for the total field $u(\mathbf{r}) = u_{(\text{in})}(\mathbf{r}) + u_{(\text{d})}(\mathbf{r})$ in 1 and $u(\mathbf{r}) = u_{(\text{d})}(\mathbf{r})$. Assuming the plane of incidence normal to the grooves, $u = E_y$ for p-waves and $u = H_y$ for s-waves, we impose the continuity of tangent components

$$\begin{aligned} u_{(\text{in})} + u_{(\text{d})1} &= u_{(\text{d})2} && \text{p- and s-waves} \\ -i \left(k_{(\text{in})x} \frac{\partial f}{\partial x} - k_{(\text{in})z} \right) u_{(\text{in})} + \frac{\partial u_{(\text{d})1}}{\partial x} \frac{\partial f}{\partial x} - \frac{\partial u_{(\text{d})1}}{\partial z} \\ &= \begin{cases} \left(\frac{\partial u_{(\text{d})2}}{\partial x} \frac{\partial f}{\partial x} - \frac{\partial u_{(\text{d})2}}{\partial z} \right) & \text{p-wave} \\ \left(\frac{\partial u_{(\text{d})2}}{\partial x} \frac{\partial f}{\partial x} - \frac{\partial u_{(\text{d})2}}{\partial z} \right) \frac{n_2^2}{n_1^2} & \text{s-wave} \end{cases} \end{aligned}$$

for $\mathbf{r} = [x, y, f(x)]$ belonging to the grating surface.

If u is a p-wave, the total field is continuous together with its gradient across the grating surface, so that it can be expanded in a Fourier series in x (Petit, 1980),

$$u(x, y, z) = \exp(-ik_{(\text{in})x}x - ik_{(\text{in})y}y) \times \sum_{m=-\infty}^{\infty} \exp(-i2\pi mx) V_m(z),$$

where the $V_m(z)$ are continuous functions to be determined. Consequently, the diffracted field is completely described by the set of functions $V_m(z)$.

For determining $V_m(z)$ we plug the above series into the wave equation and integrate over a period d , thus obtaining

$$\frac{d^2 V_m}{dz^2} + k_0^2 V_m g_0 = -k_0^2 \sum_{m \neq q} V_q g_{m-q},$$

where

$$g_{m-q}(z) = \frac{1}{d} \int_0^d \frac{\beta_q^2(x, z)}{k_0^2} \exp \left[i \frac{2\pi}{d} (m - q)x \right] dx.$$

Here

$$\beta_q^2(x, z) = k_0^2 n^2(x, z) - \left[k_{(\text{in})x} + \frac{2\pi}{d} q \right]^2 - k_{(\text{in})y}^2.$$

Now we can distinguish three regions A, B, and C, defined, respectively, by $z > f_{\text{max}}$ (region A), $f_{\text{max}} > z > f_{\text{min}}$ (region B), and $f_{\text{min}} > z$ (region C). In regions A and C we can solve for V_m by putting

$$V_m(z) = \begin{cases} R_{m1} \exp(-i\beta_{m,1}z) + \partial_m \exp(ik_{(\text{in})z}), & z > f_{\text{max}} \\ T_{m2} \exp(i\beta_{m,2}z), & z < f_{\text{min}}, \end{cases}$$

where R_{m1} and T_{m2} are respectively generally complex reflection and transmission coefficients relative to the m th-order wave.

XI. COHERENT AND INCOHERENT DIFFRACTION OPTICS

An ideal imaging system can be described mathematically as a mapping of the points of the object plane Π_o into those of the image plane Π_i . For a finite wavelength and delimited pupil a unit point source located at (x_0, y_0) produces a field distribution $K(x, y; x_0, y_0)$ called the *impulse response*, which differs from a delta function $\delta^{(2)}(x - \bar{x}, y - \bar{y})$ centered on the Gaussian image of the object having coordinates (\bar{x}, \bar{y}) . As a consequence, diffraction destroys the one-to-one correspondence between the object and the image.

The departure of K from a delta function introduces an amount of uncertainty in the reconstruction of an object through its image. This is indicated by the fact that two point sources are seen through an optical instrument as clearly separate only if their distance is larger than

a quantity W roughly coincident with the dimension of the region on Π_i where K is substantially different from zero. The parameter W , which measures the smallest dimension resolved by an instrument, is proportional to the wavelength. This explains the increasing popularity of UV sources, which have permitted the implementation of imaging systems capable of resolutions better than $0.1 \mu\text{m}$, a characteristic exploited in microelectronics for the photolithographic production of VLSI circuits.

A. Impulse Response and Point Spread Function

Let us consider a unit point source in (x_0, y_0, z_0) producing a spherical wavefront transformed by a composite lens into a wave converging toward the paraxial image point (\bar{x}, \bar{y}, z) . Using the Luneburg–Debye integral, we can express the impulse response $K(x, y; x_0, y_0)$ on Π_i as an integral extended to the wavefront of the converging wave:

$$\begin{aligned} K(x, y; x_0, y_0) &= \exp\left(-ik\frac{x_0^2 + y_0^2}{2d_0} - ik\frac{\bar{x}^2 + \bar{y}^2}{2d}\right. \\ &\quad \left.+ iv_x \bar{p} + iv_y \bar{q}\right) \times \frac{\Omega}{2\pi} \\ &\quad \times \iint P(p, q) A(p, q) \\ &\quad \times \exp[i(v_x p + v_y q)] dp dq \\ &\equiv e^{-i\Phi} \bar{K}(x - \bar{x}, y - \bar{y}), \end{aligned} \quad (37)$$

where $u_x = k(x - \bar{x}) \cdot \text{NA}$ and $u_y = k(y - \bar{y}) \cdot \text{NA}$ are the optical coordinates of the point (x, y) referred to the paraxial image (\bar{x}, \bar{y}) , NA is the numerical aperture of the lens in the image space, P is the pupil function, p and q are proportional to the optical direction cosines of the normal to the converging wavefront in the image space in such a way that $p^2 + q^2 = 1$ on the largest circle contained in the exit pupil, and d_0 and d are the distances of the object and the image from the respective principal planes.

Physically we are interested in the intensity of the field, so it is convenient to define a *point spread function* $t(x, y; x_0, y_0)$ given by

$$t(x, y; x_0, y_0) = |K(x, y; x_0, y_0) / K(\bar{x}, \bar{y}; x_0, y_0)|^2. \quad (38)$$

In particular, for *diffraction-limited instruments* with circular and square pupils, respectively, \bar{K} is

$$\bar{K} \propto 2J_1(v)/v, \quad \bar{K} \propto \frac{\sin v_x}{v_x} \frac{\sin v_y}{v_y}, \quad (39)$$

with

$$\begin{aligned} v_x &= k\text{NA}(x + Mx_0), & v_y &= k\text{NA}(y + My_0), \\ v &= (v_x^2 + v_y^2)^{1/2}. \end{aligned} \quad (40)$$

M is the magnification of the lens and J_1 the Bessel function of first order.

B. Coherent Imaging of Extended Sources

We can apply the superposition principle to calculate the image field $i(x, y, z)$ corresponding to the extended object field $o(x_0, y_0, z_0)$ defined on the plane region Σ_0 , obtaining

$$\begin{aligned} i(x, y, z) &= \iint_{\Sigma_0} K(x, y, z; x_0, y_0, z_0) o(x_0, y_0, z_0) \\ &\quad \times dx_0 dy_0. \end{aligned} \quad (41)$$

If we are interested in the intensity, we have

$$\begin{aligned} |i(x, y, z)|^2 &= \iint_{\Sigma_0} dx_0 dy_0 \iint_{\Sigma_0} dx'_0 dy'_0 \\ &\quad \times K(x, y, z; x_0, y_0, z_0) \\ &\quad \times K^*(x, y, z; x'_0, y'_0, z_0) \\ &\quad \times o(x_0, y_0, z_0). \end{aligned} \quad (42)$$

C. Optical Transfer Function

If we are interested in the intensity distribution $I(x, y, z) \propto \langle i(x, y, z) i^*(x, y, z) \rangle$, Eq. (42) gives for an isoplanatic system

$$\begin{aligned} I(x, y, z) &= \frac{1}{M^2} \iint_{\Sigma_0} I_0\left(-\frac{x_0}{M}, -\frac{y_0}{M}\right) \\ &\quad \times E(x - x_0, y - y_0) dx_0 dy_0, \end{aligned} \quad (43)$$

where M is the magnification. With Fourier transformation, this relation becomes

$$\tilde{I}(\alpha, \beta) \propto T(\alpha, \beta) \tilde{I}_0(\alpha, \beta), \quad (44)$$

where \tilde{I} is the two-dimensional Fourier transform of $I(x, y)$, and $T(\alpha, \beta)$ is the *optical transfer function* (OTF) of the system (see Fig. 12). For an isoplanatic system, $T(\alpha, \beta)$ is proportional to the Fourier transform of the point-spread function t , so in absence of aberrations

$$\begin{aligned} T(\alpha, \beta) &= \frac{\iint_{-\infty}^{+\infty} \tilde{t}(v_x, v_y) \exp(i\alpha v_x + i\beta v_y) dv_x dv_y}{\iint_{-\infty}^{+\infty} \tilde{t}(v_x, v_y) dv_x dv_y} \\ &= \frac{\iint_{-\infty}^{+\infty} P(p + \alpha/2, q + \beta/2) P(p - \alpha/2, q - \beta/2) dp dq}{\exp[-ik \iint_{-\infty}^{+\infty} P(p, q)] dp dq} \end{aligned} \quad (45)$$

by virtue of convolution and Parseval's theorems. The modulus of T is called the *modulation transfer function* (MTF).

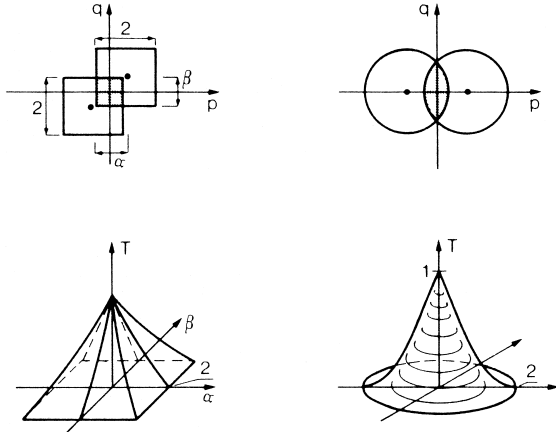


FIGURE 12 Optical transfer functions related to incoherent illumination of square (left) and round (right) pupils. T is proportional to the convolution of the pupil function [see top diagram and Eq. (45).]

For a circular pupil, the OTF reads

$$T(\omega) = \begin{cases} (2/\pi)[\arccos(\omega/2) \\ -\omega(1 - \omega^2/4)^{1/2}], & 0 \leq \omega \leq 2 \\ 0, & \text{otherwise,} \end{cases} \quad (46)$$

where $\omega = (\alpha^2 + \beta^2)^{1/2}$. Then for a circular pupil the spatial frequency ω is limited to the interval $(0, 2)$. The dimensionless frequency ω , conjugate of the optical coordinate v , is related to the *spatial frequency* f , expressed in cycles per unit length, by $f = k\omega \cdot \text{NA}$. The expression just given for the OTF indicates that the diffraction sets an upper limit to the ability of the system to resolve a bar target with a normalized spatial frequency greater than 2.

XII. RECENT DEVELOPMENTS

A. Calculation of Diffraction Integrals by Expansion of the Field in Gaussian Beams

Let us represent the field in the form

$$u(\vec{r}) = e^{-ik_0 z} A(\vec{r}),$$

where \hat{z} is an arbitrary direction and $A(\vec{r})$ satisfies the equation

$$\frac{\partial}{\partial z} A(\vec{r}) = -\frac{i}{2k_0} \left(\nabla_t^2 + \frac{\partial^2}{\partial z^2} \right) A.$$

Here ∇_t^2 stands for the transverse Laplacian.

If A is a slowly varying function of z , we can neglect the second-order derivative of it with respect to z , thus obtaining the Fock–Leontovich parabolic wave equation

$$\frac{\partial}{\partial z} A(\vec{r}) = -\frac{i}{2k_0} \nabla_t^2 A(\vec{r}).$$

Accordingly, the electromagnetic propagation is reduced to an irreversible diffusive process similar to those associated with wavefunction evolution in quantum mechanics. It is straightforward to show that the Green's function associated with the above equation coincides with the Fresnel kernel characteristic of paraxial propagation.

We look for a trial solution of the above equation by putting

$$A(\vec{\rho}, z) = \frac{1}{w(z)} f\left(\sqrt{2} \frac{\rho}{w(z)}\right) \times \exp\left[-i\left((l+m+1)\psi(z) + \frac{k}{2q(z)}\rho^2\right)\right],$$

with $\vec{\rho} = \hat{x}x + \hat{y}y$. Here $w(z)$, $q(z)$, $\psi(z)$, and f are functions to be determined by plugging the above function into the parabolic wave equation. In particular, we choose

$$f = \frac{1}{\sqrt{\pi 2^{l+m-1} l! m!}} H_l\left(\sqrt{2} \frac{x}{w}\right) H_m\left(\sqrt{2} \frac{y}{w}\right),$$

where

$$\psi(z) = \arctan \frac{2z}{kw_0^2},$$

$$w^2(z) = w_0^2(1 + z^2/z_R^2),$$

with z_R the so-called Rayleigh length (equivalent to the focal depth) of the beam, which is related to the spot size w_0 in the waist by the relation

$$z_R = \frac{\pi w_0^2}{\lambda}.$$

The complex curvature radius q is given by

$$\frac{1}{q} = \frac{1}{z + iz_R} = \frac{1}{R(z)} - i \frac{\lambda}{\pi w^2(z)},$$

with R the curvature radius of the wavefront,

$$R(z) = z + \frac{z_R^2}{z}.$$

Finally, H_m is the Hermite polynomial of order m .

If we use polar coordinates, it may be preferable to represent the field as a superposition of Gauss–Laguerre modes,

$$f = \sqrt{\frac{2}{\pi} \frac{p!}{(l+p)!}} \left(\sqrt{2} \frac{\rho}{w}\right)^l \times L_p^l\left(2 \frac{\rho^2}{w^2}\right) \cos\left(l\theta + \varepsilon \frac{\pi}{2}\right),$$

with $\varepsilon = 0, 1$; and L_p^l is the Laguerre polynomial.

These modes can be used for representing a generic field in the paraxial approximation. In particular, as we expand the field on an aperture we can immediately obtain the transmitted field by using the above expressions for propagating the single modes. For example, we can use Gaussian beams having the same curvature of the wavefront of the incoming field and a spot size depending on the illumination profile.

B. Expansion in Gauss–Laguerre Beams

If we represent the field on the output pupil of an optical system as a superposition of Gauss–Laguerre modes, we can describe the effects on the diffracted field of the limited size of the pupil and of the aberrations by introducing a matrix, i.e.,

$$\mathbf{E}_d = \mathbf{K} \cdot \mathbf{E}_i,$$

where $\mathbf{E}_{i,d}$ is a vector whose components are the generally complex amplitudes of the Gauss–Laguerre modes composing the incident (diffracted) field. \mathbf{K} is a matrix which can be represented as a product $\mathbf{K}_a \cdot \mathbf{K}_p \cdot \mathbf{K}_p$ representing the finite size of the pupil,

$$\begin{aligned} \mathbf{K}_{pp'}^{ll'\varepsilon\varepsilon'} &= \delta_{pp'}\delta_{ll'}\delta_{\varepsilon\varepsilon'} - \delta_{ll'}\delta_{\varepsilon\varepsilon'} \sqrt{\frac{p!p'!}{(l+p)!(l'+p')!}} \\ &\times e^{-u} \sum_{k=0}^{p+p'+1} \frac{u^k}{k!} \sum_{m+m'+l \geq k} (-1)^{m+m'} \\ &\times \frac{(m+m'+l)!}{m!m'!} \binom{p+l}{p-m} \binom{p'+l}{p'-m'}, \end{aligned}$$

where $u = 2a^2/w^2$, with a pupil radius and w the spot size of the Gaussian beam corresponding to it, \mathbf{K}_a accounts for the aberrations,

$$\mathbf{K}_a = \exp(ik\mathbf{W}).$$

\mathbf{W} is the aberration matrix. In particular, for spherical aberrations \mathbf{W} takes the form

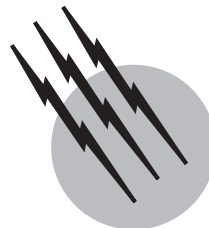
$$\begin{aligned} \mathbf{W}_{pp'\varepsilon\varepsilon'}^{ll'} &\propto \delta_{ll'}\delta_{pp'}\delta_{\varepsilon\varepsilon'} \\ &- \delta_{ll'}\delta_{\varepsilon\varepsilon'} \sqrt{\frac{p!p'!}{(l+p)!(l'+p')!}} \\ &\times \sum_{mm'} \frac{(-1)^{m+m'}}{m!m'!} \binom{p+l}{p-m} \binom{p'+l}{p'-m'} \\ &\times (m+m'+l+2)! \end{aligned}$$

SEE ALSO THE FOLLOWING ARTICLES

COLOR SCIENCE • DIFFRACTIVE OPTICAL COMPONENTS
• ELECTROMAGNETICS • WAVE PHENOMENA

BIBLIOGRAPHY

- Barakat, R. (1980). In "The Computer in Optical Research" (B. R. Frieden, ed.), pp. 35–80. Springer-Verlag, Berlin.
- Born, M., and Wolf, E. (1970). "Principles of Optics," Pergamon, Oxford.
- Durnin, J., Miceli, J. J., and Eberly, J. H. (1987). *Phys. Rev. Lett.* **58**, 1449.
- Duffieux, P. M. (1983). "The Fourier Transform and its Application to Optics," Wiley, New York.
- Gaskill, J. D. (1978). "Linear Systems, Fourier Transforms and Optics," Wiley, New York.
- Hansen, R. C. (ed.) (1981). "Geometric Theory of Diffraction," IEEE Press, New York.
- Herman, R. A., and Wiggins, T. A. (1991). *J. Opt. Soc. Am. A* **8**, 932.
- Kouyoumjan, R. G., and Pathak, P. H. (1974). *Proc. IEEE* **62**, 1448.
- Northover, F. H. (1971). "Applied Diffraction Theory," Elsevier, New York.
- Petit, R. (ed.). (1980). "Electromagnetic Theory of Gratings," Springer-Verlag, Berlin.
- Solimeno, S., Crosignani, B., Di Porto, P. (1986). "Guiding, Diffraction and Confinement of Optical Radiation," Academic Press, Orlando, FL.



Optical Isolators, Circulators

Yale Cheng

JDS Uniphase Inc.

- I. Principle of Optical Isolators
- II. Principle of Optical Circulators
- III. Reflective Design of Optical Isolators and Circulators

GLOSSARY

Birefringent material A crystalline anisotropic material having different refractive indices along different axes of the material. A birefringent material has at least two refractive indices.

Faraday effect A magneto-optic effect discovered by Michael Faraday in 1845 in which the orientation (polarization plane) of an electromagnetic wave is rotated under a magnetic field applied parallel to the propagation direction of the electromagnetic wave.

Insertion loss Energy loss during propagation through an optical element, defined by the ratio of output power at the output port of the element to the total power in decibel (dB) units.

Isolation Energy loss during reversed propagation through an optical element, defined by the ratio of output power at the input port of the element to the total power in decibel (dB) units.

Polarization Property of an electromagnetic wave describing the time varying direction and amplitude of an electric field vector.

Polarization-dependent loss (PDL) Maximum insertion loss variation of an optical device when the polarization state of a light beam launched into the device is randomly varied.

Polarization extinction ratio A parameter determines linearity of a linearly polarized beam and is defined as the power ratio at one polarization direction to that at a perpendicular polarization direction in decibel (dB) units. For example, the extinction ratio of a perfect linearly polarized light would be infinite and that of a circularly polarized light would be 0 dB.

Polarization mode dispersion (PMD) Relative time (or phase) delay introduced by an optical device between two orthogonal polarization vectors of a light beam after passing through the device.

Waveplate Also called retardation plate or phase shifter. Waveplates are made from birefringent materials for changing polarization state. The change of polarization state caused by the waveplate is reciprocal; that is, if a light beam is passed through a waveplate twice from opposite directions, the effect of the waveplate will be cancelled. Half-waveplates (introducing 180-degree phase difference between two orthogonal polarized beams) and quarter-waveplates (introducing 90-degree phase difference between two orthogonal polarized beams) are the most commonly used waveplates.

OPTICAL ISOLATORS are two-port optical devices based on the Faraday effect. There are two types of optical

isolators, polarization-dependent (also called free space) and polarization-independent (also called in-line) isolators, but the actual structures of optical isolators can vary significantly depending on materials and operating principles. The function of an optical isolator is to transmit a lightwave from the input port to the output port with a maximum intensity (minimum insertion loss), but at the same time to block any light transmission in the reversed direction (from the output to input ports). Optical circulators are at least three-port optical devices based on the same Faraday effect and are extensions of the optical isolators. Optical circulators can also be divided into polarization-dependent and polarization-independent types. An optical circulator is constructed using similar base elements to an optical isolator, but its structure design is much more complex. The function of an optical circulator is to transmit a lightwave from one port to the next sequential port with a maximum intensity, but at the same time to block any light transmission from one port to the previous port.

I. PRINCIPLE OF OPTICAL ISOLATORS

Since the invention of optical laser in 1958 by Schawlow and Townes of Bell Laboratories and the first demonstration in 1960 by Maiman, development of optical laser sources has progressed drastically. Telecommunication industries have benefited largely from this invention of optical laser, especially semiconductor lasers. Today, optical communication systems using semiconductor lasers have become the single largest information transmission media. However, one of the common problems in lasers is that the stability of a laser is very sensitive to external light feedback caused by the reflection of the laser beam. Optical isolators were initially developed to prevent any light feedback into the laser, thus improving the stability of the laser. Optical isolators have been used inside or outside a laser cavity to prevent multiple reflection of a light. Applications of optical isolators have been expanded into various fields, and currently, they are being used in optical sensors, various optical lasers including signal and pump lasers, optical amplifiers, Cable TV systems, etc.

Most of the optical isolators are designed to be functional with wavelength windows of 850, 1310, and 1550 nm for optical communication applications and with wavelength windows of 980 and 1480 nm for pumping optical amplifiers. Performances of an optical isolator are mainly determined by the insertion loss and isolation. An ideal optical isolator would have zero insertion loss and infinite isolation. There are two types of optical isolators; one is a polarization-dependent optical isolator which only functions for a linearly polarized light in forward direction, and the other is a polarization-independent

isolator which is functional for any light with any polarization in both forward and backward directions.

A. Faraday Effect

All optical isolators and circulators are based on the Faraday effect, which is a magneto-optic effect discovered by Michael Faraday in 1845. The Faraday effect is a phenomenon in which the polarization plane of an electromagnetic (light) wave is rotated in a material under a magnetic field applied parallel to the propagation direction of the lightwave. A unique feature of the Faraday effect is that the direction of the rotation is independent of the propagation direction of the light; that is, the rotation is nonreciprocal. The angle of the rotation θ is a function of the type of Faraday material, the magnetic field strength and the length of the Faraday material, and can be expressed as

$$\theta = VBL \quad (1)$$

where V is the Verdet constant of a Faraday material, B is the magnetic field strength parallel to the propagation direction of the lightwave, and L is the length of the Faraday material.

The Verdet constant is a measure of the strength of the Faraday effect in a particular material, and a large Verdet constant indicates that the material has a strong Faraday effect and is more suitable for building an optical isolator since a short length of material can be used. The Verdet constant normally varies with wavelength and temperature. The Faraday effect has been observed in various types of materials including solids, liquids, and gases. Only solid Faraday materials have so far been used in optical isolators and circulators.

Since the Verdet constant of a Faraday material is wavelength dependent, an optical isolator is typically only functional within a specific wavelength band. Depending on the operating wavelength range of the optical isolator, different Faraday materials are used in the isolator.

Rare-earth-doped glasses and rare-earth garnet crystals are the common Faraday materials used in optical isolators due to their large Verdet constant. To be specific, Terbium-doped glasses, Terbium Gallium Garnet (TGG), Yttrium Iron Garnet (YIG, $\text{Y}_3\text{Fe}_5\text{O}_{12}$), $\text{Cd}_{0.70}\text{Mn}_{0.17}\text{Hg}_{0.13}\text{Te}$, and Bismuth-substituted Iron Garnets [BIG, i.e., $\text{Gd}_{3-x}\text{Bi}_x\text{Fe}_5\text{O}_{12}$, $(\text{BiYbTb})_3\text{Fe}_5\text{O}_{12}$, $(\text{HoTbBi})_3\text{Fe}_6\text{O}_{12}$, $(\text{BiTb})_3(\text{FeGa})_5\text{O}_{12}$, etc.] are the most common materials. Terbium-doped glasses, TGG, and $\text{Cd}_{0.70}\text{Mn}_{0.17}\text{Hg}_{0.13}\text{Te}$ are frequently used for wavelengths below 1100 nm due to their low absorption and high Verdet constant at short wavelength range. $\text{Cd}_{0.70}\text{Mn}_{0.17}\text{Hg}_{0.13}\text{Te}$ is specifically developed for use at the wavelength around 980 nm. YIG and BIG are used for wavelengths above 1100 nm due to

TABLE I Typical Parameters of Various Faraday Materials at Wavelength Around 1550 nm

Materials	Insertion loss (dB)	Extinction ratio (dB)	Temperature dependence of the rotation angle (degree/°C)	Wavelength dependence of the rotation angle (degree/nm)
YIG	<0.1	>38	−0.042	−0.040
Gd _{3−x} Bi _x Fe ₅ O ₁₂	<0.1	>40	−0.1	−0.068
(BiYbTb) ₃ Fe ₅ O ₁₂	<0.1	>40	−0.045	−0.060
(HoTbBi) ₃ Fe ₅ O ₁₂	<0.1	>40	−0.06	−0.068
(BiTb) ₃ (FeGa) ₅ O ₁₂	<0.1	>40	−0.06	−0.054
(BiRE) ₃ (FeGa) ₅ O ₁₂	<0.1	>40	−0.09	−0.054

their high absorption at shorter wavelengths. Compared to YIG, the Verdet constant of the BIG is typically more than five times larger; therefore, a compact isolator can be made using the BIG crystals. All these materials usually need an external magnet to be functional as a Faraday rotator. Recently, however, a premagnetized garnet [also call latching garnet, (BiRE)₃(FeGa)₅O₁₂] crystal has been developed which eliminates the use of an external magnet to provide further potential benefit of reducing the overall size of an isolator. Table I shows the typical parameters of various Faraday materials at wavelength around 1550 nm. As can be seen, most of the materials have a similar insertion loss and extinction ratio, but the temperature and wavelength dependence of the Faraday rotation coefficient is different.

B. Principle of Polarization-Dependent (Free Space) Optical Isolators

The principle of a polarization-dependent optical isolator is shown in Fig. 1, which comprises a polarizer, an analyzer, and a Faraday rotator for providing a 45-degree polarization rotation. Direction of the analyzer is oriented 45 degrees relative to that of the polarizer.

In the forward direction (Fig. 1a), an input light is launched into the polarizer to provide a linear vertically polarized light beam (s-polarization). The polarization direction of the linearly polarized beam is then rotated −45 degrees, and the beam is passed through the analyzer without any attenuation. To minimize the loss in the forward direction, a linearly polarized light with the polarization plane matching that of the polarizer should be launched into the isolator, so that the input light can pass through all elements in the isolator with minimum energy loss. In this case, the insertion loss is determined by the absorption loss of each element, the loss due to reflection at each interface of the elements, and imperfection of the materials. These types of optical isolators are developed

for use directly with lasers taking advantage of polarized output characteristics of the laser. Polarization-dependent optical isolators are commonly placed inside a semiconductor laser package and directly coupled to the laser chip in free space. Therefore, this type of optical isolator is also called a free space isolator.

In the backward direction (Fig. 1b), when a light is launched into the isolator, most of the light is blocked by the analyzer, except for the light with a polarization direction of −45 degrees, which will pass through the analyzer. The polarization of this beam is then rotated −45 degrees again by the Faraday rotator and becomes a horizontally

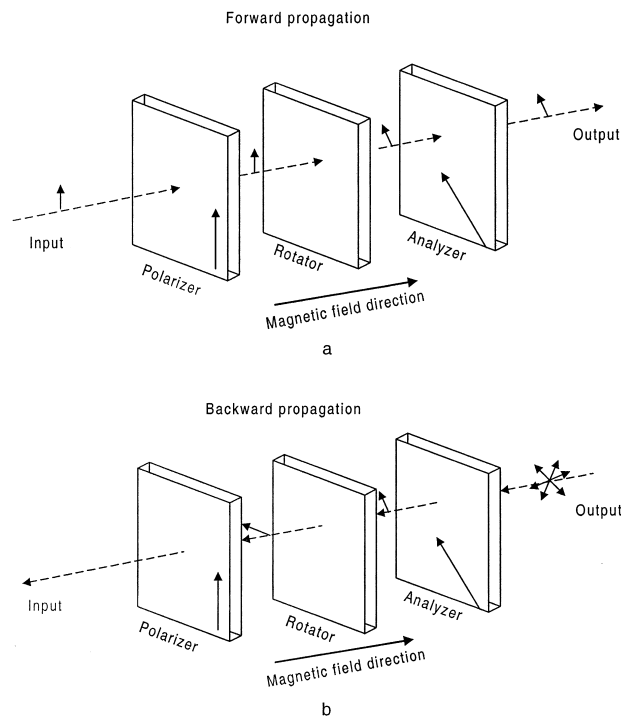


FIGURE 1 Operating principles of a single-stage, polarization-dependent optical isolator: (a) operation in the forward direction and (b) operation in the backward direction.

polarized beam (p-polarization). Since the polarization direction of this beam is 90 degrees to the direction of the polarizer, this beam is blocked by the polarizer and cannot pass through. Therefore, in the backward direction theoretically no light can pass through the optical isolator. However, in reality there is a very small amount of light passing through the isolator due to imperfection of the elements (i.e., limited extinction ratio of the polarizer and analyzer and errors in Faraday rotation angle).

Various designs of the polarization-dependent optical isolators have been proposed and developed depending on the selection of materials and structure of the polarizers and analyzers. Commonly used polarizers and analyzers are polarization prisms (polarizing beam splitter cubes, Glan-Taylor prism, etc.), plate polarizers (Polacor, LAMIPOL, etc.), and birefringent (calcite, rutile, and YVO_4) crystals.

C. Principle of Polarization-Independent (In-line) Optical Isolators

It is well known that when a light with a linear polarization plane propagates through an optical fiber, its state of polarization will not be maintained and it varies with time due to imperfection (birefringence) of the fiber. Therefore, an optical device with performances independent of the polarization variation is desirable when used in a fiber optic system. A polarization-independent optical isolator is developed to fit this need. Polarization-independent optical isolators are normally pigtailed with input and output fibers and are also called in-line-type isolators.

Operation of the polarization-independent optical isolators is based on two principles; one is based on polarization splitting and recombining, and the other is based on the anti-symmetric field conversion.

Polarization splitting and recombining based optical isolators use birefringent crystals to split an incoming beam into two beams with orthogonal polarization and then recombine the two beams after passing through a Faraday rotator. Two typical configurations of the polarization-independent optical isolators are shown in Figs. 2 and 3. The main difference is in the structure of the birefringent material. One uses wedge-shaped and the other uses plate-shaped birefringent material.

As shown in Fig. 2, two birefringent wedges with optical axes 45 degrees to each other are used for polarization splitting and combining, and a Faraday rotator for 45-degree rotation is inserted between the two wedges. In operation, in the forward direction (Fig. 2a), a light beam with an arbitrary polarization is launched to a collimating lens through an input fiber. After passing through the first birefringent wedge, the collimated beam is split into two beams with orthogonal polarization states. The optical

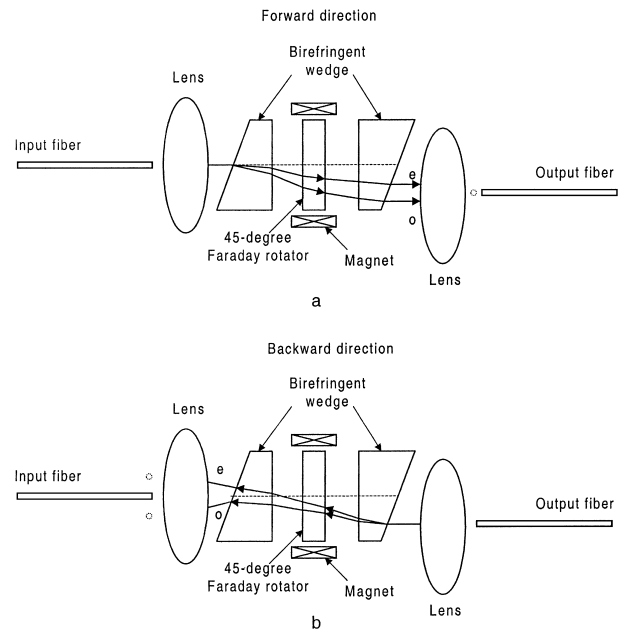


FIGURE 2 Design and principles of a birefringent wedge-type, single-stage, polarization-independent optical isolator: (a) operation in the forward direction and (b) operation in the backward direction.

axis of the first birefringent wedge is selected in such a way that the obtained two beams exit the wedge with two different angles. The Faraday rotator rotates the polarization of both beams by 45 degrees, and the polarization

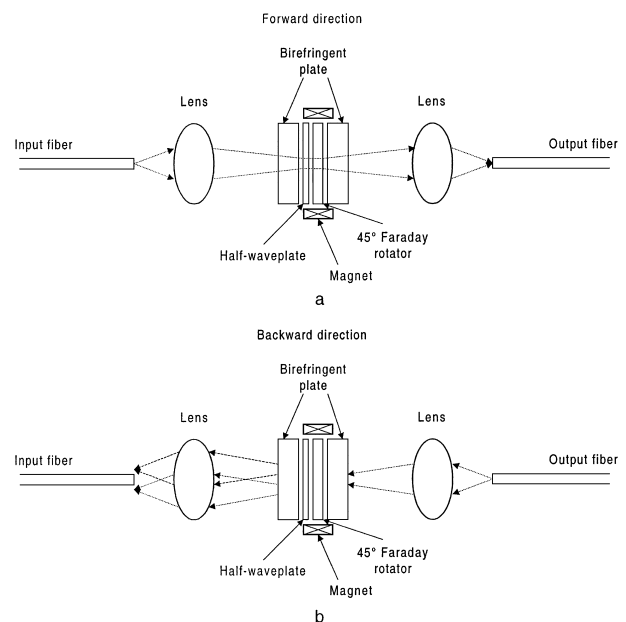


FIGURE 3 Design and principles of a parallel birefringent plate-type, single-stage, polarization-independent optical isolator: (a) operation in the forward direction and (b) operation in the backward direction.

rotated beams are passed through the second birefringent wedge. The optical axis of the second birefringent wedge is selected such that the angle between the two beams exiting the wedge becomes the same again. Finally, the two parallel beams are passed through a focusing lens and coupled to the output fiber. Since the two beams have the same angle and the displacement between the two beams is very small, the two beams can be coupled to the output fiber with very high efficiency (i.e., with very small insertion loss).

In the backward direction (Fig. 2b), when any light beam with an arbitrary polarization is launched into the output fiber, the light beam is collimated by the output lens. After passing through the output birefringent wedge, the collimated beam is split into two beams with orthogonal polarization states. The Faraday rotator rotates the polarization of both beams by 45 degrees, and the polarization rotated beams are passed through the input birefringent wedge. Due to the nonreciprocity of the Faraday rotation, the polarization states of these two beams are exchanged compared to their original states when the light beam is launched from the input port. Therefore, after passing through the input birefringent wedge, the angular difference between the two beams is further increased. The input lens converts this angular difference into a position separation, and the two beams are focused to spatial positions away from the input fiber. Therefore, the backward travelling light cannot be coupled into the input fiber. The isolation of these types of isolators depends on the angle of the wedges and the performance of the Faraday rotator. Birefringent materials such as lithium-niobate, rutile, calcite, and YVO_4 are commonly used for the wedge.

Another type of polarization splitting and combining based optical isolators is based on parallel birefringent plates. In this type of isolator, a light beam is split into two parallel beams and recombined into a single beam again after passing through a Faraday rotator. Therefore, the isolator itself is a self-complete unit, and this type of isolator is normally placed near the focal point of a lens instead of in a collimated beam. There are many different designs based on this method, and one of the designs is shown in Fig. 3. In this design, two birefringent parallel plates, a Faraday rotator and a half-waveplate, are used, and the optic axis of the half-waveplate is oriented at 22.5 degrees to the y-axis. The isolator is placed at the back focal point of the lenses. The operating principle is shown in Fig. 4, where circles show the positions of the light beams and arrows inside circles show the polarization state of each beam.

In the forward direction (Fig. 4a), the input beam is split into two beams with orthogonal polarization by the first birefringent plate. The polarization direction of the

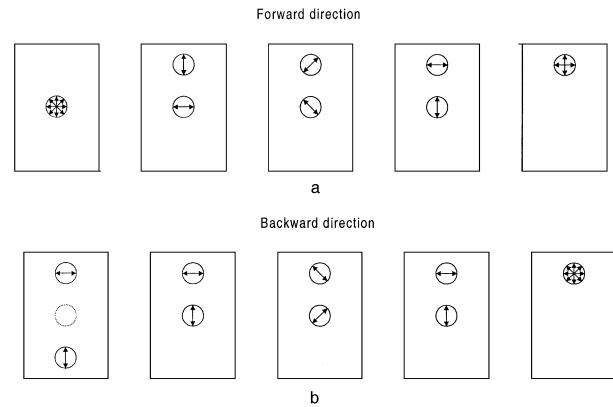


FIGURE 4 Detailed operation principles of the isolator shown in Fig. 4: (a) operation in the forward direction and (b) operation in the backward direction, where circles show beam positions and arrows in the circle show the polarization direction of the beams.

two beams is rotated by 45 degrees after passing through the half-waveplate and another 45 degrees by the Faraday rotator. Since the beam splitting direction of the second birefringent plate is designed to be the opposite of that of the first plate, the polarization rotated two beams are recombined completely by the second birefringent plate without any loss.

In the backward direction (Fig. 4b), a light beam is split into two beams with orthogonal polarization by the second birefringent plate, and their polarization directions are rotated by the Faraday rotator and the half-waveplate. Due to the nonreciprocity of the Faraday rotation, the polarization directions of the two beams reaching the first birefringent plate become opposite to the beam splitting direction of the plate; therefore, the two beams are further displaced away from the input fiber. If the displacement is large enough, any backward propagating light will not be coupled into the input fiber, thus providing good isolation. The performances of these types of isolators depend on the beam size, displacement, and characteristics of the Faraday rotator and waveplate. Birefringent materials such as rutile, calcite, and YVO_4 are commonly used for the plate due to their large birefringence.

Another type of the polarization-independent optical isolator is based on anti-symmetric field conversion (or multibeam interference). The schematic diagram of a typical structure of this type of isolator is shown in Fig. 5a. The main feature of this type of isolator is the elimination of the use of birefringent materials, which could result in potential cost reduction. In the structure shown in Fig. 5a, a collimating lens is used to collimate a light beam launched from the input fiber, and a focusing lens is used to couple the collimated beam back into the output fiber. The collimated beam is divided into two sub-beams; one sub-beam

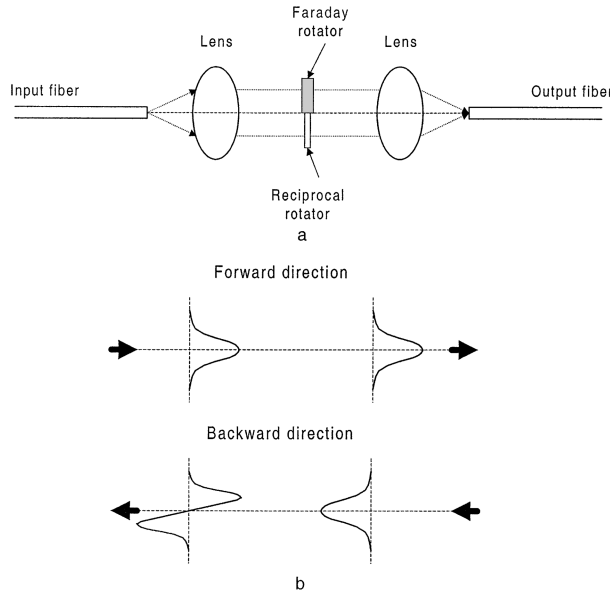


FIGURE 5 Design and principle of an anti-symmetric field conversion-type, single-stage, polarization-independent optical isolator: (a) design and (b) operation principle.

is passed through a Faraday rotator for 45-degree rotation, and the other sub-beam is passed through a reciprocal rotator such as a half-waveplate. The angle of rotation of the Faraday rotator and half-waveplate is selected so that in the forward direction the phase difference between the two sub-beams is zero and in the backward direction the phase difference is π due to the nonreciprocal Faraday rotation. Therefore, in the forward direction the two sub-beams form a symmetric field distribution at the output fiber and can be propagated in the fiber without any loss. However, in the backward direction, due to the out of phase relation between the two beams, an anti-symmetric field is formed at the end of the input fiber as shown in Fig. 5b. Since this anti-symmetric field corresponds to a higher mode of the single-mode fiber and cannot propagate through the single-mode fiber, the backward light will be leaked as cladding mode light. Therefore, these types of isolators are only functional as an optical isolator when used together with single-mode optical fibers. Diffraction loss is the main cause of the insertion loss of these types of isolators. To minimize the loss, good edge quality and large diameter beams are desirable.

The anti-symmetric field conversion isolator can also be realized in waveguide form as shown in Fig. 6. In operation, light launched into the input fiber is split into two beams and guided into two different waveguides by the Y-branching waveguide. In each waveguide, light is passed through a half-waveplate and a Faraday waveguide and then the two beams are combined by another Y-branching waveguide. An isolator can be realized by optimizing the

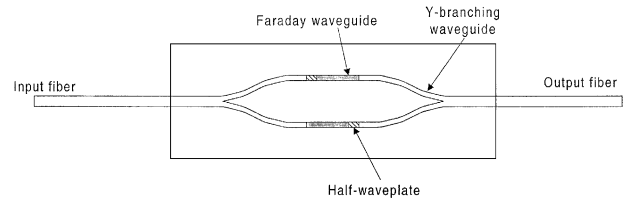


FIGURE 6 Schematic diagram of a waveguide-based, anti-symmetric field conversion-type optical isolator.

phase between two beams. The potential advantages of the waveguide device are lower manufacturing cost, compact size, and consistent performance. However, generally the coupling loss from the waveguide to the optical fiber is still high, and it is also difficult to integrate the Faraday material into the waveguide. Therefore, waveguide-type optical isolators have not yet been commercialized.

D. Performance Evaluation of Optical Isolators

Key performance parameters for all types of optical isolators are insertion loss and isolation, and in addition, for the polarization-independent optical isolators, polarization-dependent loss (PDL) and polarization mode dispersion (PMD) are important parameters as well. Performances (insertion loss and isolation) of an optical isolator can be estimated by using the Jones matrix. For example, insertion loss and isolation of a polarization-dependent isolator can be expressed as

$$\text{Insertion loss} = 10 \log \left(\frac{\eta \cos^2 \theta}{(1 + 10^{-ER})} \right) \quad (2)$$

$$\text{Isolation} = 10 \log \left(\frac{10^{-ER} \cos^2 \theta + \sin^2 \theta}{(1 + 10^{-ER})} \right) \quad (3)$$

where η is the transmittance of all materials; ER is the extinction ratio of the polarizer, analyzer, and Faraday rotator in decibel units; and θ is the rotation angle error of the Faraday rotator from the ideal rotation angle of 45 degrees. The extinction ratio is determined by material itself, and the rotation angle error is mainly caused by dimension error (length error) of Faraday rotator and environmental (temperature and wavelength) change during operation of the isolator since the Verdet constant of a Faraday material is normally a function of temperature and wavelength. Typical insertion loss and isolation of a polarization-dependent isolator is shown in Fig. 7 as a function of the Faraday rotation angle error. It is found that the insertion loss variation due to the rotation angle error is very small, and the peak isolation is determined by the extinction ratio of materials and decreases with the angle error. Therefore, to achieve certain isolation, the operating condition (wavelength range or temperature) is

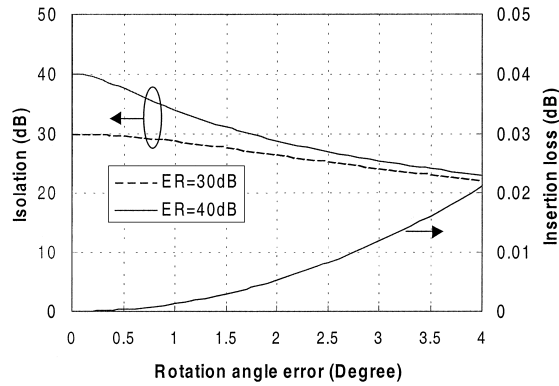


FIGURE 7 Typical isolation and insertion loss of a single-stage optical isolator as a function of rotation angle error of the Faraday rotator.

limited. For a given operating wavelength or temperature range, the achievable isolation can be estimated by using the parameters in Table I and Eq. (3).

The main causes of the insertion loss in a polarization-dependent isolator are absorption loss of materials and the Fresnel reflection loss at each interface. To achieve lower insertion loss, it is important not only to use low absorption material, but also to reduce reflection at each interface.

For polarization-independent isolators, besides the insertion loss and isolation, PDL and PMD are very important parameters. An ideal polarization-independent isolator should have zero PDL and PMD; however, in practice, PDL and PMD of a polarization-independent isolator are not zero due to assembly processes and dimensional tolerances of materials. In the polarization-independent isolators using birefringent wedges (Fig. 2), since the two orthogonal polarized beams are not combined before being launched into the output fiber, the PDL of the device is much more sensitive to the relative position among the two beams, lens, and fibers and to the environmental condition (temperature) changes. The PMD of these types of isolators is also high and in the range of 0.3 ps, however, the PMD can be compensated by utilizing additional birefringent plates to equalize the path length of the two orthogonal polarized beams. On the other hand, in the polarization-independent isolators configured with parallel plates (Figs. 3 and 4), the two orthogonal polarized beams are completely recombined before being passed through the lens and the output fiber. Therefore, its performance is much more stable, and the PMD can also be minimized easily (typically less than 0.01 ps).

The insertion loss of a polarization-independent (or inline) optical isolator is dominated by the coupling loss from input to output fibers, and the coupling loss could account for more than 70% of the overall insertion loss. Therefore, it is very important to achieve the best possible coupling between fibers. The coupling loss is generally

determined by the quality of the lens used for the coupling and is typically in the range of 0.2 dB.

E. Multistage Optical Isolators

Isolators described above are so-called single-stage isolators that have a basic configuration of a Faraday rotator sandwiched by two polarizing elements. As shown in Fig. 7, the peak isolation of a typical single-stage isolator is about 40 dB, and the overall isolation is less than 30 dB with a rotation angle error of 2 degrees, which can be easily introduced due to environmental changes. In some applications, a much higher isolation is desired. Optical isolators can be used in series (cascaded) to increase the isolation. For example, the isolation of a two-stage isolator generally is twice that of a single-stage isolator. A multistage isolator can be constructed by simply cascading multiple single-stage isolators or by optimizing configurations that use fewer materials. To construct a two-stage optical isolator, the required minimum materials are three polarizers (birefringent crystals) and two Faraday rotators, and they can be configured by inserting the Faraday rotators between the polarizers (birefringent crystals).

An example of a two-stage, polarization-dependent optical isolator is shown in Fig. 8. In the forward direction operation (Fig. 8a), an input light is launched into the polarizer to provide a linear and vertical polarized light beam (s-polarization). The polarization direction of the linearly polarized beam is then rotated -45 degrees by the first Faraday rotator, and the beam is passed through the first analyzer. After the first analyzer, the light beam is further passed through a second Faraday rotator for another -45 -degree rotation and then through a second analyzer oriented at the horizontal direction without any attenuation.

In the backward direction operation (Fig. 8b), when a light is launched into the isolator, most of the light is blocked by the output analyzer, except for the light with a horizontal polarization direction, which will pass through the analyzer. The polarization of this beam is then rotated -45 degrees again by the Faraday rotator and becomes a 45 -degree polarized beam. Since the polarization direction of this beam is 90 degrees to the direction of the first analyzer, this beam is blocked by the first analyzer and only a small amount of light can pass through due to the limited extinction ratio of the analyzer and the rotation angle error caused by the imperfection of the Faraday rotator. The polarization of this small amount of light is further rotated by the first Faraday rotator and becomes a horizontally polarized light. Since the polarization direction of this beam is 90 degrees to the direction of the input polarizer, this beam is further blocked by the input polarizer. Therefore, in the backward direction of a two-stage

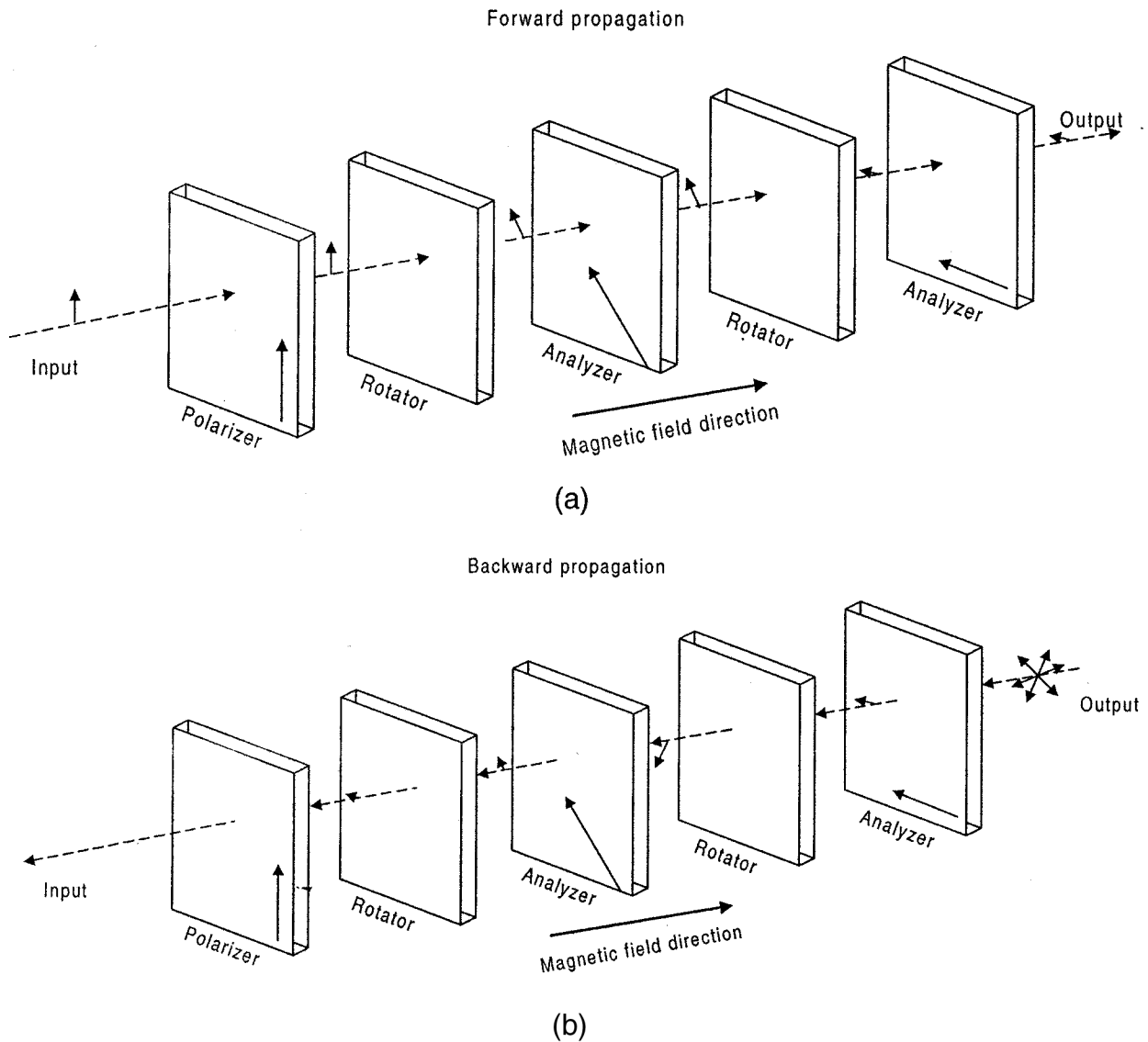


FIGURE 8 Operating principles of a two-stage, polarization-dependent optical isolator: (a) operation in the forward direction and (b) operation in the backward direction.

isolator the light is blocked twice by the polarizers, resulting in higher isolation.

The operating principle of a multistage, polarization-independent optical isolator is similar to that of the multistage, polarization-dependent isolator, except for the polarization splitting and recombining portion. Various designs of two-stage, polarization-independent isolators have been developed, and one example is shown in Fig. 9. In Fig. 9, open arrows in the blocks show the direction of beam displacement of the birefringent crystals and solid arrows show the Faraday rotation direction. The insertion loss of a multistage isolator is also higher than that of a single-stage isolator. However, the increase is minimal

in the case of the polarization-independent isolators since the dominant insertion loss in a polarization-independent isolator is the coupling loss rather than loss caused by materials.

II. PRINCIPLE OF OPTICAL CIRCULATORS

The first optical circulator was designed in the early 1980s and was primarily for use in bidirectional optical communication systems to increase the capacity of the systems. In the traditional bidirectional optical communication system, a 50/50 (3-dB) coupler, that splits a light beam into

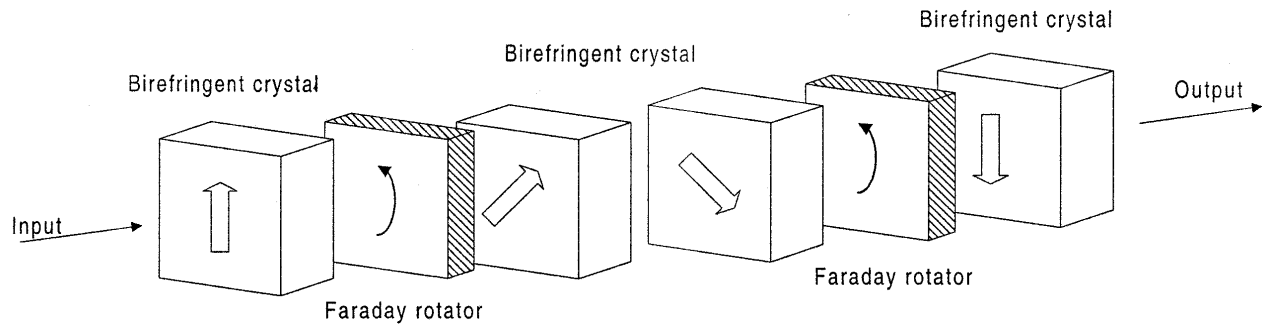
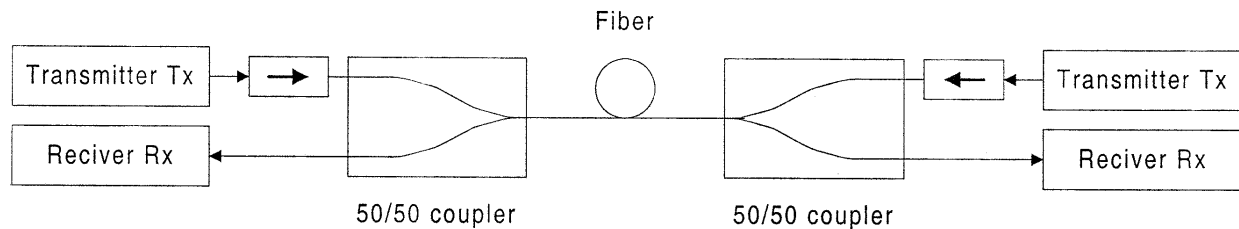


FIGURE 9 Design of a parallel, birefringent plate-type, two-stage, polarization-independent optical isolator.

two beams with equal intensity, was used to couple the transmitters and receivers as shown in Fig. 10a. However, there are two main problems with this kind of structure. One is the need for an optical isolator in the transmitters to prevent light crosstalk between the transmitters, and the other is the high insertion loss associated with the use of the 50/50 coupler, since two couplers have to be used and each has a minimum loss of 3 dB, which results in a minimum 6-dB reduction of the link budget from the system. The use of an optical circulator can solve both of the problems since an optical circulator can provide the isolation function as well as a loss of less than 3 dB (Fig. 10b). With the improvement in optical circulator designs together with the rapid growth in the optical communication

industry, especially the recent vast deployment of dense wavelength division-multiplexing (DWDM) systems, optical circulators have found many applications not only in the telecommunication field, but also in medical and energy industries as well. Optical circulators are being used in optical amplifiers, in bidirectional optical systems, in optical DWDM systems as an Add/Drop device or demultiplexing device together with the fiber Bragg gratings, and in medical instruments for improving imaging, etc.

Similar to optical isolators, optical circulators can also be divided into two categories, that is, polarization-dependent and -independent optical circulators. However, since the initial application of an optical circulator is intended for use with optical fibers, the majority of optical



(a)

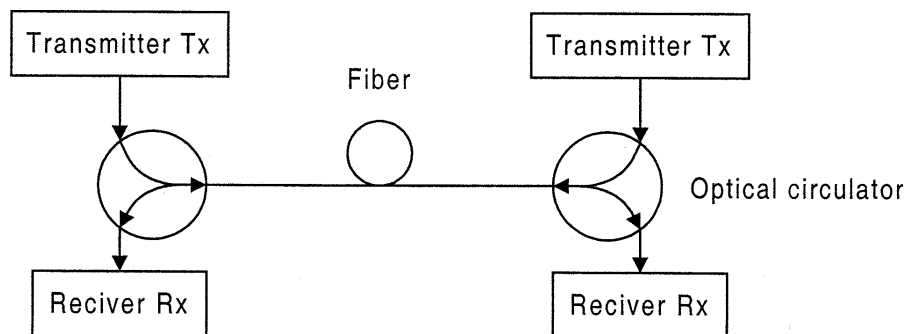


FIGURE 10 A typical bidirectional optical communication system using 50/50 couplers.

circulators are designed for polarization-independent operation, and there is little effort on designing polarization-dependent circulators. Optical circulators can be divided into two groups based on their function: one is the full circulator, in which light passes through all ports in a complete circle (i.e., light from the last port is transmitted back to the first port), and the other is the quasi-circulator, in which light passes through all ports sequentially but light from the last port is blocked and cannot be transmitted to the first port. For example, in the case of a full three-port circulator, light passes through from port 1 to port 2, port 2 to port 3, and port 3 back to port 1. However, in a quasi-three-port circulator, light passes through from port 1 to port 2 and port 2 to port 3, but any light from port 3 will be blocked and cannot be propagated back to port 1. In most of the applications, only a quasi-circulator is required.

Key performance parameters of optical circulators include insertion loss, isolation, PDL, PMD, and return loss. Although performances of optical circulators depend on the particular design and materials, typically an insertion loss of 1 dB, an isolation of more than 40 dB, a PDL of less than 0.1 dB, and a PMD of less than 0.1 ps are achievable regardless of the design. The main differences among various designs are performance stability, size, and cost. The isolation of an optical circulator can also be increased by utilization of multiple stages. However, unlike optical isolators, a multistage circulator cannot be achieved by simply cascading multiple circulators. A multistage circulator has to be designed differently from the beginning. Performances of an optical circulator can also be estimated by using the Jones matrix and Eqs. (2) and (3).

Polarization beam splitters based on dielectric coatings were used in early stage circulators. However, the isolation of those optical circulators was relatively low due to the limited extinction ratio (around 20 dB) of the polarization beam splitters. Various designs using birefringent crystals have been proposed to increase the isolation. However, the use of the birefringent crystals generally results in the increase of size and cost of the circulator. Extensive development efforts have been concentrated on improvement of various designs. To construct a polarization-independent optical circulator, a minimum of three functional elements is required, that is, a polarization beam splitting and combining element, a nonreciprocal polarization rotation element, and a beam shifting element.

A. Principle of a Polarization-Independent Quasi-circulator

There are many designs of polarization-independent optical quasi-circulators. However, most of them follow a common design concept, which uses the three basic elements mentioned above. One example of two-stage, three-

port, polarization-independent quasi-circulators is shown in Fig. 11, where (a) shows the schematic diagram of the design and (b) shows the operating principle. In this particular design, birefringent crystals are used as both the polarization splitting and combining element and the beam shifting element, and Faraday rotators and half-waveplates are used as the polarization rotation element. In the operation from port 1 to port 2, a light beam launched into port 1 is first split into two beams with orthogonal polarization, then the upper beam receives a 45-degree clockwise rotation by a half-waveplate and the lower beam receives a 45-degree counterclockwise rotation by another half-waveplate, so that after passing through the waveplates the polarization of the two beams becomes the same. Polarization of the two beams is further rotated by the Faraday rotator and becomes vertically polarized beams so that they match the ordinary state of the beam shifting birefringent crystal and are unaffected rotator. Both beams become vertically polarized so that therefore, the positions of the two beams are unchanged and can be recombined by the crystal and coupled into port 2 after passing through the second Faraday rotator and waveplates (Fig. 11b).

In the operation from port 2 to port 3, a light beam launched into port 2 is first split into two beams with orthogonal polarization, then the upper beam receives a 45-degree clockwise rotation by a half-waveplate and the lower beam receives a 45-degree counterclockwise rotation by another half-waveplate, so that after passing through the waveplates the polarization of the two beams becomes the same. Polarization of the two beams is further rotated by the Faraday rotator and becomes horizontally polarized beams due to the nonreciprocal rotation of the Faraday rotator, which matches the extraordinary ray of the beam shifting birefringent crystal; therefore, both beams are shifted horizontally to new locations. After passing through the Faraday rotator and waveplates, the two beams are then recombined by the crystal at a location away from port 1 and, therefore, can be coupled to an additional port (port 3, Fig. 11b). Multiport circulators with more than three ports can be configured using the same design by extending ports laterally. There are many modified designs by utilizing different combinations of polarization rotators and birefringent crystals. An advantage of these types of designs is its high isolation. However, the size of the device is relatively large and the cost is also higher due to the use of a large number of different elements.

B. Principle of a Polarization-Independent Full Circulator

The first polarization-independent optical circulator was designed using only polarization beam splitters and polarization rotators and also featured full circulation function.

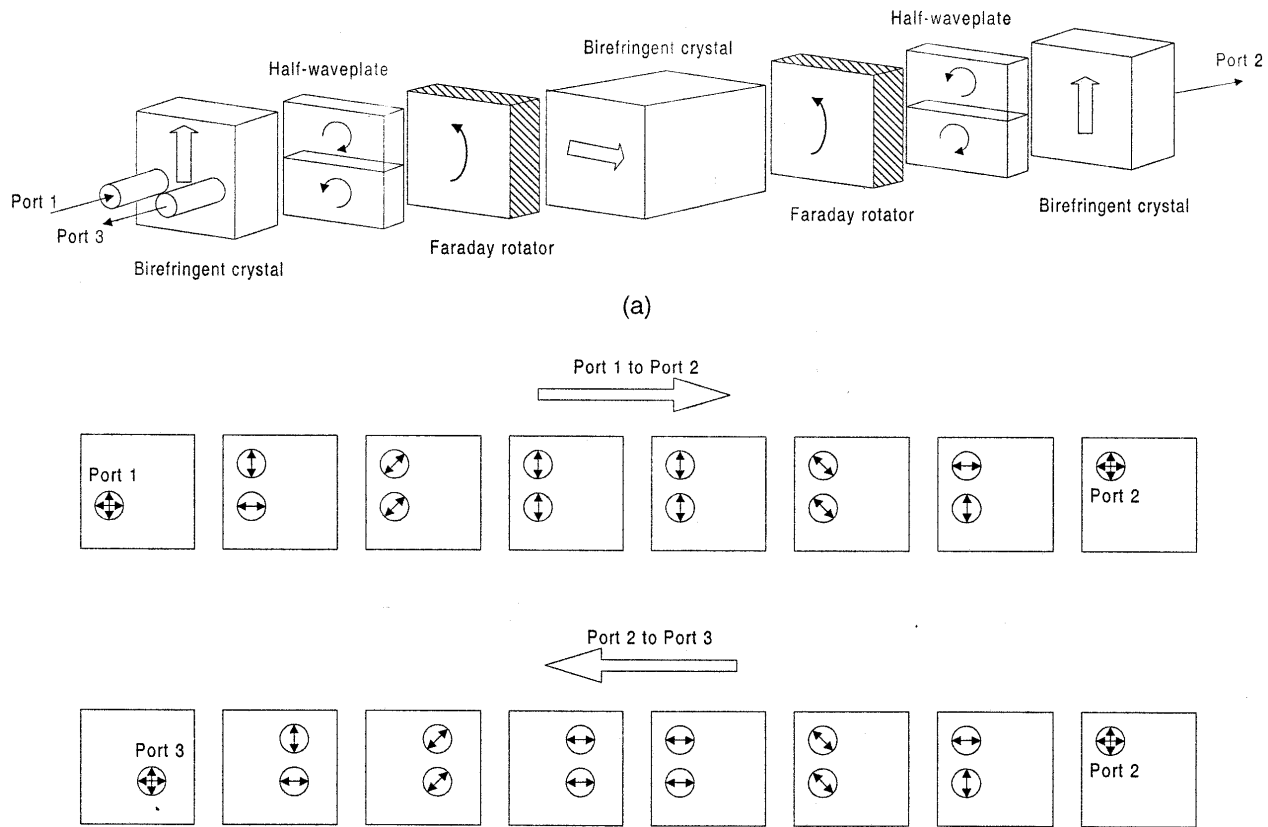


FIGURE 11 Design and principle of a quasi-circulator using birefringent crystals: (a) design and (b) operating principle where circles show beam positions and arrows in the circle show polarization direction of the beams.

The schematic diagram of that circulator is shown in Fig. 12. In operation, light launched into port 1 is split by a polarization beam splitter into two beams with orthogonal polarization, and one beam is reflected by a mirror. The two beams are then passed through a Faraday rotator and a half-waveplate and receive a 90-degree polarization rotation. The two beams are then recombined into port 2 by another polarization beam splitter. Similarly, light launched into port 2 is split and passed through the same rotators. However, due to the nonreciprocal rotation of the Faraday rotator, the polarization of two beams remains unchanged. Therefore, the two beams are recombined by the polarization beam splitter at the location different than port 1 and are coupled into port 3. The same principle applies to the propagation from port 3 to port 4 and port 4 to port 1. This design is very simple and easy to make, but it has a drawback of low isolation due to low extinction ratio of the beam splitter, and also the maximum achievable port numbers are limited to four. It is also difficult to construct a multistage circulator using this design. To increase the isolation various improvements have been proposed, but all of them involve the use of birefringent crystals and become very complicated.

C. Design and Principle of Compact Optical Circulators

Although optical circulators bring significant design advantages in the optical communication system, it was difficult to realize an optical circulator that meets the performance and reliability requirements of the optical communication industry due to the complicity of its design and the vast usage of different crystals. Its relatively high cost and large size also limited its wide application. However, in recent years the popularity of the internet has significantly increased the demand for more bandwidth, resulting in the rapid growth in optical communication systems. Demand for development of low cost and compact optical circulators has been stronger than ever.

In the conventional optical circulator, light from each fiber is first collimated and then launched into the circulator. Therefore, there are limitations of how close each port can be positioned, which in turn limits the minimum size of crystals, since the length of the crystals is directly proportional to the distance between two ports. The efforts of developing low cost and compact optical circulators have

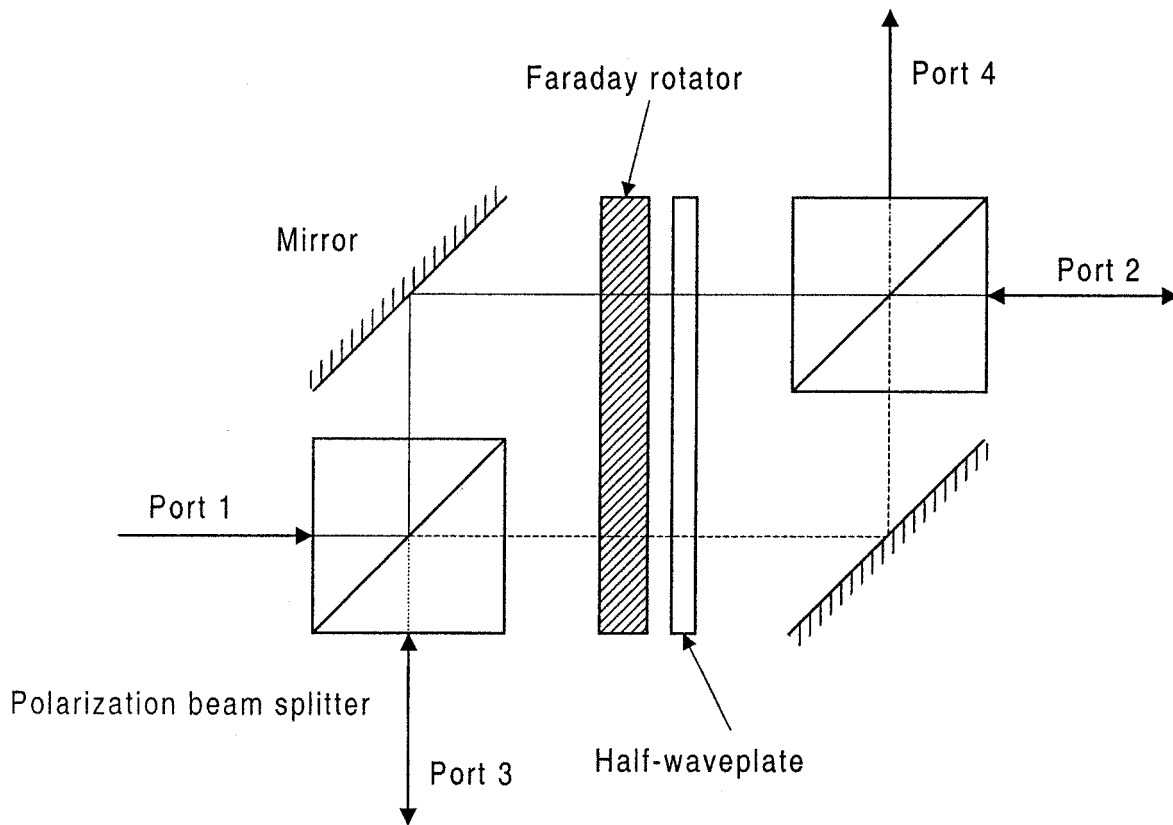


FIGURE 12 Schematic diagram of a polarization-independent full circulator.

been concentrated on significant reduction of material usage and simplifying optical fiber coupling. To reduce the crystal size, the two fibers must be as close as possible? Two types of low cost and compact circulators have been commercialized. One uses the angle deviation method, and the other uses the noncollimating displacement method. Both of them use the dual fiber collimating method for the fiber coupling in which two fibers can be placed in contact.

In the angle deviation method shown in Fig. 13, light from each port is still collimated first and the circulator is still placed in the collimated beam. However, instead of physically displacing a beam by using a birefringent crystal as in the conventional circulator design, a Wollaston prism is used to introduce different angles for different po-

larization states of the beam. During the operation, when a light is launched into port 1, it is collimated and then split into two beams with orthogonal polarization. The polarization state of the two beams is properly rotated to the same direction. After passing through the Wollaston prism, the two beams receive a first angle deviation and are then combined and coupled into port 2. In the operation of port 1 to port 2, the light launched into port 2 is split into two beams and properly polarization rotated. After passing through the Wollaston prism, the two beams receive another angle deviation which is the same angle as the first one but in the opposite direction due to the nonreciprocal rotation of the Faraday rotator. The collimating lens converts this angle into a displacement at a place different than port 1, and by properly designing the Wollaston prism and

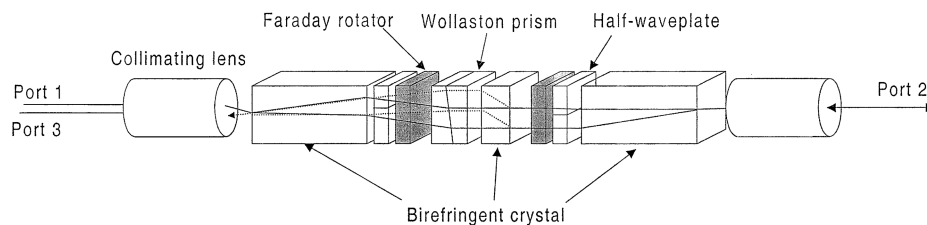


FIGURE 13 Schematic diagram of an angle deviation-type, compact, polarization-independent optical circulator.

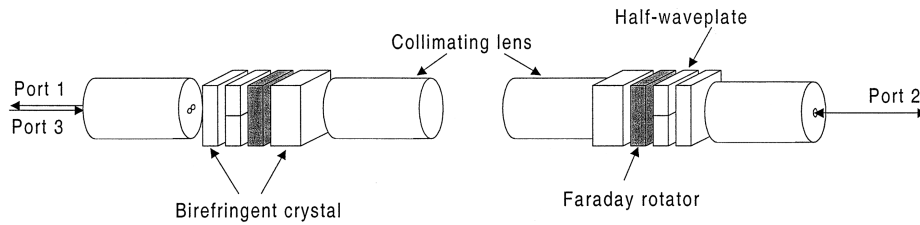


FIGURE 14 Schematic diagram of a noncollimating beam displacement-type, compact, polarization-independent optical circulator.

the fiber separation between the port 1 and port 3, the light launched into port 2 can be efficiently coupled into port 3.

In the noncollimating displacement method shown in Fig. 14, the circulator design itself is exactly the same as that of the conventional circulator, except that the circulator is placed in a noncollimating beam instead of a collimating beam. In the conventional design, in order to couple a light beam into different fibers, a minimum displacement equal to the diameter of a lens is required, which is generally a couple of millimeters. However, in the noncollimating displacement, fibers can be placed side by side and the minimum displacement required is only equivalent to the diameter of the fiber, which is only $125\ \mu\text{m}$. Therefore, by placing the circulator into the noncollimated beam, material reduction can be achieved by a factor of more than 10. During the operation, light launched into port 1 is split into two beams and their polarization directions are properly rotated. The two beams are then directed to the collimating lens without position displacement. After passing through another lens the two beams are polarization rotated again and recombined and coupled into port 2. Similarly, light launched into port 2 is split and polarization rotated. However, due to the nonreciprocal rotation of the Faraday rotator, this time beams are displaced in position for an amount equal to the separation between port 1 and port 3. Therefore, after passing through the lenses and rotators, the beams are combined and coupled into port 3.

Since the displacement angle of commonly available birefringent crystals is very close to the numerical aperture (NA) of a normal single-mode fiber, to achieve a complete beam separation for the diverging beam from the fiber it is necessary to reduce the NA of the fiber. Therefore, in the noncollimating circulator design, mode-field expanded fibers are typically used and the mode-field expansion is typically achieved by heating the fiber to a high temperature, hence causing thermal diffusion of the core dopant material into the cladding.

D. Other Optical Circulator Designs

Optical circulators can also be realized using the anti-symmetric field conversion or using polarization-main-

taining (PM) fiber splitters. The principles are very similar to that of isolators, except for the utilization of the 180-degree out of phase property of the couplers. A typical design of the anti-symmetric field conversion-type circulator is shown in Fig. 15 in a waveguide form. An optical circulator can also be formed by replacing the X-branching waveguide in Fig. 15 with a PM fused splitter, which is made by fusing two PM fibers together and by adjusting the coupling between two fibers so that the light beams with orthogonal polarization are propagated into different fibers.

III. REFLECTIVE DESIGN OF OPTICAL ISOLATORS AND CIRCULATORS

Optical isolators and circulators mentioned in Sections I and II are so-called transmissive devices; that is, the light is propagated along one direction and the input and output ports are on the opposite side of the devices. To further simplify design and reduce materials usage and cost, a reflective design concept has been introduced based on the fact that most of the transmissive designs have an image plane and all elements are symmetric in respect of the image plane. In the reflective design, all ports of the device are coming out from one side, further providing the advantage of easy installation in the applications. Typical designs of single-stage and two-stage reflective isolators are shown in Figs. 16 and 17, respectively. The operating principles of the reflective isolators are the same as

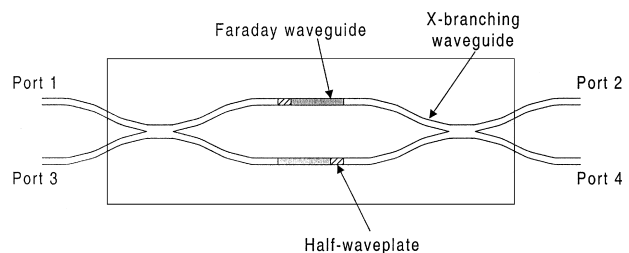


FIGURE 15 Schematic diagram of an anti-symmetric field conversion, waveguide-type, polarization-independent optical circulator.

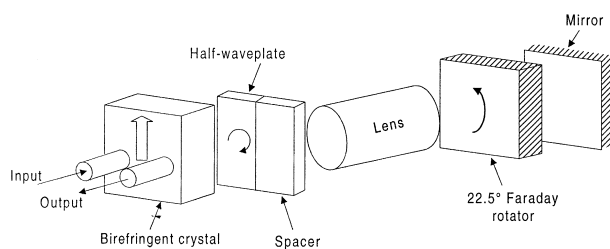


FIGURE 16 Schematic diagram of a reflective, single-stage, polarization-independent isolator.

those of the transmissive designs, except for the sharing of common elements and the utilization of a mirror for beam folding. In the single-stage reflective isolator, a Faraday rotator for 22.5-degree rotation is used instead of the common 45-degree Faraday rotator due to the fact that a beam is passed through the Faraday rotator twice, providing a total rotation of 45 degrees. With the reflective design, the material usage can be reduced almost by a factor of 2.

A typical design of a reflective optical circulator is shown in Fig. 18, where a 90-degree polarization reflector, consisting of a quarter-waveplate and a mirror, is used for

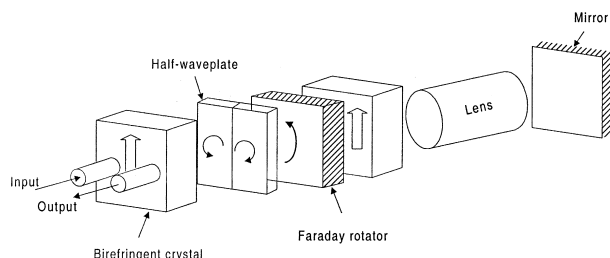


FIGURE 17 Schematic diagram of a reflective, two-stage, polarization-independent isolator.

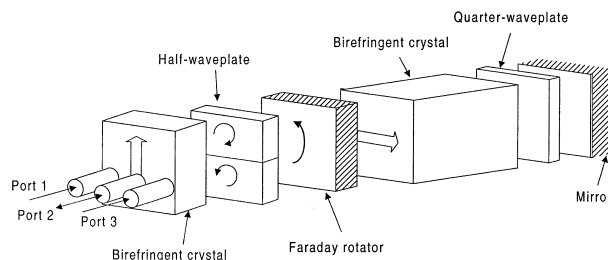


FIGURE 18 Schematic diagram of a reflective, polarization-independent optical circulator.

the beam folding. The operating principle of a reflective circulator is very similar to that of a transmissive circulator, except for the beam folding by the 90-degree polarization reflector. Compact circulators can also be designed in the reflective form and are currently being developed for commercialization.

SEE ALSO THE FOLLOWING ARTICLES

LASERS, SEMICONDUCTOR • OPTICAL AMPLIFIERS (SEMICONDUCTOR) • OPTICAL FIBER COMMUNICATIONS • OPTICAL WAVEGUIDES AND WAVEGUIDE DEVICES

BIBLIOGRAPHY

- Cheng, Y. (1996). "Passive hybrid optical components for fiber amplifiers," *N. Fiber Op. Eng. Conf. Proc.* **3**, 553–558.
- Huard, S. (1997). "Polarization of Light," Wiley, New York.
- Kashima, N. (1995). "Passive Optical Components for Optical Fiber Transmission," Artech House, Boston.
- Zhang, B., and Lu, L. (1998). "Isolators protect fiberoptic systems and optical amplifiers," *Laser Focus World* **November**, 147–152.



Optical Waveguides (Planar) and Waveguide Devices

Katsunari Okamoto

NTT Photonics Laboratories

- I. Waveguide Fabrication by Flame Hydrolysis Deposition
- II. $N \times N$ Star Coupler
- III. Arrayed Waveguide Grating
- IV. Optical Add/Drop Multiplexer
- V. $N \times N$ Matrix Switches
- VI. Lattice-Form Programmable Dispersion Equalizers
- VII. Hybrid Integration Technology Using PLC Platforms
- VIII. Conclusion

GLOSSARY

DFB (distributed feedback) laser Laser diode in which light feedback for oscillation is realized by diffraction grating.

DSF (dispersion shifted fiber) Single-mode optical fiber whose zero-dispersion wavelength is shifted from the normal $1.3 \hat{=} \frac{1}{4} \mu\text{m}$ to $1.55\text{--}1 \hat{=} \frac{1}{4} \mu\text{m}$ low-loss wavelength region.

FHD (flame hydrolysis deposition) Doped silica glass deposition method in which fine glass particles are produced in oxy-hydrogen torch.

LPCVD (low-pressure chemical vapor deposition) Doped silica glass deposition method in which glass

particles are produced by low pressure chemical vapor deposition.

PECVD (plasma-enhanced chemical vapor deposition)

Doped silica glass deposition method in which gases used for depositing films are dissociated by electron impact in a glow discharge plasma.

PLC (planar lightwave circuits) Optical waveguides fabricated on planar substrate.

RIE (reactive ion etching) Etching method in which chemically active ions are accelerated in electric field and used for perpendicular etching.

SOA (semiconductor optical amplifier) Laser diode in which light feedback is suppressed to achieve single-pass, high-speed light signal amplification.

TE (transversal electric) mode Electromagnetic field distribution in which electric field vector lies in the plane perpendicular to the propagation direction.

TM (transversal magnetic) mode Electromagnetic field distribution in which magnetic field vector lies in the plane perpendicular to the propagation direction.

WDM (wavelength division multiplexing) Optical transmission scheme in which multiple of lightwave signals are carried with different color of light.

PLANAR LIGHTWAVE CIRCUITS (PLCs) are constructed using various kinds of waveguide components as shown in Fig. 1. Figures 1(a)–(f) show various kinds of core ridge structures which are normally placed on the substrate material having a lower refractive index than that of the core. In order to reduce contamination and radiation loss, these core ridge structures are covered by the over-cladding material. Refractive-index of the over-cladding is also lower than that of the core. Bent and S-shaped waveguides are used to change the direction of the lightwaves. The tapered waveguide is used to change the width of the waveguide; Y-branch and crossing waveguides are used for splitting and combining lightwaves; and the directional coupler is used for the interference effect.

Silica-based PLCs are fabricated with various kinds of technologies. The substrate material is either silica or silicon. Several typical fabrication technologies are (1) flame hydrolysis deposition (FHD), (2) low-pressure chemical vapor deposition (LPCVD), and (3) plasma-enhanced chemical vapor deposition (PECVD). FHD will be de-

scribed in the following section. LPCVD involves steps like the ones used in microelectronics silicon processing. The waveguide is formed on a substrate by low-pressure chemical vapor deposition of an undoped thick buffer layer ($\sim 15 \mu\text{m}$), on which a core layer is deposited. PECVD is widely used in silicon integrated circuit applications. Unlike LPCVD, the gases used for depositing films are dissociated by electron impact in a glow discharge plasma. This allows deposition of relatively dense oxide films at low temperatures ($< 300^\circ\text{C}$). The mechanism of film formation in this technique is distinctly different from FHD and LPCVD, due to the highly energetic nature of the charged particles which impinge on the substrates and the growing films.

I. WAVEGUIDE FABRICATION BY FLAME HYDROLYSIS DEPOSITION

Planar lightwave circuits using silica-based optical waveguides are fabricated on silicon or silica substrate by a combination of FHD and reactive ion etching (RIE). In this technique, fine glass particles are produced in the oxy-hydrogen flame and deposited on substrates. After depositing under-cladding and core glass layers, the wafer is heated to high temperature ($> 1000^\circ\text{C}$) for consolidation. The circuit pattern is fabricated by photolithography and reactive ion etching. Then core ridge structures are covered with an over-cladding layer and consolidated again. The most prominent feature of the planar waveguides fabricated by FHD and RIE is their simple and well-defined waveguide structures (Okamoto, 2000). This allows us

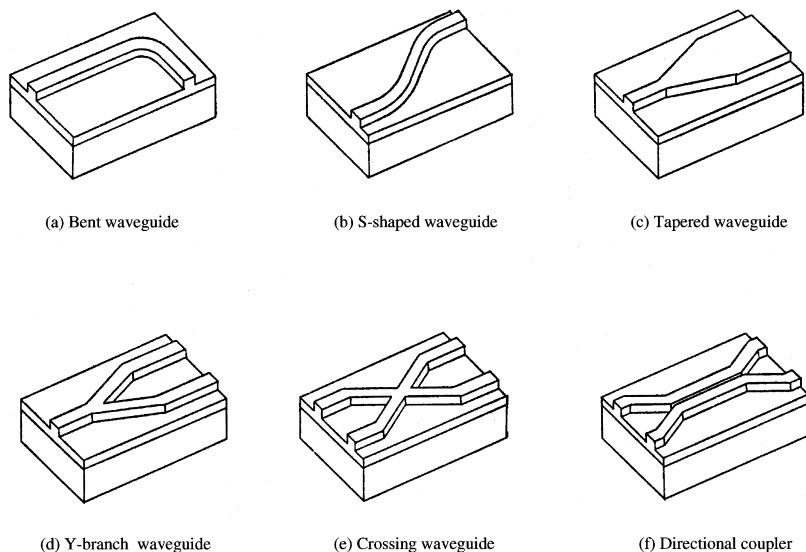


FIGURE 1 Various kinds of planar waveguide components.

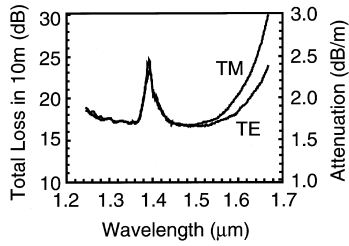


FIGURE 2 Loss spectra of 10-m long waveguide.

to fabricate multibeam or multistage interference devices such as arrayed-waveguide gratings and lattice-form programmable dispersion equalizers.

Figure 2 shows loss spectra of a 10-m-long waveguide with $\Delta = 0.45\%$ index difference ($R = 15$ mm) (Hida *et al.*, 1995). The higher loss for TM mode (electric field vector is perpendicular to the waveguide plane) maybe due to the waveguide wall roughness caused by the RIE etching process. However, the mode conversion from TE to TM or vice versa was less than -20 dB in 10-m-long waveguides. Various kinds of waveguides are utilized depending on the circuit configurations. Table I summarizes the waveguide parameters and propagation characteristics of four kinds of waveguides. The propagation losses of low- Δ and medium- Δ waveguides are about 0.01 dB/cm and those of high- Δ and super high- Δ waveguides are about 0.04~0.07 dB/cm, respectively. The low- Δ waveguides are superior to the high- Δ waveguides in terms of fiber coupling losses with the standard single-mode fibers. On the other hand, the minimum bending radii for high- Δ waveguides are much smaller than those for low- Δ waveguides. Therefore, high- Δ waveguides are indispensable for constructing highly integrated and large-scale optical circuits such as $N \times N$ star couplers, arrayed-waveguide grating multiplexers, and dispersion equalizers.

TABLE I Waveguide Parameters and Propagation Characteristics

	Low- Δ	Medium- Δ	High- Δ	Super high- Δ
Index difference (%)	0.3	0.45	0.75	1.5~2.0
Core size (μm)	8×8	7×7	6×6	4.5×4.5 $\sim 3 \times 3$
Loss (dB/cm)	<0.01	0.02	0.04	0.07
Coupling Loss ^a (dB/point)	<0.1	0.1	0.4	2.0
Bending radius ^b (mm)	25	15	5	2

^a Coupling loss with standard single-mode fiber.

^b Bending radius at which bending loss in a 90° arc is 0.1 dB.

II. $N \times N$ STAR COUPLER

$N \times N$ star couplers are quite useful in high-speed, multiple-access optical networks, since they evenly distribute the input signal among many receivers and make possible the interconnection between them. A free space type integrated-optic star couplers, in which a slab waveguide region is located between the fan-shaped input and output channel waveguide arrays, are quite advantageous for constructing large-scale $N \times N$ star couplers (Dragone *et al.*, 1989; Okamoto *et al.*, 1991). Figure 3 shows the schematic configuration of the $N \times N$ star coupler. The input power, from any one of the N channel waveguides in the input array, is radiated to the slab (free space) region and it is received by the output array. In order to get the uniform power distribution into N output waveguides, the radiation pattern at the output side slab-array interface should be uniform over a sector of N waveguides. Since the radiation pattern is the Fraunhofer pattern (Fourier transform) of the field profile at the input side slab-array interface, proper sidelobes must be produced by the mode coupling from the excited input waveguide to neighboring guides so as to make a rectangular field pattern. Therefore, the dummy waveguides are necessary for the marginal guides to guarantee the equal coupling condition as the central guides. Figure 4 shows the splitting loss histogram of the fabricated 64×64 star coupler measured at $\lambda = 1.55 \mu\text{m}$ (Okamoto *et al.*, 1992). The essential splitting loss when light is evenly distributed into 64 output waveguides is 18.1 dB. Therefore, the average of the excess losses is 2.5 dB and the standard deviation of the splitting losses is 0.8 dB, respectively. Among the 2.5-dB additional loss, the inevitable imperfection (theoretical) loss is about 1.5 dB, the propagation loss is about 0.6 dB, and the coupling loss with a single-mode fiber is

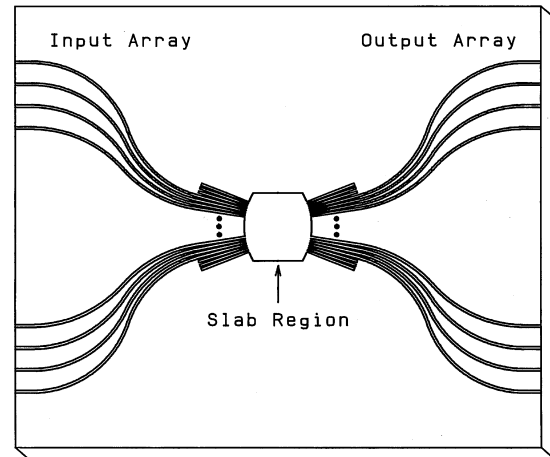


FIGURE 3 Schematic configuration of $N \times N$ star coupler.

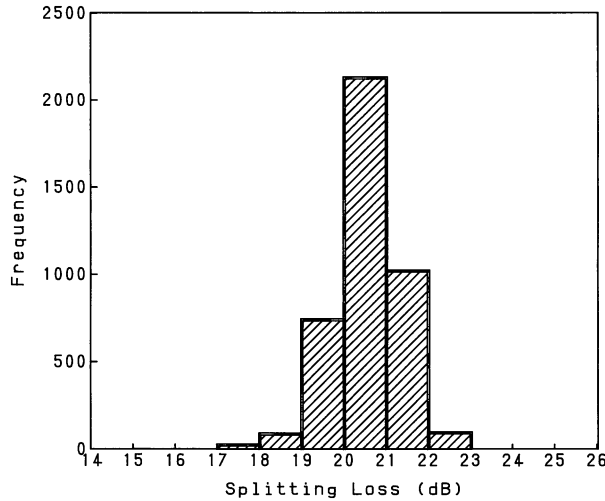


FIGURE 4 Splitting loss histogram of 64×64 star coupler.

0.2 dB/facet, respectively. Various kinds of star couplers ranging from 8×8 to 256×256 have been fabricated.

III. ARRAYED WAVEGUIDE GRATING

A. Principle of Operation and Fundamental Characteristics

An $N \times N$ arrayed-waveguide grating (AWG) multiplexer is very attractive in optical WDM networks since it is capable of increasing the aggregate transmission capacity of a single-strand optical fiber (Smit, 1988; Takahashi *et al.*, 1990). AWG consists of input/output waveguides, two focusing slab regions, and a phase-array of multiple channel waveguides with the constant path length difference ΔL between neighboring waveguides (Fig. 5). The

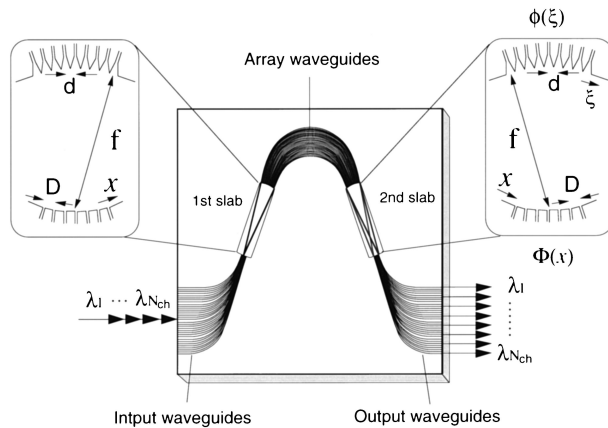


FIGURE 5 Schematic configuration of arrayed-waveguide grating multiplexer.

input light is radiated to the first slab and then excites the arrayed channel waveguides. After traveling through the arrayed waveguides, the light beams constructively interfere into one focal point in the second slab. The location of this focal point depends on the signal wavelength since the relative phase delay in each waveguide is given by $\Delta L/\lambda$. The dispersion of the focal position x with respect to the wavelength λ is given by

$$\frac{\Delta x}{\Delta \lambda} = \frac{f m}{n_s d}, \quad (1)$$

where f is the focal length of the converging slab, n_s is the effective index in the slab region, d is the pitch of the channel waveguides at their exits. The diffraction order m is given by $m = n_c \Delta L/\lambda$, where n_c is the effective index of the channel waveguide. It is known from Eq. (1) that the multiple WDM signals are dispersed and focused simultaneously to each prescribed position.

Various kinds of multiplexers ranging from 50-nm spacing 8-channel AWG to 25-GHz spacing AWG have been fabricated (Okamoto, 1999; Okamoto *et al.*, 1995, 1996). Figure 6 shows the measured loss spectra of 16-channel–100-GHz spacing AWG when light from tunable laser is coupled into center input port #8. In this case, the total number of channels of AWG is 32. Crosstalks, which is defined by the light leakage at the center of the neighboring channel, are about -40 dB. Figure 7 shows loss spectra of 64-channel–50-GHz multiplexer. Crosstalks of less than -35 dB have been achieved even in this narrow channel spacing AWG. By the improvement of the fabrication technology and the optimization of the waveguide configurations, good crosstalk characteristics have also been achieved in 128-channel–25-GHz spacing AWG as shown in Fig. 8. The average crosstalks to the neighboring channels are about -20 dB.

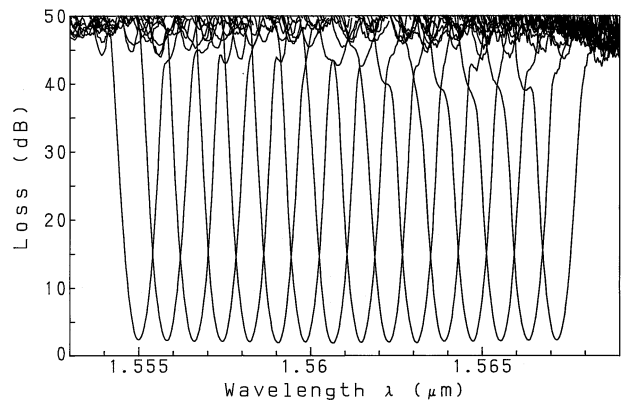


FIGURE 6 Demultiplexing properties of 16-channel–100-GHz spacing AWG.

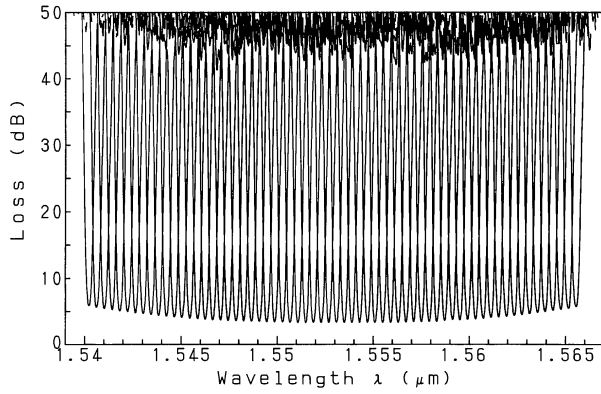


FIGURE 7 Demultiplexing properties of 64-channel-50-GHz multiplexer.

B. Flat Spectral Response AWG

Since the dispersion of the focal position x with respect to the wavelength λ is almost constant, the transmission loss of normal AWG monotonically increases around the center wavelength of each channel. This places tight restrictions on the wavelength tolerance of laser diodes and requires accurate temperature control for both AWGs and laser diodes. Moreover, since optical signals are transmitted through several filters in the WDM ring/bus networks, the cumulative passband width of each channel becomes much narrower than that of the single-stage AWG filter. Therefore, flattened and broadened spectral responses are required for AWG multiplexers. Several approaches have been proposed to flatten the pass bands of AWGs (Amersfoort *et al.*, 1996; Okamoto and Yamada, 1995; Trouchet *et al.*, 1997). Among them a flat-response AWG multiplexer having parabolic waveguide horns in the input waveguides (Okamoto and Sugita, 1996) shows quite satisfactory filter characteristics. Figure 9 shows the enlarged view of the interface between (a) input wave-

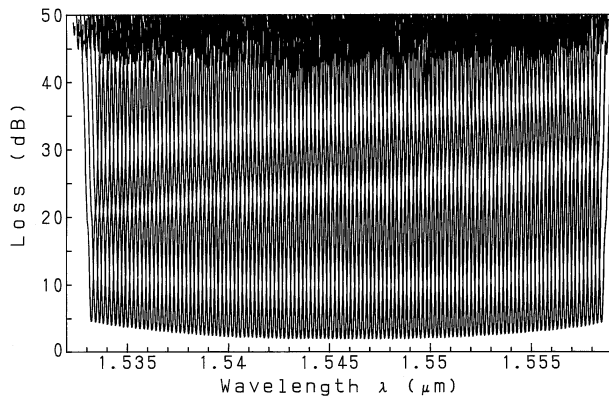


FIGURE 8 Demultiplexing properties of 128-channel-25-GHz spacing AWG.

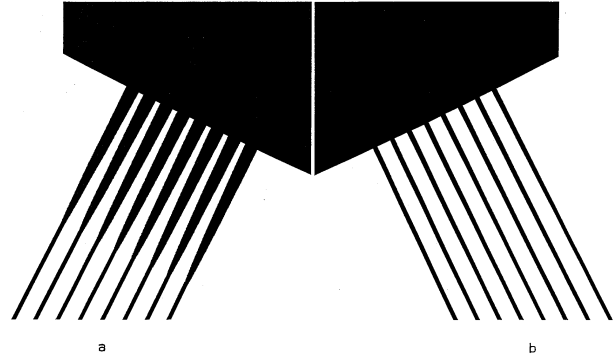


FIGURE 9 Enlarged view of the interface between (a) input waveguides and first slab and (b) second slab and in flat response 8-channel AWG.

guides and first slab and (b) second slab and output waveguides, respectively. The width of the parabolic horn along the propagation direction z is given by (Burns *et al.*, 1977)

$$W(z) = \sqrt{2\alpha\lambda_g z + (2a)^2}, \quad (2)$$

where α is a constant less than unity, λ_g is the wavelength in the guide ($\lambda_g = \lambda/n_{eff}$), and $2a$ is the core width of the channel waveguide (see inset of Fig. 10), respectively. At the proper horn length $z = \ell$ less than the collimator length, a slightly double-peaked intensity distribution can be obtained as shown in Fig. 10. A broadened and sharp-falling optical intensity profile is obtainable by the parabolic waveguide horn, which is quite advantageous for achieving a wide passband without deteriorating the nearest neighbor crosstalk characteristics. The broadened and double-peaked field is imaged onto the entrance of an output waveguide having a normal core width. The overlap integral of the refocused field with the local normal mode of

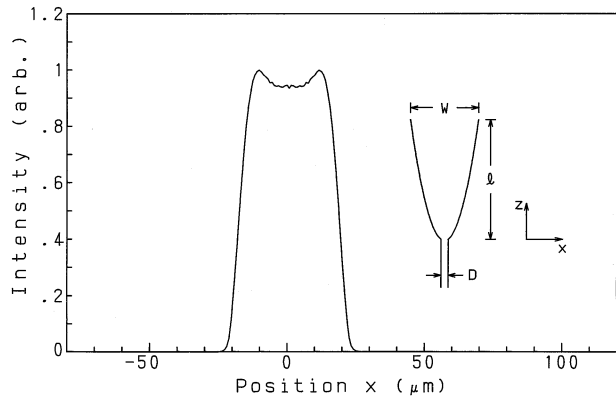


FIGURE 10 Intensity profile calculated by the beam propagation method for the parabolic horn with $W = 40 \mu\text{m}$ and $\ell = 800 \mu\text{m}$. Inset shows the schematic configuration of parabolic waveguide horn.

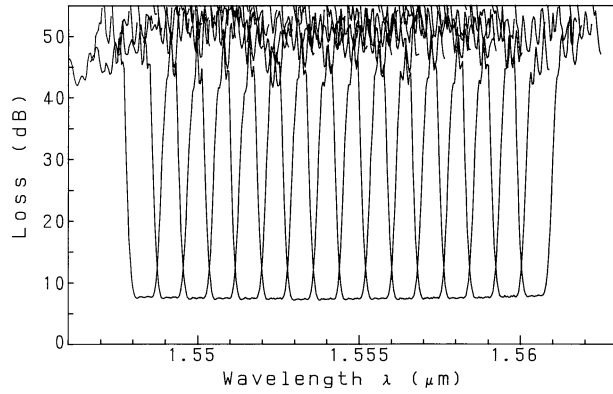


FIGURE 11 Demultiplexing properties of 16-channel-100-GHz spacing AWG having parabolic horns with $W=40\ \mu\text{m}$ and $z=800\ \mu\text{m}$.

the output waveguide gives the flattened spectral response of AWG. Figure 11 shows the demultiplexing properties of 16-channel-100-GHz spacing AWG having parabolic horns with $W=40\ \mu\text{m}$ and $z=800\ \mu\text{m}$. The crosstalks to the neighboring channels are less than $-35\ \text{dB}$ and the on-chip loss is about $7.0\ \text{dB}$. The average 1-, 3-, and 20-dB bandwidths are 86.4, 100.6, and 143.3 GHz, respectively.

It has been confirmed that in order to obtain a flat spectral response, it is necessary to produce the rectangular electric field profile at the focal plane (interface between the second slab and output waveguides). Since the electric field profile in the focal plane is the Fourier transform of the field in the array output aperture (interface between the array waveguide and second slab), such a rectangular field profile could be generated when the electric field at the array output aperture obeys a $\sin(\xi)/\xi$ distribution, where ξ is measured along the array output aperture (Okamoto and Yamada, 1995). Figure 12 shows the electric field amplitude and relative phase delays (excess phase value added to $i \times \Delta L$ where i denotes the i -th array waveguide) in

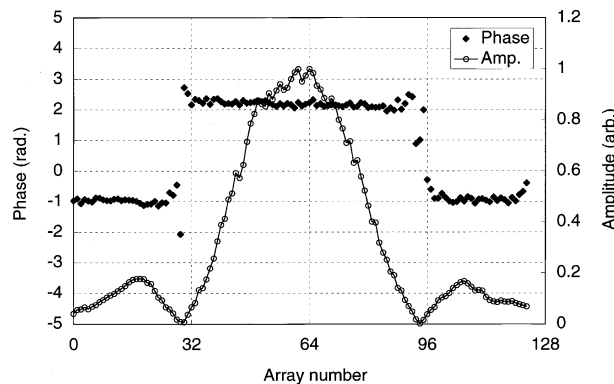


FIGURE 12 Electric field amplitude and relative phase delays in the sinc-type flat response AWG.

the sinc-type flat response AWG measured by the low coherence Fourier transform spectroscopy (Takada *et al.*, 1994). Sinc-shaped electric field amplitude distribution is realized by introducing an additional loss to each array waveguide. The excess path length differences of π for the negative sinc values are realized by the additional path length $\delta\ell = \lambda/2n_c$ to the corresponding array waveguides. The crosstalk and flat passband characteristics are almost the same as those of parabola-type AWGs.

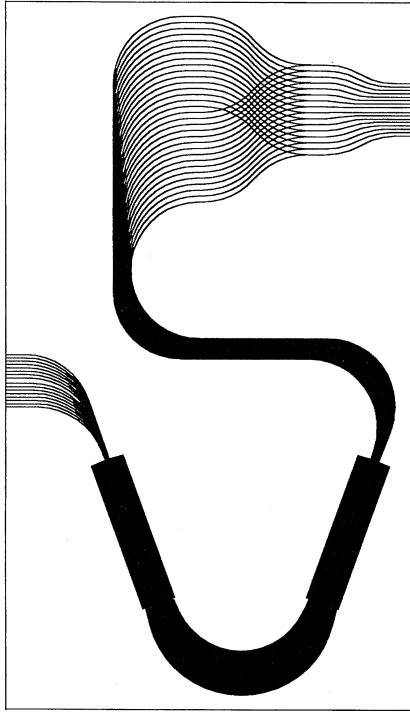
C. Uniform-Loss and Cyclic-Frequency (ULCF) AWG

In principle, $N \times N$ signal interconnection can be achieved in AWG when free spectral range (FSR) of AWG is N times the channel spacing. Here FSR is given by

$$FSR = \frac{n_c v_0}{N_c m}, \quad (3)$$

where v_0 is the center frequency of WDM signals. Generally light beams with three different diffraction orders of $m-1$, m , and $m+1$ are utilized to achieve $N \times N$ interconnections (Dragone *et al.*, 1991). The cyclic property provides an important additional functionality as compared to simple multiplexers or demultiplexers and plays a key role in more complex devices as add/drop multiplexers and wavelength switches. However, such interconnectivity cannot always be realized with the conventional AWGs. The typical diffraction order of 32-channel AWG with 100-GHz channel spacing is $m=60$. Since FSR is inversely proportional to m , substantial pass frequency mismatch is brought by the difference between three FSRs. Also insertion losses of AWG for peripheral input and output ports are $2 \sim 3\ \text{dB}$ higher than those for central ports as shown in Figs. 7 and 8. These noncyclic frequency characteristics and loss nonuniformity in conventional AWGs are main obstacles which prevent the development of practical $N \times N$ routing networks.

Novel 32×32 arrayed-waveguide grating having uniform-loss and cyclic-frequency characteristics is proposed and fabricated to solve the problems in conventional AWGs. Figure 13 shows the schematic configuration of uniform-loss and cyclic-frequency (ULCF) arrayed-waveguide grating (Okamoto *et al.*, 1997). It consists of an 80-channel AWG multiplexer with 100-GHz spacing and 32 optical combiners which are connected to 72 output waveguides of the multiplexer. The arc length of slab is $f=24.55\ \text{mm}$ and the number of array waveguides is 300 having the constant path length difference $\Delta L=24.6\ \mu\text{m}$ between neighboring waveguides. The diffraction order is $m=23$ which gives a free spectral range of $FSR=8\ \text{THz}$. In the input side, 32 waveguides ranging from #25 to #56 are used for the input waveguides

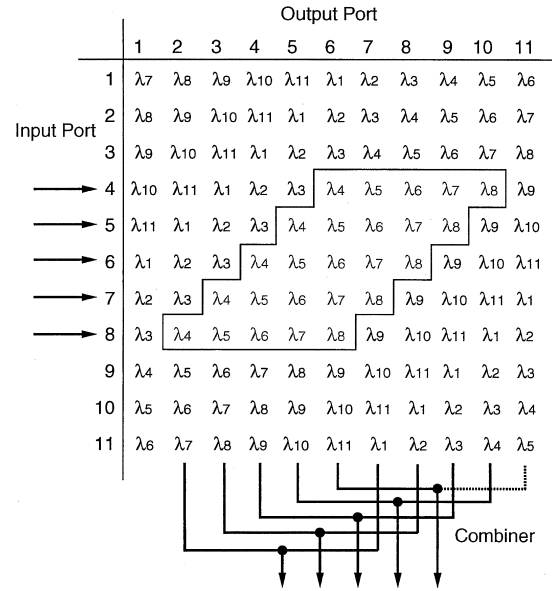


Configuration of ULCF 16x16 AWG

FIGURE 13 Schematic configuration of uniform-loss and cyclic-frequency AWG.

so as to secure the uniform loss characteristics. In the output side, two waveguides ($(i + 8)$ -th and $(i + 40)$ -th waveguide for $i = 1 \sim 32$) are combined through waveguide intersection and a multimode interference (MMI) coupler (Veerman *et al.*, 1992) to make one output port. Since the peripheral output ports are not used, uniform loss characteristics are obtained. Figure 14 shows the principle of how the ULCF arrayed-waveguide grating is constructed. In this example a 5×5 ULCF AWG is fabricated from an 11×11 original AWG. Figures 15(a) and (b) show measured histograms of insertion losses and channel pass frequency deviations for 32-input/32-output combinations in the uniform-loss and cyclic-frequency AWG. The average and standard deviation of on-chip losses are $\alpha_{av} = 6.6$ dB and $\sigma_{loss} = 0.25$ dB. The standard deviation of center pass frequencies (wavelength) is $\sigma_{freq} = 3.9$ GHz ($\sigma_{wl} = 0.025$ nm).

An example of all optical $N \times N$ ($N = 5$) interconnection systems using an arrayed-waveguide grating as a router is shown in Fig. 16. Figure 16(a) shows the physical topology between AWG router and N nodes. Some routes may be used for connection to other networks. Based on the interconnectivity of AWG, the resulting logical connectivity patterns become $N \times N$ star network as shown in



Operational Principle of ULCF 4x4 AWG

FIGURE 14 Operational principle of the ULCF arrayed-waveguide grating.

Fig. 16(b). All the nodes can communicate with each other at the same time, thus enabling N^2 optical connections simultaneously. Signals can be freely routed by changing the carrier wavelength of the signal. When combined with the wavelength conversion lasers, this AWG router can construct signal routing networks without using space-division optical switches.

D. Athermal AWG

Temperature sensitivity of the pass wavelength (frequency) in the silica-based AWG is about $d\lambda/dT = 1.2 \times 10^{-2}$ (nm/deg) ($dv/dT = -1.5$ (GHz/deg)), which is mainly determined by the temperature dependence of silica glass itself ($dn_c/dT = 1.1 \times 10^{-5}$ (1/deg)). The AWG multiplexer should be temperature controlled with a heater or a Peltier cooler to stabilize the channel wavelengths. This requires the constant power consumption of a few watts and a lot of equipment for the temperature control. AWG configuration to achieve an athermal (temperature insensitive) operation over $0 \sim 85^\circ\text{C}$ temperature range is reported (Inoue *et al.*, 1997). Figure 17 shows a schematic configuration of athermal AWG. The temperature-dependent optical path difference in silica waveguides is compensated with a trapezoidal groove filled with silicone adhesive which has a negative thermal coefficient. Since the pass wavelength is given by $\lambda_0 = n_c \Delta L / m$, the optical path length difference

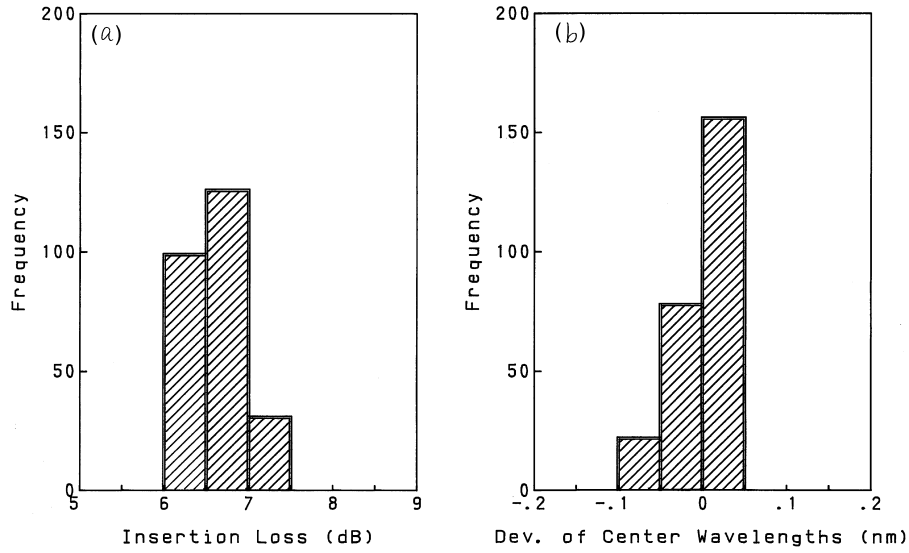


FIGURE 15 (a) Measured insertion loss histograms for entire 32×32 input/output combinations in the uniform-loss and cyclic-frequency AWG. (b) Measured histograms of channel center frequency deviations for entire 32×32 input/output combinations in the uniform-loss and cyclic-frequency AWG.

$n_c \Delta L$ should be made insensitive to temperature. Therefore the groove is designed to satisfy the following conditions:

$$n_c \Delta L = n_c \Delta \ell + \hat{n}_c \Delta \hat{\ell}, \quad (4)$$

and

$$\frac{d(n_c \Delta L)}{dT} = \frac{dn_c}{dT} \Delta \ell + \frac{d\hat{n}_c}{dT} \Delta \hat{\ell} = 0, \quad (5)$$

where \hat{n}_c is the refractive index of silicone and $\Delta \ell$ and $\Delta \hat{\ell}$ are the path length differences of silica waveguides and the silicone region. Equation (4) is a condition to sat-

isfy the AWG specifications, and Eq. (5) is the athermal condition, respectively. The temperature sensitivity of silicone is $d\hat{n}_c/dT = -37 \times 10^{-5}$ (1/deg). Therefore the path length difference of silicone is $\Delta \hat{\ell} \cong \Delta \ell / 37$. Figure 18 shows temperature dependencies of pass wavelengths in conventional and athermal AWGs. The temperature-dependent wavelength change has been reduced from 0.95 to 0.05 nm in the $0 \sim 85^\circ\text{C}$ range. The excess loss caused by the groove is about 2 dB which is mainly a diffraction loss in the groove. The insertion loss caused by diffraction loss can be reduced by segmenting a single trapezoidal silicone region into multiple groove

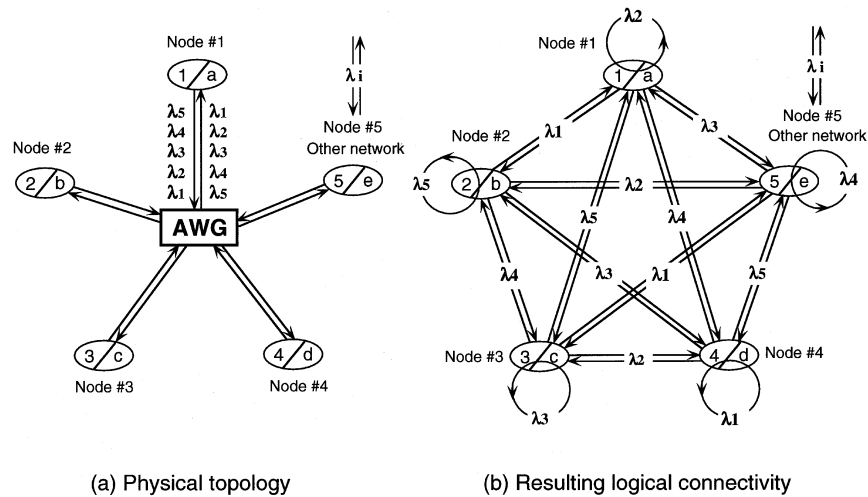


FIGURE 16 All optical $N \times N$ interconnection system using ULCF AWG as a router. (a) Physical topology, (b) Resulting logical connectivity.

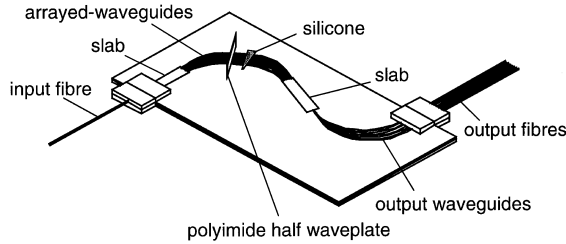


FIGURE 17 Schematic configuration of athermal AWG.

regions (Kaneko *et al.*, 1999). In the segmented groove regions, the light beam is periodically focused. Therefore insertion loss is reduced to about 0.4 dB. During the heat cycle test between -40 and 85°C , the loss change is smaller than 0.2 dB. Furthermore, the channel wavelength change was less than 0.02 nm in a long-term test over 5000 hr at 75°C and 90% relative humidity.

E. Phase Error Compensation of AWG

Crosstalk improvement is the major concern for the AWG multiplexers, especially for narrow channel spacing AWGs and $N \times N$ AWG routers. Crosstalks to other channels are caused by the sidelobes of the focused beam in the second slab region. These sidelobes are mainly attributed to the phase fluctuations of the total electric field profile at the output side array-slab interface since the focused beam profile is the Fourier transform of the electric field in the array waveguides. The phase errors are caused by the nonuniformity of effective-index and/or core geometry in the arrayed-waveguide region. Phase errors in the AWGs are measured by using Fourier transform spectroscopy (Okamoto *et al.*, 1997). In order to improve the crosstalk characteristics of AWGs, a phase error compen-

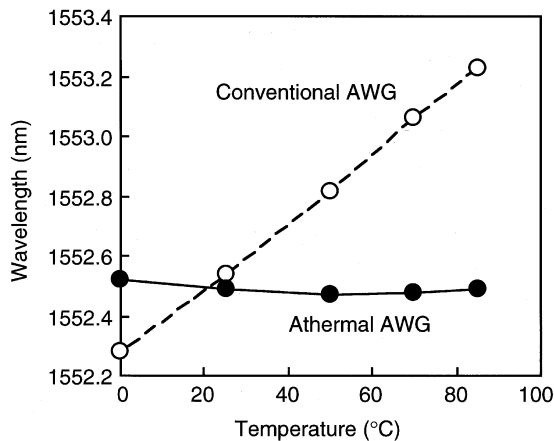


FIGURE 18 Temperature dependences of pass wavelengths.

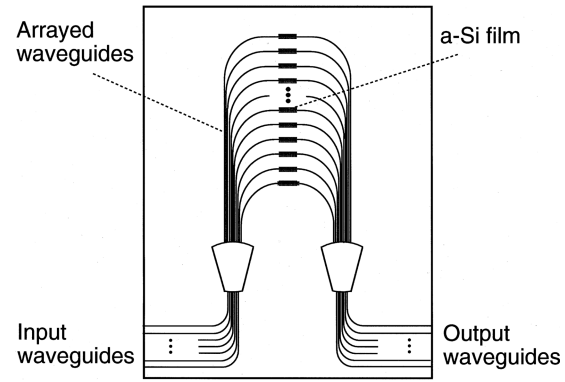


FIGURE 19 Configuration of phase-compensated AWG with a Si film.

sation experiment is carried out using 16-channel-10-GHz spacing ultra narrow AWG filter (Yamada *et al.*, 1996). Figure 19 shows the configuration of phase-compensated AWG using an a-Si film for phase trimming. The path length difference of AWG is $\Delta L = 1271 \mu\text{m}$ in 64 array waveguides. Since the array-waveguide region occupies a large area, accumulated phase errors become quite large. Therefore the crosstalk of AWG without phase error compensation was about -8 dB. A Si stress-applying film was deposited on top of the over-cladding of each array waveguide and the amount of photoelastic effect was adjusted by the Ar ion laser trimming (evaporation) of an Si film length. Figure 20 shows phase-error distributions before and after phase compensation. The peak-to-peak fluctuation is reduced to about ± 2 degrees over the 60 mm average path length. The effective-index fluctuations are reduced to within 1.2×10^{-7} . Figure 21 shows the measured demultiplexing properties of phase compensated 16-channel-10-GHz spacing AWG. Crosstalks to the neighboring and all other channels are less than -31 dB.

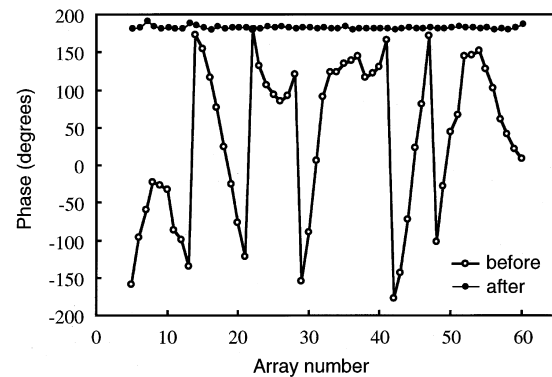


FIGURE 20 Phase error distributions for the TE mode before and after phase compensation.

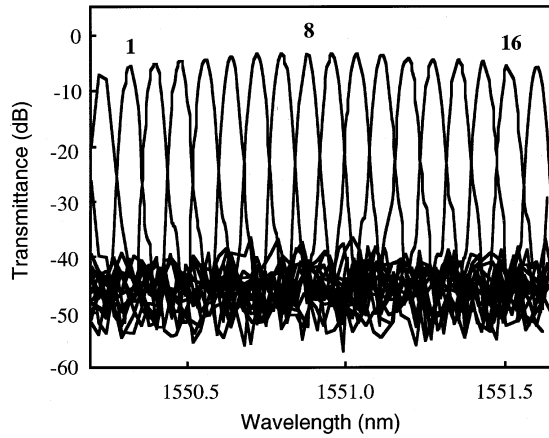


FIGURE 21 Transmittance of phase-compensated 16-channel–10-GHz spacing AWG.

IV. OPTICAL ADD/DROP MULTIPLEXER

An optical add/drop multiplexer (ADM) is a device that gives simultaneous access to all wavelength channels in WDM communication systems. A novel integrated-optic ADM is fabricated and basic functions of individually routing 16 different wavelength channels with 100-GHz channel spacing were demonstrated (Okamoto, Okuno *et al.*, 1996). The waveguide configuration of a 16-channel optical ADM is shown in Fig. 22. It consists of four arrayed-waveguide gratings and 16 double-gate thermo-

optic (TO) switches. Four AWGs are allocated crossing their slab regions with each other. These AWGs have the same grating parameters; they are the channel spacing of 100 GHz and the free spectral range of 3300 GHz (26.4 nm) at a 1.55- μm region. Equally spaced WDM signals, $\lambda_1, \lambda_2, \dots, \lambda_{16}$, which are coupled to the main input port (add port) in Fig. 22 are first demultiplexed by the AWG₁ (AWG₂) and then 16 signals are introduced into the left-hand side arms (right-hand side arms) of double-gate TO switches. The cross angles of the intersecting waveguides are designed to be larger than 30° so as to make the crosstalk and insertion loss negligible. When the double-gate switch is “off,” the demultiplexed light by AWG₁ (AWG₂) goes to the cross arm and is multiplexed again by the AWG₃ (AWG₄). On the other hand, if double-gate switch is “on” state the demultiplexed light by AWG₁ (AWG₂) goes to the through arm and is multiplexed by the AWG₄ (AWG₃). Therefore, any specific wavelength signal can be extracted from the main output port and led to the drop port by changing the corresponding switch condition. A signal at the same wavelength as that of the dropped component can be added to the main output port when it is coupled into an add port (Fig. 22). When TO switches $SW_2, SW_4, SW_6, SW_7, SW_9, SW_{12}, SW_{13}$, and SW_{15} , for example, are turned to “on” the selected signals $\lambda_2, \lambda_4, \lambda_6, \lambda_7, \lambda_9, \lambda_{12}, \lambda_{13}$ and λ_{15} are extracted from the main output port (solid line) and led to the drop port (dotted line) as shown in Fig. 23. The on-off crosstalk is smaller than -30 dB with the on-chip losses of 8 ~ 10 dB.

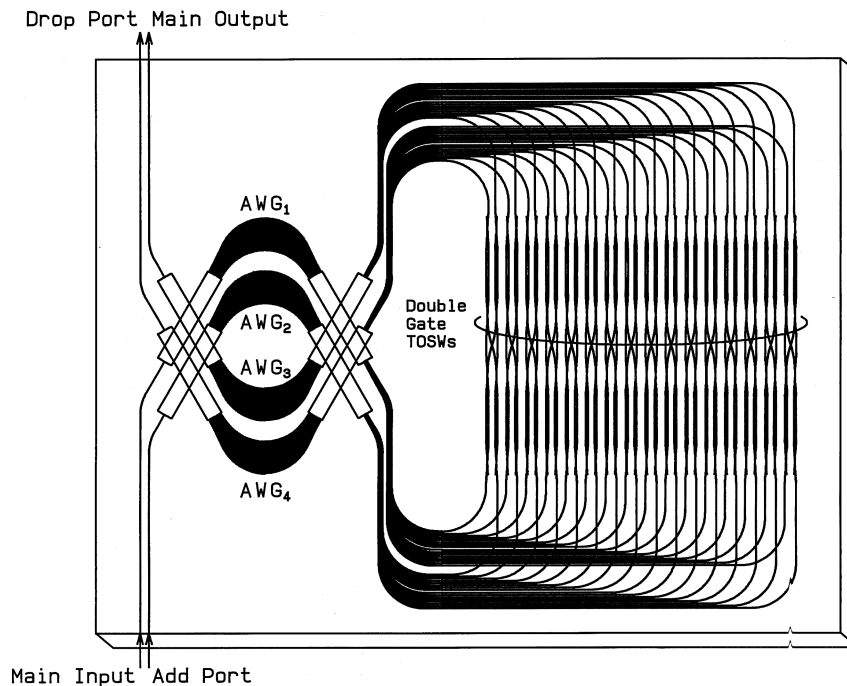


FIGURE 22 Waveguide configuration of 16-channel optical ADM with double-gate TO switches.

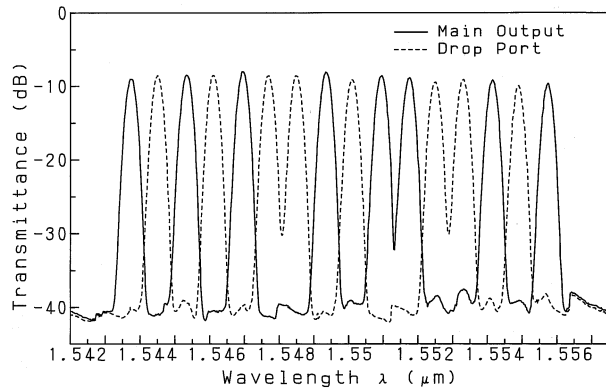


FIGURE 23 Transmission spectra from main input port to main output port and drop port when TO witches SW_2 , SW_4 , SW_6 , SW_7 , SW_9 , SW_{12} , SW_{13} , and SW_{15} are “on.”

Though the electric power necessary to drive double-gate switch becomes two times larger than the conventional TO switch, the power consumption itself can be reduced to almost $1/5 \sim 1/2$ when bridge-suspended phase shifters are utilized. The present optical ADM can transport all input signals to the succeeding stages without inherent power losses. Therefore, these ADMs are very attractive for all optical WDM routing systems and allow the network to be transparent to signal formats and bit rates.

V. $N \times N$ MATRIX SWITCHES

Space-division optical switches are one of the indispensable optical devices for the reconfigurable interconnects in cross-connect systems, fiber-optic subscriber line connectors, and photonic intermodule connectors (Ito *et al.*, 1992). Figure 24 shows the logical arrangement of a 16×16 strictly nonblocking matrix switch with a path-independent insertion loss (PI-Loss) configuration (Goh *et al.*, 1998). This arrangement is quite advantageous for reducing total circuit length since it requires only N switching stages to construct the $N \times N$ switch.

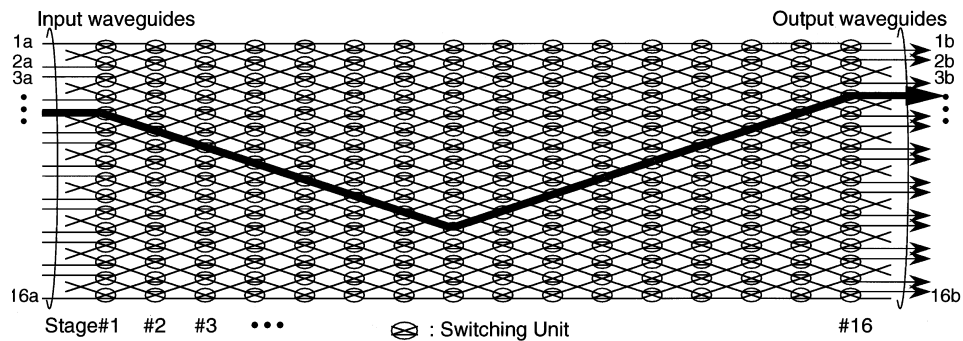


FIGURE 24 Logical arrangement of a 16×16 matrix switch.

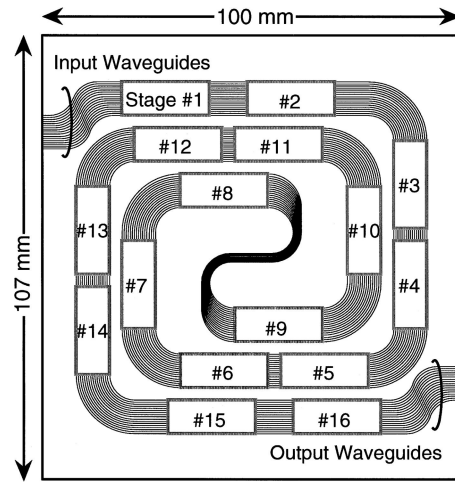


FIGURE 25 Circuit layout of the 16×16 matrix switch.

The switching unit consists of double-gate switch which has been described in Fig. 22. The circuit layout of the 16×16 matrix switch is shown in Fig. 25. Sixteen switching stages are allocated along the serpentine waveguides. There are 16 switching units in one switching stage. The total circuit length is 66 cm. Figures 26(a) and (b) show measured insertion losses and extinction ratios for all 256 input and output connection patterns. The insertion loss ranges from 6.0 to 8.0 dB, with an average of 6.6 dB. The extinction ratio ranges from 45 to 67 dB, with an average of 55 dB. The total electric power for operating the 16×16 matrix switch is about 17 W (1.06 W for each switching unit).

VI. LATTICE-FORM PROGRAMMABLE DISPERSION EQUALIZERS

The transmission distance in optical fiber communications has been greatly increased by the development of erbium-doped fiber amplifiers. Consequently, the main factor limiting the maximum repeater span is now the fiber

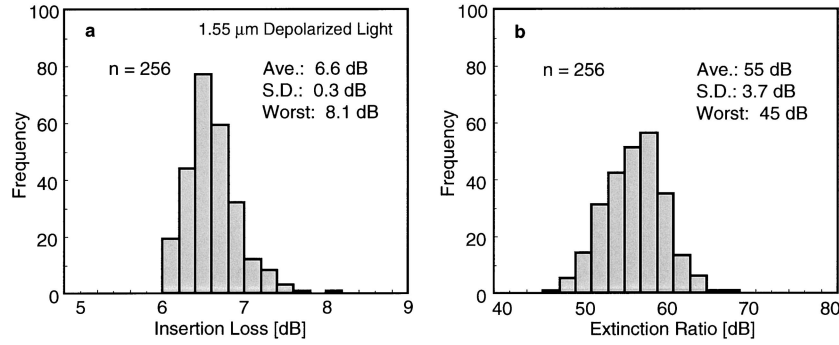


FIGURE 26 Measured (a) insertion losses and (b) extinction ratios for all 256 input and output connection patterns.

chromatic dispersion. Several techniques have been reported to compensate for the delay distortion in optical stages (Gnauck *et al.*, 1993; Hill *et al.*, 1994; Okamoto *et al.*, 1993; Takiguchi *et al.*, 1994, 1995; Vengsarkar *et al.*, 1993). An advantage of the PLC optical delay equalizer (Takiguchi *et al.*, 1994, 1995) is that variable group-delay characteristics can be achieved by the phase control of silica waveguides. The basic configuration of the PLC delay equalizer is shown in Fig. 27. It consists of N ($= 8$) asymmetrical Mach-Zehnder interferometers and $N + 1$ ($= 9$) tunable couplers, which are cascaded alternately in series. The cross-port transfer function of the optical circuit is expressed by a Fourier series as

$$H(z) = \sum_{k=0}^N a_k z^{-k+N/2}, \quad (6)$$

where z denotes $\exp(j2\pi\nu\Delta t)$ (ν : optical frequency, $\Delta t = n_c \Delta L / c$: unit delay time difference in asymmetrical MZ interferometer) and a_k is the complex expansion

coefficient. The circuit design procedures are as follows. First the equalizer transfer function to be realized is expressed by the analytical function. Then coefficient a'_k s are determined by expanding the analytical function into a Fourier series. Finally coupling ratio (ϕ_i) and phase shift value (θ_i) in each stage of lattice filter are determined by the filter synthesis method (Jinguji, 1995).

In ultrahigh speed optical fiber transmission systems (> 100 Gbit/s), the effect of the higher-order dispersion (third-order dispersion or dispersion slope) in the dispersion shifted fiber (DSF) is one of the major factors limiting the transmission distance (Kawanishi *et al.*, 1995). Programmable dispersion equalizers can be designed so as to compensate for the higher-order dispersion of DSFs. Figure 28(a) shows the measured power transmittance and relative delay time of the PLC higher-order dispersion equalizer (Takiguchi *et al.*, 1996). The dispersion slope of the equalizer is calculated to be -15.8 ps/nm². Figure 28(b) is the relative delay of the 300-km DSF.

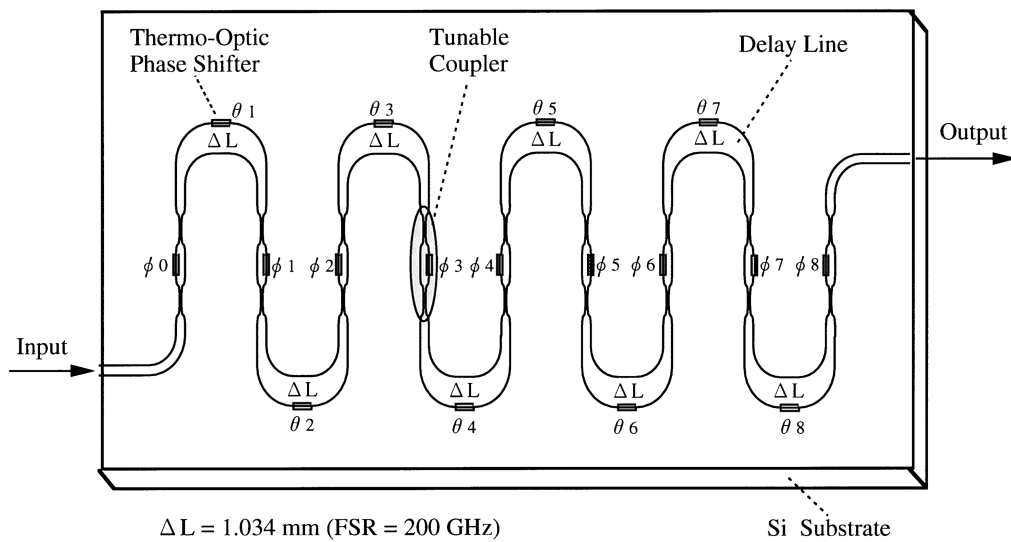
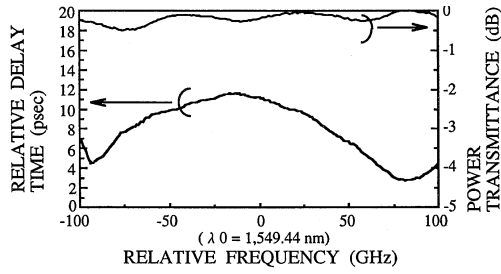
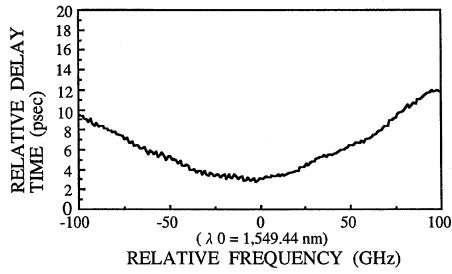


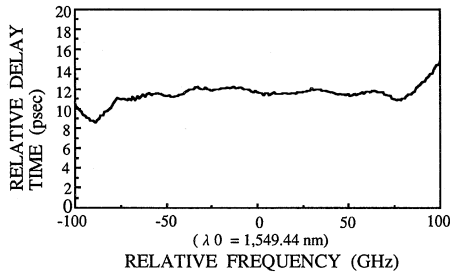
FIGURE 27 Basic configuration of the PLC equalizer.



(a) Higher-Order Dispersion Equalizer



(b) 300-km Dispersion Shifted Fiber



(c) 300-km DSF + Equalizer

FIGURE 28 Relative delay times of (a) PLC higher-order dispersion equalizer, (b) 300-km DSF and (c) 300-km DSF + equalizer.

The dispersion slope of DSF is $0.05 \sim 0.06$ ps/nm²/km. Therefore, the equalizer can compensate the higher-order dispersion of ~ 300 -km DSF. Figure 28(c) shows the relative delay time of 300-km DSF cascaded with the equalizer. The positive dispersion slope of the DSF is almost completely compensated by the PLC equalizer. A time-division (200 Gbit/s) multiplexed transmission experiment using a dispersion slope equalizer has been carried out over a 100-km fiber length (Takiguchi, Kawanishi, Takara, Kamatani *et al.*, 1996). The pulse distortion caused by the dispersion slope was almost completely recovered, and the power penalty was improved by more than 4 dB.

VII. HYBRID INTEGRATION TECHNOLOGY USING PLC PLATFORMS

It is widely recognized that optical hybrid integration is potentially a key technology for fabricating advanced

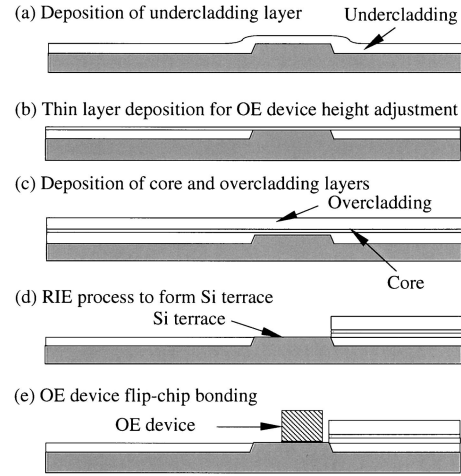


FIGURE 29 PLC platform fabrication process.

integrated optical devices. A silica-based waveguide on an Si substrate is a promising candidate for the hybrid integration platform since high-performance PLCs have already been fabricated using silica-based waveguides and Si has highly stable mechanical and thermal properties which make it suitable as an optical bench. Figure 29 shows the PLC platform fabrication process (Yamada *et al.*, 1993). First, a thick under-cladding is deposited on a Si substrate with a terraced region using FHD, and then the surface of the substrate is flattened by mechanical polishing. To minimize the optical coupling loss between the optoelectronics (OE) device on the terrace and optical waveguide, a thin layer is deposited on the polished substrate surface. The thickness of the layer corresponds to the height of the active region of the OE device on the terrace. Then, a core layer is deposited and patterned into a core ridge by RIE. The core ridge is then covered by the over-cladding layer. Finally, RIE is used to form the Si terrace for the OE devices on the PLC and the terrace surface is exposed. The relative positions of the core and Si terrace surface are determined precisely because the terrace acts as an etch-stop layer during the RIE process. As a result, Si terrace functions as both a high-precision alignment plane

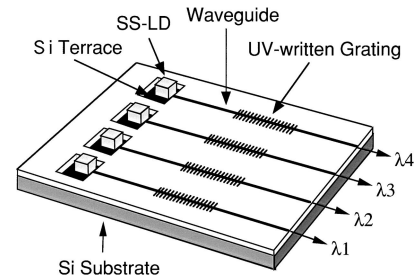


FIGURE 30 Configuration of hybrid integrated multi-wavelength external cavity laser.

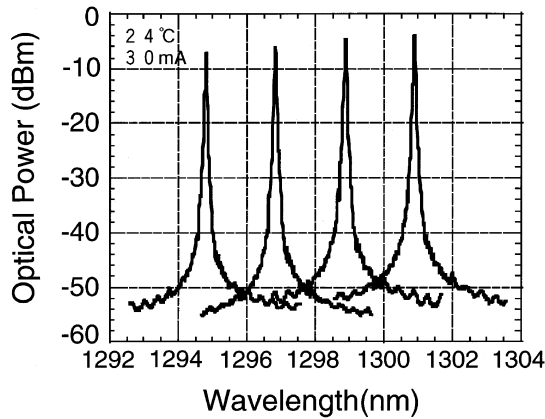


FIGURE 31 Oscillation spectra of 4-wavelength laser.

and a heat sink when the OE device is flip-chip bonded on the terrace.

PLC platform technology has also been utilized in fabrication of a hybrid integrated external cavity laser (Tanaka *et al.*, 1996). Figure 30 shows the configuration of multi-wavelength external cavity laser with a Uv written grating (Tanaka *et al.*, 1997). Bragg gratings with a 2-nm wavelength interval are written into each waveguide by ArF excimer laser irradiation through phasemasks. Figure 31 shows the measured output spectra. Each laser operates in a single longitudinal mode with a side-mode suppression of 40 dB. The temperature sensitivity of the oscillation frequency is -1.7 GHz/deg, which is one eighth of the DFB lasers. A four-channel simultaneous modulation experiment has been successfully carried out at 2.5 Gbit/s (Takahashi *et al.*, 1997). A temperature-stable multiwavelength source will play an important role in WDM transmission and access network systems.

Semiconductor optical amplifier (SOA) gate switches having spot-size converters on both facets have been successfully hybrid-integrated on PLC platforms to construct

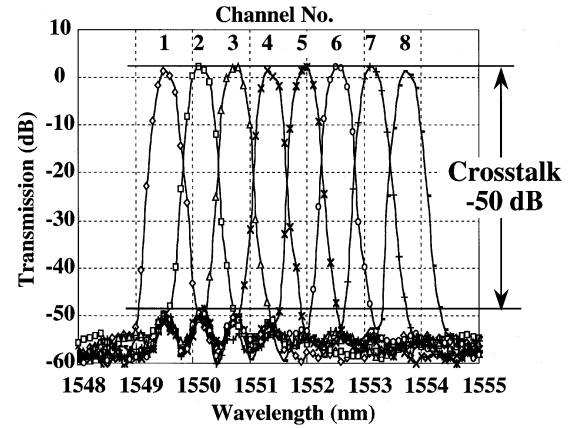


FIGURE 33 Optical transmission spectra of the wavelength channel selector when only one SOA gate switch is activated successively.

high-speed wavelength channel selectors and 4×4 optical matrix switches (Kato *et al.*, 1998; Ogawa *et al.*, 1998). Figure 32 shows the configuration of 8-channel optical wavelength selector module. It consists of two AWG chips with 75-GHz channel spacing and hybrid integrated SOA gate array chip. It selects and picks up any wavelength channel from multiplexed signals by activating the corresponding SOA gate switch. Three PLC chips are directly attached to each other using Uv curable adhesive. The length of the SOA gate switch is $1200 \mu\text{m}$ and their separation is $400 \mu\text{m}$. The coupling loss between the SOA and PLC waveguides ranges from 3.9 to 4.9 dB. Figure 33 shows the optical transmission spectra of the wavelength channel selector when only one SOA gate switch is activated successively. The SOA injection current is 50 mA for all SOAs. The peak transmittances have 1 ~ 3 dB gains; there are 16 ~ 19 dB total chip losses, and fiber coupling losses, are compensated by SOA gains. The crosstalk is less than -50 dB and the

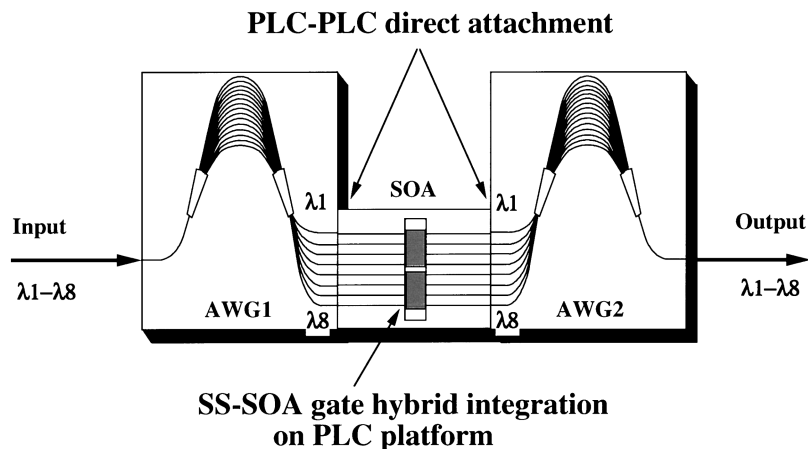


FIGURE 32 Configuration of 8-channel optical wavelength selector module.

polarization-dependent loss is smaller than 1.4 dB, respectively. In the high-speed switching experiments, the rise and fall time is confirmed to be less than 1 nsec.

VIII. CONCLUSION

Although silica-based waveguides are simple circuit elements, various functional devices are fabricated by utilizing spatial multibeam or temporal multistage interference effects such as arrayed-waveguide grating multiplexers and lattice-form programmable filters. Hybrid integration technologies will further enable us to realize much more functional and high-speed devices. The PLC technologies supported by continuous improvements in waveguide fabrication, circuit design, and device packaging will further proceed to a higher level of integration of optics and electronics aiming at the next generation of telecommunication systems.

SEE ALSO THE FOLLOWING ARTICLES

DIFFRACTIVE OPTICAL COMPONENTS • INTEGRATED CIRCUIT PACKAGING • MULTIPLEXING • OPTICAL FIBER COMMUNICATIONS

BIBLIOGRAPHY

- Amersfoort, M. R., Soole, J. B. D., LeBlanc, H. P., Andreadakis, N. C., Rajhei, A., and Caneau, C. (1996). "Passband broadening of integrated arrayed waveguide filters using multimode interference couplers," *Electron. Lett.* **32**, 449–451.
- Burns, W. K., Milton A. F., and Lee, A. B. (1977). "Optical waveguide parabolic coupling horns," *Appl. Phys. Lett.* **30**, 28–30.
- Dragone, C., Henry, C. H., Kaminow, I. P., and Kistler, R. C. (1989). "Efficient multichannel integrated optics star coupler on silicon," *IEEE Photonics Tech. Lett.* **1**, 241–243.
- Dragone, C., Edwards, C. A., and Kistler, R. C. (1991). "Integrated optics $N \times N$ multiplexer on silicon," *Photon. Tech. Lett.* **3**, 896–899.
- Goh, T., Yasu, M., Hattori, K., Himeno, A., Okuno, M., and Ohmori, Y. (1998). "Low-loss and high-extinction-ratio silica-based strictly non-blocking 16×16 thermo-optic matrix switch," *IEEE Photonics Tech. Lett.* **10**, 810–812.
- Gnauck, A. H., Jopson, R. M., and Derosier, R. M. (1993). "10-Gb/s 360-km transmission over dispersive fiber using midsystem spectral inversion," *IEEE Photonics Tech. Lett.* **5**, 663–666.
- Hida, Y., Hibino, Y., Okazaki, H., and Ohmori, Y. (1995). "10-m-long Silica-Based Waveguide with a Loss of 1.7 dB/m," IPR'95 ITHC6, Dana Point, CA.
- Hill, K. O., Theriault, S., Malo, B., Bilodeau, F., Kitagawa, T., Johnson, D. C., Albert, J., Takiguchi, K., Kataoka, T., and Hagimoto, K. (1994). "Chirped in-fiber Bragg grating dispersion compensators; linearization of dispersion characteristics and demonstration of dispersion compensation in 100 km, 10 Gbit/s optical fiber link," *Electron. Lett.* **30**, 1755–1756.
- Inoue, Y., Kaneko, A., Hanawa, F., Takahashi, H., Hattori, K., and Sumida, S. (1997). "Athermal silica-based arrayed-waveguide grating multiplexer," *Electron. Lett.* **33**, 1945–1946.
- Ito, T., Himeno, A., Takada, K., and Kato, K. (1992). "Photonic inter-module connector using silica-based optical switches," Proc. GLOBE-COM'92, Orlando, FL, pp. 187–191.
- Jinguji, K. (1995). "Synthesis of coherent two-port lattice-form optical delay-line circuit," *IEEE J. Lightwave Tech.* **13**, 73–82.
- Kato, T., Sasaki, J., Shimoda, T., Hatakeyama, H., Tamanuki, T., Yamaguchi, M., Titamura, M., and Itoh, M. (1998). "10 Gb/s Photonic Cell Switching with Hybrid 4×4 Optical Matrix Switch Module on Silica Based Planar Waveguide Platform," OFC'98 postdeadline paper PDP3, San Jose, CA.
- Kaneko, A., Kamei, S., Inoue, Y., Takahashi, H., and A. Sugita, (1999). "Athermal silica-based arrayed-waveguide grating (AWG) multiplexers with new low loss groove design," OFC-IOOC '99 TuO1, pp. 204–206, San Diego, CA.
- Kawanishi, S., Takara, H., Morioka, T., Kamatani, O., and Saruwatari, M. (1995). "200 Gbit/s, 100 km time-division-multiplexed optical transmission using supercontinuum pulses with prescaled PLL timing extraction and all-optical demultiplexing," *Electron. Lett.* **31**, 816–817.
- Okamoto, K. (1999). "Recent progress of integrated optics planar light-wave circuits," *Opt. Quantum Electron.* **31**, 107–129.
- Okamoto, K. "Fundamentals of Optical Waveguides," Academic Press.
- Ogawa, I., Ebisawa, F., Yoshimoto, N., Takiguchi, K., Hanawa, F., Hashimoto, T., Sugita, A., Yanagisawa, M., Inoue, Y., Yamada, Y., Tohmori, Y., Mino, S., Ito, T., Magari, K., Kawaguchi, Y., Himeno, A., and Kato, K. (1998). "Loss-less Hybrid Integrated 8-ch Optical Wavelength Selector Module Using PLC Platform and PLC-PLC Direct Attachment Techniques," OFC'98 postdeadline paper PDP4, San Jose, CA.
- Okamoto, K., Hibino, Y., and Ishii, M. (1993). "Guided-wave optical equalizer with α -power chirped grating," *IEEE J. Lightwave Tech.* **11**, 1325–1330.
- Okamoto, K., Hasegawa, H., Ishida, O., Himeno, A., and Ohmori, Y. (1997). " 32×32 arrayed-waveguide grating multiplexer with uniform loss and cyclic frequency characteristics," *Electron. Lett.* **33**, 1865–1866.
- Okamoto, K., Moriwaki, K., and Suzuki, S. (1995). "Fabrication of 64×64 arrayed-waveguide grating multiplexer on silicon," *Electron. Lett.* **31**, 184–185.
- Okamoto, K., Okuno, M., Himeno, A., and Ohmori, Y. (1996). "16-channel optical Add/Drop multiplexer consisting of arrayed-waveguide gratings and double-gate switches," *Electron. Lett.* **32**, 1471–1472.
- Okamoto, K., Okazaki, H., Ohmori, Y., and Kato, K. (1992). "Fabrication of large scale integrated-optic $N \times N$ star couplers," *IEEE Photonics Tech. Lett.* **4**, 1032–1035.
- Okamoto, K., and Sugita, A. (1996). "Flat spectral response arrayed-waveguide grating multiplexer with parabolic waveguide horns," *Electron. Lett.* **32**, 1661–1662.
- Okamoto, K., Syuto, K., Takahashi, H., and Ohmori, Y. (1996). "Fabrication of 128-channel arrayed-waveguide grating multiplexer with a 25-GHz channel spacing," *Electron. Lett.* **32**, 1474–1476.
- Okamoto, K., Takahashi, H., Suzuki, S., Sugita, A., and Ohmori, Y. (1991). "Design and fabrication of integrated-optic 8×8 star coupler," *Electron. Lett.* **27**, 774–775.
- Okamoto, K., and Yamada, H. (1995). "Arrayed-waveguide grating multiplexer with flat spectral response," *Opt. Lett.* **20**, 43–45.
- Smit, M. K. (1988). "New focusing and dispersive planar component based on an optical phased array," *Electron. Lett.* **24**, 385–386.

- Takahashi, H., Suzuki, S., Kato, K., and Nishi, I. (1990). "Arrayed-waveguide grating for wavelength division multi/demultiplexer with nanometer resolution," *Electron. Lett.* **26**, 87–88.
- Takada, K., Inoue, Y., Yamada, H., and Horiguchi, M. (1994). "Measurement of phase error distributions in silica-based arrayed-waveguide grating multiplexers by using Fourier transform spectroscopy," *Electron. Lett.* **30**, 1671–1672.
- Takahashi, H., Tanaka, T., Akahori, Y., Hashimoto, T., Yamada, Y., and Itaya, Y. (1997). "A 2.5 Gb/s, 4-Channel Multiwavelength Light Source Composed of Uv Written Waveguide Gratings and Laser Diodes Integrated on Si," IOOC-ECOC '97 pp. 355–358, Edinburgh, UK.
- Takiguchi, K., Jinguji, K., Okamoto, K., and Ohmori, Y. (1995). "Dispersion compensation using a variable group-delay dispersion equalizer," *Electron. Lett.* **31**, 2192–2193.
- Takiguchi, K., Kawanishi, S., Takara, H., Okamoto, K., Jinguji, K., and Ohmori, Y. (1996). "Higher order dispersion equalizer of dispersion shifted fiber using a lattice-form programmable optical filter," *Electron. Lett.* **32**, 755–757.
- Takiguchi, K., Kawanishi, S., Takara, H., Kamatani, O., Uchiyama, K., Himeno, A., and Jinguji, K. (1996). "Dispersion slope equalizing experiment using planar lightwave circuit for 200 Gbit/s time-division-multiplexed transmission," *Electron. Lett.* **32**, 2083–2084.
- Takiguchi, K., Okamoto, K., Suzuki, S., and Ohmori, Y. (1994). "Planar lightwave circuit optical dispersion equalizer," *IEEE Photonics Tech. Lett.* **6**, 86–88.
- Tanaka, T., Takahashi, H., Hashimoto, T., Yamada, Y., and Itaya, Y. (1997). "Fabrication of Hybrid Integrated 4-Wavelength Laser Composed of Uv Written Waveguide Gratings and Laser Diodes," OECC '97 10D3-3, Seoul, Korea.
- Trouchet, D., Beguin, A., Boek, H., Prel, C., Lermiaux, C., and Maschmeyer, R. O. (1997). "Passband flattening of PHASAR WDM using input and output star couplers designed with two focal points," *Proc. OFC '97 ThM7*, Dallas, TX.
- Tanaka, T., Takahashi, H., Oguma, M., Hashimoto, T., Hibino, Y., Yamada, Y., Itaya, Y., Albert, J., and Hill, K. O. (1996). "Integrated external cavity laser composed of spot-size converted LD and uv written grating in silica waveguide on Si," *Electron. Lett.* **32**, 1202–1203.
- Veerman, F. B., Schalkwijk, P. J., Pennings, E. C. M., Smit, M. K., and Verbeek, B. H. (1992). "An optical passive 3-dB TMI-coupler with reduced fabrication tolerance sensitivity," *J. Lightwave Tech.* **10**, 306–311.
- Vengsarkar, A. M., Miller, A. E., and Reed, W. A. (1993). "Highly efficient single-mode fiber for broadband dispersion compensation," OFC '93 Postdeadline paper PD13, San Jose, CA.
- Yamada, H., Takada, K., Inoue, Y., Ohmori, Y., and Mitachi, S. (1996). "Statically-phase-compensated 10 GHz-spaced arrayed-waveguide grating," *Electron. Lett.* **32**, 1580–1582.
- Yamada, Y., Takagi, A., Ogawa, I., Kawachi, M., and Kobayashi, M. (1993). "Silica-based optical waveguide on terraced silicon substrate as hybrid integration platform," *Electron. Lett.* **29**, 444–445.



Photographic Processes and Materials

P. S. Vincett

Xerox Research Centre of Canada, Ltd.

M. R. V. Sahyun

University of Wisconsin-Eau Claire

- I. The Imaging Chain
- II. Silver Halide (AgX)—Black-and-White
- III. Silver Halide—Color
- IV. One-Step and Other Related Processes
- V. Photopolymers
- VI. Digital Photography
- VII. Nonimpact Printing Technologies (NIP)
- VIII. Conclusions

GLOSSARY

Amplification Extent to which the light-sensitivity of a photographic process exceeds that of a hypothetical process in which light *directly* produces or destroys the light-absorbing material of an image. Amplification usually occurs in the development step.

Density In photography, the ability of a material to absorb light, usually as a result of exposure and development. Accurately, density refers to $\log_{10}(I_0/I)$, where I_0 and I are the light incident on, and transmitted by, the material. Often referred to as optical density.

Development Process of making the latent image visible by providing appropriate energy or matter to the imaging material; this couples to the latent image differently from the rest of the material.

Emulsion As (incorrectly) used in photography, the active (light-sensitive) layer of the photographic material.

Fixing Removal of the light-sensitivity of the emulsion after exposure (usually after development and usually by removing the light-sensitive material).

Gamma A contrast index, measuring of the rate at which the final image density increases (or decreases) as a function of the exposure. More accurately, gamma is the slope of the nearly linear region that usually occurs in the plot of density versus the logarithm of exposure known as the “characteristic curve.”

Latent image Invisible image-precursor formed by the imaging light; it makes a region of the material developable.

Latitude Ability of a material to provide an acceptable

image when various imaging parameters, particularly the exposure, are varied.

Oxidation Reverse of reduction, particularly when an element or compound combines with oxygen.

Photoconductor Insulator or semiinsulator whose electrical conductivity is changed (usually increased) when light is shone on it.

Photon Minimum packet of light, having some of the characteristics of a discrete particle.

Pixel A picture element. In most digital imaging processes the image is divided up into a multiplicity of pixels; within each pixel the image is assumed *de facto* to display uniform characteristics, e.g., color, tone, etc. These characteristics are sequentially digitized on a pixel-by-pixel basis to form the record of the image.

Positive/negative and positive/negative-working Positive working process or material renders the lighter and darker objects of an original lighter and darker, respectively, in the image. A negative-working system does the opposite. A positive is an image whose lighter and darker areas correspond to the lighter and darker areas of a typical original, such as a scene or a document; a negative is the opposite.

Reciprocity failure Failure of a photographic response to be dependent only on the integrated illumination energy but not on the rate at which the energy is supplied. There is then no longer a reciprocal relationship between the exposure intensity and exposure time for a given response.

Reduction Chemical process of providing one or more electrons to the metal ion in a chemical compound so as to reduce its oxidation state, often returning it to the metal by removing nonmetallic elements. For example, a silver halide can be reduced to silver.

Sharpness Used in this article to include the various factors related to image definition of a material. These include acutance, a measure of the faithfulness with which the edge of an object is rendered, and resolution, a measure of the ability of the material to keep closely spaced objects separate in the image. These factors are related but do not correlate perfectly.

Speed A measure of the photographic sensitivity of a material after processing. It is usually expressed as the reciprocal of the exposure required to achieve a given density, and is normally expressed in ISO units.

Silver halide Compounds or mixed compounds of silver with chlorine, bromine, and/or iodine.

FOR OUR PURPOSES, photography may be defined as the creation of a chemical or physical spatial-image pattern by the action of a corresponding pattern of visible light (or other radiation) on a sensitive material; the pattern

is either visible and more or less permanent or, much more commonly, is made so by a separate development process applied to the pattern.

The present article discusses the processes and materials used in photography, with an emphasis on the fundamental mechanisms, but does not describe the corresponding hardware (cameras, etc.), except where this is necessary to understand the process. The creative aspects of photography are also not discussed. While silver halide materials are still important, much more attention must now be given to newer technologies that have come of age at the start of a new century.

I. THE IMAGING CHAIN

Photographic materials and processes are used in the imaging chain, that sequence of operations by which an image is: (1) captured; (2) processed; (3) stored; (4) communicated; and (5) displayed.

An image for the purpose of this entry may be defined as a two-dimensional array of information, organized according to its own internal system of logic. Such an image may be, among other things, a picture of reality, a document, or an abstract work of art.

A. Image Capture

Image capture may involve photochemical or digital or analog electronic technology. Traditional photography is an example of the former; digital photography is an example of the latter. Image capture usually occurs in some type of camera, which optically focuses the image onto the capture device or medium. In astronomical photography, the entire telescope serves as the camera, while in medical radiography, the image is captured initially by a phosphor screen which converts x-radiation to light, then it is recaptured by a piece of photographic film in contact with the screen.

Increasingly, however, various sorts of scanners serve as image capture devices. The original image, e.g., document or art work, is scanned with a light source focused to a spot, and the absorption of that light by the original image is sensed by a light measuring device, e.g., photomultiplier. As a result the light-absorbing characteristics of the original are converted to an electronic record on a spot-by-spot (pixel) basis. Likewise where photographic film or paper may be the capture (or display) medium, the image may be created on a spot-by-spot basis, using a modulated laser beam to "write" the image on the sensitive medium. Image capture technologies are a major focus of this article, and as will be seen below, many of them also double as display technologies.

The quantity of light required for exposure of a photographic system may be expressed in terms of ergs/cm^2 ; this is numerically equal to mJ/m^2 . The higher the number, the less sensitive the material. For comparison, a *very* rough rule of thumb would be that arithmetic ISO sensitivities (ASA) are equal to the reciprocal of the exposure in ergs/cm^2 required to obtain slight development; reasonably complete development, however, may need roughly 10–100 times more exposure. (This correlation is a rough one, since ISO sensitivities are most appropriate for white-light sensitivity, while ergs/cm^2 numbers usually refer to the exposure required at the most sensitive wavelength.) Sensitive AgX emulsions (say ISO 1000–10,000) thus require exposures for *slight development* in the region of 10^{-3} – 10^{-4} ergs/cm^2 while the very highest sharpness AgX materials need as much as 10^3 ergs/cm^2 . A typical xerographic office-copying photoreceptor needs about 0.5 and 5 ergs/cm^2 for slight and complete development, respectively, and the XDM electrophotographic film system is presently several times more sensitive. A hypothetical system that formed one molecule of an efficiently absorbing dye for every photon absorbed but had no amplification and no sources of inefficiency would need energy in the 10^5 ergs/cm^2 region. Many insensitive systems with little or no amplification are actually in the 10^6 – 10^7 ergs/cm^2 range.

B. Image Processing

Image processing involves the manipulation of the imagery as recorded for a variety of purposes. These may include enhancement of the esthetic qualities of the image, analysis of the information content of the image, and compression, i.e., reduction of the amount of data which has to be stored or communicated in order to represent the image. A modicum of image processing is built into all photographic materials and processes, owing to their limited response characteristics, as well as the way in which they are employed. Modern materials have been designed with this realization in mind, and their response characteristics have been carefully engineered to match the requirements of their intended applications. Thus different color films are available for exposure under daylight and under incandescent lighting, in order to obtain a photograph in either case in which the image colors seem “natural.”

Much more sophisticated image processing can be carried out with digital imagery. It is now not uncommon for photographs on film to be digitized using a scanner; the now digital image is then processed to optimize contrast, color balance, and sharpness. The processed image can then be output using a digital printing technology to obtain a picture indistinguishable from a photograph produced entirely by photochemical means but with im-

proved tone and color reproduction characteristics. This sort of image processing, which occurs in a hybrid, i.e., not completely photochemical and not completely digital, imaging chain, is becoming increasingly commonplace in consumer photofinishing. In this case the customers may receive not only negatives and prints of their pictures, but digital files on floppy or compact disks as well. A common file format used at the time of this writing for these disks, e.g., in the Fuji digital photofinishing system, is JPEG, which involves a high degree of image compression. Thus the information content of a representative 35-mm negative in digital terms is about 20 MB, greater than a frame of HDTV, while the JPEG file may be typically 40 kB, i.e., nearly three orders of magnitude less information than the 35-mm negative and an order of magnitude less than a frame of conventional (NTSC) broadcast television. The latter image is, however, quite suitable for computer monitor display or small prints, but not for enlargement to, say, a 16×20 inch print, which should be realizable from any good 35-mm negative. The JPEG format is currently popular with consumers who wish to transmit their photographs over the internet, since file sizes are suitable for transmission over low bandwidth networks in reasonable times. The highly compressed file format also allows the entire output of a 36 exposure roll of 35-mm film to be stored on one floppy disk. Since most image processing is carried out digitally, details of the technologies involved are beyond the scope of this entry.

C. Image Storage and Communication

Image storage and communication may occur at any stage of the imaging chain. Thus, in the illustration above, the camera-exposed negatives, the digital record of the enhanced and compressed imagery on a computer disk, and the final display prints are all examples of image storage, each occurring at a different point in the chain. Image communication may involve physical transfer of the image record, whether in a “hard” (photographic negative or print) or “soft” form (recorded on computer memory media), or electronic transmission of digitized imagery. In the latter case, owing the large amounts of information that images represent, image compression is particularly critical, especially when low bandwidth channels are used. Image communication, like image storage, may occur at any point in the imaging chain.

D. Image Display

Like storage media, displays may be “soft” or “hard.” Computer monitors and projection systems are examples of the former. Photographic prints and transparencies, whether created by optical or digital means, are examples

of the latter. Hard display technologies will be the other main focus of this article.

A special case of image display is represented by stereoscopic photography. Steroscopic effects can be obtained by making two photographs from viewpoints separated by about the distance between the eyes, and presenting them independently to the two eyes, for example, by using two independent viewing boxes for transparencies. Alternatively, monochrome images can be printed in two different colors and viewed through spectacles of similar colors. Color stereo pairs can be placed over one another in polarizing layers and viewed with polarizing spectacles, or the stereo pair (or more than two images) can be split into narrow interlaced strips and viewed more or less independently through a carefully aligned grid of microlenses placed on the composite (lenticular photography).

II. SILVER HALIDE (AgX)—BLACK-AND-WHITE

Most commercially available photographic materials employ silver halide (AgX) “emulsions” as the light-sensitive element. A black-and-white (b/w) “emulsion,” which is actually a solid dispersion, consists mainly of small AgX crystals (“grains”) and a protective matrix, usually gelatin (Fig. 1). The end product of exposure and development of the AgX crystals is an image whose dark areas usually consist of a large number of small grains of silver (Ag). Exposure, in a camera or other suitable device, gives rise to an initial, minute, invisible latent image (LI) in some or all of the exposed crystals, the LI consisting typically of just a few atoms of Ag per crystal. Subsequent development with chemical reducing agents can convert such crystals to a mass of opaque metallic Ag, thus providing amplification factors up to around one billion (10^9).

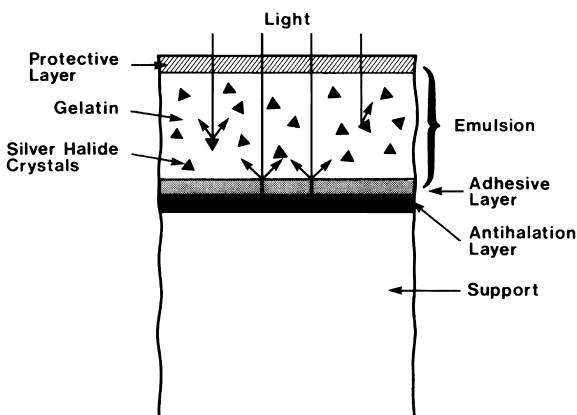


FIGURE 1 Schematic cross section of a typical black-and-white AgX emulsion.

The visible darkening is determined by the exposure and the extent of development, and depends on the number of crystals developed in a particular area and their size. Non-exposed crystals are left largely unreduced. Typically, the unchanged AgX is then dissolved away (“fixing”); then the emulsion is washed to remove the remaining chemicals and is dried. Since more density is thus usually produced where the original was lighter, the process is usually negative-working. A positive can then be made by exposing another AgX emulsion (often on paper) through the negative and again developing, fixing, and washing. Some materials, however, give rise to a direct positive, and others give a positive by modified so-called “reversal” processing. In conventional color processes, the developer which is oxidized in the process of reducing the AgX is generally used to form or destroy organic dyes, and the Ag itself is usually removed. Instant color processes typically rely on the oxidation of the developer to cause changes in the mobility of dyes, which then drift (diffuse) to and are captured by a receiver layer.

A. Conventional B/W Silver Halide Emulsions

Typical b/w emulsions contain 30–40% by weight (about 7–12% by volume) of AgX in gelatin; smaller amounts of other compounds are added for specific purposes as described later. The AgX crystals are typically mixed ones, usually containing silver chlorobromide, iodochloride, or iodobromide, depending on the intended use of the emulsion. The proportion of chloride is higher where rapid processing (reflecting its higher solubility) is important and its lower sensitivity can be tolerated, as in printing papers; more sensitive emulsions use iodobromide with about 1–10% iodide.

1. AgX Crystals

To a first approximation, the AgX crystals act as independent units during exposure and development, and so their size and size-distribution exert a profound influence on the photographic properties of the emulsion. An entire crystal can be developed from one LI speck, which must contain a certain minimum number of Ag atoms; this number (and the number of absorbed photons required to produce it) depends very little on crystal size, at least up to a point. Other things being equal, the number of photons absorbed by a given crystal during a given exposure increases proportionately with its size; the number of incident photons per unit area required to form the minimum LI therefore falls rapidly as the crystal size increases, so the sensitivity of the material increases. Its sharpness will also fall as the crystal size is increased, although sharpness is also determined by light scattering within the emulsion, which

depends on factors such as its thickness but which can be controlled to a certain extent by use of dyes, called acutance dyes, added to the gelatin phase. Furthermore, as the spread of crystal sizes is increased, the variation of sensitivity among crystals increases; other things being equal, the gamma of the emulsion is thus decreased and its exposure latitude and tonal range increased. Other factors reducing the gamma include crystal-to-crystal variations in the number of absorbed photons required to render the crystal developable, changes in the intensity of the light as it passes through the emulsion (in which the AgX is usually distributed randomly), and to a small extent even statistical variations in light absorption caused by the random arrival of the photons themselves. This latter factor becomes of major importance in the case of direct exposure of AgX materials to high-energy radiation, such as X-rays. The gamma is also affected by the development conditions. The mean crystal size in high-sensitivity emulsions is in the region of $1\text{ }\mu\text{m}$, while for an ultrahigh-sharpness emulsion the mean may be about $0.05\text{ }\mu\text{m}$. It is important to appreciate that there is always a reciprocal relationship between image quality (graininess, sharpness) and photographic speed (sensitivity). Over the past decade technological advances have allowed the image quality previously typical of ISO 100 films to be realized in ISO 400 products. Much room for further improvement in this trade-off remains.

The developed image from a given crystal, however, may be significantly larger (frequently 2–10 times) than the crystal itself. The spread of crystal sizes is typically of the same order as the mean size. A high-sensitivity emulsion is typically about $7\text{ }\mu\text{m}$ thick, and it may contain roughly 5×10^8 crystals per cm^2 .

Silver halide grains are typically cubic, with six faces, or octahedral, with eight faces, in habit. Crystals of mixed habit (cubo-octahedral) with 14 faces may also be obtained. The silver ion concentration, usually expressed as pAg ($= -\log_{10}[\text{Ag}^+]$) is manipulated during emulsion production (below) to select the desired grain habit. Grain growth controlling additives, commonly organic compounds, may be used in small amounts during emulsion production to further enhance the crystallographic selectivity of the process. Intrinsically, cubic grains are more photosensitive than octahedral ones. After chemical and spectral sensitization, however, the situation is often reversed.

Most modern photographic materials utilize core-shell grains. In these crystals, a core of one halide composition is formed, and then, in a second stage of emulsion production, a shell of differing halide composition is deposited, completely encapsulating the original grains. The iodide content of the core phase is usually higher than for the shell phase. Recently, experimental emulsions have been

disclosed in which the core is an inert salt, e.g., CaF_2 , so that only the shell is radiation sensitive. One aim of this work has been to produce emulsions with high light sensitivity but reduced sensitivity to cosmic radiation, for extraterrestrial application. They have reportedly been used successfully in the Russian space program.

Recently, a major change in emulsion-making technology, the T-grain (T for tabular) system was introduced and used for the first time in the Ferrania ISO 800 color negative-working film, and more recently in other films to improve sensitivity/sharpness. The AgX crystals are made flat and thin (i.e., tablet-shaped) in contrast to the more symmetrical crystals of conventional emulsions. They tend to lie parallel to the layers. The increased surface-to-volume ratio enables better dye sensitization than heretofore, because such sensitization is limited by the quantity of dye that can be incorporated into no more than a monolayer on the crystals. Thus, increased sensitivity can be attained without increasing the “depth” of the crystals, which would increase the thickness of the emulsion and degrade the image sharpness via light scattering. Furthermore, the tendency for the crystals to lie parallel to the layers also tends to cut down on lateral light scattering. Furthermore, the shape of the crystals improves the covering power of b/w images, and the high surface area improves processing speed. Dye-sensitization in the blue can also produce substantial sensitivity advantages with these crystals.

2. Gelatin

Gelatin is a high-molecular-weight natural protein made by the hydrolysis of collagen, which in turn is extracted from the hides and bones of animals, particularly cows. When sufficiently concentrated aqueous solutions of gelatin are cooled, they set reversibly to form gels, from which the water can subsequently be removed; this layer can later be penetrated and swelled by developing solutions, although the swelling can be prevented from becoming excessive, especially at the higher processing temperatures used for rapid development, by the addition of hardeners to induce permanent crosslinks between the gelatin chains. The gelatin acts as a mechanical support and also fulfills various other functions: it attaches strongly to the AgX crystals during emulsion-making, permitting their growth while preventing coagulation; during development, it permits chemicals to penetrate the emulsion, but it retards the reduction of the AgX, thus facilitating the selective development of the exposed crystals; it also stabilizes the LI against oxidation (without gelatin, the LI can decay in seconds in air).

The gelatin must be purified as far as possible before use, but controlled amounts of trace compounds may then

be added. Gelatin absorbs water vapor, so the stability of the emulsion depends on the relative humidity, and fungus and bacteria growth may occur under appropriate conditions. Phenolic antifungal compounds are usually added to photographic emulsions.

3. Emulsion Production and Additives

The emulsion is formed by mixing solutions of a soluble silver salt and of the appropriate halide(s) in the presence of gelatin. Fine particles of the AgX are precipitated. The precipitation is followed by a ripening stage, during which larger particles grow at the expense of smaller ones, and the average particle size increases; this occurs principally because the solubility of small particles is larger than that of large ones, owing to the greater surface energy of the former. If the composition of the bulk solution is held fairly constant via continuous addition of the original solutions (double-jetting), a quasi-equilibrium situation can arise and all the AgX crystals grow to a similar size. If one solution (usually the silver salt) is added continuously to the other (single-jetting), however, conditions change throughout the precipitation, a range of crystal sizes is formed, and the final emulsion will have a lower gamma. After removal of by-products, the emulsion is treated with minute amounts of sensitizers and other additives such as coating agents, antistatic materials, hardeners, plasticizers, stabilizers, disinfectants (to prevent bacterial growth), and development modifiers; it is then coated onto a base and dried. A more sensitive layer may be coated over one or more less sensitive ones, so that reasonable exposure latitude is obtained; alternatively, the high- and low-sensitivity components may be blended. Subsequent aging leads to significant sensitivity changes for at least several days; in fact, it is said that one cannot predict the exact characteristics of a given emulsion in advance of this aging process.

4. Substrates

The most common film substrates (bases) are cellulose triacetate and polyester, typically a few thousandths of an inch thick; the latter is used principally when dimensional stability or high strength are important. An adhesive subcoating is used between the emulsion and the base, and an anticurl layer may be placed on the rear of the base. Protection is often provided against halation (reflection at the base/air interface of light that has passed through the emulsion, leading to a spurious image where it reenters the emulsion). This can be done by coating one side of the base with a light-absorbing material or by incorporating it in the emulsion; it is removed or bleached during processing of the film. If an increase in background den-

sity is acceptable, a gray base may be used instead. Blue-tinted polyester film base is normally used for medical diagnostic films, mostly on the basis of tradition. An anti-static layer and an abrasion-resistant topcoat may also be incorporated.

A paper substrate (with high wet-strength and appropriate chemical inertness) is generally used when transparency is not required. Barium sulfate (baryta) in a gelatin binder may be applied first to provide a reasonably smooth surface and to improve opacity and brightness, and colorless blue-fluorescent brighteners may also be added to the paper to improve "whiteness." More recently, "resin-coated" (RC) paper, which is coated on both sides with a thin layer of polyethylene, has gained wide acceptance, primarily because it absorbs far less of the developing solutions and thus reduces the processing time, making machine processing rapid and easy. There is also much less tendency for the paper to tear when wet; furthermore, the curl of the paper can be controlled by the polyethylene, the dimensional stability is improved, and desired gloss levels can be obtained without separate postdevelopment processes. Pigments, such as titanium dioxide, in the polyethylene can be used to control the background color and brightness and may also improve the image sharpness. Finally, when extreme dimensional stability or flatness is required, e.g., in holography, glass-plate substrates are used.

B. Latent Image (LI) Formation

On exposure of the emulsion to low-intensity light such as that transmitted by a camera lens, an LI consisting of a very minute cluster of elemental Ag, containing just a few Ag atoms, is formed in some or all of the exposed AgX crystals. This LI is generally located on the surface of the crystals; although it may contain only 1 part in 10^7 – 10^{10} of a crystal's Ag, it renders the whole crystal developable. That is, during subsequent chemical treatment, reduction of the crystal to metallic Ag (commencing at the Ag speck) is significantly faster than it is in the absence of the LI. The number of photons that must be absorbed by a crystal to make it developable is on the order of 6–15 for sensitive emulsions, although a significant portion of the crystals may be made developable by four photons each; when certain special sensitizers are used under special conditions, every crystal may be made developable by only two or three absorbed photons each. However, the sensitivity of AgX is limited by the fairly low absorption (typically of order 10%) of incident photons by each small crystal.

Remarkably, the most basic mechanisms of LI formation are still controversial and the subject of ongoing research. We shall describe what is perhaps the most

widely held theory. First, however, we must give some background.

1. AgX Electronic Structure

Silver halides are ionic crystals consisting of a regular cubic lattice of Ag and halide ions together with a small proportion of defects, such as Ag ions that have been displaced from their regular lattice position to another “interstitial” position (the Ag ions are much smaller than the halide ions), and the corresponding vacancy in the lattice. Although the lattice itself is fairly rigid, such interstitials and vacancies are fairly mobile, because they can jump through the crystal, one lattice spacing at a time, without other ions having to move substantially to make room for them. These defects are more likely to occur at the crystal surface. Since the interstitial silver ions, but not their vacancies, can diffuse into the crystal itself, the surface acquires a negative charge, called the space charge, with an associated electric field. When ions are brought together in a crystal, their sharply defined electronic energy levels, corresponding to the energy states of the outer-shell electrons, are broadened by the interactions between them into a series of bands; electrons within these bands are not necessarily each localized on one ion, and may belong instead to the crystal as a whole.

In AgX, all the energy levels of the highest occupied energy band (the valence band, VB) are filled with electrons. Although the electrons are delocalized, any increase of one electron's momentum under the action of an electric field must be balanced by a reduction of momentum of another one, because on spare energy states are available. Thus, in a perfect crystal, essentially no electronic conduction occurs. The next higher energy band (the conduction band, CB) is essentially empty in a perfect AgX crystal, being too far (2.5 eV or so) above the VB for significant thermally excited transitions to take place from the VB. In a perfect crystal, there are no energy levels in the gap between the VB and CB, but in real microcrystals there are a few such levels associated with crystal defects. These defects states are critical to photosensitivity. It has been shown experimentally that the purest, most perfect crystals of AgX are virtually insensitive to light.

When a photon is absorbed by an AgX, there is a high probability of formation of an “electron-hole” pair: an electron is raised from the VB to the nearly empty CB, in which it can move more freely, while the “hole” is a way of conceptualizing the behavior of the previously full VB from which the one electron has been removed. It may be thought of as electron-like particle with a positive charge, and is sometimes, in fact, called a “defect electron.” Since the depleted band of electrons responds to an electric field, for example, like the full band *less* one electron,

and since the full band gives no net response, the overall effect is like that of *minus one* electron moving in the VB. Such electron-hole formation occurs in all photoconductors upon light absorption, but in most cases the electrons and holes recombine rapidly in the absence of an external electric field; the space charge may provide this field. In AgX, however, the combination of electrons with Ag ions to form Ag atoms can be at least comparable to recombination. Electron-ion combination cannot occur if both the electron and ion are mobile, since the total energy of the ion-electron system would make them separate again, nor can it occur if the Ag ion is in a lattice position, since there is not enough energy to overcome the force holding the ion in the lattice. However, it can happen if the electron or ion is at a trap; an electron-trap, for example, is a position in the crystal lattice (usually associated with a defect) where the electron's energy is lower than it is when free and where it may consequently reside for a time. One reason that iodide is included in most commercial emulsions is that iodide ions or ion clusters in an AgX lattice trap holes, thus interfering with the energy-wasting recombination process.

2. Step-by-Step Mechanism

One of the most important points that any theory of the LI formation process must explain is how it is that photons that are absorbed through out each AgX crystal become concentrated, in the sense that the Ag formed from them becomes concentrated in one or just a few LI specks. The most widely accepted theory is the “step-by-step” mechanism, which had its origins with Gurney and Mott as long ago as 1938. The electron is considered to move through the crystal in the CB until it reaches a trap, where it remains for long enough that a mobile Ag ion can move to it because of the attraction of their opposite charges, and combine with it to form an Ag atom. Once this happens, the same thing can occur again at the same place, until a tiny speck of Ag is built up to form an LI. In fact, a stable Ag nucleus of more than a certain size is itself a deep electron trap, so beyond this size further electrons are even more likely to be trapped at the growing nucleus rather than elsewhere. Thus, the mobility of the electrons and their subsequent trapping provide the concentration mechanism. In practice, the all-important trap or “sensitivity center” is usually associated with areas of the surface of the crystal that have been chemically sensitized as discussed below, so that usually the LI is formed primarily at the surface. In fact, unsensitized crystals are fairly insensitive, especially when conventional developers (which do not dissolve the AgX and therefore act primarily on surface LIs) are used. Of course, there may be more than one such trap per crystal, although (since one LI center can render a whole crystal developable) it is clearly more

efficient to have as few LIs as possible. If too many electrons are produced simultaneously, the ionic step of the stepwise growth mechanism doesn't have time to occur, so each electron gets trapped at different trapping sites, and the resulting silver atoms are not concentrated into a useful LI. This effect gives rise to high-intensity reciprocity failure, usually observed under conditions where exposure times are $1\ \mu\text{s}$ or shorter. It is critical to overcome high-intensity reciprocity failure for certain technological applications of photography, e.g., exposure in very rapid laser scanners or high-speed scientific photography using ultrashort (nanosecond) exposures.

It is important to note that one Ag atom, formed as above, is not completely stable against redissociation to the ion and electron, and in fact there is a minimum size (probably two atoms) of the Ag speck for stability (as opposed to developability which requires further growth). Thus, the probability of formation of an LI depends not only on the total number of photons absorbed but also on the time between absorption of successive photons. If the time is too long, the chance of two Ag atoms being formed at the same place at the same time drops, because the first atom dissociates before the second one can form. This gives rise to low-intensity reciprocity failure, an effect well known to anyone who has attempted photography in very low light levels, especially with color films, and particularly to photographers of astronomical objects. In addition, especially in the absence of chemical sensitization or in the presence of oxygen or moisture, electron-hole recombination is a significant factor limiting the photographic sensitivity.

3. Chemical Sensitization

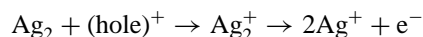
In practical emulsions, the sensitivity and reciprocity behavior are usually greatly increased by sensitizers, compounds that are added during emulsion-making and are adsorbed on the surface of the crystals. "Chemical sensitization" produces a major decrease in the number of absorbed photons needed to make a crystal developable and can greatly decrease low-intensity reciprocity failure. Care must be exercised, however, since too much sensitization can cause the formation of chemical nuclei that can act like LI during development, giving rise to fog (spurious density in nonexposed areas).

The three most common forms of chemical sensitization are based on reaction of the AgX with very small quantities of sulfur compounds, salts of gold or other noble metals, and reducing agents. The premier commercial technique for high sensitivity is a combination of the first two. The sulfur compounds react with the AgX to form some silver sulfide on the surface of the crystals. While catalytic mechanisms have been proposed, and enhanced hole trap-

ping is also possible, it is more commonly believed that sulfur sensitization acts by formation of electron traps, thus diminishing recombination, perhaps by increasing the electron trap depth of an interstitial Ag ion. Moreover, it may stabilize single Ag atoms somewhat against redissociation and leads to less formation of LI in the interior of the crystal where most developers cannot reach it. If too much sulfur is added, the sensitivity decreases again; this has been attributed to formation of multiple LI centers, which causes an inefficient use of photons.

Gold sensitization causes gold to become part of the LI (by formation of a mixed silver-gold cluster) and also part of the sensitivity specks, by formation of silver-gold sulfide. This decreases the size at which the LI becomes developable: a three-atom (or possibly two-atom) gold-containing cluster is developable (e.g., Ag_2Au or AgAu_2), compared with four Ag atoms. Gold sensitization may also improve the stability of Ag atoms and increase the depth of associated electron traps.

Reduction sensitization is produced by heating the emulsion with reducing agents such as dimethylamine borane; they are believed to act by forming very small Ag clusters. While these could act as electron traps and thus constitute LI nuclei, it is believed more likely that they capture photogenerated holes, reducing the hole concentration and hence recombination. Indeed, if a photon produces an electron and a hole, and the hole attacks a reduction center consisting of two Ag atoms, the Ag_2^+ first formed may dissociate spontaneously into two Ag ions and another electron:



The net result is two electrons for one photon. An active area of current AgX research is the development of new sensitizing reagents which can produce this same effect without the drawbacks of reduction sensitization, which normally include elevated fog levels and incompatibility with gold sensitization.

For emulsions designed for high-intensity, short-duration exposure, e.g., in laser scanners, or where ultrahigh contrast is required, e.g., the graphic arts, it is now usual practice to dope the crystals with trivalent metal ions, such as Ir^{+3} , Rh^{+3} , or even Fe^{+3} . These ions provide deep electron traps within the bulk of the AgX grains. Over time, after completion of exposure, these deep traps thermally release the electrons to the conduction band of the crystal, thus simulating a lower intensity, longer duration exposure, and overcoming high-intensity reciprocity failure.

It should be kept in mind that new mechanisms of LI formation and development, either combinations of existing proposals or new ones such as a mechanism based on the redox potential of the AgX crystals, are still being introduced.

C. Spectral Sensitization and Hypersensitization

1. Spectral Sensitization

A pure or chemically sensitized AgX emulsion absorbs significantly (and thus produces electron–hole pairs) only up to the wavelength of its absorption edge, corresponding to the gap between the VB and the CB. For silver bromide and silver chloride, the absorption falls to very low levels by about 490 nm (with a slight extension for iodobromide) and 420 nm, respectively, leading to sensitivity only in the blue region. To extend the sensitivity to other regions, emulsions are often doped with dyes that absorb in the wavelength region desired. Emulsions sensitive from the blue to the green–yellow are often termed orthochromatic, while those with reasonably constant sensitivity throughout the visible region (~ 400 – 700 nm) are said to be panchromatic. Good sensitivity out to the 900-nm region of the near infrared can be obtained, and some response to around 1300 nm. (Beyond this, either the stability of the emulsions becomes very low or indirect methods must be used, e.g., detection by an electrical infrared cell, which then modulates a scanning light source.)

Sensitizing dyes are usually from the polymethine class. The three major subclasses of these dyes are the cyanines (the most common), the merocyanines, and the oxonols. These dyes have high absorption coefficients over narrow bandwidths, a wide range of structures and absorption wavelengths, and are easily purified. Generally, the sensitizing dye molecules are effective only when in close contact with the AgX crystals. The absorption of appropriate visible (or near-infrared) radiation by the dye results in the promotion of an electron from the highest filled electronic level of the dye molecule to one of the lowest unoccupied levels. As explained below, the result is an electron in the CB of the AgX that then undergoes the same processes we have described for directly produced photoelectrons to form LI. While this effect increases with increasing light absorption by a given dye, the dye may also desensitize the emulsion by acting on some of the secondary processes of LI formation; since the increase in absorption can become slower than the increase in desensitization, there is an optimum dye concentration, which generally corresponds to less than complete molecular monolayer coverage of the crystals. This is very important, since it limits the sensitivity obtainable. Because of this, dye-enhanced sensitivities tend to be proportional to the area of the crystals, while intrinsic sensitivities (which depend on bulk absorption by the AgX) tend to be proportional to crystal volume. Both proportionalities tend to saturate at about $1\ \mu\text{m}$ crystal diameter: above this, imperfections in the crystals tend to stop them from acting as one unit.

2. Spectral Sensitization Mechanisms

The mechanism of sensitization is generally believed to be direct transfer of an electron (Fig. 2a) from the excited electronic state of the dye into the CB of the AgX (by quantum-mechanical tunneling through the energy barrier between them and, if necessary, by slight thermal excitation), which may or may not be followed by regeneration of the dye by transfer of an electron from a defect state in the AgX bandgap, leaving a trapped hole. As would therefore be expected, there is a reasonably good correlation between the energy levels of otherwise similar dyes and their sensitizing efficiency. Another possible sensitizing mechanism is indirect transfer of the excited electron by an energy transfer mechanism: here the excitation energy of the excited dye is transferred by resonance to a defect state lying in the bandgap of the AgX, promoting an electron into the CB as before (Fig. 2b). This mechanism has been shown to operate almost certainly under special circumstances but is generally assumed to be of secondary importance in practical situations.

Because of local variations in the relevant energy levels, sensitizing dyes could also desensitize when their excited energy level is below the CB of the AgX, causing them to behave as an electron trap (Fig. 2c); this would possibly allow the electron to combine with oxygen and prevent it from taking part in the formation of metallic Ag. Also, a dye may trap a hole (Fig. 2d) by transferring an electron

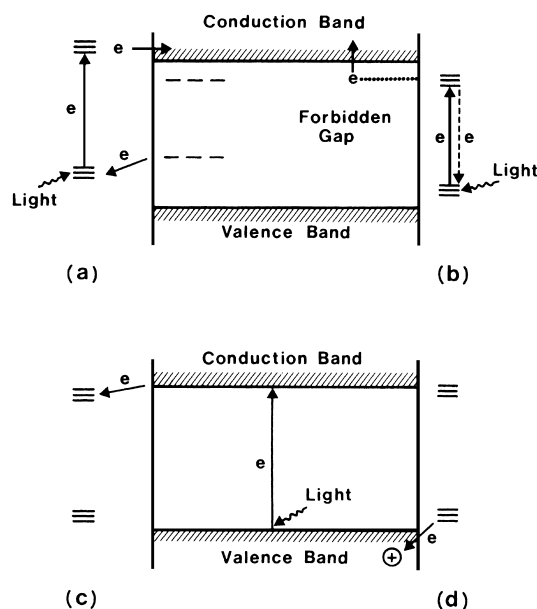


FIGURE 2 Energy levels of AgX and of color-sensitizing dyes, showing sensitization by (a) electron- and (b) energy-transfer, and desensitization by (c) electron trapping and (d) hole trapping. “e” represents electron pathways.

from its ground state to a hole in the VB, or (if the ground state of the dye is near or below the VB) an electron may be injected from the VB to the photoexcited dye, creating a hole in the VB. Such trapped or free holes may then recombine with photoelectrons. On the other hand, an additional electron may be captured by the excited from the VB or from a supersensitizing dye. (See below.) This may change the relative energy levels of the dye and this AgX enough to allow electron transfer.

3. Supersensitization

A second sensitizing dye or related compound adsorbed on the crystal may have a major effect on dye sensitization, and the improvement may be greater than merely an additivity of sensitizing effects. Such “supersensitization” can arise from increased absorption caused by a change in the structure of the dyes on the surface, but the more important form involves an actual increase of the efficiency of transfer of electrons to the CB. A number of theories have been proposed for this. For example, transfer of an electron from the second dye to the first, so that the excited electron and the corresponding hole are separated, may reduce the possibility of recombination and may also change the energy levels favorably. The supersensitizing dye may also act as a discontinuity in the dye layer, trapping the excitation long enough at a given point to give time for transfer to the crystal. In addition, the supersensitizer may provide an energetically favorable site for transfer because of electrostatic interactions with the first dye.

4. Hypersensitization

For applications such as low-light-level astronomical photography, hypersensitization procedures are often applied to the emulsion, which is exposed shortly thereafter to avoid the formation of fog. Such procedures can greatly reduce low-intensity reciprocity failure as well as increasing overall sensitivity, and can allow the recording of images from light that is one-millionth of the intensity required to be visible to the naked eye. Hydrogen, for example, has been used as a reduction hypersensitizer on sulfurplus-gold sensitized emulsions. In favorable circumstances, *all* the AgX crystals have then been made developable by absorption of only two or three photons each. This is to be compared with about 6–15 photons for 50% developability in typical sulfur-plus-gold sensitized silver bromide emulsions, about 20–50 for sulfur sensitization alone, and about 150–300 for the unsensitized emulsion. Part, but not all, of this sensitivity increase is due to initial vacuum-outgassing of the emulsion, since oxygen or moisture causes significant desensitization of AgX emulsions, especially at low exposures. Another significant factor in the increase of

sensitivity by hypersensitization may be the formation of hole traps, including Ag₂, reducing the electron–hole recombination rate. Other forms of hypersensitization include soaking the emulsion in water, ammonia, or silver nitrate solution. These treatments are believed to work by mechanisms involving increased Ag ion concentration in the vicinity of the crystals. The sensitivity gains can be spectacular: one infrared emulsion increased in sensitivity by a factor of 18,000.

5. UV and X-Ray Sensitivity

For decreasing wavelengths below that of the absorption edge of the AgX, the practical limitation on sensitivity is absorption by the gelatin, and below 200 nm special emulsions must be used; they may be very low in gelatin, may have AgX crystals projecting above the top surface of the gelatin, or may contain materials that give longer-wavelength fluorescence. Gelatin becomes transparent again in the soft X-ray region. The electrons first liberated by an X-ray photon (or by a high-energy particle) have too great a velocity to be trapped and to combine with Ag ions. Their collisions with ions in the crystal, however, release secondary electrons of decreasing velocity. For energies of the order of tens of kilovolts, the secondary electrons do not leave the crystal, but it becomes developable. Higher energies may expose other crystals. However, the sensitivity in terms of *incident* energy or incident photons is low, because of the weak absorption of X-rays in the emulsion and the inefficient use of their energy. AgX emulsions can be used directly for some X-ray applications, but when a minimum X-ray dose is required (as in medical applications), the film (emulsion-coated on both sides) is usually sandwiched between two phosphor screens that absorb the X-rays efficiently and fluoresce in color regions to which the emulsion is highly sensitive. Alternatively, the screen can be photographed using a camera, or the energy stored in a special screen can later be made to fluoresce. For short-wavelength X-rays, as used, for example, in various industrial processes, a metal screen such as lead may be used in place of the phosphor. Medical films need high gamma, because the contrast of the X-ray image is low; various contrast-enhancement techniques can also be used.

D. Development

During development, exposed crystals are chemically reduced to Ag much more rapidly than nonexposed crystals, while the developing agent is oxidized. Since more silver is usually produced at exposed crystals, typical b/w development is negative-working. Dye image formation, which is necessary for color photography, depends on the

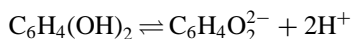
oxidation of the developing agent and will be discussed in a later section.

1. Developers

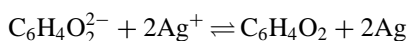
Most of the known developers are organic compounds, although some use certain nonmetallic inorganics or metal ions that can be converted to a higher oxidation state by Ag. Most are aromatic, usually derivatives of benzene, and are typically phenols or amines (or simple derivatives) with at least two hydroxy groups or two amine groups, or one of each arranged ortho or para to each other (adjacent or opposite, respectively) on the benzene ring. Examples of such developers are hydroquinone (1,4-dihydroxybenzene), *N*-methyl-*p*-aminophenol sulfate (Elon, Metol), Amiodl (the dihydrochloride of 2,4-diaminophenol), and derivatives of *p*-phenylene-diamine. Other important organic developers are the reductones (including ascorbic acid) and the pyrazolidones (including 1-phenyl-3-pyrazolidone or Phenidone). Combinations of hydroquinone with Metol and with Phenidone are particularly important.

In addition, practical developers usually contain (1) an alkali (commonly a hydroxide or carbonate of sodium or potassium, or sodium metaborate or tetraborate) to adjust and maintain the hydrogen ion concentration (acidity or pH) to a value appropriate for sufficient development activity; (2) an easily oxidized sulfite, to react with oxidized developing agent (maintaining the reaction kinetics reasonably constant and removing unwanted colored products that may cause staining) and to act as a preservative against oxidation of the developer by air; and (3) one or more antifoggants, such as potassium bromide and certain organic heterocyclics, that retard the development of nonexposed regions of the emulsion more than that of exposed regions and thus improve the discrimination between the two.

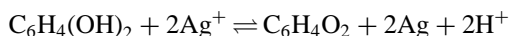
The sensitivity of many developers to pH is very important, since in some situations it can be used to "turn on" the development when required. Many developers have one or more ionizable hydroxy groups and become far more active when this group is ionized. For hydroquinone, for example, the equilibrium between nonionized and ionized forms, and the subsequent reduction of the silver ions of the AgX, can be represented thus:



and



or, overall,



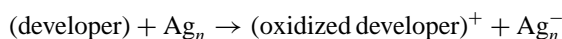
The large concentration of OH^- ions present in an alkaline solution tends to remove the H^+ ions (forming water), thus pushing the equilibria toward the right and favoring development.

2. Types of Development

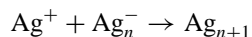
If, as is most common, essentially all the Ag for the growing image comes from the AgX crystals, the process is often called "chemical" development; most of such development is from the "direct" process in which the Ag ions are reduced directly at the Ag/AgX interface. If the developer contains a soluble silver salt to supply Ag to the image, this is called "physical" development, a misnomer since the process is still chemical. Physical development as such is of little practical importance, but "solution physical development," in which Ag ions from exposed or nonexposed AgX crystals pass into solution and are then reduced elsewhere at an LI or a growing visible image, is very important in some processes and often plays some part in normal "chemical" development. The balance between these various processes can be changed by the concentration of AgX solvents, such as sodium sulfite or sodium thiosulfate, in the developer.

3. Mechanisms

In conventional negative-working development, reduction of the exposed crystal starts at the one or more LI nuclei on the crystal surface and proceeds until the entire crystal is reduced or the developer is removed. The electron microscope reveals, however, that the final black Ag image usually does not have the same shape as the crystal, but rather is a tangled mass of 15- to 25-nm-diameter filaments of Ag, roughly resembling a piece of steel wool. Very fine crystals may give only one or a few filaments each; such images tend to be somewhat colored, because of wavelength-dependent light scattering. Exposed crystals are reduced to Ag much more rapidly than nonexposed crystals, because the LI nuclei in the former act as catalytic centers for reduction of further AgX (i.e., it is easier to form Ag where Ag already exists). The new Ag in turn accelerates the reduction of further Ag ions, and the process of reduction thus accelerates as an autocatalytic reaction. The presence of Ag or LI has been shown to reduce the activation energy for the reduction process by several kilocalories per mole, but the exact mechanism of the catalysis is not fully understood. An electrochemical mechanism has been proposed, in which the growing Ag nucleus (or the LI) acts as an electrode to permit both an anodic partial process, in which the nucleus (Ag_n) gains an electron:



and a cathodic partial process, in which Ag ions are neutralized at the nucleus:



A second possible mechanism is adsorption catalysis, in which one or more of the reactants is first adsorbed by the Ag LI or at the Ag/AgX interface. An intermediary complex is formed between the developing agent and the Ag ion. This complex interacts with the Ag (or the adsorption site on the Ag surface may form part of the complex), and the complex is distorted, lowering the activation energy required for the reduction reaction. Reaction then occurs by electronic rearrangement, and the complex decomposes to Ag and oxidized developer. Of course, adsorption could also be involved in the electrochemical mechanism. Adsorption by Ag, AgX, or both has actually been demonstrated for many developers. It is possible that the mechanism changes from one to the other during development, especially since a very small LI might not be able to act as an electrode in the early states. It has also been suggested that, in the very early stages of development, developers may be able to transfer an electron indirectly to the conduction band of the AgX, giving rise to growth of the LI just as if the electron had been produced by light absorption.

The formation of Ag filaments is also not fully understood, but elongation of the initial spherical nuclei may start when new Ag atoms do not have time to diffuse far enough to maintain a spherical form. This elongation may be extended by the electrostatic repulsion of electrons in the Ag to the ends of the spheroid, giving rise to a substantial electric field in the AgX and concentrating further growth at the ends. Physical development often gives rise to much less filament formation than chemical development.

4. Rate of Development

The rate of development (and of subsequent processing) of the exposed crystals is an important practical consideration and is influenced by the rate of the primary processes, the rate of swelling of the gelatin when it is immersed in the developer, and the rates of diffusion of the processing chemicals in the swelled gelatin. Rates are increased by agitation of the solutions and by increase of temperature. Too high a temperature, however, can cause excessive swelling of the gelatin and consequent physical damage to the emulsion; this and subsequent damage can be decreased by hardening (cross-linking) the gelatin during manufacture or processing. Alternatively, a high concentration of a salt that restrains swelling, such as a sulfate, can be added. Certain development accelerators can increase the rate of development, sometimes by neutralizing the negative charge found on AgX crystals to permit the easier access of negatively charged developer ions.

Some combinations of developers (e.g., hydroquinone with Metol or Phenidone) show superadditive rates of development. In these cases, one agent is active as the undissociated molecule or as the singly charged ion and is therefore not too affected by the charge on the crystal; it also typically has a chemical affinity for the AgX and/or the Ag and adsorbs on it. This material therefore rapidly initiates development, forming a fairly stable oxidation product that probably remains adsorbed. The second developer is doubly or triply charged in its most active form when used alone but can rapidly reduce the oxidation product, regenerating the first developer and removing any tendency for the oxidation product to retard development. In effect, the first developer finds the LI and acts as a "wire" for electrons from the second developer.

The effective sensitivity of AgX emulsions can be increased by "pushing" (developing for longer times, at higher temperatures, or with more active developers). However, such processes are limited by fog formation, by changes in the characteristics of the emulsion (for example, in the gamma), and by large increases in the graininess of the image.

5. Fine-Grain Development

When high magnification will be used to view or print the finished image, high sharpness is required, and special "fine-grain" developers can be used. These contain AgX solvents such as sodium sulfite and use low to medium rates of development, so that solution physical development is significant. Fine-grain development results in a decrease of the average size of the Ag grains, perhaps by separating several LI centers from the AgX crystal and allowing them to grow separately by solution physical development; more importantly, there is a decrease in the size and probability of formation of randomly distributed clumps of Ag grains.

Fine grain is particularly important for microfilm applications, in which documents or computer output displays are photographed at reduction ratios of typically about 25–50; this is still (and will probably remain) a widely used and inexpensive way to store large quantities of data. In newer micrographic systems, documents are digitized, then copied onto film at high reduction using either a laser scanner or CRT exposure device. Such systems may be preferable to digital storage both because the microfilm has a higher data storage density than most digital memories, and because the resulting record is eye-readable. The question of storing an image of the whole document versus storing just the digital data corresponding to the information contained by the document, receives a different answer depending on the content and character of the document, e.g., whether it is written in alphanumeric symbols

or Kanji characters. Photothermographic films (below) are increasingly used for such reduced format document storage applications.

E. Post-Development Processing

1. Fixing, Washing, and Drying

After development, other processes are usually employed to render the image stable, and rinsing and washing steps may be used at various points. Preferably the image is first rinsed in water or in a stop bath or clearing bath to lower the pH in the emulsion (make it less alkaline) and to remove developer chemicals, including oxidized developer, which can cause staining, especially in color materials. The material is then immersed in a “fixing” bath to dissolve out undeveloped AgX. This bath usually contains sodium thiosulfate (“hypo”), or ammonium thiosulfate (“ammonium hypo”) for faster results. It may be buffered acidic to neutralize any alkaline developer carried into it, in which case a preservative (e.g., bisulfite) can also be added to retard decomposition of the thiosulfate to bisulfite and sulfur. In recent years, with the increasing price of Ag and more emphasis on pollution control, there has been increasing interest in recovering the Ag from the fixing bath, which can also increase the bath’s useful life. Finally, washing removes fixing-bath chemicals and soluble Ag ion complexes, and the material is dried as uniformly as possible.

Drying the washed film or paper using warmed air is also a critical step in processing. Overly rapid or severe drying conditions can lead to physical deformation of the emulsion layer or reticulation. Modern rapid processes, as used for medical X-ray films or in minilabs, require drying to occur as rapidly as possible. Making the emulsion layers as thin as possible, along with chemical hardening of the layers (cross-linking) to preclude excessive swelling, enable more rapid drying without image deformation or distortion. A special problem in drying is posed by holographic plates, where dimensional stability of the image is critical, but very thick emulsions have to be employed. The deformation of the layer which occurs even on careful, normal warm air drying causes disastrous reduction in the brightness of the recorded hologram. To this end, an additional washing step, in which the water from the first wash is replaced with alcohol, is now employed in processing full color display holograms. The alcohol reduces the swelling of the layer uniformly, and it can be removed from the layer under milder drying conditions.

2. Image Stability and Modification

If residual thiosulfate is not removed sufficiently, the image will be more or less unstable. In fact, the likely keeping

(archival) properties of a given Ag image are determined by measuring this residue. If washing is incomplete, thiosulfate or its decomposition products can react with finely divided Ag, or silver thiosulfate complexes already adsorbed to the image Ag can decompose to give yellow or brown silver sulfide; thiosulfate can also enhance atmospheric oxidation of the Ag. Complete washing of a film can take as long as 30 min, and porous-paper prints (as opposed to those using the more recent resin-coated papers described earlier) can take much longer. Washing can be speeded-up by initially using salt solutions to displace the thiosulfate. Even after proper fixing and washing, Ag images can be physically or chemically unstable under some environmental conditions.

Various chemical treatments can be applied to the final image to change its color, or to increase or reduce its density, or to improve stability (e.g., by converting the image Ag to the sulfide or selenide). Conversion of the silver image to silver sulfide, often using thiourea as the reagent, is known as sepia toning. Historically, conversion of the image silver to gold, platinum, or palladium was an important procedure, both to render the image more environmentally stable, i.e., inhibit image degradation on storage, and for the artistic effects obtainable thereby. These methods fell into disuse in the second half of the twentieth century, when automated, standardized processing was emphasized. Recently they have been revived, largely because of artistic considerations. Density-enhancement techniques, such as using Ag-catalyzed dye formation reactions, may become more important in the future with the growing need to conserve finite supplies of Ag. In special circumstances, Ag images have even been rendered radioactive, and an image with enhanced density has been formed by the absorption of radioactive emanations in an adjacent emulsion or glass plate. Various retouching methods may also be used to remove or obscure unwanted features; airbrushing, for example, applies a fine localized chemical spray.

F. Obtaining a Positive

Negative-working development is used in most b/w work, often followed by exposure of another emulsion (e.g., an AgX printing paper) through the negative to reverse the sign again and form a positive. Papers of various “hardness” (gammas) can be used to correct nonoptimum density range in the negative or to achieve desired esthetic effects. Variable-gamma papers contain a mixture of emulsions having different gammas and different color sensitivities; the color of the printing light controls their effective gamma.

For some applications, such as b/w motion picture film, it is useful to produce a positive at once. This is known,

somewhat confusingly, as reversal processing (since the sign of the image is reversed during processing), although the process is in fact positive-working. Most commonly, the initial unfixed negative Ag image is removed by a solvent (e.g., a solution of potassium dichromate and sulfuric acid, which converts the Ag to the soluble sulfate). Some of the remaining AgX may also be dissolved away (with sodium thiosulfate, for example) to reduce the ultimate density of the final image and provide clear highlights. A cleaning bath (sodium or potassium metabisulfite) stops the bleaching action and prevents staining. The resulting image pattern of AgX is exposed uniformly to light or to the action of fogging (reducing) materials such as stannous compounds or amine boranes; it is then developed in a conventional developer and may also be fixed.

III. SILVER HALIDE—COLOR

A. Basic Principles

Present-day color films are tricolor systems: they can reproduce most natural colors, within a limited color-saturation range, by adding together suitable proportions of the three subtractive primary colors, yellow, magenta, and cyan in the form of dyes. Tricolor systems can also work by combining the actual primary lights; this is known as additive synthesis and is used in color TV for example. Most color photographs, however, are subtractive: white light (which contains all colors) passes successively through three single-color image layers, each of which absorbs some or all of one of the primary colors: the colors of these dye images are cyan (which absorbs red and transmits blue and green), yellow (which absorbs blue and transmits red and green, which is why yellow filters are often used in b/w photography to darken the rendering of the sky), and magenta (which absorbs green and transmits red and blue, appearing as a purplish pink). In principle, then, appropriate optical densities in the three layers can produce any proportion of the three additive primaries (blue, green, and red) in the transmitted light. With transparencies, the colored light is transmitted directly to the eye, while in photographic prints the light must pass through the image twice and is reflected from the underlying substrate, usually paper.

Recent research on the nature of color vision suggests that more accurate color reproduction in a photographic system should be achieved using four “primary” colors. Some recent high-quality color negative films from the Fuji Photo Film Co., Ltd., have been reported to exhibit maximum response at four different wavelengths. In addition to response maxima corresponding to blue, green, and red light, a blue-green response has been introduced.

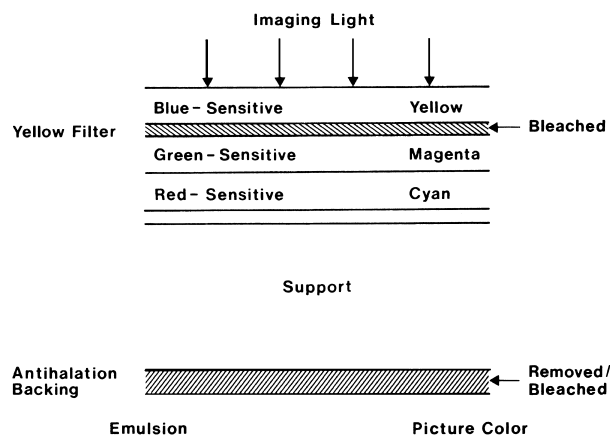


FIGURE 3 Schematic illustrations of the main layers in a typical chromogenic color emulsion. The color sensitivity of the active layers is shown on the left-hand side, and the colors formed after development are given on the right-hand side.

Practical color photographic materials are usually multilayer structures with separate but superimposed red-, green-, and blue-sensitive AgX layers. Most color development processes are chromogenic: the dye images are formed during development by reactions between the oxidation products of the developer and colorcouplers; the latter are special chemicals (which may or may not be colored themselves) contained either in the developing solution or within the layers of the emulsion. It is important to note that these color-couplers and the dyes resulting from their reaction with the oxidized developer are quite separate and distinct from the dyes used to sensitize the AgX to the appropriate exposure colors. A simplified chromogenic color film is shown schematically in Fig. 3, together with the colors that can form in each layer. A cyan image develops where the red-sensitive layer is exposed, a yellow image where the blue-sensitive layer is exposed, and a magenta image where the green-sensitive layer is exposed.

The developing agents used for chromogenic development are generally members of the *p*-phenylenediamine class. Quinonediimines formed by oxidation of the developer react with the couplers, which are often precursors of indophenol or azomethine dyes. Although removal of the oxidation products of the developer by sulfite is not desired in chromogenic development, a little sulfite is usually still added to protect the developer from oxidation by air. As in b/w development, alkali and antifoggant materials are also present.

Thus, with direct development, a pure blue area (for example) in the original scene causes development only in the blue-sensitive layer, yielding a yellow area in the final image; this absorbs blue, so the image is a negative of the original, analogous to a b/w negative. Similarly, a white area in the original scene causes development and dye formation in all the layers, yielding black.

Reversal (positive-working) color materials for slides and motion pictures usually employ b/w (nonchromogenic) development of exposed crystals, followed by reexposure of the remaining AgX and chromogenic development. In the above example, the blue-sensitive layer is then not available for chromogenic development but the other two are; cyan and magenta therefore develop. Only blue is transmitted by both of them, so a positive results. In the area exposed to white light, none of the layers is available for chromogenic development, leaving a clear area; this appears white when viewed normally.

The yellow filter of colloidal Ag (which is bleached during processing) protects the underlying layers (Fig. 3), which are intrinsically sensitive to blue as well as to the colors to which they are sensitized, from blue light transmitted by the blue-sensitive layer. It is crucial that the various layers stay separate and do not undergo unwanted interlayer diffusion of some of the chemicals; this was one of the key problems that delayed the introduction of color films. The dyes used are never perfect matches to the primaries (especially given the other criteria, such as stability and chemical properties, that they must fulfill), and exposure and development errors can cause additional color distortion. Nevertheless, in practice, a reasonable representation of most originals can be obtained. This is done by paying particular attention, both in selecting the dyes and during the photofinishing process, to various colors to which the eye is especially critical, particularly skin-tones.

B. Reversal Processes

Kodachrome, introduced in 1935, was the first commercially successful color film, although some earlier additive systems with poor brightness and sensitivity had been marketed. It forms reversal (positive) dye transparency images in successive color development steps, using soluble couplers in the developer. Typical development steps with such materials are as follows. After b/w (nonchromogenic) development of exposed crystals, the remaining AgX is uniformly exposed to red light to make the remaining parts of the red-sensitive layer developable; this layer is then developed (using a cyan-forming coupler in the developer) to form the positive cyan image. The emulsion is then exposed to blue light to make the remaining parts of the blue-sensitive layer developable, and this layer is developed to form the yellow image. Finally, the emulsion is exposed to white light (or to a chemical fogging compound) to make the remaining AgX (in the green-sensitive layer) developable, and the corresponding magenta image is developed. (Any unused AgX in the red- or blue-sensitive layers is removed by nonchromogenic “auxiliary” development immediately after the respective color development steps, to prevent color contamination by

magenta in the yellow and cyan layers.) Processing finishes with bleaching (to remove the Ag), fixing, washing, and drying steps. Water rinses, stop baths, and gelatin-hardening baths may be used at appropriate points in the process.

The bleach step, which may be combined with the fixing step in a so-called “blix,” is important to both the economics of color photography and its sustainability. It enables essentially all the silver to be recovered with high efficiency for recycling and reuse in photographic films and papers. Silver recovery is typically accomplished by passing spent bleach and fix baths through a column packed with steel wool, on which the dissolved silver plates out in metallic form. Electrochemical recovery processes are also employed. At the same time the bleach bath has been a major source of environmental problems for the photographic industry. Historically, Fe^{+3} compounds, including ferricyanides, have been used as bleach reagents, rendering photographic waste highly toxic even after silver recovery. Recent technological advances have allowed introduction of more benign bleaching reagents, e.g., peroxides. The final washing step is now often replaced by a stabilization treatment using special chemicals; this greatly cuts down on the water used and the effluent produced, which is especially important for small, self-contained “minilab” processing systems.

Processing of reversal emulsions is simpler if colorless couplers are incorporated right into the emulsion rather than being provided by the developer. This, of course, must be done in such a way that there is little interlayer diffusion of the couplers, or the “wrong” colors will be produced. Materials of this type include Agfachrome, Fujichrome and Ektachrome: such a material was first achieved by Agfa in 1936. By convention, consumer reversal films for slides are named by all major manufacturers with the suffix “chrome.” Films whose trade names end in the suffix “color” are by the same convention negative working.

After nonchromogenic development and reexposure to white light (or to a fogging compound), the three dye images are developed in a single step, followed by bleaching of the silver, fixing, etc. An additional step may be included to stabilize the unused couplers. Interlayer diffusion of the couplers may be prevented (1) by attaching long-chain substituents to the molecules, thus immobilizing them, (2) by dispersing them in oily globules which are immobile in the gelatin, or (3) by reacting them with other materials to form longchain molecules (“polymers”); a dispersion of very small droplets of such polymers in water (a latex) is incorporated into the film, providing an immobile distribution of couplers. Incorporated coupler films tend to have a somewhat poorer sensitivity/sharpness trade-off than the Kodachrome system, although this is of real importance only in high-magnification situations.

C. Negative Processes

Most color paper prints, and some transparent images (e.g., motion picture release “prints”), are made by first forming a negative image on a camera film and then exposing a negative-working color paper or film through the original negative. Negative-working color materials work by principles similar to the reversal process just described, with the couplers in the emulsion, but yield a negative dye image by chromogenic development of the original exposed image, followed by bleaching of the Ag, fixing, etc. The couplers in negative-working camera films are often not colorless but are colored red or yellow or both. This color is destroyed when the couplers react to form the image dye, but the unreacted couplers remain as positive images and compensate for unwanted absorptions of the image dyes. Such couplers give the negative the familiar overall red or orange tint.

The “print” emulsion work by principles similar to the camera films, although they usually have a different emulsion sequence. In papers, the sequence is reversed compared with camera films; the yellow layer is then at the bottom, where the mottle of the underlying paper has less effect, because of the low visual contrast of the yellow layer. Negative-working *films* for producing professional motion picture release “prints” from the camera negative, on the other hand, usually have the magenta layer on top; this reduces the visual effect of light scattering in the lower layers of the film: the eye is most sensitive where the magenta absorbs, so that this layer carries most of the sharpness information. The color sensitivity of print materials may be different from those of the camera films, to compensate for the characteristics of the original negative dyes. Also, the intrinsic (nonsensitized) blue sensitivity of the cyan and magenta layers must be very small (and that of the yellow layer very high) with this emulsion order, since the cyan and magenta layers are not protected by a blue-absorbing layer.

Reversal color papers also exist and may be used to duplicate an existing paper print or to make a paper print from a transparency; alternatively, a negative-working “internegative” can be made and a normal printing paper used.

D. Further Details

In practice, there may be a large number of layers in color emulsions. In addition to antihalation, abrasion-resistant, and UV-absorbing layers, for example, gelatin interlayers can be used to cut down interlayer interactions and may contain chemical scavengers to prevent unwanted interlayer migration of oxidized developer; also, as with b/w emulsions, there may be more than one light-sensitive layer for each color. With all color films, good exposure latitude is more difficult to obtain than with b/w emul-

sions, because too low a gamma can reduce the saturation of the colors obtained.

Because the materials controlling the color-sensitivity of a layer in a color film and its color formation process are quite independent, one or more layers can deliberately be made to form a “wrong” color; some films, especially those sensitive to infrared, use such false color for special applications. For example, they may distinguish visually between items of similar color which reflect infrared differently, such as vegetation and camouflage paint, or live and dead vegetation. (Infrared film also tends to “cut through” haze, because of its low infrared scattering.)

Some b/w emulsions (e.g., Ilford’s XP1 400 and Agfa’s Vario-XL) also use chromogenic dye chemistry. The Ag can be recovered, and there is less visible grain and light scattering than with an Ag image. A wide usable density range is possible; overexposed areas do not become completely clogged with massive silver deposits and thus remain capable of discrimination between dark and still-darker areas. As a result, in effect, the film can be treated as having user-variable sensitivity, even from frame to frame; the appropriate corrections can be made at the final printing stage. Excellent sensitivity-sharpness characteristics and tonal range can be achieved.

Remarkable improvements have been made to color emulsions and processes over the years. Many of the individual processing steps have been combined, and the times required have been much reduced, partly by using high processing temperatures while still maintaining adequate gelatin hardness. For example, Kodachrome development, originally needing 28 steps over $3\frac{1}{2}$ hr, has been reduced to 16 steps over about $\frac{1}{2}$ hr; Ektacolor, originally involving seven solutions and 50 min wet time, can now be developed with two solutions and 3 min wet time. Dyes with more ideal color absorption characteristics and much better stability have been introduced, and a UV-absorbing layer may be used. Nevertheless, even more than with b/w materials, it is still important to store color images at moderate temperatures and humidities, and as far as possible in the dark. Sharpness has been improved by a variety of methods. One of the most important was to reduce light scattering by reducing the thickness of the layers. The development of T-grain technology has been especially significant in this regard. This also involved better control over the mechanics of coating and use of the minimum amount of imaging materials consistent with good photographic sensitivity. The quantity of couplers and of AgX was reduced by using couplers of maximum efficiency and improving the sensitivity of the AgX. In a recent refinement of the latex coupler-immobilization technique mentioned above, many couplers were linked to each “anchor” of the polymer, reducing the required volume of coupler and hence the thickness of the emulsion and the light

scattering. The softening of sharp image edges has also been minimized by using edge-enhancement techniques. For example, certain couplers can release a development-inhibiting compound during the coupling reaction; since this inhibitor tends to migrate laterally, it is less concentrated at the edges of any object in the image, thus allowing enhanced development there.

E. Dye Destruction Systems

A different method of forming tricolor images uses dye destruction instead of dye creation. The emulsion starts with the final dyes present; these are destroyed during development (by reactions involving or catalyzed by Ag) in proportion to the Ag image formed. Straightforward processing is positive-working, but more elaborate methods (analogous to normal reversal processing) yield a negative. Such dye-bleach materials (e.g., Cibachrome) are claimed to use far more stable dyes (generally of the azo class) than those produced chromogenically. The dyes also act as a barrier to light scattering during exposure, thus reducing the loss of sharpness in the image, but also significantly reducing the sensitivity. These materials are more suited to production of display prints and transparencies, e.g., from color slides, than for in-camera applications.

IV. ONE-STEP AND OTHER RELATED PROCESSES

One-step photography (which is one-step only from the point of view of the user) is based on technology originally developed in the 1930s by the Agfa Corporation for office copying application, e.g., so-called “photostat” systems. Its in-camera application for the consumer and scientific photography markets is now normally associated with the name of the Polaroid Corporation, although other systems have been developed. Such photography depends on diffusion processes (the tendency of matter to move so as to even-out concentration differences), and these processes are also the basis of other forms of photography. One-step photography has a great advantage in applications that need quick access to a finished image, coupled with high sensitivity (e.g., in amateur photography, exposure testing, and scientific photography). However, sharpness is often not high (because diffusion takes place over significant distances, leading to sideways “spread”), duplication is difficult, and manufacturing is not inexpensive.

A. One-Step Black-and-White Photography

In the Polaroid one-step b/w systems, a solvent (e.g., sodium thiosulfate) for the AgX is included in the

developing solution, which is also made viscous (“thick”) by incorporation of soluble plastics. A small quantity of this solution is contained in an impervious pod in the film pack. After exposure, the emulsion (the “donor”) and a “receiver” sheet containing development nuclei (particles that can catalyze the reduction of silver, analogous to LIs in conventional materials) are drawn together through a pair of pressure rollers; this action ruptures the pod and releases the processing reagent, which spreads between the sheets. AgX that is not developed in the emulsion dissolves in the processing solution, and the Ag ion complex so formed diffuses to the receiver where the Ag ions are reduced by solution physical development at the development nuclei. The process is positive-working, since the amount of Ag complex reaching the receiver varies inversely with the amount of AgX developed in the donor. After an appropriate time, the donor and receiver sheets are stripped apart, and a stabilizing coating is often applied to the surface of the positive print. The negative on the donor may be discarded or, with appropriate films, used after washing as a negative for making more prints. Integral films, in which no peel-apart step occurs, are now available for b/w use; they work using techniques similar to the color systems discussed later.

Silver formed by solution physical development tends not to form filaments, so the original images tended to be brownish (e.g., sepia) rather than black: this is because the small, fairly compact Ag particles scattered blue light preferentially. Special techniques are therefore used to make the color more neutral, by making the particles grow in small regions where they clump together to form larger aggregates of approximately 100 nm in diameter, still much smaller than the original AgX crystals and of a very high covering power. This can be done by precipitating metallic sulfide catalytic nuclei in the interstices (holes) of a very fine colloidal silica suspension which comprises the coating of the receiver sheet.

Most of the developing reagent is removed by adherence to the negative when it is stripped from the positive, but some reagents are left in the very thin ($\sim 0.4 \mu\text{m}$) images formed in colloidal silica; moreover, the image is susceptible to abrasion and to attack by the atmosphere. Such images are therefore coated with an acidic aqueous solution of a film-forming plastic; the acid tends to deactivate the developer, and the coating action washes away the reagent and provides the required protection. Coaterless Polaroid emulsions deposit the Ag in depth within a receiving layer of regenerated cellulose, where it is protected from abrasion and atmospheric attack. The developing reagents employed have colorless oxidation products and leave no residues, and image stabilizers and immobilized polymeric acids (to neutralize excess alkali) may be incorporated into the receiver.

Naturally, the various processes occurring during development must be properly balanced for this system to work. However, as a result of the tendency to remove any AgX that is not quickly developed, very active developers can be used without serious fogging and (since they are sealed in the pod prior to use) without attack by air; because of this, and perhaps because these images are rarely enlarged, very high effective sensitivities (even ISO 20,000) can be attained. The high covering power (high density) of the small Ag image particles assists the attainment of such sensitivities and also leads to a reduction in the quantity of Ag necessary in the original emulsion.

Another important application of this photochemical technology is the direct production of short-run litho-plates, by photographing the document to be printed. In a typical material for this use (3M Camera Plate) the emulsion layer is coated over the silica-based receiving layer on the film base. After exposure, the plate (film) is developed in a high pH developing solution incorporating sodium thiosulfate. As in one-step photography, a negative image forms in the emulsion layer and a positive image forms in the silica receiving layer. Processing is completed by washing off the emulsion layer with warm water and bleaching the silver, which correspondingly degrades the silica layer image-wise. This decomposition of the silica layer exposes the ink-receptive support, while the remaining silica in the background areas remains water-receptive. The resulting plate can be printed on a conventional off-set press. Other products from other manufacturers, while corresponding in general to the principle of operation of the Camera Plate, rely on the ink receptivity of the silver image itself, rather than employing a bleaching step.

B. Color One-Step Processes

1. Polacolor™ Process

In the Polacolor process, the donor consists of three AgX emulsion layers separated by spacer layers, each associated with a layer that contains an actual dye (not a coupler) complementary in color to that of the light to which the emulsion layer is sensitive. Each dye molecule is chemically linked to a hydroquinone developer molecule to form a “dye-developer.” That is, there are six active layers: (1) a blue-sensitive layer backed by (2) a layer having a yellow dye-developer compound, (3) a green-sensitive layer backed by (4) a magenta dye-developer, and (5) a red-sensitive layer backed by (6) a cyan dye-developer. In addition to ultimately forming part of the image, each dye serves during exposure as an antihalation layer for the layer above it and as a color filter for those below it. The dye-developer molecules are insoluble in acid media but are made soluble by the alkaline processing solution

(released from the ruptured pod as the film pack passes through a set of rollers at the start of development); this ionizes the hydroquinone portion of the molecule, making it an active developing agent. Dye-developer molecules that react with exposed AgX are oxidized to the quinone, which is fairly insoluble and may react with gelatin, thus remaining in the donor layer. Unreacted dye-developer molecules diffuse into the receiver layer and are immobilized by reaction with compounds contained in it, thus forming the image. This is then stripped apart from the donor. The receiver also contains an immobile polymeric acid to neutralize the alkali of the processing solution (once it has passed through a “timing” spacer layer between it and the acid) and to stabilize the dye image.

In order that the oxidation state of the developer shall not affect the color of the dye part of the molecule, the two are separated by “insulating” (nonconjugated) chemical links. The efficiency of development can be increased by adding small auxiliary “messenger” developers of greater alkali-solubility and greater diffusion rate than the dye-developers. Such a developer reaches the AgX crystals quickly to initiate development, and its oxidation product can then oxidize and immobilize the slower moving dye developer, regenerating the messenger. Phenidone and Metol were among the first messengers used; later systems included substituted hydroquinones, alone or combined with Phenidone. The timing of the processes is critical in systems of this kind: for example, development in the individual layers must be near complete before the dyes diffuse significantly, otherwise the “wrong” dye could develop the crystals; this timing is assisted by the spacer layers, and steps are also taken to deactivate the emulsions after a specific time.

2. The SX-70 System and Its Refinements

The Polaroid SX-70 system was the first fully integral one-step color film, in the sense that no strip-apart step is involved. The original system is shown in Fig. 4. After exposure, the film is ejected from the camera through rollers to rupture the pod containing the processing fluid. The very high optical density of the new reagent layer so formed, due to the presence of white TiO₂ and of opacifying (light-blocking) dyes, protects the emulsion from further exposure. The development process is in principle broadly similar to the Polacolor system, except that the donor and receiver are integral and remain so, the dyes which remain in the donor being hidden behind the TiO₂ layer; the donor as well as the receiver must therefore be stable after completion of the image formation process. The opacifying dyes (of the phthalein pH indicator class) of the processing fluid are colored only in alkaline media and become colorless as the alkali is finally neutralized by

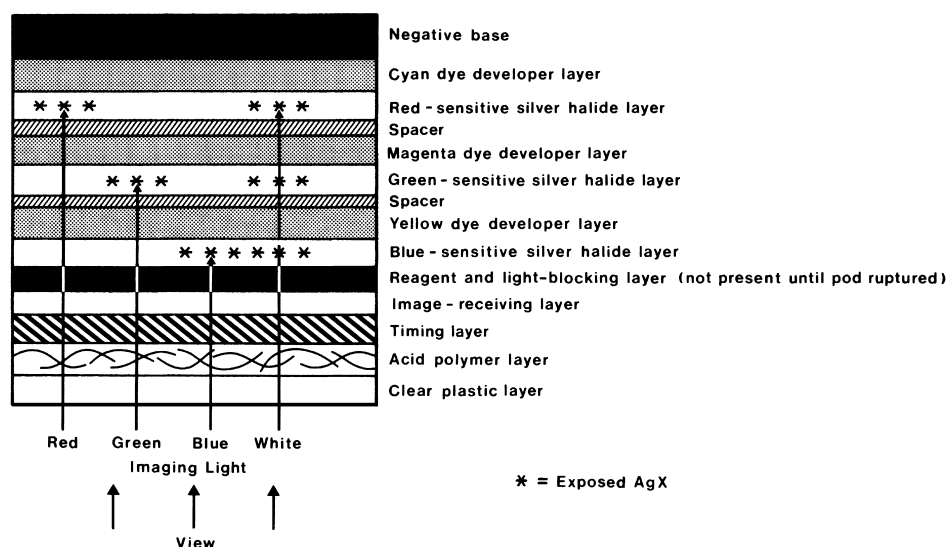


FIGURE 4 Schematic cross section of the Polaroid SX-70 integral one-step color film. Imaging and eventual viewing take place from the same side of the film.

the polymeric acid, leaving the TiO_2 to provide a white background for the image.

The Time Zero system, which provides more rapid image development against a white background, uses a reformulated structure similar in principle to that of the original SX-70 system. A new clearing layer was added between the image receiving layer and the reactants. This decolorizes the dyes immediately adjacent to the image receiving layer, without destroying elsewhere the required opacifying protection for the reagent layers. Time Zero (and later Polaroid integral films) incorporate a quarter-wavelength antireflection coating on the viewing surface; this reduces flare (light scattering) during exposure and improves light transmission efficiency in exposure and viewing. Type 600 film is analogous to Time Zero but uses new high speed emulsions and other refinements.

C. Thermally Processed Silver Materials

Thermally processed silver films and papers, commonly known as Dry Silver, and originally developed by 3M Co., are somewhat similar in principle to conventional negative-working AgX except that development simply uses heat. The emulsion contains AgX crystals and a much larger quantity of organic silver compounds that are insensitive to visible light; also present are a developing agent, toners, additives to promote stability, and spectral sensitizing dyes for the AgX. The image-forming mixture is dispersed in a synthetic plastic binder, usually polyvinyl butyral, in place of gelatin. The organic silver compounds are silver salts of long-chain organic acids, e.g., behenic

acid, and the AgX must be mixed with them on an intimate basis. It can be formed *in situ* by partial conversion of the silver behenate to AgX using halide salts, or the silver behenate can be prepared in the presence of preformed AgX microcrystals. Very small AgX crystals, typically $0.05\mu\text{m}$, are used in this application, normally without chemical sensitization, which is reportedly ineffective in this system. This implies that LI formation mechanisms may be greatly altered from those in conventional AgX emulsions. The AgX normally comprises only a small fraction of the total silver in a thermally processed film, usually less than 5%. Higher sensitivity might be obtained at higher levels of AgX, but it must be kept in mind that these films and papers are not fixed during processing, so that the AgX is still present after thermal development. Higher levels of AgX would lead to noticeable photochemical darkening of this residual AgX in the imaged materials when handled, viewed, or stored under ordinary room light.

The developers are much less active than those used for wet processing, but on being heated to $115\text{--}120^\circ\text{C}$ for a few seconds, they reduce the organic silver compounds, but not AgX itself, to silver in those regions where an LI catalyzes the process. The mechanism of development is thought to be somewhat analogous to solution physical development of conventional AgX films. The toners are actually organic silver complexing agents, e.g., phthalazine, phthalimide, etc., and on heating the film, they convert the silver behenate to mobile silver complexes which diffuse to the growing LI center. Complex formation produces behenic acid as a by-product; this compound is molten at the development temperature and is thought to provide the solvent both for mobility of the complexes and for the

development reaction. On development of a photothermographic film the AgX component is neither reduced nor fixed. It remains light sensitive. In principle this means that on handling the films under normal lighting conditions the AgX could darken to create a background stain in the image, called print-out. This was a major problem with early photothermographic films. Through careful study of the photochemical mechanism of these films, as well as computer modeling, the research organizations of the major manufacturing companies have learned to minimize the amount of AgX required to provide adequate photosensitivity for their applications. Thus over the life cycle of these products the fraction of total silver present as AgX has been reduced from about 25% to below 5%. This level is low enough that print-out is no longer a serious problem with photothermographic films and papers.

Dry Silver materials are easy and quick to develop, but their sensitivity, shelf-life, and image stability have historically been drawbacks, because the processes occurring during development have some tendency to occur during storage, either before or after exposure and development. Significant advances have been made in the stability and sensitivity of these materials in the last few years, and products based on this technology are offered by most major photographic manufacturers. They have good sharpness and can be made sensitive, not only throughout the visible region, but in the near infrared as well. They have achieved market acceptance in such areas as computer-generated microfilm and microfiche, with CRT or laser exposure, image setting for the graphic arts, enlarging papers, e.g., for military uses in environments where water-based processing is precluded, and recording papers, e.g., for logging oil well drilling operations. The most successful current use is in medical diagnostic imaging (Kodak DryView™). Here digital data from a variety of diagnostic modalities, e.g., CAT scan, MRI, ultrasound, digital radiography, etc., can be fed to a common laser scanner, in which a near-infrared laser writes the full-size diagnostic images on the thermally processed film. The Dry Silver emulsion is, of course, coated on a blue-tinted polyester base for this application, so that the resulting image looks like a traditional medical radiograph.

Formulations similar to those used in Dry Silver but without the AgX component are light stable but can be imaged thermally. The silver behenate develops to form a black silver image where (and only where) heat is applied. Materials of this type were originally intended for reflex exposure. Here the film is sandwiched with a document and exposed to a source of infrared radiation. The ink on the document absorbs the infrared and converts it to heat, which initiates development of the film, where it is in contact with the ink image on the paper. Overhead projector transparencies, sometimes called “view graphs,”

have long been made by this process (3M ThermoFax™), originally developed for office copying. Nowadays it is more common to expose such thermally imageable films either with a thermal print head or an infrared laser. In addition to the overhead transparencies on film, paper-based materials of this type are used for on-demand generation of bar codes, e.g., airline luggage labels and point-of-sale receipts. More will be said about thermal imaging below.

V. PHOTOPOLYMERS

A. Principles and Technology

A polymer is a large molecule built up by the repeated linking of smaller units called monomers. “Plastics” are polymers processed into useful objects. In the broad sense, photopolymerization processes may be defined as those involving a light-induced change of molecular weight in a polymeric system. As such, they may involve actual photoinduced polymerization (polymerization of a monomer on exposure to light), and/or photocrosslinking (linking of polymer chains to form higher-molecular-weight, three-dimensional networks with greatly decreased solubility), or photodegradation of a polymer to lower molecular weight fragments via the breaking of chemical bonds. In addition, photopolymers may work by a change of properties, such as solubility, occurring without major change in molecular weight, either because of light-induced changes in some or all of the repeat units or because of a change in an additive incorporated in the polymer. Principal photo imaging applications for photopolymers are for printing plate fabrication, holographic displays, and electronic manufacture.

1. Processes

Very often an initiator such as benzoyl peroxide is added to a monomer to start polymerization. In photopolymerization, this initiator is activated by light to form a material such as a free radical or an ion radical; this starts a chain reaction of polymerization, the radical reaction center propagating by successive addition of large numbers of monomer molecules. Photoreducible dyes may be used to extend the sensitivity from the UV into the visible; the dye in its excited state reacts with a weak reducing agent to generate free radicals. The chain reaction, which can be controlled and inhibited to achieve good shelf-life, builds an amplification step into the process, limited by chain termination arising from unwanted reactions or from the ends of the growing chain locally running out of monomer. Photocrosslinking also gives amplification, because large areas can be cross-linked with comparatively

few links; chain reactions can also occur. Amplification can be achieved with photodegradation if, for example, a degradation product initiates another reaction or if a polymer chain, cleaved at one point, is unstable and then undergoes partial or complete depolymerization.

Nonpolymeric materials, apart from initiators, may undergo light-induced changes that may have a variety of effects on the polymers with which they are in contact. The new material may, for example, take part in reactions such as cross-linking, or it may be a catalyst (e.g., an acid) for polymerization, depolymerization, changes to some or all of the repeat units, and so on; alternatively, it may inhibit the solubility of a polymer more or less than the original material.

2. Materials

Actual photopolymerization materials are typically in the form of a monomer liquid, alone or dispersed in a film-forming polymer; in the latter case, the photopolymerization of the monomer also causes linkages to the existing polymer and thus changes the mixture's solubility and other properties. In photocrosslinking, the reaction may occur via excitation of groups in the polymer or of an additive that can link to two chains. An example of the former is the cinnamoyl group which undergoes cyclodimerization to form a cyclobutane ring between chains containing the group. Cross-linking effectively converts thermoplastic polymers to thermosets, i.e., obviates their heat softenableity. Azide additives photodecompose to nitrogen and imidogen radicals or nitrenes which can cross-link polymer chains, such as cyclized rubbers, that contain residual bond-unsaturation or other suitable groups.

A well-known older system uses bichromate ions as the light-sensitive element; a change in oxidation state of the chromium during exposure produces ions that cross-link many natural and synthetic polymers, including proteins, gelatin, albumen, certain glues, shellac, and poly(vinyl alcohol). However, the sensitivity of bichromate systems is low and depends on a number of factors that are difficult to control, and their use has declined greatly although they are still used in holography. A common photopolymer system employs a diazo-oxide dissolved in a phenol-formaldehyde-type acidic polymer. The diazo-oxide inhibits the dissolution of the polymer in aqueous alkali solutions, but light-exposure converts the inhibitor to an acid, permitting the polymer to dissolve. This system is most commonly used as a photoresist in electronic manufacturing applications.

An important class of initiators for photoinitiated polymerization are the 'onium salts, primarily the diaryliodonium and triaryl- or arylalkylsulfonium salts. One advantage of these initiators is that on photolysis these

compounds may form both free radicals, which initiate polymerization of unsaturated monomers, e.g., acrylates, and cations, which yield strong Lewis or Brønsted acid, which can, in turn, catalyze cationic polymerizations, e.g., of epoxides. In addition to photoinitiated polymerization, the radicals produced by onium salt photolysis can be used to activate crosslinking processes. Likewise the photogenerated acid can be used to catalyze various photodepolymerization processes. Not only do these initiators have broad applicability to a variety of polymerization mechanisms, but their photolysis may be spectrally sensitized. The mechanism of this sensitization is by electron transfer from photoexcited dye to the onium salt. This sensitization pathway is essentially the same as that described above for spectral sensitization of the silver halides.

Another category of photoinitiators of growing importance comprises transition metal complexes. Titanocenes used to photocure the polymers used in manufacture of CDs and various other optical disks in the Phillips Corp. process are examples of this class. Various metal carbonyl compounds are also useful in this manner. In general they cannot be spectrally sensitized and require UV sources, e.g., mercury lamps, excimer lamps and lasers, as light sources. These compounds photodecompose to eject one or more ligands attached to the central metal atom, leaving it free to bind to the unsaturated monomer, e.g. acrylate. The metal complexed monomer adds to another monomer molecule to initiate the polymerization thermally. Thus it is proper to term these polymerization reactions as photocatalyzed rather than photoinitiated. Another application of the metal complex initiators is in fabrication of electronic circuits; ionic photoinitiators such as the onium salts may not be useful in these cases because of their effects on the electronic properties of the resulting polymer.

3. Development

The development of photopolymers generally makes use of one of the changes in properties that occurs on exposure, most typically by washing away the more soluble regions, or more recently by removing one component by vacuum processes, such as etching with reactive ions. Other techniques include transfer of the softer regions to another support; dusting-on of a pigment (using differences in tackiness) followed by transfer of the pigmented areas to another support; delamination or peeling apart of a film sandwich containing the photopolymer (using changes in the adhesive properties); using the change in hardness of a polymer shell to control release of encapsulated dyes or dye precursors inside; and making use of changes in the transmissibility of inks or chemicals (which, for example, can permit color-coupling reactions to proceed imagewise, analogous to AgX chromogenic

materials) or of differences in light scattering. In many cases, as we shall see, the visibility of a photopolymer “image” is not the prime concern, although the material may be specially pigmented where it would not otherwise be visible. When nonexposed material is not removed, for example in the case of certain light-scattering images, appropriate materials may be included in the formulation so that the primary reactants can be desensitized to form nonreactive compounds after exposure, for example, by heating or by irradiation with light of a wavelength below that used for exposure.

B. Applications

1. Holography

Holography involves formation of an interference pattern between coherent, i.e., laser, light beams, one of which is modulated by the object being imaged. The resulting pattern, which must be recorded with very high resolution and low noise, can be reconstructed with either a laser beam or a point source to create a three-dimensional image of the original object. Its practice is described in more detail in a separate entry in this *Encyclopedia*. Although holography has been practiced for over thirty years, since the advent of lasers, it is only in the past decade that it has become really significant commercially and not just an esoteric curiosity. Applications include various security applications, e.g., credit cards, drivers' licenses, etc., industrial nondestructive testing, image display for virtual reality, and “heads-up” displays for military aircraft. Similar displays are beginning to find their way into production automobiles. Holography is increasing in acceptance in the visual arts community as well. An illustration of the latter is an exhibition of the sculpture of Michaelangelo, held at a major North American art museum (Musée des Beaux-Arts de Montréal), in which the works which for reasons of size or safety could not be brought to the museum, were displayed in the form of full-size holograms. Holography continues to hold out great promise for high-density information storage, especially for random access memories for optical computing.

Holographic display media are usually coated on glass plates, as dimensional stability is very critical and generally cannot be achieved on other substrates. Principal materials for holographic display media are silver halide plates, dichromated gelatin, both of which have been described previously, and photopolymers. Two modes of holography may utilize photopolymers: (1) volume holograms, in which the interference patterns are recorded as refractive image differences between monomer and polymer in the bulk of the photopolymer layer, and (2) relief holograms, which exploit the volume difference between

monomer and polymer, i.e., exploit the shrinkage of the layer where photoinitiated polymerization has occurred, so that the interference pattern is represented by a relief pattern of “hills” and “valleys” on the surface of the photopolymer layer. The dimensional scale of these surface features may, of course, be comparable to the wavelength of light.

An interesting feature of these relief holograms is the ease with which they can be mass reproduced. To this end the original hologram can be coated with a metal layer, usually nickel, by evaporation or electroless plating, i.e., solution physical development, and then used as a stamper to emboss a complementary pattern into a heat-softened thermoplastic foil. The replica in the thermoplastic can be rendered permanent by a photocrosslinking step which converts the heat-softenable thermoplastic into an unmeltable thermoset plastic. Such images are familiar on many credit and identification cards. The process of their fabrication, sometimes called microreplication, is, of course, essentially the same as that used in mass production of CDs, DVDs, etc.

2. Printing Plate Applications

For printing plates, photopolymers may be applied to the baseplate during manufacture, or they may be coated by the user. A wide variety of photopolymerization, photocrosslinking, and other materials is available. Removal of the more soluble part of the image after exposure, leaving the underlying material, can give rise to a more-or-less planar plate, part of which is hydrophilic (water loving) and part oleophilic (oil loving) as needed for lithography. With thicker polymer layers the relief pattern needed by letterpress printing may be created.

3. Electronic Applications

During the microfabrication of typical electronic integrated circuits, a resist is coated onto an Si/SiO₂ wafer and exposed to radiation through a patterned mask. (The mask itself was conventionally made via photographic reduction processes onto high-sharpness AgX film, although electron-beam methods are now widely used.) After solvent development, the mask pattern formed by the resist is left on the wafer. Subsequent chemical etching of the SiO₂ along the pattern and removal of the remaining resist in a strong solvent generates a pattern of bared Si on the wafer, into which dopants and the like can be diffused. The whole process may be repeated several times.

Photopolymer resists may be applied from a solvent using techniques such as spin coating, in which the substrate is rapidly rotated as the solution is applied, leading to a uniform layer roughly 1 μm in thickness.

Very high sharpness is obtainable with these systems, and they are used to photofabricate lines down to the 1–2 μm range. Positive-working photoresists (in which the irradiated regions are dissolved) commonly employ the diazo–oxide/phenol–formaldehyde solubility-inhibition system described previously, while negative-working resists have generally used photocrosslinking, often employing the cyclized rubber/bisazide system. To reduce the effects of diffraction on the sharpness of the lines obtainable, photopolymers may be exposed with deep-UV (toward 200 nm), electron beams, X-rays, or ion beams. Experimentally, objects down to about 10 nm have been made with electron beams. For these applications, different resists are needed. In the case of positive-working deep-UV resists, for example, chain scission of poly(methyl methacrylate) has been widely used. It is desirable to eliminate the wet-development step after exposure of the resist. Systems have been reported in which the light exposure degraded the polymer to volatile products; another approach is to use a dry-development technique such as plasma etching. Multiple-layer systems are being used in which a thin, high-sharpness resist acts (after development) as an *in situ* mask for a thicker one that smooths-over steps in the substrate, and optically bleachable overlayers with nonlinear optical properties can enhance contrast and hence sharpness. Inorganic resists are under investigation, since they may give higher effective sharpness as well as being plasma-developable.

A major trend in photoresist technology is the development of resists which function at increasingly short wavelengths, available with excimer lamps and lasers. The driving force for this development is the fact that resolution with which the master image may be registered on the resist is inversely related to the wavelength of light used in the process. Increased resolution is in turn driven by the desire of the computer industry to put more functions per unit area onto a microcircuit chip (Moore's Law). Patterning of photoresists at 194 nm is now common practice. A major issue in the development of photoresists for such short wavelength is their transparency; many organic materials absorb so strongly at these wavelengths that only the surface would get exposed, and the bulk of the photoresist layer wouldn't capture the image.

Resists for lower-sharpness applications, such as printed-circuit boards, use much thicker emulsions than those for fine-line applications, and dry, self-supporting, adhesive photopolymer films are available.

VI. DIGITAL PHOTOGRAPHY

The principal competition to traditional photography for imaging applications, whether in the form of conventional

silver halide, one-step, photothermography, or photopolymer technology is digital photography. Digital video, for example, has replaced analog video, which, in turn, has replaced motion picture photography in consumer, scientific, and journalism applications, and even in some commercial motion picture production. This has transpired mostly within the past fifteen years. The growing trend to imaging applications of digital technologies reflects two factors: (1) the ability to interface these technologies with other digital capabilities, primarily the internet; and (2) the applicability of Moore's Law to image-capture devices. Moore's Law, promulgated in the mid-1960s, states that the number of solid-state devices comprising a single microcircuit chip doubles every eighteen months. If each of these elements corresponds to a single pixel in an image-capture device, it can be seen that over the above-referenced fifteen-year interval, the number of pixels comprising a single digital image, i.e., the amount of image information which can be captured digitally, has increased by three orders of magnitude. Thus it is now possible for users to realize the convenience, speed, and flexibility of digital photography, all the while being able to capture amounts of image data comparable to a traditional photograph.

The most common image capture device in digital image capture is the charge-coupled device array (CCD). At this writing, 5 megapixel CCD arrays have just become available in consumer products which are approaching being competitive, in terms of the amount of image information captured, to AgX film in the 35-mm format. The two factors that ultimately determine the competitiveness of a digital technology, whether for image capture, storage or display, with its photochemical counterpart, are cost, convenience, and image quality.

A. Image Quality

Various metrics are used to quantify the quality of an image.

1. Resolution

Resolution, or resolving power, is a figure of merit related to sharpness insofar as it measures the fineness of detail which an image process can capture. From the point of view of digital photography it represents the number of image receiving or display elements, i.e., pixels, per unit area of the imaging device. It is inappropriate to compare AgX films, which are analog image-capture devices, with digital image-capture devices, as the definition of resolution has to be different in these two cases. Instead it is more appropriate to compare them on the more common basis of the information content of an image. Thus a 5 megapixel CCD array can capture approximately 5 MB of information for one image, if each element is capable

of discriminating 256 gray levels. Psychometric studies have shown that this number of gray levels can provide an acceptable simulation of continuous tone in an image. The 5 MB would correspond to the amount of information captured by on a 35-mm color negative, if the film were capable of 40 line pairs per millimeter resolution. Even very high speed commercial films currently on the market have at least twice this resolving power, but, unlike digital devices, may yield reduces force scale on their resolution limit is approached.

2. Modulation Transfer Function (MTF)

MTF and the related metrics, edge transfer function and contrast transfer function, are mathematical functions, rather than a figure of merit, and describe the sharpness of an edge in an image. It must be understood that the MTF of a displayed image, whether captured digitally or photochemically, reflects not only the MTF of the capture device but all the image processing which may have occurred between capture and display, i.e., the whole of the imaging chain is involved in determining the quality of the final displayed image.

3. Granularity

Granularity is a measure of the noise content of an image. The term comes from the fact that in conventional photography a high noise content image appears grainy to the viewer. Zero granularity is, of course, impossible. Consider a finite number of photons falling on an array of detectors. They will be randomly distributed among the detectors, even if the exposure is uniform. The distribution can be described by Poisson statistics. Thus if N photons on average impact each detector, the standard deviation in the exposure experienced by the detectors in the array will be $N^{0.5}$. This variability gives rise to what is known as “shot noise” or “quantum mottle.”

4. Detective Quantum Efficiency (DQE)

DQE is an overall figure of merit for an imaging device or process, which incorporates photographic speed, contrast (gamma), and granularity. DQE is an embodiment of the speed-image quality trade-off. It is defined as

$$DQE = (S/N)_{out}^2 / (S/N)_{in}^2$$

where $(S/N)_{out}$ is the signal-to-noise ratio of the output data stream, or recorded image, and $(S/N)_{in}$ is the input signal-to-noise ratio. It can be understood as a measure of how efficiently an image-capture device uses the photons reflected by the subject into the image-capture device to record the image. The input signal-to-noise ratio is always finite because of the shot noise.

State-of-the-art AgX films and state-of-the-art CCD devices both are characterized by DQEs of 8–10%. The upper limit of the DQE for AgX materials is 50%, because the AgX grains function as coincidence counters. From a consideration of the mechanism of LI formation, above, it can be seen that at least two photons must be absorbed by one grain within a very short period of time, usually less than one second, for the grain to behave as exposed. This feature is actually an advantage of the AgX process as it limits unwanted, accidental “exposure” by heat, cosmic rays, not to mention airport security systems, and is critical to the shelf-life of AgX films. The upper limit for most CCD arrays of current design is 33%, because each device in the array is sensitive to only one of the three primary colors (red, blue, or green) owing to the presence of an array of color filters over the array of solid-state detector devices. This is again intentional, as the design enables the three color records of the image to be collected, processed, and displayed separately, using only a single CCD array.

B. Image Capture Devices

Digital photography is differentiated from conventional photography by the use of solid-state, electronic image-capture devices, instead of photographic film. The camera used for digital photography may look the same as that used for film photography. Only the solid-state device replaces the film at the plane of focus, and the power supply for the electronics replaces the motor drive of the camera. Because the optical components of the cameras for both kinds of photography are essentially the same, many of the same manufacturers are active in supplying the marketplace with both digital and film cameras.

The most commonly used image-capture device in a digital camera is an array of CCDs. The schematic of a small portion of such an array is shown in Fig. 5. Each individual device element, i.e., each pixel, comprises a potential well in the doped silicon substrate. It is linearly connected with other such elements in series. Photoelectrons created by exposure of the pixel, on a more or less one-for-one basis, are trapped in the well associated with that pixel, then shifted to the well of the next pixel in the series by sequential application of a voltage to the electrodes (A, B and C in Fig. 5) adjacent to the well. The electric fields associated with these voltages effectively create artificial potential wells under the electrodes. The packets of photoelectrons ultimately emerge sequentially into an external circuit, in the inverse order to the physical arrangement of the pixels relative to the external circuit connection (drain).

An important feature of practical CCD arrays is an array of microlenses over the electronic devices. These focus all the photons incident on the individual pixel onto the

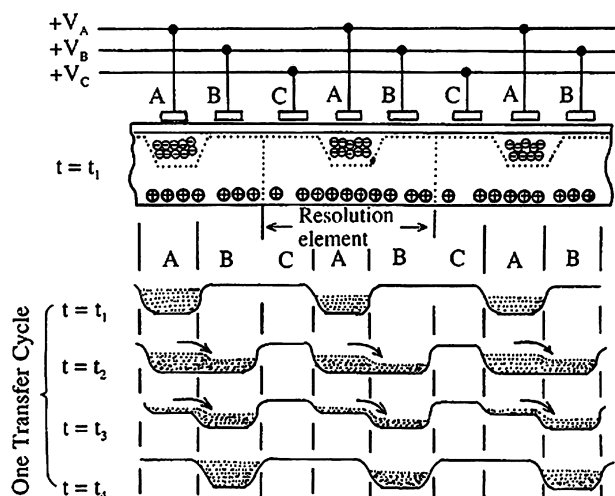


FIGURE 5 Operational schematic for a CCD array. One charge transfer cycle is shown; three such cycles are required to move the photogenerated charge packet from the potential well of one device element to that of the next. [From Tani, T. (1998). *J. Imaging Sci. Technol.* **42**, 1–14; by permission of the copyright holder, the Society for Imaging Science and Technology.]

photoactive area of the silicon chip, which has been doped to create a p-n junction, at which the electrons can be photogenerated.

The size of each element, i.e., pixel, in a 1 megapixel CCD array is typically $9\ \mu\text{m}$ to accommodate the array on a 1 cm square chip. Creation of a 5 megapixel array thus requires that the individual pixel size be reduced to $4\ \mu\text{m}$. As can be seen from Fig. 5, these individual elements are fairly complex. Clearly meeting the challenge of increasing the number of pixels in an array, so as to make digital photography increasingly competitive with AgX technology, requires advances in microelectronics fabrication technology. From Moore's Law, however, we might expect the number of pixels available at a given cost to quadruple, i.e., pixel size to decrease by 50%, over a period of three years, as has already happened between 1998 and 2001. Another significant characteristic of digital image-capture devices, especially CCDs, is the linearity of their response. This means that the photocurrent representing a single pixel is directly proportional to the number of photons which were absorbed by the sensor in that pixel. By contrast AgX films exhibit a logarithmic response, more like that of the human eye. Thus the dynamic range, i.e., the range of light to dark tones which they can reproduce, of the digital devices tends to be less than that of traditional films.

Even at its present level of elaboration, however, photography with digital cameras has become popular with many consumers who are willing to tolerate the lower image quality compared to AgX in return for easy interfacing

of their cameras with home computers and the internet. It is already well-established in the photojournalism community. In the marketplace digital photography represents a particularly significant threat to one-step photography, insofar as the digitally recorded pictures are instantly available with the appropriate hardware.

Conventional cameras utilize the film itself for image storage. In digital cameras, temporary image storage is required until the images are transferred to a more permanent digital record, e.g., CD or hard drive. Most consumer cameras utilize solid-state memory cards, so-called flash cards or Fuji Photo's SmartMedia, for this application. Miniaturized floppy disks and disk drives have also been developed for in-camera use.

C. Digital Photofinishing

1. State-of-the-Art

Currently the best features of AgX and digital photography can be combined in a hybrid imaging chain in which the image is captured on conventional color negative film. The information on the negative is digitized either by scanning it using a scanning laser or, more usually, by imaging it onto a CCD array. The digital image is then processed by computer, as required for color correction, sharpness, tonal scale adjustment, and image format. With respect to the latter variable, marketing studies conducted by the Eastman Kodak Co. have indicated that, given the choice, consumers would like to crop approximately 40% of the photos they take. While, of course, home computers offer them this capability, PC-based image processing still has to become a lot easier for it to be attractive to a majority of amateur picture takers.

In digital photofinishing the processed image is output onto a print medium, which may be a silver halide paper or may utilize any of several nonimpact printing technologies. Of the latter, ink-jet printing, see below, seems to be preferred for so-called all-digital minilabs or microlabs, where the entire digital photofinishing process may be carried out in a customer-accessible location, e.g., drug store. For centralized, large-scale photofinishing operations, printing onto silver halide color papers is preferred. The AgX print material may be exposed using a CRT display, a light-emitting diode (LED) array, a liquid crystal light valve array, or a set of scanned lasers. Prototypes of such systems have been around for nearly twenty years (LaserColorTM), but only recently have data transfer rates sufficient for the speed of mass-market photofinishing been commonplace. In addition, recent advances in CCD array manufacture, enabling high megapixel CCD arrays for image capture, enable enough information to be digitized from the negatives for digital

photofinishing to provide any size print of which the negative is capable.

For home computer users, the processed images can be returned to the customers in digital form on the internet (PhotoWorld in Europe, PictureVision™ and AOL's "You've got Pictures" in North America, etc.), for printing via their PCs. Again ink-jet technology is likely to be used for obtaining the final display prints, owing to the widespread distribution of these printers in homes and offices. This option is expected to become more popular with increasing penetration of broad bandwidth digital subscriber line (DSL) technology. It goes without saying that the input to a digital photofinishing system may also be the data file corresponding to an image captured by a digital camera, as well as a scanned silver halide negative. In the former case the scanning step is obviated.

2. Advanced Photo System (APS)

The replacement of traditional optical photofinishing with digital photofinishing to form a hybrid imaging chain represents an opportunity for photographic film manufacturers to create new films tailored to the requirements of these systems. From one point of view, the new films may be simpler in construction than present ones. With all optical photofinishing, contrast control, color correction, edge enhancement, etc., must be accomplished in the film using elaborate chemistries. (See above.) Most of these forms of imaging processing, necessary to yield a pleasing final print, are done digitally in the hybrid systems, so the chemistries that enable analog image processing can be left out of the film. The realization of these possibilities are still in the future as of this writing, however.

The first step toward development of films engineered for digital photofinishing is represented by the Advanced Photography System (APS or Advantix™). This is the first new film format in over a decade, and has been established as the result of collaboration among the major photographic manufacturers worldwide. The actual frame area on the APS film is somewhat smaller than a conventional 35-mm frame and uses the 1.85:1 format of HDTV. It has obviously been anticipated that consumers will use to view their photos on their HDTV sets which should be capable of displaying their full photographic quality. Advances in emulsion technology enable equal or superior overall image quality to be obtained with these products, and it is commonplace for photofinishing labs to deliver to consumers larger prints from their APS films than from 35 mm, even with traditional optical printing of the negatives. Since it has to be compatible with conventional optical photofinishing as well as new digital technology, APS still incorporates all the now familiar chemistry for *in situ* analog image processing. Cameras designed for

APS format films represent the fastest growing segment of the consumer camera market at this writing.

VII. NONIMPACT PRINTING TECHNOLOGIES (NIP)

During the past decade the family of technologies which comprise nonimpact printing, i.e., the creation of images on paper or other substrate without use of a conventional printing press, has grown to be more important than either silver halide or digital photography, in commercial terms. Its further growth promises to be even more dramatic, given the possibilities these technologies offer for on-demand printing, including direct computer-to-press scenarios, which are currently being implemented by companies such as NexPress and Iris Graphics in the United States, Xeikon in Europe, and Indigo in Israel. This development will put digital image output into direct competition with traditional modes of high-volume printing, such as web and sheet offset, and letterpress. To place developments in NIP in the context of photographic materials and processes, it must be realized that NIP essentially represents a merger of traditional photography with the graphic arts, joining photography, office copying, microfilm, medical imaging, and printing into one technological family.

The principal technologies included under the term nonimpact printing are (1) electrophotography, also known as xerography and as electrostatic printing, (2) ink-jet printing, and (3) thermal imaging. Of these, only electrophotography normally offers an image-capture capability, though certain versions of thermal imaging can, in principle, enable image capture, e.g., as an infrared pattern. Examples of this sort of thermal imaging were discussed above in Section IV.C.

A. Electrophotography

Electrophotography is based on the formation of images by the movement of matter, usually charged particles, in electric fields. Most commonly, the fields are reduced or destroyed in exposed areas via the action of a photoconductive insulator. Amplification of at least $\sim 10^5$ is achievable, because previously stored charge at high voltages can be discharged by photons that each carry the energy equivalent of only a few volts, and comparatively large particles can be moved by the attraction of charges equivalent to relatively few electrons. In fact, in typical situations, it is possible to deposit about 15 million pigment molecules per photon of exposure. The exposure required for complete development by typical, efficient electrophotographic processes is thus in the several ergs/cm²

region, easily good enough for direct photocopying applications (a lens-coupled camera-exposure situation in which fairly high illumination intensities can be supplied) and also for a number of more conventional photographic applications e.g., micrographics. Modern office copying increasingly involves digital technology. The original document is scanned with an array of light-sensing devices, e.g., a CCD array, or using a scanning diode laser. These data are digitized and the image is processed as required, e.g., for reduction, enlargement, or contrast control. The digital image is then written on the photoconductor belt (see below) in an electrophotographic engine using a scanning, modulated diode laser. This latter step is identical to that employed in laser printing for computer output. A number of machines are now on the market which can function both as copiers and as laser printers, or even as fax machines as well. This multifunctionality indicates the power of the imaging chain concept, insofar as the same engine is used to output image data, whether captured by a fax machine, a scanner, or created in a computer.

The most important electrophotographic process is transfer xerography, which is the basis of the majority of modern office copiers and laser printers. A photoconductive material, which is the active part of the photoreceptor (P/R), is electrically charged with ions produced from the air and is exposed to the imaging light; in exposed areas, the voltage across the P/R is largely discharged, leaving an imagewise pattern of charge on its surface. This is the latent image. Development is achieved by the differential attraction to the charged or uncharged areas of charged pigmented plastic toner particles; the toner is usually a dry powder and the pigment is commonly carbon black. The resulting toner image is then transferred to the final substrate, usually paper, to which it is fixed. After discharge and cleaning steps, the P/R is ready for reuse for as many as 300,000–500,000 cycles, or even more.

Electrophotographic development is rapid and usually involves no liquids. With reusable P/Rs, the main consumables are plain paper and a toner system, so the materials costs per unit area are extremely low by photographic standards, although not by commercial printing standards. The final paper image is very stable. Electrophotographic *films* can be handled in the light before use and can be updated after development and use; their shelf-life and image stability can be very good, because the sensitivity to light is “gated” by application and later removal of the electric field. Although the sharpness of many electrophotographic processes is fairly low, as is appropriate for document copying applications and the like, there is nothing fundamental about this, and excellent sharpness and pictorial quality can be attained under appropriate cir-

cumstances. Both positive- and negative-working systems are available. Most electrophotographic processes show little, if any, reciprocity failure.

1. Transfer Xerography: Equipment and Processes

The basic steps in transfer xerography are illustrated in Fig. 6. They are generally carried out around the periphery of a continuously rotating P/R drum or flexible belt. The charging step most commonly employs a corotron, a device with one or several thin wires largely surrounded by a grounded shield; the wires are held at a high enough voltage (typically 2–10 kV) to ionize the air in the vicinity of the wire. Because the field is high only in the vicinity of the thin wire, no sparking occurs, but ion multiplication processes near the wire form a stream of positive or negative ions, such as hydrated protons and CO_3^- , some of which flow to and thus charge the nearby P/R. Since most of the current flows to the shield, this tends to stabilize the process against uncontrolled variations in factors such as the corotron-P/R distance. “Scorotrons” use a control grid between the wire and the P/R, somewhat analogous to a triode: the ion flow tends to stop when the P/R surface voltage is similar to that connected to the grid. In other cases, a corotron may consist of a grid of wires, so that it is essentially transparent when placed (out of focus) in an imaging system.

“Dicorotrons” are similar to corotrons, except that the wire is usually coated with an insulator such as glass; an AC, instead of DC, voltage is applied to the wire, and the desired polarity of ions is driven to the P/R by a bias applied to the shield. The coating on the dicorotron wire reduces its susceptibility to dirt, and the charging is generally more uniform; disadvantages include cost, increased ozone generation, and susceptibility to insulator damage. Charging can also be carried out using pins: the high electric field at the sharp points of the pins produces ionization; uniform charging over the full width of the P/R is accomplished by using many pins spaced a few millimeters apart, often in several linear arrays. Pin charging usually uses a scorotron-type control grid. Advantages of pin charging include high charging speed, uniformity, rigidity (there is no wire to vibrate), and ease of cleaning; much less ozone can be produced if no shield is used.

Recently, comparatively low-voltage charging (about 300 V) has been done in some machines by pressing a conductive rubber roller against the P/R; the roller is coated with a thin insulator and held at the desired voltage. This rather simple approach would presumably not work for the higher voltages needed for much xerography, since there would presumably be uncontrolled air breakdown as the roller approached and separated from the P/R; excellent

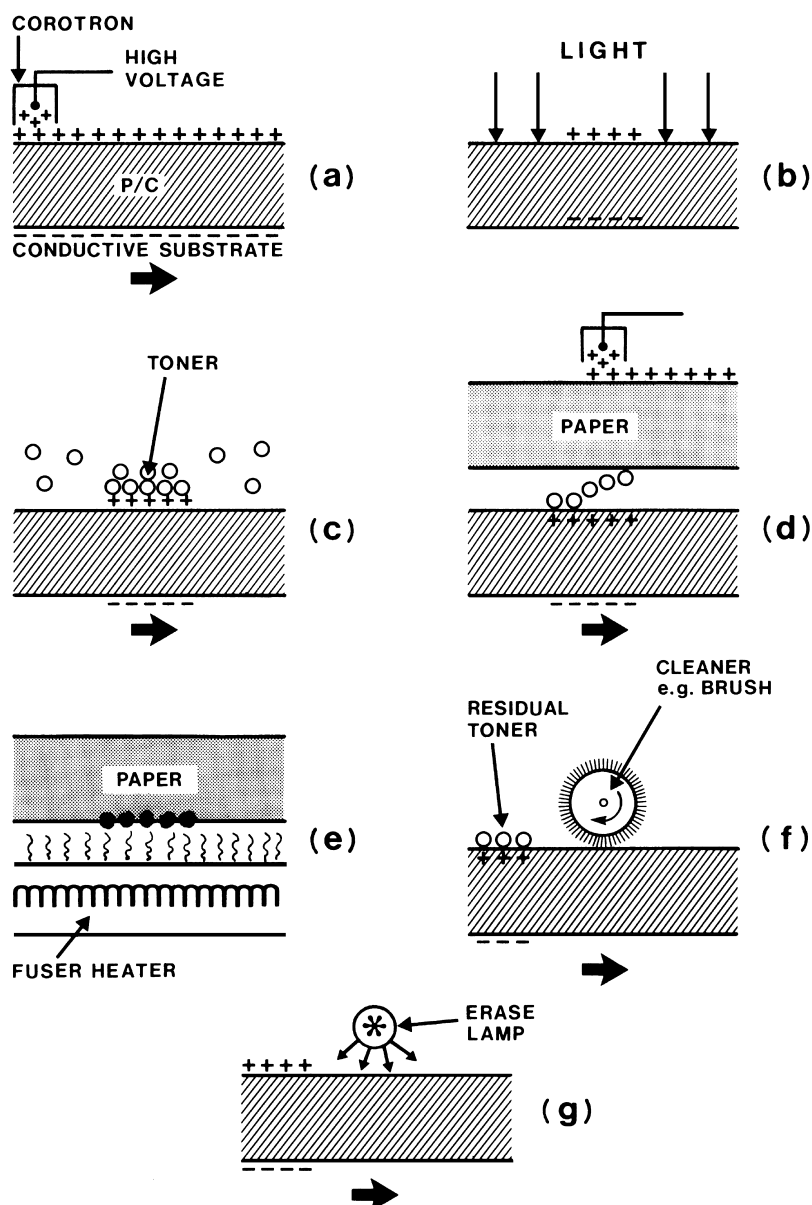


FIGURE 6 The basic steps in transfer xerography: (a) charging; (b) light exposure; (c) toning; (d) toner transfer to paper; (e) toner fusing; (f) P/R cleaning; (g) P/R discharge.

contact between the roller and the P/R must also be needed, and P/R contamination may be an issue.

The ions are reasonably stable on the surface of appropriate P/Rs, such as noncrystalline (glassy) selenium, amorphous silicon, or various thermoplastics, and the resulting field across the P/R is reasonably constant for times typically on the order of seconds. The surface charge induces an equal and opposite charge in the grounded base electrode to maintain overall charge neutrality.

Illumination of the document to be copied typically involves slit-scanning, in which only a thin strip of the document is illuminated and imaged onto the P/R at any in-

stant; the P/R moves (often on a revolving drum) at a rate similar to the scanning lamp. A long thin mirror typically moves with the lamp and collects light from the illuminated portion of the document; other mirrors in the optical system are also typically moved to keep constancy in the distance between the illuminated portion of the paper and the P/R. Some systems move the platen on which the paper is placed, rather than the lamp. The position of the imaged "slit" is most often kept constant in space as the P/R moves past it, or it may be scanned a small distance parallel to the P/R surface motion (but in the opposite direction) in order to minimize the length of P/R (and thus the time)

wasted as the scanning lamp “flies back” between scans; for a given process speed, this can increase the copy rate by 30% or more.

Slit scanning is optically inefficient, since only a small portion of the document is illuminated at any one time, and it can produce slight distortion of the copy. High speed machines therefore often use flash exposure of the whole document at once; to keep all the image in focus requires that the exposed portion of the P/R be flat, or that a curved platen be used.

A lens may be employed to transmit light from the paper to the P/R. In addition to the lens-coupled copying of existing documents, however, xerography is now widely used in electronic printing devices including digital copiers. In such applications the “original” is not a light pattern, but a stream of electronic signals. Usually, a laser beam is scanned repeatedly in one dimension back and forth across the surface of the P/R, using a rapidly spinning multifaceted mirror, while the P/R moves slowly at right-angles to the scan direction. The laser beam is turned on and off (“modulated”) by the electronic “bit stream” as necessary to generate a pattern on the P/R corresponding to the print desired.

Other techniques of coupling electronic input into xerographic systems are becoming important; particularly significant are “image bar” systems, which include active linear displays (e.g., LED arrays) and passive displays (e.g., liquid crystal arrays), usually of the same length as the P/R. The passive displays are used as a line of “shutters” in conjunction with a continuous light source such as a fluorescent tube. CRTs have also been used to expose the P/R. One of the advantages of image bars (in addition to factors such as mechanical simplicity and cost) is that they can often be used with space-conserving SELFOC® lenses; these consist of several hundred individual gradient-index rod lenses arranged in one to three rows, a page-width in length. Since the focal length of each small lens can be small, the total size of the optical system can also be kept small. SELFOC arrays are used in some copiers too.

In all cases, exposure of the P/R to light produces electron-hole pairs, somewhat analogous to the first step in the AgX latent image formation process; rather than causing chemical changes, these then separate under the action of the applied field. With a selenium P/R, for example, the light is absorbed near the surface, but if the initial ionic corotron charge is positive, the holes can travel right to the bottom of the selenium in a fraction of a second, with little deep trapping. The electrons need move only a short distance to neutralize the surface charge. Thus, virtually total discharge is achieved in the exposed areas.

To achieve a photographically positive image when the ionic charge is positive, negatively charged toner particles must be attracted as uniformly as possible to the nondis-

charged P/R areas. However, in the absence of another grounded surface, and if the charged area is not small, a toner particle near the surface of the P/R “sees” as much negative charge on the electrode as positive on the surface, even though the toner may be nearer to the positive charge: the effect of an “infinitely” extended sheet of uniform charge does not depend on the distance. However, if the toner is near an edge of a charged area on the P/R, the charge does not appear infinite in extent, and then the toner’s closer proximity to the positive charge leads to the intuitively expected attractive force. Expressing the situation another way, the electric fields associated with the surface charges extend almost entirely to the opposite charges induced in the conducting substrate; only near the edges of a charged area do significant fields extend above the P/R surface where a charged toner particle can experience the resulting force. Figure 7a shows this situation, expressed in terms of the direction and density of the lines of electrical force that a negatively charged toner tends to follow toward a large positively charged area.

In these circumstances, then, toner is attracted primarily to the edges of charged areas, and this leads to the well-known edge effects of some xerographic processes: thin objects such as typewritten characters are well reproduced, but solid black areas are not. This is useful up to a point in enhancing the reproduction of documents and has been used to great effect in xeroradiography, in which medical X-rays are detected xerographically; images that would otherwise have very low density variations can be greatly enhanced by the edge effects. The effect also confers process latitude, since development is purely a function of voltage differences between adjacent areas rather than depending on absolute voltages.

For most applications, however, it is desirable to reduce these edge effects or enhance solid-area development. Screening techniques can assist in this regard. Usually, however, a development electrode is placed in close proximity to the P/R (Fig. 7b): negative charges are then induced on it at the expense of the P/R substrate electrode, and the toner now sees a net force toward the P/R

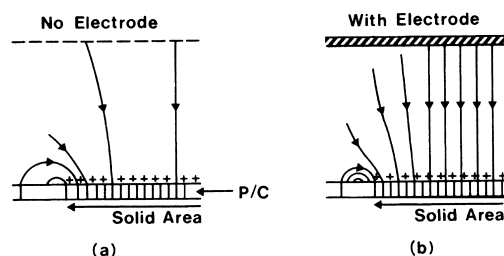


FIGURE 7 Illustration of the direction and density of the electrical lines of force “seen” by a negatively charged toner near the edge of a uniform positively charged area, (a) without and (b) with a development electrode.

even in solid areas. The difficulty that must be overcome with this is that for maximum effect the spacing between this electrode and the P/R must be small, so the electrode can interfere with access of the toner. This sort of physical image processing is obviated in digital copiers and laser printing machines, where the image data are processed electronically to allow development of an output image of the desired characteristics.

The toning process can be made negative-working by using a development electrode biased with the same sign of voltage as the original charge on the P/R. Alternatively, if only small lines are to be developed (using edge fields), then simply reversing the charge of the toner is sufficient, because the edge fields reverse direction (Fig. 7) a small distance from a charged area.

After the P/R has been toned, the toner is transferred to the paper (or other substrate such as a coated polyester for transparencies), usually assisted by a bias field applied by a corotron to the rear of the paper. An electrically biased roller can be used in place of the corotron, and this can prevent ionization effects occurring on separation of the P/R and the paper; the electrostatic image on the P/R can then be reused to produce more than 100 copies, although this is rarely done commercially. Steps can be taken to suppress the transfer of some the toner on the P/R, thus reducing background (toner deposition in areas of the paper that should be toner-free). To facilitate paper separation from the P/R, a second corotron may discharge the paper after transfer charge has been applied. After the paper is separated from the P/R, the toner is made to flow partly into the paper to fix it in position. (Note the somewhat different usage of the term “fix” in electrophotography and conventional photography.) This may be done by “fusing” (melting) the toner via contact with a heated roll which has a fluorocarbon or silicone release surface or a thin surface layer of a suitable liquid; alternatively, the toner may have a release agent in it. Image gloss can be controlled by the hardness of the roller, and hence its intimacy of contact with the paper. Compliant rollers may be used, which increase the contact time with the toned image and thus heat it for longer; sometimes, the lifetime of such rollers is limited, however. An extreme version of this concept replaces the roller with a heated fusing belt of silicone rubber or the like; this can be kept in contact with toned image for still longer and also reduces energy consumption.

Other systems employ radiant heat (steady heat from coils or lamps, or flash heating from a flash lamp); this has the advantage of reducing the copier warm-up time and eliminating the energy used to keep a roll hot, but pressure is not available to assist bonding to the paper. Cold-pressure alone may be employed, but the high pressures required may give the paper an objectionably glossy

appearance and feel; also, a careful balance must be struck in terms of toner softness, so that toners of this kind may tend not to stick well enough to the paper or to stick too well to the P/R, making cleaning difficult. Cold pressure can be assisted by auxiliary flash heating. Solvent vapor has also been used for fixing.

Residual charge on the P/R is removed by uniform illumination. The last toner residues are neutralized by corotron driven by an alternating voltage and are removed by mechanical means such as mechanical brushes, wiper blades, disposable sheets, or developers working in reverse (often with an electrical bias); removal of the toner may be assisted with a vacuum, and solid lubricant can be incorporated in the toner to assist transfer and cleaning.

2. Xerographic Materials

A wide range of P/R materials has been used in transfer xerography. Noncrystalline selenium (produced by vacuum deposition onto a cylindrical drum held at a special temperature) has been widely used but responds primarily to blue light; green illumination can “weight” the response toward longer wavelengths well enough for many applications, but nevertheless the inability to copy pure blue lines (which seem like white to the selenium) is well known. Alloys of selenium with tellurium can give panchromatic response, although conduction in the dark is much increased; one way around this is to use a thin layer of such an alloy in conjunction with a thick layer of selenium. Compounds of selenium with arsenic also improve panchromaticity; in addition, they provide better stability against crystallization, since the arsenic effectively cross-links the selenium chains. The P/R is often deposited on a thin barrier layer to prevent unwanted charge injection from the base electrode, which is typically aluminum.

Dispersions of photoconductive pigments such as cadmium sulfide or zinc oxide in insulating (usually plastic) binders may be used, provided that either the volume loading of pigment is high enough for good charge transport, or a charge-transport material (see below) is used in the insulator; many such materials can be dye-sensitized, analogous to AgX. Special organic plastic layers have been used as the P/C, a well-known example being the charge-transfer complex between poly(*N*-vinylcarbazole) and trinitrofluorenone; this has good photoconductivity and color response and transports both electrons and holes. Recently there has also been some use of noncrystalline hydrogenated silicon P/Rs. These can be made with a very hard and durable surface and is very photoactive but may be expensive to produce.

It is often desirable to have a very stable, abrasion-resistant material, such as a tough plastic, at the top surface

of the P/R, to protect a softer photoconductive material from abrasion, attack by ions, and so on. Of course, a simple overcoating would tend to trap charge, although a corotron erase step may prevent this; alternatively, photodielectric processes, one of which was described above, may be used. In addition, however, better flexibility of the P/R is very desirable; flexible P/R belts (typically on aluminized polyester) can be then used in place of the conventional P/R drum. Such belts need not move in a circular path and can be made flat in the exposure region. This enables optically efficient “flash” exposure as explained previously. Furthermore, the various process stations can be placed more conveniently with a flexible belt, high copy speed is easier since several images can be made for each rotation of the P/R without requiring a very large drum, and the belt can be bent round a sharp corner to assist separation of the paper after toner transfer. (Alternatively, belts have been rolled up and placed inside a drum, around which part of the belt was wrapped for use; when this part of the belt was worn out, new material was advanced from inside the drum.) Such flexibility is hard to obtain with wholly inorganic P/R materials, although selenium on special flexible nickel belts has been used. It is somewhat easier with organics, but the requirement both to photogenerate and to transport charge carriers is difficult to achieve with a single-layer organic. Finally, it is often important to improve on the dark conductivity of single-layer P/Rs.

For all these reasons, a wide range of P/Rs nowadays separates the functions of photoactivity and charge transport by using multiple layers. The thin photoactive layer may be inorganic, for example, trigonal selenium (which has better color response than noncrystalline selenium) in a plastic dispersion, or an alloy of selenium, tellurium, and arsenic; very often, however, organic P/Cs (such as phthalocyanines or thiopyryliums) with or without sensitization are used, usually also in plastic dispersions. A much thicker overlayer of a tough plastic, such as a polycarbonate, provides the charge transport and hence the insulating ability and voltage contrast (the contrast is proportional to the thickness for a given charge). Photogenerated carriers are injected into the overlayer and are transported therein by virtue of a special charge-transport material dissolved in the layer. In the case of a negatively charged P/R, this charge transport material needs an unusually small amount of energy to lose an electron (i.e., it has a low “ionization potential”) and can therefore easily donate an electron to neutralize the hole left in the photosensitive layer after photoexcitation; after this, the charged molecule can accept another electron from a neighboring molecule, and so on. The initial photogenerated hole is thus effectively transported through the insulator to neutralize the negative charge at the surface.

Near-infrared-active materials such as phthalocyanines are becoming especially important in materials of this type, because they permit the use of infrared diode lasers (which are compact, inexpensive, and easily modulated) for electronic printing applications; such materials can have quantum efficiencies (electron-hole pairs produced per photon absorbed) near unity and near-infrared sensitivities (for 50% discharge) of about 1 erg/cm^2 . There are worldwide efforts underway, however, to produce usable visible-light diode lasers.

Any P/R must meet a number of requirements, some of them in conflict. The sensitivity must be high; this requires good absorption of the appropriate light, high efficiency of electron-hole generation and transport (which tends to imply high fields), and reasonably low surface charge (which is proportional to V/d , where V is the charging voltage and d the P/R thickness). V must be high enough for good development. The thickness d must not become so large that the carriers start to be deeply trapped before they traverse the whole P/R thickness, and economic and adhesive factors speak for lower values of d . The field in the P/R must not be so large that breakdown or excessive dark conduction occurs. In addition, the P/R must cycle reproducibly many times without excessive trapped charge (which can prevent proper photodischarge or can effectively increase the P/R dark conductivity because of an increase in field at the electrode), without other increases in conductivity, and without mechanical damage or physical changes such as crystallization.

Typical charging voltages are in the region of 1000 V, although excellent results have been obtained with only a hundred volts or so; selenium-based P/Rs are usually about $50 \mu\text{m}$ thick while organic layers tend to be about $10\text{--}25 \mu\text{m}$. A typical selenium-based P/R therefore acquires an initial surface charge density of about 10^{-7} C/cm^2 or the equivalent of about 7×10^{11} electronic charges/ cm^2 . Substantial photodischarge with an efficiency of, say, 0.7 electron-hole pairs per absorbed photon therefore requires about 10^{12} photons/ cm^2 . If each photon is blue and therefore carries about 3 eV or about 5×10^{-12} ergs, the required exposure for full development is about 5 ergs/ cm^2 . The charge density in the nonexposed areas is enough to attract a dense black layer of toner particles.

Typical dry xerographic toners are plastic particles (often styrene-acrylic copolymers) typically about $10\text{--}15 \mu\text{m}$ in size, colored black with about 5–10% of submicrometer carbon black particles, or (for the color processes discussed later) colored with cyan, magenta, or yellow colorants. The function of the plastic is to physically bond to the paper after the fusing step. Early toners were charged by passing insulating particles through electrically active nozzles, thus forming a charged “powder cloud.” This process can give excellent continuous-tone images and can

allow a close-spaced development electrode but is difficult to control, can tend to give background, and is difficult to scale to high process speeds.

Most modern dry toners are charged via contact electrification caused by contact with relatively massive carrier particles: if different materials are brought into intimate contact, they tend to charge each other via electron flow under the influence of their different electronic energy levels. Charge control additives may also be used to control such charging. Since the toner particles, which usually constitute about 1–3% of the total mass of these “two-component” development systems, are attracted to the carrier as well as to the P/R, there is less tendency for them to be deposited in uncharged areas, and if they are deposited there they tend to be cleaned up later by reattachment to other carriers. In the classical technique of cascade development, this mixture is poured over the exposed P/R, and toner is knocked loose, redistributing itself between the carrier and charged areas of the P/R. Inverted cascade development uses a similar concept, but the carrier is cascaded down the development electrode, with the P/R above it: gravity thus acts against transfer in this case. With insulating carriers, such as small glass beads, edge development is prominent, because the size of the carrier beads (typically 250–600 μm) makes it impractical to bring a development electrode close to the P/R.

Conductive carrier beads, however, can act as their own development electrode. A widely used technique which can employ either conductive or insulating carriers is the magnetic brush system. A magnet carries a mixture of fine toner with coarse (typically 100 μm but sometimes down to about 20 μm) magnetized carrier beads; the beads may be of ferrite or steel, usually overcoated with an insulator such as PTFE to confer the required charging properties and release the toner easily. The system is usually in the form of a soft brushlike structure (Fig. 8) which is moved against the P/R. The brush electrode may be electrically biased to compensate for incomplete P/R discharge or to change the gamma of the process. By adjustment of the conductivity of the developer mixture (“conductive magnetic brush” or “insulating magnetic brush”), either edge-contrast or solid-area development can be favored; control of conductivity through the life of a developer may also be less of an issue with insulating systems. The magnetic brush system offers good process control, is capable of high process speed, gives low background, can be fitted into a smaller space than cascade systems, and does not depend on gravity, giving greater process freedom; it does, however, require sophisticated control of the materials package if it is to be stable over many copies. For high-speed systems, the (insulating) developer may be forced through a nip between the developer and a flexible P/R which is partly wrapped around it, thus forming more of

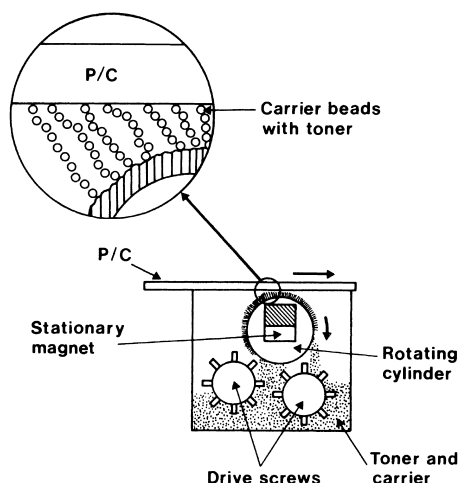


FIGURE 8 Schematic illustration of the magnetic brush development process. For simplicity a flat P/R is assumed.

a mat than a brush; the agitation so produced allows toner (and charge) to be replaced quickly at the end of the brush, preventing “starvation” effects; several cylinders may be used in sequence to obtain more complete development without excessive background. Agitation of the developer can also be achieved by rotating the magnet within the cylinder. In some systems, the size of the carrier and that of the toners is comparable; this gives a very soft brush and good fine-line rendition. It is possible to reduce machine dirt and increase developer life by using magnetic toners; this reduces the possibility of toners becoming detached from the carrier when they should not.

Single-component development systems dispense with the carrier and may use toner particles that themselves carry magnetic materials. Such toner may be applied, for example, from a magnetic roller. Single-component developers have the advantages that there is no carrier to wear out and no necessity to keep a reasonably constant balance of toner and carrier, there may be less P/R abrasion, and the development unit may be simpler. The toner-bearing surface is usually still moved against the P/R to give a good supply of toner. Alternatively, “touch-down” or “impression” development may be used: a toner layer, typically on a soft donor surface (such as a conductive brush fabric or an elastomer layer) is pressed into contact with the P/R without shearing motion. Close spacing between the roller or sheet and the P/R can provide solid-area coverage, and an electrical bias may be applied. However, close control of the P/R voltages is required, and care must be taken to avoid background.

Single-component systems may use relatively conductive toner, so that the charged P/R itself induces charge separation in them, giving the required attraction to the P/R. However, toner transfer to the paper can then be more difficult, since the conductive toner may exchange charge with

the paper and be repelled; moreover, the conductivity may be humidity-dependent. Other systems contact-charge an insulating single-component toner by contacting it to a blade held at a suitable voltage or by applying it to the P/R from a roller of an appropriate insulating material (which essentially constitutes one large carrier).

One interesting insulating single-component magnetic system first forms a uniform layer of blade-charged toner on a magnetic roller. The toner then jumps a 300- μm gap to the P/R, propelled by a combination of an AC and DC bias between the rotating roller and the P/R. As the distance between any given region of the toner layer and the drum reduces and then increases again, toner transfers, oscillates between the two, and finally is back-transferred from the light areas. Presumably, the AC approach tends to reduce the chance of toner staying on an area where it should not be. A somewhat similar AC concept has been used to improve the efficiency and quality of an effectively two-component development system. Further disadvantages of various single-component systems include the inability to use magnetic toners in color systems (because of the dark color of the magnetic material), and limited charging blade and donor roll life.

Toner plastics are chosen to melt at the lowest temperature possible (since fuser power and warm-up time are important practical considerations) but must not stick to each other, to the carrier during storage, or in the photocopier prior to use, or too well to the P/R. Tacky toner encapsulated in a higher-melting coat may be a way around this conflict.

A different form of development uses micrometer or submicrometer pigment or pigmented plastic particles in a kerosene-type liquid. The particles may become charged by selective absorption of ionic species that are added as charge-control agents; reasonable stability against aggregation is then achieved by the resulting electrostatic repulsion. Alternatively, long-chain molecules may be attached to the particles to charge them and to form a diffuse shell around them to keep them apart. Very sharp images can be obtained with these small particles, which are carried very close to the latent image by the liquid and thus sample the field very close to it. (Dry toners in this size range are difficult to handle and may be a health risk in some size ranges.) Heat fixing may be avoided, since the image may fix itself as the liquid dries, especially if a soluble film-forming plastic is contained in the liquid. Major disadvantages include the tendency of conventional toner suspensions to aggregate and settle out, and carry-out of liquid on the paper, often with absorption of toner-bearing liquid in the background areas. Moreover, resulable P/Rs can be hard to clean when such developers are used, and some P/Rs are incompatible with the toner liquid over long periods. Fine lines may be smeared as the liquid image is transferred to

paper. Special paper may be used to prevent undue spread of the liquid image, using a blotting-paper effect.

Various ways of avoiding or reducing the liquid carry-out problem have been used. For example, a counter-rotating roller may be brought close to the P/R to remove excess liquid. Alternatively, much of the liquid can be squeezed from the toner suspension by means of a charge applied to it on the P/R; at the transfer stage, the toner can be made to jump to the paper under the influence of a field (the required separation may be maintained by large particles in the liquid), and enough solvent is transferred that no fixing is necessary. Another system applies the liquid toner on a sponge-rubber roller covered with a thin insulating mesh. This is pushed against the P/R and is thus squeezed. Excess toner and liquid are removed by the sponge as it expands on separating from the P/R. Yet another approach is to use a P/R that is not wetted by the developer except with the aid of the field from the surface charge. In a different liquid-development method, conductive ink was electrostatically drawn out of the recesses of a finely patterned applicator onto the P/R only in charged, nonexposed areas.

Liquid toner technology has enabled several direct-xerographic microfilms to be introduced. To obtain high sharpness, these are developed with liquid toners, and no transfer step is used. Since the P/R is visible in the final image, it is typically an organic P/C in a transparent plastic with dye-sensitizers also present. Such films are charged, exposed, toned, and dried, and the toner is fused onto the P/R, usually by heat. Drying usually employs an air-knife (a long thin jet of air), a vacuum, or warm air. As with photoplastic films, the need to keep the P/R reasonably transparent dictates typical exposures in the 100 ergs/cm² range. A material using a crystalline solid P/R layer such as cadmium sulfide has also been described.

Handling liquid toners and keeping them stable is inconvenient, but processing is rapid; the sharpness can be very high with fine toners, and color response can be reasonably good. Like photoplastic materials, these films can be handled in reasonable light prior to charging, are more or less stable before and after imaging, and can be updated after use by addition of new images next to old ones. The lifetime of the charge and of the latent image is not very long, and camera systems must be designed with this in mind. Positive- and negative-working systems can be designed using appropriate ionic-charging and toner-charge polarities, and (especially with a development bias) the gamma can be made low enough for continuous tone reproduction.

3. Direct Xerography

In direct xerography, the P/R is factory-coated onto the image substrate, usually paper, and remains there after imaging. There is no requirement for image transfer or P/R

cleaning and recycling, which significantly reduces the mechanical complexity of the copier. However, the paper no longer appears “plain” (apart from its greasy feel, it cannot be perfectly white if it is to absorb visible light), and its cost is significantly increased. Typically, zinc oxide, mixed with a roughly equal volume of a binder and with appropriate dye sensitizers, is coated to a thickness of about 12–15 μm and used as the P/C. The paper is usually treated to make it reliably conductive (reasonably independent of humidity), so as to form the required underlying electrode. Development may be with dry or liquid toners, and repeated exposures through appropriate filters can give color images. Direct xerography with no charging step has also been reported from development of the small “Dember” voltages arising from independent migration of photogenerated electrons and holes in the absence of a field.

4. Color Xerography

Full-color transfer xerographic systems most often use three or four full charge, exposure, development, transfer, and cleaning passes of the photoreceptor for each print; the exposures are through appropriate color filters, and development is done with three or four different colored toners. The (subtractive) color principles are similar to those discussed in Section II.A, with the fourth color (if used) usually being black, to improve the depth and color-neutrality of black areas. These systems need careful registration of the images during the transfer process.

Other systems produce three or four toned images on the P/R before transfer, but this is not simple, since the second and subsequent charge expose cycles are not onto a clean P/R; this approach is easier with color printing (rather than copying) since each exposure can then be performed with an infrared laser, which may not be absorbed by the previously deposited toners, as long as black is done last. In another approach, the three or four color images are transferred to a reusable intermediate sheet after each image is produced; the final image is then transferred to paper. Presumably, the intermediate is less sensitive to size variations (especially those caused by temperature and humidity changes, or by stretching) than is paper, thus making the registration process easier. Colored slide originals can be used with color copiers by projection of the image onto a special flat “Fresnel lens” placed on the document platen to redirect the projected light toward the copier lens.

Xerography is typically a fairly high gamma process, which is appropriate for applications like document copying. The high gamma arises principally because the discharge of most P/Rs occurs over a fairly narrow range of light exposures, the resulting sharply defined voltages then being converted nearly proportionally to density with most modern developers. For pictorial applications, how-

ever, screening techniques (Section IV.A) (nowadays often done by manipulating the image electronically or, sometimes, by building a screen right into the P/R) can reduce the effective gamma to levels acceptable for pictorial purposes. Color fidelity can be controlled either with special P/Rs and multiexposure processes or nowadays by appropriate correction of electronic image input. Electronic systems can also correct for various process variations and imperfections. Given the power of electronic color correction, screening, and other image manipulation processes, even color *copying* is commonly done nowadays by scanning the paper input electronically (using a CCD array and colored filters), manipulating it electronically, and then printing it xerographically.

Electronic correction systems, and the smaller dry toners which are now becoming available (in the 5–8 μm range), together with smaller carriers and laser spots, permit very good dry-toner color copying and printing without the graininess that large toners can cause. The trend to smaller toners, with consequent improvements to image quality, may well continue, if the difficulty of releasing small toners from the P/R and the carrier can be overcome. Paper smoothness is also important in obtaining a smooth looking image; of course, smooth (even coated) paper is more acceptable for pictorial than for document applications. Extremely accurate scanning systems are also important for high quality color.

Meanwhile, a significant trend in liquid toners for color xerography is toward larger (up to about 2 μm) particles or agglomerates; these are easier to transfer from the P/R (there is no need to release such toners from a solid carrier, of course), particularly if special particle shapes are used which touch the P/R only at a few points. These larger agglomerates are still smaller than present-day solid toners and are small enough to give very good image sharpness. Other advantages of liquid color toners over their larger solid counterparts include better toner transparency, and hence image brightness, and less mass to fuse to the paper: four color images can contain so much solid toner that it is difficult to fuse it without drying and curling the paper. Solvent carry-out is a disadvantage. The images available with liquid developed xerographic systems are now good enough that several systems have been developed for color proofing (getting a quick indication that a particular commercial print run will give acceptable quality); this is very demanding of quality and reproducibility.

Very high xerographic speeds are now being achieved for some commercial printing applications: for example, a recent introduction using liquid toning, while not yet of pictorial quality, has a process speed in the vicinity of 300 ft/min. Issues which have to be dealt with in such devices include the high charging rate needed, solvent carry-out, and attack of the P/R by the liquid toner. Tandem

xerography (in which the three or four colors are automatically laser-printed in succession on separate machines) is an important recent innovation which increases the effective color printing speed, analogous to the tandem systems of conventional commercial printing.

5. Other Electrostatic Processes

A traditional difficulty with electrophotography is its failure to reproduce uniform gray areas of an image. This problem occurs because electric fields which attract toner particles tend to be maximized where there is a sharp differential in surface potential on the electrophotographic receptor, i.e., at edges in the image, even when development electrodes are employed (see Fig. 7). Use of digital inputs to electrostatic printing engines in a hybrid imaging chain provides a means to overcome this difficulty, insofar as the image is represented as an array of pixels. Every pixel boundary, even in nominally uniform areas of the image, represents an edge, and develops accordingly. If the pixels are small enough, the developed area appears uniform to the viewer. Two approaches to pixel-wise generation of an electrostatic latent image have recently been developed.

In the digital micromirror (DMD) technology developed at Texas Instruments Corp., the optical input to a xerographic P/R is divided into pixels. The digital image is composed on an array of micromirrors, each of which represents a pixel of the image. Micromirrors which are “on” reflect light from a light source to the P/R, while micromirrors which are “off” do not. The fabrication of a micromirror array is of special interest. In one version, developed originally by the Parkin-Elenor Corp., a flexible sheet of a suitable plastic, e.g., polyester or MylarTM, is coated with a reflective metal layer, which also serves as an electrode. The metallized, reflective sheet is stretched over a honeycomb support, behind which is an array of electrodes. Each aperture in the honeycomb corresponds to a pixel and has its own electrode, not in contact with the plastic sheet. Each element of the honeycomb array is thus a microcapacitor, insulated from its neighbors. Application of a voltage to one of these electrodes charges the capacitor and leads to deformation of the flexible plastic mirror surface, thus deflecting the reflected light at this point away from the P/R, i.e., turning the pixel “off.” The entire image is composed in this manner on the micromirror array and reflected, in part or as a whole, onto the charged P/R, which then forms an electrostatic latent image and is developed in the usual way. Alternatively, of course, the micromirror array can be comprised of a large number of individual, electromechanically tipped, bistable mirrors. This approach is actually used on currently commercial DMDs.

The second approach, developed at the NexPress Corp., does not use a P/R. Instead the latent image is formed on a receptor which comprises an array of metallized domains, or rectangular islands, on the surface of an insulating layer, coated over a continuous, grounded electrode. The islands are isolated from each other and represent individual pixels. The array can be fabricated by usual microlithographic methods used in the microelectronics industry. The receptor is charged by contact with flexible metal “fingers” which are themselves addressed by the digital data stream. Thus when a given finger is “on” in response to a “1” in the data stream, charge is transferred to the metallic island on the receptor which happens to be under the finger at that time; the receptor is constantly moving, so, shortly, another island is under the finger, which, by then, may be “on” or “off” in response to the data stream. The resulting charge pattern on the receptor can then be developed by the usual methods of dry or transfer xerography.

6. Electronic Paper

Flexible, paper-like displays represent another opportunity for image output or display based on electrostatic technology. Unlike electrostatic printing where the xerographic engine is used to create a toner image that can be transferred to paper, in electronic paper displays all the active components are included in the paper itself, known as electronic paper or e-paper. Furthermore in e-paper the displayed information isn’t static but can be updated on demand. Some applications envisioned for e-paper in the 21st century include: (1) a newspaper that wirelessly updates itself throughout the day; (2) hand-held computer displays that are thin and flexible enough to be rolled up and placed in pocket or purse; (3) wallpaper that can be changed to match other changes in decor; (4) business cards that can be updated with new job assignment or contact information; and (5) point-of-sale advertising displays that be updated with current price information and sale announcements. The last of these is the only one which has been realized commercially, at the time of this writing, using either GyriconTM technology, originally developed at the Xerox Corp., or E-Ink technology, developed at Lucent Technologies’ Bell Labs and at MIT, but the others are all under active development.

The possibility of e-paper is a consequence of the development of thin transparent electronic circuitry based on conducting and semiconducting synthetic polymers. The development of these materials was recognized with the Nobel Prize in 2000. The Gyricon approach involves a thin, transparent layer of silicone rubber, impregnated with a high concentration of polyethylene spheres, typically 50–100 μm in diameter. Each sphere is black on one side and white on the other; in addition, the black

and white sides are oppositely charged, electrostatically. Each sphere floats in its own spherical, oil-filled cavity in the rubber layer. The rubber layer is sandwiched between transparent electrodes, which are patterned so as to divide the sheet into an array of isolated, individually addressable pixels. When a voltage is applied between the electrodes of a pixel the electric field causes the spheres in the rubber layer in that pixel to rotate so as to display their black or white sides, depending on the sign of the applied voltage. Electronic notepaper, in which the charge is applied across the silicone rubber layer using an electrically charged pencil, instead of a top surface electrode, has also been demonstrated by Gyricon.

In the E-Ink approach, the paper incorporating the organic electronic circuitry is coated with a water-based polymer containing a high concentration of clear plastic microcapsules, similar in size to the Gyricon microspheres. These microcapsules contain charged titanium dioxide particles suspended in a hydrocarbon oil that is dyed black. When the voltage is applied to an individual pixel of the coated paper, the white particles move to one side of the microcapsule or the other, depending on the sign of the voltage. If they move to the side toward the viewer, the paper appears white at this point; if they move to the side away from the viewer it appears black.

The complex patterning of the electrodes and their driving electronics, to allow each pixel of the electronic paper to be addressed individually, depends on another imaging technology, microcontact printing, which involves using a very high resolution (non scale) rubber stamp with a long-chain alkanethiol as the ink. This ink adheres image-wise to a thin gold layer, the electrode precursor, formed by vacuum deposition or by electroless plating onto the paper, actually a plastic sheet, substrate. The latter technique is akin to solution physical development, described above in connection with conventional photography. Regardless of the origin of the gold layer, the alkanethiol ink strongly chemisorbs to it and acts as a resist, which is impenetrable to etching chemistry. The etching solution dissolves away the gold where it is unprotected, however; the alkanethiol can then be removed with an appropriate solvent to leave the patterned electrode. Developers of the e-paper technologies agree that this high resolution patterning technology has been critical to the feasibility of this sort of image display.

B. Ink-Jet Printing

Ink jet describes a form of dot-matrix printing in which small droplets of ink are ejected from a small aperture in the imaging device and directed to a specified position of the imaging medium, e.g., the paper, to print the image. There is no specific image-capture technology as-

sociated with ink-jet printing, but the input used to control the ink-jetting device is usually a digital data stream. Ink-jet printing has the capability of creating full color images using, usually, four inks: yellow, magenta, cyan, and black. This four color strategy is analogous to that used in conventional letterpress and offset printing. Unlike the impact printing technologies used in printing presses, ink-jet printing has a limited capability for continuous tone: either the number of drops or the size of the drops reaching a given pixel may be controlled. Another major advantage of ink-jet printing, compared to electrostatic color printing, above, or color thermography, below, is that the registration problem is simplified. In-ink jet printing the four image colors are deposited simultaneously. In the other systems, the four colors are put down sequentially, which means that very high precision is required in repositioning the imaging medium in the printing engine for each cycle, so that the colors overlap perfectly. This registration, as it is called, is not trivial from the viewpoint of the mechanical engineering required.

There are two basic types of ink-jet printers: continuous ink jet and drop-on-demand. In the former a stream of drops is continuously produced, and droplets are deflected electrostatically as required to the print medium to create the image. Ink from unused droplets is collected and recycled. This strategy has been employed in many high-end printers designed to produce images for the high quality commercial market, e.g., by Iris Graphics. In the latter type of printer, the drops are generated only as required. This strategy is used primarily in desktop printers for the office and consumer markets, e.g., Canon BubbleJet™ or Hewlett-Packard ThinkJet™.

1. Printer Types

The design of the print head in a continuous ink-jet system is shown in Fig. 9. In this example of a binary deflection system the drops are generated continuously by a piezoelectric device and pass between a pair of charging electrodes which, depending on the signal applied,

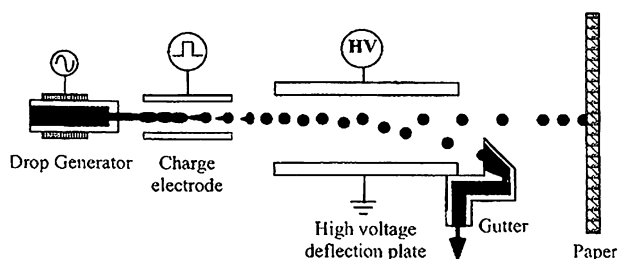


FIGURE 9 A binary deflection continuous ink jet print head. [From Le, H. P. (1998). *J. Imaging Sci. Technol.* **42**, 49–62; by permission of the copyright holder, the Society for Imaging Science and Technology.]

may or may not apply a charge to a particular drop (or set of drops). A deflection plate deflects the charged drops electrostatically to the gutter, where the ink is collected for reuse, while the uncharged drops proceed to the paper. There are also similarly functioning multiple deflection systems, wherein the charge applied to the droplets is variable; the angle of deflection of the droplets on passing through the deflection plate is likewise variable. The drops can accordingly be deflected to different spots on the paper. Both types of devices are used for industrial coding, marking, and labeling applications. The multiple deflection continuous ink-jet method is particularly suitable for high-quality, continuous tone images, e.g., in the graphic arts market. Very large images, up to billboard size, can be produced using continuous ink-jet technology.

Drop-on-demand ink-jet print heads are less complex and are commonly used in printers for the office and home markets. Evolution of the drop-on-demand technology to increasingly small droplet sizes (currently picoliters) has enabled extremely high quality to be obtained with these printers, as well, though not with as high speed as the continuous ink-jet printers. Most drop-on-demand ink jet printers utilize either a thermal or a piezoelectric driver. The thermal ink jet is considered to be the most reliable method on the market at this time. A diagram of a representative thermal ink jet is shown in Fig. 10. The device involves a resistive heating element which turns an electrical pulse into a heat pulse. The heat pulse causes the ink to boil where it is in contact with the element, and the bubble-forces a droplet of the ink out of the orifice and impels it onto the paper. It is this bubble-forming feature of the print head that inspired one manufacturer of printers utilizing this method (Canon) to adopt the trade name BubbleJet. Several other manufacturers, e.g., Lexmark and Hewlett-Packard, however, utilize this printing method with only minor variations.

In the piezoelectric method, a piezoceramic deforms in response to the electronic signal and forces ink out of the orifice. Different manufacturers of these devices utilize different geometries, giving rise to so-called squeeze, bend, push, and shear modes of operation. Squeeze mode print heads are, by now, more a matter of historical interest, and shear mode devices have not yet made a significant commercial impact. Bend and push mode devices, illustrated in Fig. 11, are both commonplace designs in ink-jet printers offered to office, home, and commercial printing markets by a number of manufacturers, including Dataproducts, Epson, Sharp, Tektronix, and Xerox.

2. Materials

It is axiomatic in the paper industry that there is no such thing as "plain" paper. All papers are designed and pro-

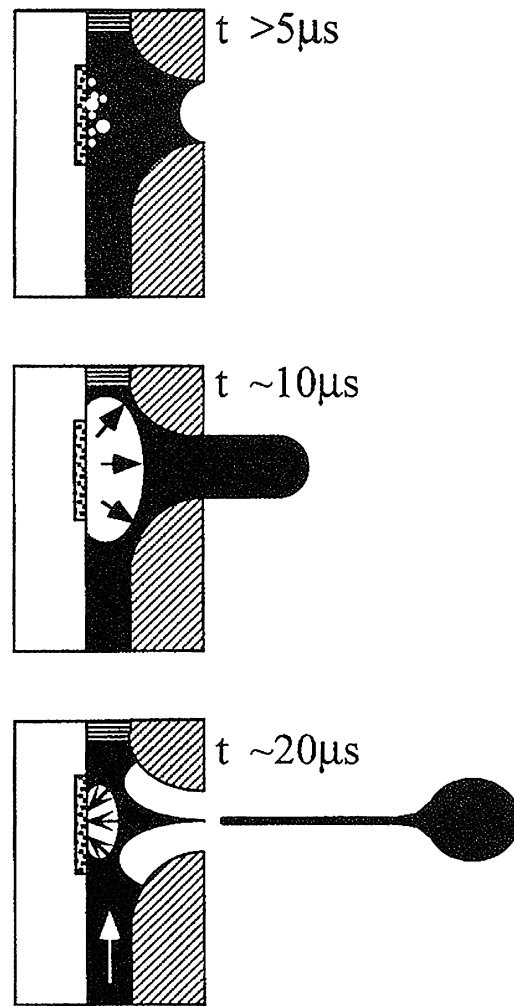


FIGURE 10 A thermal ink jet print head, showing the evolution of bubble formed from the ink heated by the resistive heating element. [From Le, H. P. (1998). *J. Imaging Sci. Technol.* **42**, 49–62; by permission of the copyright holder, the Society for Imaging Science and Technology.]

duced for specific product applications, and incorporate unique features, e.g., coatings, which enable these applications. Yet printer manufacturers offer products which they advertise as designed to print on "plain" paper. Water- or solvent-based inks, once transferred to the paper, dry by a combination of evaporation and penetration of the water or solvent into the paper, leaving the colorant on or near the surface. The challenge in the design of ink-jet papers is to facilitate the absorption of the solvent into the paper while minimizing the spread of the ink in the surface layers of the paper. The development of specialty coatings which facilitate this effect has been driven recently by the trend to higher resolution print-head technology and smaller single drop volumes, offering 1440 dots per inch in the home and office markets. As discussed above in connection with

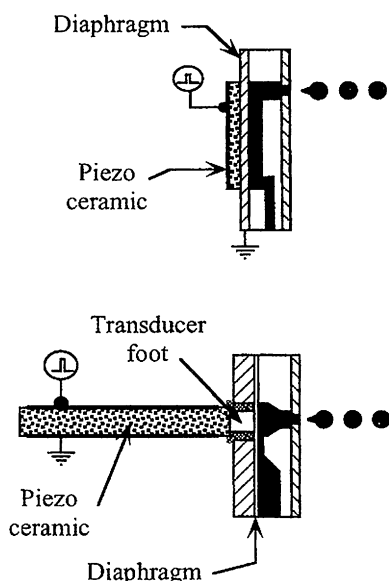


FIGURE 11 Two geometries of piezoelectric driven drop-on-demand print heads: bend-mode (top) and push-mode (bottom). (From Le, H. P. (1998). *J. Imaging Sci. Technol.* **42**, 49–62; by permission of the copyright holder, the Society for Imaging Science and Technology.)

image quality, such resolving powers are necessary to the achievement of photographic quality in digital imaging. One manufacturer (Canon) has commercialized a printing system which applies an appropriate coating to the paper as the first step in the printing process, to allow effective ink-jet printing. Colloidal silica is a principal ingredient in the coatings used on ink-jet papers, and also on films designed for making transparencies on ink-jet printers. It dries to yield a porous, hydrophilic coating, which the ink can penetrate quickly, avoiding surface smearing while the solvent is eliminated by evaporation and wicking into the paper substrate.

The most critical material in ink-jet printing is, of course, the ink. The rheology of the ink not only affects the appearance and quality of the printed image, it also determines the drop ejection characteristics of the print head. Thus a specific ink, or set of inks, is usually designed for a specific print-head design. Water-based inks are usually used for drop-on-demand home and office printers, largely to avoid the problem of solvent vapors in the workplace environment. Water is also the material of choice for thermal ink-jet print heads because of its bubble-forming ability. Because water penetrates so easily into paper, both laterally and in depth, very poor resolution and optical density may result from the use of water-based inks on uncoated papers. Therefore solvent-based inks are used in continuous ink-jet printers for the highest quality commercial printing applications. Solvent-based inks must also be

used when substrates which do not allow absorption or penetration are involved, e.g., glass, metal, or many plastics. Ink-jet colorants may be dyes or pigments. Most current inks employ dyes. Pigments, however, offer several advantages. Unlike dye molecules, the pigment particles are trapped as deposited in the pores of the surface coating on the substrate, or remain adhered to the surface fibers of an uncoated paper. This leads to more intense coloration, as well as higher resolution in the image. The pigments also tend to be more inert chemically and less prone to fading than typical dyes, thus providing more stable and durable printed images. On the down side, however, pigments are more expensive than dyes, and the dispersion of the pigment particles in the ink vehicle may not be stable, leading to clogging of the ink-jet orifices.

A recent development ink-jet printing technology is the use of phase change or so-called hot melt inks. Instead of water or a solvent, in these inks the carrier is a wax which is solid at room temperature. It is heated in a reservoir in the printing engine to provide a supply of liquid ink for the drop-on-demand print heads, which usually employ the piezoelectric method. When the ink is transferred to the substrate, it cools to room temperature and solidifies immediately. To obtain good adhesion to, e.g., paper, the ink drop is then fused in place using a pressure fusing roller, just as is done with toner images in electrophotography.

Among emerging applications, an offset printing press (digital color press) based on ink-jet printing with phase change inks has been demonstrated by Tektronix (non Xerox copy.). In this machine, the molten ink is used to print the image onto a heated, silicone oil coated aluminum roller, which is maintained above the glass transition, i.e., softening, temperature of the ink, but below its melting point. The printed image is thus stable on this roller, but can be transferred quickly and easily under pressure to the paper, which has also been preheated. This press provides very high speed printing on demand from digital image data. Still another approach to volume printing using phase change inks involves formation of an oleophilic image (the solidified ink) on a hydrophilic, usually aluminum, substrate. The resulting composition can then be run as a printing plate on a conventional offset press. This strategy enables direct generation of a printing plate from a digital data stream, i.e., computer-to-plate. (See below for another example of computer-to-plate.)

C. Thermal Printing

There are two main approaches to thermal printing, direct thermal and transfer thermography. Some examples of direct thermal printing were discussed, above, in the section on thermally developed silver systems. Current applications for this technology are for the most part

in the low-end market, i.e., low cost, low image quality applications such as labels, bar codes, etc. Most applications for higher quality thermal printing involve transfer thermography. Up until the 1990s thermal print heads were the principal means of image input for both direct and transfer thermography. When digital photography was first introduced to the marketplace in the early 1980s with the Sony MAVICA camera, transfer thermography was the technology of choice for image display. However, with its higher resolution, greater dynamic range (more gray levels), and lower cost, ink-jet printing has now completely taken over this application. More recently new materials have been developed for laser thermography, whereby the image can be written on the imaging material with a scanned near-infrared laser, using the thermal energy of the laser beam, not any photochemical effects.

1. Transfer Thermography

The thermal print heads used in transfer thermography convert an electrical signal into heat by heating of an electrically resistive element in contact with the imaging medium. These resistive elements are of three types: (1) thick film resistors, which are relatively inexpensive; (2) thin film resistors, which are produced by sputtering, are capable of the highest resolution, and, being expensive, are usually protected by an anti-abrasion layer to prevent their being worn down by the print medium; and (3) silicon resistors, in which the individual resistive elements are created in a silicon monolith by conventional microelectronic fabrication technology. The resolution achievable in thermal printing is, however, not limited ultimately by the print head fabrication technology, but by thermal diffusion in the imaging medium itself, i.e., spot spread. At the peak of its market penetration in the early 1990s, a resolving power of 400 dots per inch was standard for color transfer thermography, though premium engines, costing over \$ 10,000, and their media were capable of, e.g., 720 dots per inch. Compare these figures with the current standards of 1440 and 2880 dots per inch in low cost (under \$ 100) home and office ink-jet printers. Transfer thermography continues to be popular, however, for some specialized applications such as color proofing in the graphic arts and textile printing.

There are two principal technologies for transfer thermography, mass transfer, sometimes called wax transfer, and sublimation printing. Both strategies require a donor sheet, one for each primary color to be printed. The donor sheet(s) passes over the print head in contact with the receptor, i.e., the paper. In the former technique, the heat melts a wax-based ink coated on the donor, and causes it to transfer to the receptor; in terms of ink for-

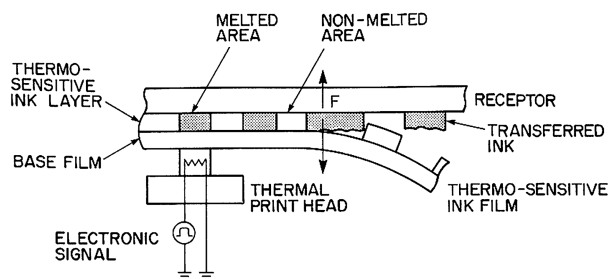


FIGURE 12 Thermal mass transfer printing process. [From Ohno, T. (1981). *J. Appl. Photogr. Eng.* 7, 171; by permission of the copyright holder, the Society for Imaging Science and Technology].

mulations these inks are not dissimilar to those used in phase change ink-jet printing, and colorants may be either dyes or pigments. The former are to be preferred both for cost and transparency, where colors are to be overlaid; the latter offers better light fastness and more intense coloration. The method is illustrated in Fig. 12. Sublimation printing differs in that only the colorant, in this case a volatile dye, is transferred by evaporation from a relatively heat-insensitive donor sheet. An important distinction between wax transfer thermal printing and sublimation printing is that the former is binary, transfer is an all-or-nothing affair, while the latter is, in principle, capable of continuous tone, i.e., longer or more intense heat pulses transfer more dye than shorter or less intense ones. In any digital printing scheme, however, true continuous tone is not achieved. Owing to the digital input to the print heads, color is transferred to the receptor in quantized levels (gray levels). If there are enough levels and the density differential between levels is small enough, the image is perceived by the viewer as exhibiting continuously variable tones. Psychometric studies have led the imaging industry to accept 256 gray levels (8 bits per pixel) as visually equivalent to continuous tone; these studies were initially carried out using sublimation transfer media.

Additional disadvantages of transfer thermography are its inefficiency, and the image registration problem. The donor is the expensive consumable in transfer thermography. As can be seen from Fig. 12, where the ink isn't transferred to the receptor it remains on the donor, which must be discarded after one use. Usually that means that most of the ink or sublimable dye of each color ends up being wasted. It is also necessary to apply the primary colors sequentially, much as is done on a conventional printing press. This means that the physical position of the donor and print heads for each step must correspond precisely with that for the donor and print heads in the preceding and following steps. Another factor limiting resolution in transfer thermography is thus the mechanical precision with which this registration, as it is called, can be achieved.

2. Laser Thermography

In dye sublimation transfer printing by laser heating, it is possible to obtain high resolution and continuous tone, because a modulated and focused near-infrared laser beam is used as the heat source. As the same laser can be used as the heat source for transfer of each primary color, only the donor sheet, not the heat source or receptor, need move during the sequential deposition of the colors, so the registration problem is also obviated.

For laser thermographic imaging media, sensitivity is usually expressed in terms of the exposure energy required to produce full modulation of the medium, i.e., maximum density on the image. A good benchmark for the energy requirement of typical laser thermographic materials is of the order of 250 mJ/cm^2 , which is about six orders of magnitude greater than representative AgX media. For producing a $20 \times 30 \text{ cm}$ image in one minute on a material requiring this exposure energy, a scanning laser power of 2.5 W would be required. In recent years, diode lasers, which can be electronically modulated at the speeds required for such an imaging system, have become cost effective, but typically these devices do not offer more than 100 MW power. In practice, it is therefore necessary to combine optically the output of an array of several diodes in one spot in order to image these materials. Several different schemes, involving conventional and fiber optics, are used to this end in the commercial printers for laser thermography, and their availability has driven the development and commercialization of this technology.

In laser thermography a slightly more complicated donor sheet is required, however, as shown in Fig. 13. In addition to the supporting film base and the ink layer incorporating the volatile dye, a light-absorbing layer must be interspersed. The laser light impinging on this layer is ab-

sorbed and converted to heat, which then effects the evaporation of the dye. Carbon black is a usual constituent of the light-absorbing layer, though dyes which selectively absorb the laser wavelength, usually between 750 and 900 nm , have also been disclosed, e.g., titanil phthalocyanine for use with an 825-nm laser. Typical laser spot sizes are of the order of $3\text{--}20 \mu\text{m}$, and resolving powers of over 8000 dots per inch have been reported for this method. In order to obtain continuous tone, e.g., for photo printing, it is generally preferable to vary input energy to an individual pixel by pulse width modulation, i.e., varying the length of time the laser is on while exposing that pixel. Attempting to modulate the laser intensity typically leads to overheating of the donor; in this case not only does the dye transfer, but the polymer binder in which the dye is dissolved also decomposes. Above the onset of binder decomposition, termed ablation, the whole layer is blown apart and essentially all the dye in it is transferred to the receiver, with concomitant loss of control over the image density on the receptor. It should be realized that, if we think of the image on the receptor as a positive, then the donor sheet bears the corresponding negative image.

Laser thermal imaging media have found an important niche in the graphic arts, although their application in digital photofinishing has also been discussed. The Kodak Approval color proofing system is based on sublimation transfer. It functions essentially as shown in Fig. 13, but has a couple of unique features. First of all the four primary colors are sequentially transferred to an intermediate receptor. The dye image is transferred from the intermediate support to the final print stock. Secondly, the donor is separated from the receptor by a close-packed layer of spacer beads, typically from $4\text{--}20 \mu\text{m}$ in diameter, between each donor and receptor. These beads both prevent the donor and receptor from sticking to each other, and they thermally isolate the receptor from the donor, thereby preventing the receptor from acting as a heat sink and wasting laser energy.

Laser ablation of the donor layer can, however, be exploited to produce the image. In the Kodak Professional Direct Image film, no receptor is used. Intense laser exposure simply blows away the dyed polymer layer corresponding to the donor in Fig. 13. Like wax transfer thermography, the imaging process is binary, and a high contrast dye image results. As noted above, this process produces an image which is the negative of that normally produced by transfer thermography. Resolution of up to 1800 dots per inch has been demonstrated for this technology. Laser ablation imaging tends to require very high laser energies, in order to produce the requisite exposures of several hundred mJ/cm^2 .

Several strategies have been developed to utilize thermal imaging in computer-to-plate applications. The major

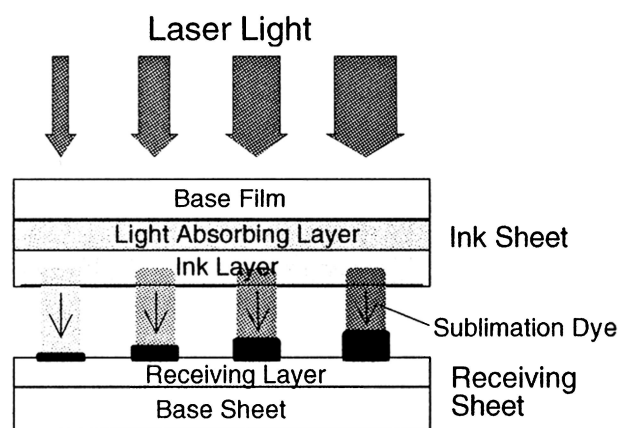


FIGURE 13 Schematic of laser dye thermal transfer printing. [From Kitamura, T. *et al.* (1999). Proceedings of IS&T's 1999 PICs Conference, IS&T, Springfield, VA; by permission of the copyright holder, the Society for Imaging Science and Technology].

attraction of computer-to-plate from the point of the printing industry is the simplification in prepress workflow that results from being able to take the digital data stream representing the material to be printed and creating the necessary press-ready plates directly therefrom. In one approach (Presstek) amplification is obtained, and laser power requirements reduced by coating an energetic polymer, e.g., nitrocellulose, as a thin film under the light-absorbing layer of a construction that is otherwise like the donor sheet in Fig. 13. The idea behind this approach is to use the chemical energy in such materials to assist the laser. When the laser heats the spot above some threshold temperature, the polymer decomposes explosively by a self-oxidizing mechanism, which produces gaseous products. Rapid expansion of the gas bubble physically disrupts the light-absorbing and imaging layers, allowing their easy removal. Amplification factors obtained by this strategy are typically three or four, compared to standard laser ablation imaging materials. For printing-plate applications, the exposed support is ink receptive. The imaging layer (corresponding to the ink layer in Fig. 13) is a silicone polymer which is ink repellent. The resulting imaged construction can be put on a dry offset press and used as a printing plate.

VIII. CONCLUSIONS

Despite its enormous achievements and its major role in our lives, the field of photographic processes and materials is very far from being a static one. Over the past decade we have seen it evolve from the traditional, photochemistry-based technology into the field of imaging, in which various digital technologies exist side-by-side with the photochemical and photophysical ones. Most major periodical publications and conference series which address picture taking and picture making technology have substituted the term "Imaging" for "Photography" in their titles over the past decade. Imaging technologies have provided a unified technological base for photography, office copying, microfilm, medical imaging, and the graphic arts.

Scenarios for the future, which were envisioned in this entry in the previous edition of the *Encyclopedia of Physical Science and Technology*, have nearly all become realities. Important engines driving the evolution of photography into imaging have been the explosive expansion of the graphics capabilities of PCs, which has made image processing and desktop publishing widely accessible, and the image communication capability of the internet. With Moore's Law continuing to apply to PCs and the rapid penetration of broad bandwidth digital communication technologies for the internet, it is apparent that these trends will continue, with even more profound impacts on the practice of imaging. It is important to realize, however,

that there is major room for improvement in all the currently existing technologies, as described above, and the historically important ones of silver halide photography and electrophotography have shown themselves to be especially robust and resilient in adapting to the new world of imaging.

Despite the convenience, speed, and flexibility of digital photography, the attractiveness of AgX as an art form promises to keep this technology alive for generations to come. It is not perhaps too great a stretch of the imagination to see the future of traditional photography in terms of the revitalization and liberation experienced by painting and drawing when photography first appeared in the 19th century.

SEE ALSO THE FOLLOWING ARTICLES

COLOR SCIENCE • IMAGE PROCESSING • METAL CLUSTER CHEMISTRY • OPTICAL INFORMATION PROCESSING • PHOTOCHROMIC GLASSES • PHOSPHORS • POLYMERS, ELECTRONIC PROPERTIES • POLYMERS, PHOTORESPONSIVE (IN ELECTRONIC APPLICATIONS) • SMART PIXELS

BIBLIOGRAPHY

- Allen, J. B. (1999). Comparison of the single pixel development of DMD and laser exposure modules in electrographic printing. *J. Imaging Sci. Technol.* **43**, 309.
- Anon., ed. (2000). Proceedings of IS&T/SPSTJ's International Symposium on Silver, IS&T, Springfield, VA.
- Bartscher, G., Cormier, S. O., Hauptmann, G., Kostyk, D., Lyness, R., Morgenwick, F., Preissig, K.-U., Rohde, D., and Schein, L. (2001). Single pixel development in a new electrographic printing system. *J. Imaging Sci. Technol.* **45**, in press.
- Carroll, B. H., Higgins, G. C., and James, T. H. (1980). "Introduction to Photographic Theory. The Silver Halide Process," Wiley (Interscience), New York.
- Chernov, S., and Zakharov, V. (2001). "Additional Treatment of Photomaterials," MAXPress, Moscow (in English).
- Dagani, R. (2001). Here comes Paper 2.0, *Chem. Eng. News*, pp. 40–44; <http://pubs.acs.org/cen>.
- Dainty, J. C., and Shaw, R. (1974). "Image Science," Academic Press, London.
- DeBoer, C. (1998). Laser thermal media: the new graphic arts paradigm. *J. Imaging Sci. Technol.* **42**, 63–69.
- Dessauer, J. H., and Clark, H. E., eds. (1965). "Xerography and Related Processes," Focal Press, London.
- DeVoe, R. J., Olofson, P. M., and Sahyun, M. R. V. (1992). Photochemistry and photophysics of 'onium salts. *Adv. Photochem.* **17**, 313–356.
- Jacobson, R. E. (1983). *J. Photographic Sci.* **31**, 1–12.
- Jacobson, R. E., Ray, S. F., and Attridge, G. G. (1988). "The Manual of Photography," 8th ed., Focal Press, London.
- James, T. H. ed. (1977). "The Theory of the Photographic Process," 4th ed., Macmillan, New York.

- Kitamura, T., Kinoshita, M., and Hoshino, K. (1999). A High Definition Continuous Tone Color Image in Dye Thermal Transfer Printing by Laser Heating, Proceedings of IS&T's 1999 PICs Conference, IS&T, Springfield, VA.
- Koulikov, S. G., and Dlott, D. D. (2000). Effects of Energetic polymers on laser photothermal imaging materials. *J. Imaging Sci. Technol.* **44**, 111–119.
- Le, H. P. (1998). Progress and trends in ink jet printing technology. *J. Imaging Sci. Technol.* **42**, 49–62.
- Marchetti, A. P., and Eachus, R. S. (1992). Photochemistry and photo-physics of the silver halides. *Adv. Photochem.* **17**, 145–216.
- Ohno, S., Fujii, T., Oka, K., and Kato, N. (1998). Image quality of digital photography Prints-I. *J. Imaging Sci. Technol.* **42**, 269–275; Ohno, S., Takakura, M., and Kato, N. (2000). Image quality of digital photography Prints-II. *J. Imaging Sci. Technol.* **44**, 51–60.
- Ohno, T. (1981). High speed thermal ink transfer recording and its applications. *J. Appl. Photogr. Eng.* **7**, 171.
- SahyM. R. V. (1998). Thermally developable photographic materials (TDPM). *J. Imaging Sci. Technol.* **42**, 23–30.
- Schaffert, R. M. (1975). "Electrophotography," Focal Press, London.
- Schulze, D. (2000). Photofinishing report. *IS&T Reporter*, **15**, 12–20; <http://www.imaging.org>.
- Shaw, R. (1999). Digital Photography: the Influence of CCD Pixel Size on Imaging Performance, Proceedings of IS&T's 1999 PICs Conference, IS&T, Springfield, VA.
- Sturge, J. M., Walworth, V. K., and Shepp. A., eds. (1989). "Imaging Processes and Materials, Neblette's Eighth Edition," van Nostrand Reinhold, New York.
- Tani, T. (1995). "Photographic Sensitivity," Oxford University Press, Oxford, U.K.
- Tani, T. (1998). Progress and future prospects of silver halide photography compared with digital imaging. *J. Imaging Sci. Technol.* **42**, 1–14.
- Weigl, J. W. (1977). *Angewandte Chemie*, International ed. in English **16**, 374–392.



Photonic Bandgap Materials

Sajeev John
Ovidiu Toader

University of Toronto

Kurt Busch

Universität Karlsruhe

- I. Introduction
- II. Photonic Bandgap (PBG) Formation
- III. Two-Dimensional PBG Materials
- IV. Two-Dimensional Photonic Crystal Devices
- V. Synthesis of Three-Dimensional PBG Materials
- VI. Quantum and Nonlinear Optics in Three-Dimensional PBG Materials
- VII. Tunable PBG Materials
- VIII. Outlook

GLOSSARY

Bandgap A range of frequencies or energies where propagating modes of a wave are absent.

Light localization A state in which light of a particular frequency is completely confined to a small, finite region of space and cannot propagate away except through some nonlinear interaction.

PBG material A non-light-absorbing material which contains a bandgap for electromagnetic waves propagating in all directions.

Photonic Crystal A non-light-absorbing material with a refractive index which exhibits periodic modulation in at least two orthogonal spatial directions.

Photonics The science of moulding the flow of light.

LIGHT IN CERTAIN engineered dielectric microstructures can flow in a way similar to electrical currents in

semiconductor chips. These microstructures represent a new frontier in the field of optics. They provide a foundation for the development of novel micro-photonic devices and the integration of such devices into an optical microchip.

I. INTRODUCTION

Electromagnetism is the fundamental mediator of all interactions in atomic physics and condensed matter physics, in other words, the force that governs the structure of ordinary matter. In a novel class of engineered dielectric materials known as photonic bandgap (PBG) materials, a fundamentally new electromagnetic effect can be realized. This phenomenon is the localization of light ([John, 1984](#); [Anderson, 1985](#)) and it may prove central to the utilization of light waves for information and communication technologies.

In the 19th century, James Clerk Maxwell deduced a precise and elegant mathematical description of the propagation of light. Maxwell's theory of electromagnetic wave propagation was shortly thereafter tested and verified by Heinrich Hertz. There began the age of wireless communication. This allowed ships at sea to communicate with land. Today, the same basic discovery provides us with the use of radio, television, and mobile phones. Despite the widespread knowledge and use of Maxwell's theory of electromagnetism, it has only recently been recognized that light waves not only propagate, they can also be trapped.

The invention of the laser in the latter half of the 20th century was another milestone in the science and technology of light. Laser light allows us to probe the structure of matter with unprecedented accuracy, it provides medical practitioners with a cutting tool sharper than any surgeon's knife, and allows us to modulate communication signals for high-speed data transfer along the Internet. Today, using laser light, undersea fiber optic cables carry enormous amounts of voice communications and data with such clarity that a pin drop in Karlsruhe, Germany, can be heard clearly in Toronto, Canada.

Optical fibers are replacing electrical wires in shorter distance communications such as local access networks and computer-to-computer communications. In a completely seamless network, communications between nearby computer chips and even within a single computer chip would take place with tiny beams of laser light rather than electricity. Optical computers of this type may be faster and support neural architectures (circuit interconnections resembling that of a human brain), unlike their electronic counterparts, which are restricted in architecture due to electrical cross-talk between nearby wires. That is to say, electrical current in one wire can disturb electrical signals passing through nearby wires. Laser beams, on the other hand, can coexist without disrupting one another.

The trapping and micromolding of light flow needed for the applications suggested above requires materials which can scatter light much more strongly than any naturally occurring material. We experience multiple light scattering when it becomes dark on a cloudy day. Light from the sun scatters many times from water droplets, following a tortuous diffusion path before reaching the ground. The distance the light travels within the cloud before it is scattered into a random direction is called the mean free path. The effect of multiple scattering is that the amount of light transmitted through the cloud is reduced by a factor of the ratio of the cloud thickness to the mean free path. The rest comes back out the other side, which is why clouds appear white. Multiple light scattering also takes place in human tissue. Here the transport

mean free path for light of 1 μm wavelength is about 1 mm.

Neither clouds nor human tissue can scatter light sufficiently strongly to localize light. For this to happen a collection of microscopic dielectric structures that scatter light 1000 times more strongly than human tissue is required. In this case the transport mean free path becomes as short as the wavelength of light itself. If, in addition to the strong resonant scattering of the individual dielectric particles, there is a periodic arrangement of the scatterers, then pathways for light propagation over specific frequency intervals can be completely removed. The removal of pathways over all directions over a band of frequencies is referred to as the creation of a photonic bandgap (PBG). Dielectric microstructures that exhibit this effect are called PBG materials.

In electronic microcircuits, electrical currents are guided by thin metal wires. Electrons are bound within the cross section of the wire by the so-called work function (confining potential) of the metal. As a result, electrical currents follow the path prescribed by the wire without escaping into the background. The situation is very different for optical waves. Although optical fibers guide light over long distances, microcircuits of light based on fibers do not exist. This is because empty space is already an ideal conductor of light waves. The light in an optical fiber can easily escape into the background electromagnetic modes of empty space if the fiber is bent or distorted on a microscopic scale. PBG materials remove this problem by removing all the background electromagnetic modes over the relevant band of frequencies. Light paths can be created inside a PBG material in the form of engineered waveguide channels. The PBG localizes the light and prevents it from escaping the optical microcircuit.

The question of whether light can be localized can be posed in the form of an analogy between Maxwell's equations for electromagnetic wave propagation and Schrödinger's equation for electrons propagating in a scattering potential. Consider a monochromatic electromagnetic wave of frequency ω propagating in a medium whose dielectric constant varies from point to point in space as

$$\varepsilon(x) = \varepsilon_0 + \varepsilon_{\text{fluct}}(x). \quad (1)$$

Here, ε_0 is the average part of the dielectric constant and $\varepsilon_{\text{fluct}}(x)$ represents the part of the dielectric constant that varies from point to point in space. We assume that the dielectric microstructure does not absorb the light and that the total dielectric constant is everywhere real and positive. The wave equation for such an optical field is given by

$$-\nabla^2 \vec{E} + \vec{\nabla}(\vec{\nabla} \cdot \vec{E}) - \frac{\omega^2}{c^2} \varepsilon_{\text{fluct}}(x) \vec{E} = \frac{\omega^2}{c^2} \varepsilon_0 \vec{E}. \quad (2)$$

Equation (2) has been written in a form which resembles the quantum mechanical Schrödinger equation. Here, the first two terms are the analogs of the “kinetic energy” terms in Schrödinger’s equation, $(\omega^2/c^2)\varepsilon_{\text{fluct}}(x)$ plays the role of a scattering potential, and $(\omega^2/c^2)\varepsilon_0$ is analogous to an energy eigenvalue. The Maxwell wave equation (2) describes a vector wave field as opposed to complex scalar wavefunction in Schrödinger’s equation.

The essential subtlety of light localization is apparent from Eq. (2). Unlike electrons, which can be trapped and localized for negative eigenvalues (bound states) in corresponding negative-energy potential wells, the overall positivity of the dielectric constant (1) leads to the constraint that the energy eigenvalue, $(\omega^2/c^2)\varepsilon_0$ is always greater than the highest of the potential barriers presented by $(\omega^2/c^2)\varepsilon_{\text{fluct}}(x)$. The occurrence of bound states of the electromagnetic wave field in this spectral range requires extremely specialized and carefully engineered dielectric materials.

II. PHOTONIC BANDGAP (PBG) FORMATION

PBG materials were predicted theoretically as a means to realize the localization and trapping of light in a bulk material over a band of frequencies (John, 1987). A direct corollary of this principle is the complete inhibition of spontaneous emission (Yablonovitch, 1987; Bykov, 1975) from an atom, molecule, or electron–hole pair excitation that is placed within the PBG material. Indeed, if the emission frequency from the atom lies within the PBG, the photon that would normally be emitted forms a bound state to the atom. Nearly all of the novel consequences of PBG materials are a direct consequence of these remarkable effects. Unlike optical confinement of a single resonance mode in a high-quality (Q factor) optical cavity, localized electromagnetic modes in a bulk PBG material are completely decoupled from the vacuum modes of free space. Unlike well-known layered dielectric structures (including Fabry–Perot resonators and distributed feedback laser cavities) which may confine light in one spatial dimension, the PBG material facilitates coherent localization of light in all spatial directions. This unique combination of light localization and the complete control of radiative dynamics distinguishes PBG materials from any previously studied optical system.

Photonic bandgap formation can be understood as a synergistic interplay between two distinct resonance scattering mechanisms. One is the microscopic scattering resonance from the dielectric material contained in a single unit cell of the photonic crystal. A simple illustration of

this is provided (Fig. 1) by the scattering of a wave from a square well potential. When one-half of the optical wavelength fits into the width of the well, the transmission of light from left to right is maximum and the least light is reflected. When one-quarter of a wavelength fits into the width of the well, the least amount of light is transmitted and the maximum amount of light is reflected. This quarter-wavelength condition is a simple example of the condition for a microscopic scattering resonance. The second resonance is the macroscopic resonance from the geometrical arrangement of repeating unit cells in the dielectric microstructure. If there is a periodic arrangement of unit cells, this is called Bragg scattering. This occurs whenever the spacing between adjacent unit cells is an integer multiple of half of the optical wavelength. Photonic bandgap formation is facilitated if the geometrical parameters of the photonic crystal are chosen so that both the microscopic and macroscopic resonances occur at precisely the same wavelength. In addition, both of these scattering mechanisms must individually be quite strong. In practice, this means that the underlying solid material must have a very high refractive index (typically about 3.0 or higher) and at the same time exhibit negligible absorption or extinction of the light (less than 1 dB of attenuation over 1 cm). These conditions on the scattering strength, geometry, and purity of the dielectric material severely restrict the set of engineered dielectrics that exhibit a PBG. Candidate materials for the PBG backbone include silicon, germanium, gallium arsenide, and indium phosphide.

Scalar waves (sound waves and other compression waves in solids) readily exhibit complete three-dimensional gaps for simple structures such as a face-centered cubic (FCC) lattice of spherical scatterers. In contrast, electromagnetic waves involve three-component electric and magnetic field vectors. This leads to much more restrictive conditions on the dielectric constant, the solid volume fraction, and the microstructure connectivity for the formation of a PBG. One very widely studied class of PBG materials is that based on a diamond lattice of dielectric scatterers. For example, a diamond lattice of nonoverlapping spheres of very high dielectric constant exhibits a PBG for which electromagnetic wave propagation is completely forbidden over a narrow frequency band. An “inverse diamond” lattice of overlapping air spheres in a high-refractive index background (Fig. 2) exhibits a much larger PBG. In the microwave regime, other diamond-like structures obtained by drilling cylindrical holes in a bulk dielectric material (with refractive index of 3.5) have been demonstrated to exhibit bandgap-to-center frequency ratios as large as 20% (Yablonovitch *et al.*, 1991). Since then, numerous structures amenable to layer-by-layer fabrication have been

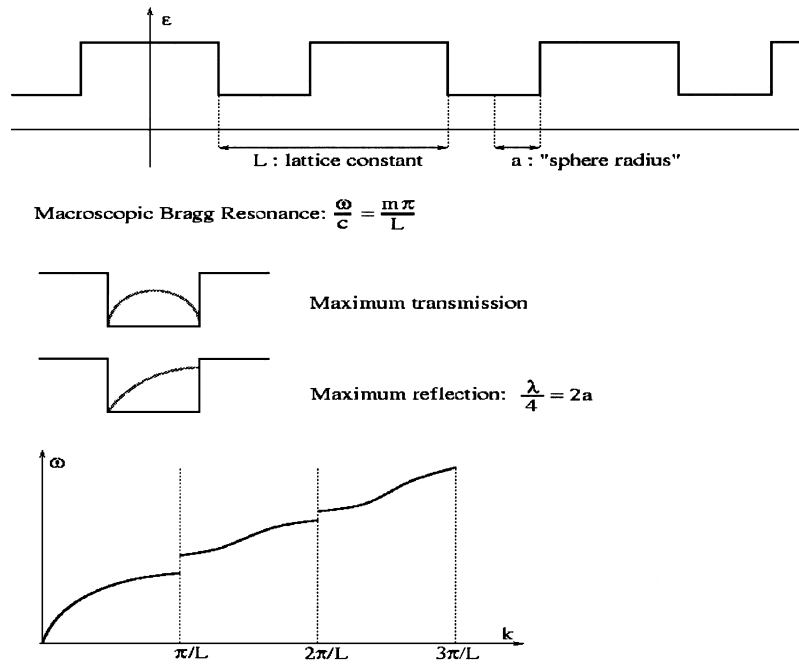


FIGURE 1 PBG formation can be regarded as the synergetic interplay between two distinct resonance scattering mechanisms. The first is the “macroscopic” Bragg resonance from a periodic array of scatterers. This leads to electromagnetic stop gaps when the wave propagates in the direction of periodic modulation when an integer number $m = 1, 2, 3, \dots$, of half-wavelengths coincides with the lattice spacing L of the dielectric microstructure. The second is a “microscopic” scattering resonance from a single unit cell of the material. In the illustration, this (maximum backscattering) occurs when precisely one-quarter of the wavelength coincides with the diameter $2a$ of a single dielectric well of refractive index n . PBG formation is enhanced by choosing the materials parameters a , L , and n such that both the macroscopic and microscopic resonances occur at the same frequency.

suggested, the most notable being the “woodpile” structure (e.g., [Ho et al., 1993](#)). A number of structures have already been fabricated with PBGs in the range of millimeter waves. The diamond structure and its cousins constitute a family of PBG materials which are characterized by a large and complete PBG between the second and third bands (fundamental gap) in the photonic band structure. The FCC lattice structure, on the other hand, does not exhibit a complete PBG between the second and third bands. Instead it has a small PBG between the eighth and ninth bands. While it has been a common belief that only the diamond structure is associated with large and complete PBGs, very recently a new class of microstructures based on a tetragonal lattice has been discovered that exhibits a large PBG between the fourth and fifth electromagnetic dispersion bands ([Toader and John, 2001](#)). An illustration of such a lattice consisting of square spiral posts is shown in [Fig. 3](#). When the posts are made of silicon, this structure exhibits a 15% PBG. The corresponding electromagnetic density of states is shown in [Fig. 4](#). This structure is amenable to micro-fabrication using a technique called glancing angle deposition (GLAD) ([Robbie and Brett, 1996](#); [Rollbbee et al., 1998](#)).

Simple applications of PBG materials can be found in the microwave- to millimeter-wave range. For instance, an antenna mounted on a conventional dielectric substrate emits the majority of its radiation into the substrate itself. If the substrate is engineered into the form of a PBG material with a gap at the radiation frequency, the losses can be minimized, leading to highly directional transmitters. Other applications include optical filters with tailor-made characteristics and cladding material for preventing losses in waveguide structures that contain bends or junctions. Nevertheless, it is for visible and near-infrared frequencies where PBG materials are likely to have their most important impact. For example, applications in telecommunications require the fabrication of PBG materials with a gap centered in the 1.3- to 1.5- μm range, where there is minimal absorption of light in silica-based optical fibers. These applications include the design of ultracompact lasers that emit coherent light with almost no pumping threshold, all-optical switching fabrics for routing data along the Internet, optical switches with an on-off cycle time of less than one trillionth of a second, and all-optical transistors.

PBG materials also represent a new frontier in fundamental aspects of quantum and nonlinear optics. While

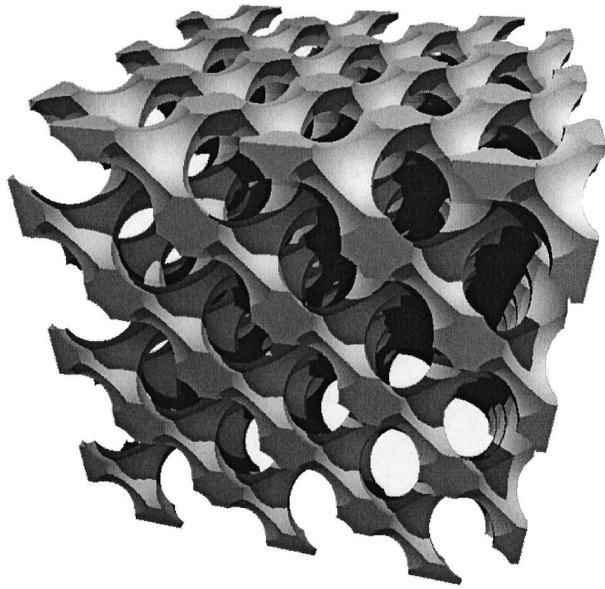


FIGURE 2 The “inverted diamond” structure was one of the first prototype structures predicted to exhibit a large and robust 3D PBG. It consists of an overlapping array of air spheres arranged in a diamond lattice. This structure can be mimicked by drilling an array of criss-crossing cylindrical holes into a bulk dielectric. The solid backbone consists of a high-refractive index material such as silicon leading to a 3D PBG as large as 27% of the center frequency. The minimum refractive index (of the backbone) for the emergence of a PBG is roughly 2.0. Practical difficulties in the synthesis of this structure have motivated simpler but closely related designs such as the “woodpile” structure (see Fig. 9).

linear wave propagation is absent in the gap of a PBG material, nonlinear propagation effects can still occur. When the backbone of the PBG exhibits a nonlinear refractive index (which depends on the intensity of the electromagnetic wave field) certain high-intensity “light bullets” (solitary waves) can pass through the material even at frequencies within the gap. In addition, PBG materials exhibit novel quantum optical features, related to the drastic alteration of the photon density of states (DOS). A vanishing DOS leads to bound photon–atom states, suppressed spontaneous emission (Yablonovitch, 1987; Bykov, 1975), and strong localization of photons (John, 1987). Localization of photons implies that emission of light from an initially excited atom occurs in a way that is very different from that in ordinary vacuum. In a PBG material, the atom has a long-time memory of the fact that it was optically excited at an earlier time. One consequence of such memory and intrinsic feedback is that lasing can occur at a photonic band edge without recourse to a standard laser cavity involving a pair of mirrors. Other novel phenomena predicted to occur in a PBG material include (i) collective switching of two-level atoms from ground to ex-

cited state with low-intensity applied laser fields, leading to all-optical transistor action (John and Quang, 1997), (ii) single-atom memory effects for possible quantum computer applications (Quang *et al.*, 1997), and (iii) low-threshold and other anomalous nonlinear optical response (John and Quang, 1996).

III. TWO-DIMENSIONAL PBG MATERIALS

In many applications, the polarization state of guided light can be fixed in a particular direction and only the passive optical guiding characteristics of a PBG material come into play. Two-dimensional (2D) periodic microstructures are often sufficient for such applications. For 2D periodic dielectrics, advanced planar microstructuring techniques borrowed from semiconductor technology can greatly simplify the fabrication process. Such structures are referred to as photonic crystals exhibiting a 2D PBG. The “aspect ratio” of a 2D PBG material is defined as the ratio of the sample depth (vertical direction) to the lattice constant (transverse direction).

High-quality photonic crystals with aspect ratios of up to 5:1 can be manufactured through plasma etching or electron beam lithography techniques (e.g., Painter *et al.*, 1999). These are sometimes referred to as membrane structures. Alternatively, high-aspect ratio 2D PBG structures can be manufactured by photoelectrochemically growing ordered macropores into a silicon wafer (Grüning *et al.*, 1996). With these 2D PBG materials aspect ratios of 200:1 and 2D bandgaps centered at wavelengths at 1.3–1.5 μm have been achieved (Schilling *et al.*, 2001).

The quality of 2D PBG materials synthesized using the photoelectrochemical process is shown in Fig. 5. Here, we depict a bar of macroporous silicon consisting of 22 pore layers with a lattice constant of 1.5 μm . This structure exhibits a 2D PBG centered near 3 μm . That is to say there is a band of frequencies over which light polarized with either the electric field (E-polarization) or magnetic field (H-polarization) parallel to the pores cannot propagate through the material. During the fabrication process of a 2D PBG material, light paths within the gap can be engineered through the introduction of defects. For instance, if a single pore is modified or left out altogether (by placing an etch mask over silicon wafer), an optical microcavity is formed and leads to a localized mode of light inside the PBG. A chain of such point defects can act as a linear waveguide channel and facilitate very sharp waveguide bends. It can also provide ultrasmall beamsplitters, Mach–Zehnder interferometers, and functional microoptical elements such as wavelength add-drop filters (Fan *et al.*, 1998). Some defect structures are shown in Fig. 6.

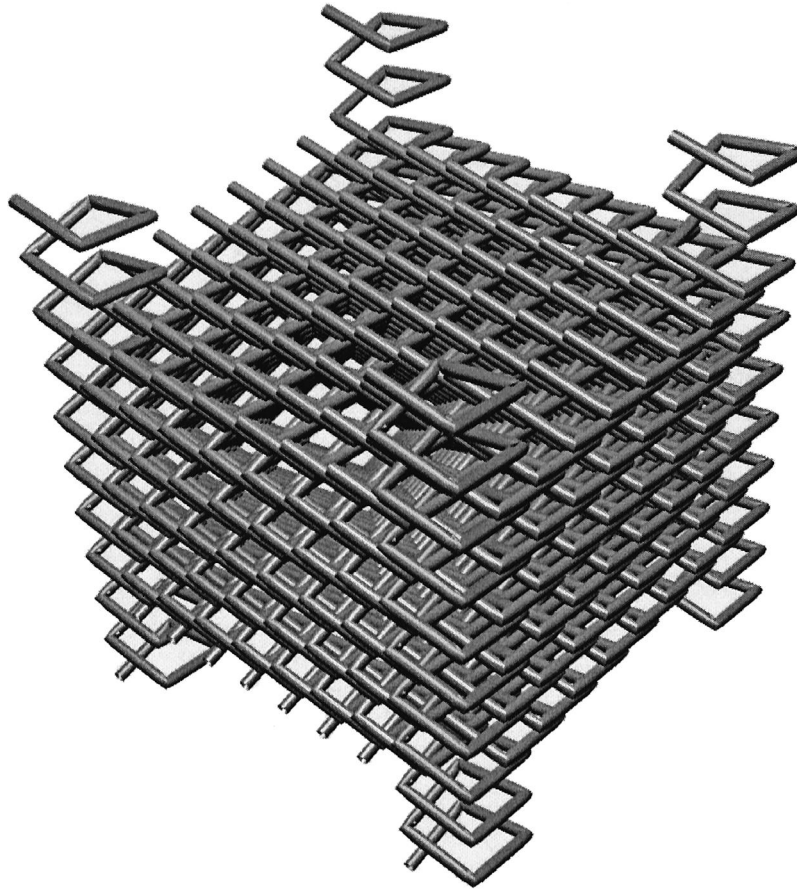


FIGURE 3 The tetragonal lattice of square spiral posts exhibits a complete 3D PBG and can be synthesized using a glancing angle deposition (GLAD) method. This chiral structure consists of slightly overlapping square spiral posts grown on a 2D substrate that is initially seeded with a square lattice of growth centers. Computer-controlled motion of the substrate leads to spiraling growth of the posts. A large and robust 3D PBG emerges between the fourth and fifth bands of the photon dispersion. The “inverse structure” consisting of air posts in a solid background exhibits an even larger 3D PBG.

IV. TWO-DIMENSIONAL PHOTONIC CRYSTAL DEVICES

An add-drop filter for a wavelength division multiplexed (WDM) communication system is depicted in [Fig. 7](#). Here, light from an optical fiber carrying many different frequencies, F_1, F_2, \dots , is inserted into a 2D PBG structure by means of a waveguide channel (missing row of holes). The frequencies F_1, F_2, \dots lie within the 2D PBG and cannot escape from the waveguide channel except in places where the periodicity of the background pores is disrupted by means of defects. For example, a hole that is larger than all the other background holes can act as a resonator which picks out a particular frequency, say F_1 , from the waveguide channel, while allowing other frequencies to propagate freely along the waveguide. Channel drop tunneling through localized defect modes with more sophisticated geometries has already been designed ([Fan et al.,](#)

[1998](#)) and tested ([Noda et al., 2000](#)). A collection of such defects could serve to pick up a band of frequencies and route them to a specified destination.

A number of prototype “active” devices based on 2D PBG heterostructures have been designed and tested. For example, 2D photonic crystal microlasers already rival the best available microcavity lasers in size and performance. For a 2D PBG the localization of light and control of spontaneous emission from excited two-level systems is incomplete. Nevertheless, low-threshold lasing from ultracompact devices has been demonstrated ([Painter et al., 1999](#); [Imada et al., 1999](#)).

[Figure 8](#) depicts two distinct types of microlasers. The first ([Fig. 8a](#)) is a “band edge microlaser” in which light emission from electrically injected electron-hole pairs in a multiple quantum well array occurs near the band edge of a 2D PBG. It has been predicted ([John and Quang, 1994](#)) that strong feedback and memory effects accompany collective

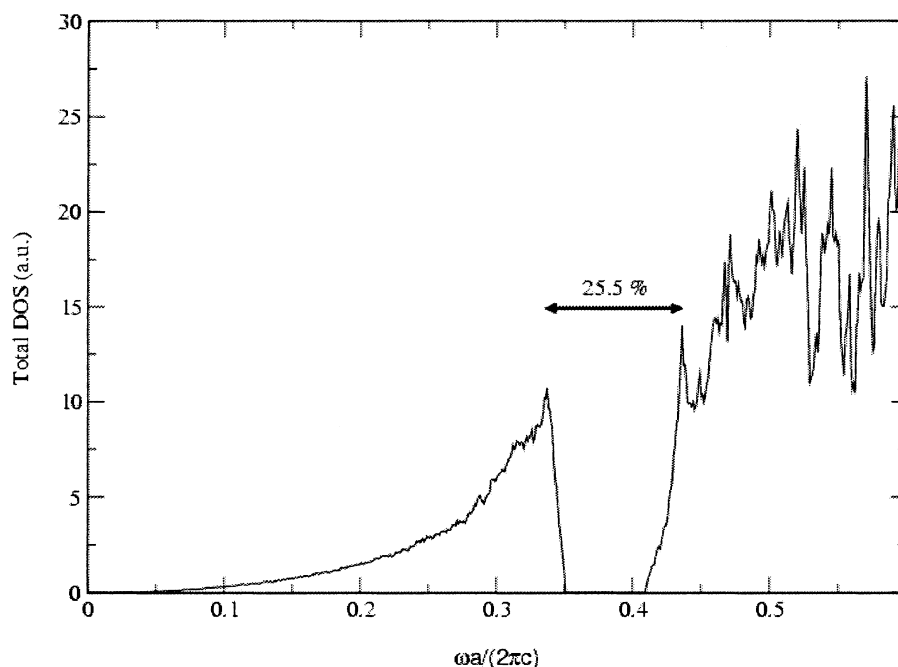


FIGURE 4 Electromagnetic density of states as a function of frequency for the tetragonal square spiral structure of Fig. 3. The total 3D density of states (DOS) vanishes over an interval of 15% of the center frequency between the fourth and fifth photon dispersion bands. In this interval, the material is “emptier than vacuum” in the sense that even the zero-point (quantum) fluctuations of the electromagnetic field have been eliminated. The size of the (larger) pseudogap over which the DOS differs significantly from that of ordinary vacuum is roughly 25% of the center frequency.

light emission near the photonic band edge. Near a true 3D photonic band edge, this would lead to lasing without a conventional optical cavity. A precursor to this effect is seen in the 2D band edge microlaser in which lasing occurs preferentially at the 2D photonic band edge (Imada *et al.*, 1999); even though the emission from the active region has a broad frequency distribution.

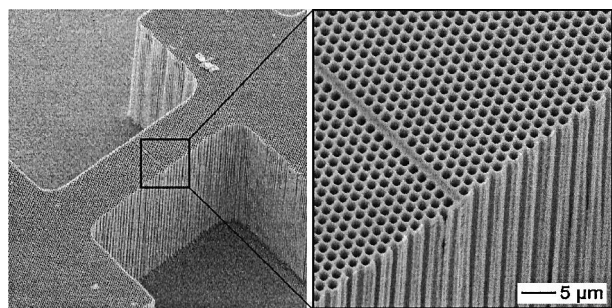


FIGURE 5 Laterally structured sample of macroporous silicon material with an incorporated defect line. The H-like structure facilitates the positioning of a fiber for the coupling in and out of light. The lattice constant is $1.5\ \mu\text{m}$ and the pore height is $100\ \mu\text{m}$. For a suitably focused beam, the structure represents a truly 2D PBG material with a band gap around $3\ \mu\text{m}$. [Courtesy of Ulrich Gösele, Max-Planck Institute for Microstructure Physics, Halle, Germany.]

The second type of microlaser (Fig. 8b) utilizes a localized state defect mode as a laser cavity (Painter *et al.*, 1999). Here, the localized electromagnetic mode is associated with a missing hole in the 2D triangular lattice. This particular structure has been proposed as the “world’s smallest microlaser” with a cavity volume of $0.03\ \mu\text{m}^3$. Spontaneous emission from electron–hole pair recombination in the multiple quantum well active region occurs preferentially into the localized state. Since the photonic crystal is two-dimensional, spontaneous emission is not exclusive to the lasing mode. This results in a finite pumping threshold before lasing occurs.

V. SYNTHESIS OF THREE-DIMENSIONAL PBG MATERIALS

While 2D PBG materials can confine light in two spatial dimensions, 3D PBG materials facilitate complete localization of light and can facilitate complete inhibition of spontaneous emission of light from atoms, molecules, and other excitations. If the transition frequency from such an atom lies within a 3D PBG, the photon that would normally be emitted and escape from the atom forms a bound state to the atom. Such feedback effects have important

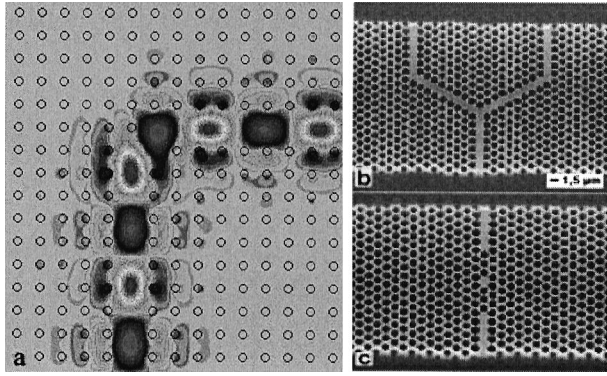


FIGURE 6 (a) Sharp bend waveguide channel in a 2D photonic crystal. The colors show the propagation of an electromagnetic mode around the bend with no reflection or scattering losses. [Courtesy of J. D. Joannopoulos, Massachusetts Institute of Technology.] (b, c) Electron micrographs of other defect structures realized in macroporous silicon with a lattice constant of 1.5 μm . The split waveguide in (b) may be used as an optical interferometer and the air holes within the waveguide channel (c) can be used as Bragg mirrors to isolate a resonator cavity within the waveguide. [Courtesy of Max-Planck Institute for Microstructure Physics, Halle, Germany.]

consequences on laser action from a collection of atoms. Indeed lasing may occur near a photonic band edge even without the need for mirrors as in a conventional laser cavity.

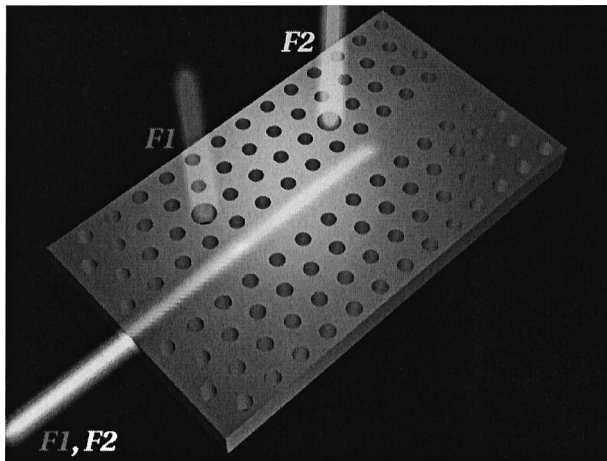
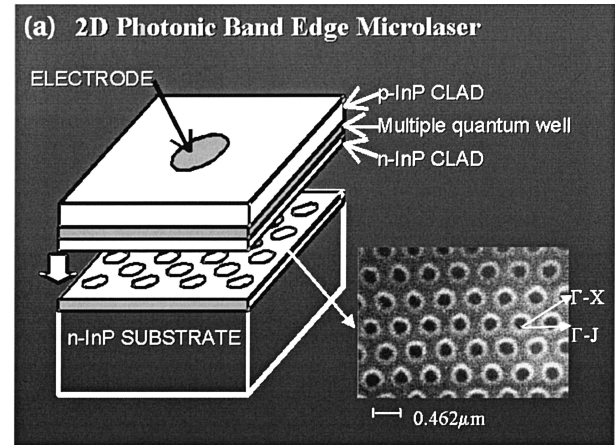


FIGURE 7 Add-drop filter for a dense wavelength division multiplexed optical communication system. Multiple streams of data carried at different frequencies F1, F2, etc. (yellow) enter the optical microchip from an external optical fiber and are carried through a waveguide channel (missing row of pores). Data streams at frequency F1 (red) and F2 (green) tunnel into localized defect modes and are routed to different destinations. The frequency of the drop filter is defined by the defect pore diameter, which is different from the pore diameter of the background photonic crystal.



(b) Defect Mode Photonic Crystal Microlaser

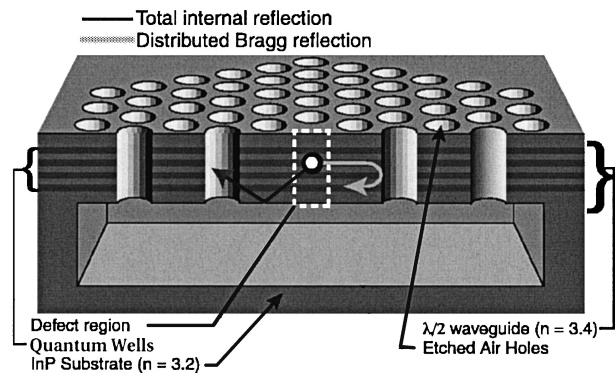


FIGURE 8 Architectures for 2D photonic crystal microlasers. (a) The band edge microlaser utilizes the unique feedback and memory effects associated with a photonic band edge and stimulated emission (arising from electron-hole recombination) from the multiple quantum well active region occurs preferentially at the band edge. There is no defect mode engineered in the 2D PBG. [Courtesy of S. Noda, Kyoto University.] (b) Defect mode microlaser requires the engineering of a localized state of light within the 2D PBG. This is created through a missing pore in the 2D photonic crystal. Stimulated emission from the multiple quantum well active region occurs preferentially into the localized mode. [Courtesy of Axel Scherer, California Institute of Technology.]

A. Layer-by-Layer Structures

The “woodpile” structure (Ho *et al.*, 1993) represents a three-dimensional PBG material that lends itself to layer-by-layer fabrication. It resembles (see Fig. 9) a criss-crossed stack of wooden logs, where in each layer the logs are in parallel orientation to each other. To fabricate one layer of the stack, a SiO_2 layer is grown on a substrate wafer, then patterned and etched. Next, the resulting trenches are filled with a high-index material such as silicon or GaAs and the surface of the wafer is polished in order to allow the next SiO_2 layer to be grown. The logs of

Integrated Optical Circuitry On a 3D PBG Microchip

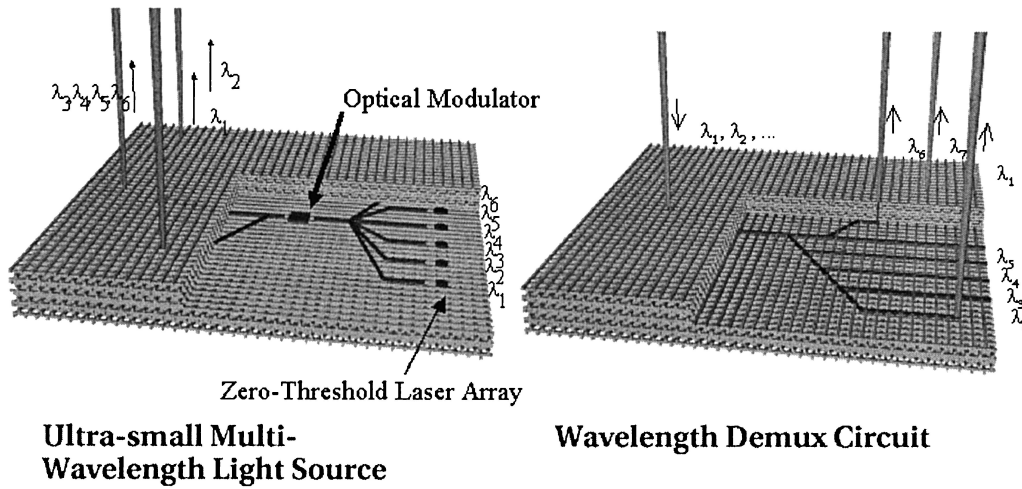


FIGURE 9 An artist's conception of a 3D PBG woodpile structure into which a microlaser array and demultiplexing (DEMUX) circuit have been integrated. [Courtesy of S. Noda, Kyoto University, Japan.]

second nearest layers are displaced midway between the logs of the original layer. As a consequence, four layers are necessary to obtain one unit cell in the stacking direction. In a final step, the SiO_2 is removed through a selective etching process leaving behind the high-index logs.

Recent work reports the successful fabrication of such a layer-by-layer PBG material made from silicon with a PBG around $1.5 \mu\text{m}$ (Lin and Fleming, 1999). However, this structure consisted of only five layers in the stacking direction. Instead of depositing successively more layers, wafer-bonding technology may be applied to single-layer substrates (Noda *et al.*, 1999). Bonding together two single-layer substrates and subsequent removal of the upper substrate results in a double-layer structure. The ensuing technique is multiplicative but tedious and expensive. This type of complex microlithography has led to the successful fabrication of an eight-layer structure (two unit cells) in the stacking direction.

B. Self-Organizing Structures

In three dimensions a number of large-scale self-assembling periodic structures exist. These include colloidal systems (Tarhan and Watson, 1996) and artificial opals (Vlasov *et al.*, 1997). Unfortunately, these readily available materials do not satisfy the necessary criteria of high index contrast and correct network topology to produce a complete PBG. Theoretical studies, however, indicate the possibility of a complete PBG in closely related structures. Face-centered cubic lattices consisting

of low-dielectric inclusions in a connected high-dielectric network (henceforth called inverse structures) can exhibit small PBGs. The recipe for producing inverse structures from artificial opals is to infiltrate them with a high-dielectric material such as silicon (Blanco *et al.*, 2000) and subsequently to etch out the SiO_2 spheres, leaving behind a connected network of high-dielectric material with filling ratios of about 26% by volume. Such a “Swiss cheese structures” with air voids in a silicon backbone is displayed in Fig. 10. This large-scale inverse opal PBG material exhibits a complete 5% PBG relative to its center frequency at $1.5 \mu\text{m}$ (Blanco *et al.*, 2000). The etching out of the SiO_2 provides the necessary dielectric contrast for the emergence of a complete 3D PBG. Moreover, the presence of air voids rather than solid SiO_2 may allow the injection of atomic vapors with which quantum optical experiments can be carried out. It also facilitates the infiltration of optically anisotropic materials such as nematic liquid crystals for the realization of electrooptic tuning effects and enables the infiltration of active materials such as conjugated polymers and dyes for laser applications.

VI. QUANTUM AND NONLINEAR OPTICS IN THREE-DIMENSIONAL PBG MATERIALS

Photonic bandgap materials represent a fundamentally new paradigm for low-threshold nonlinear optical

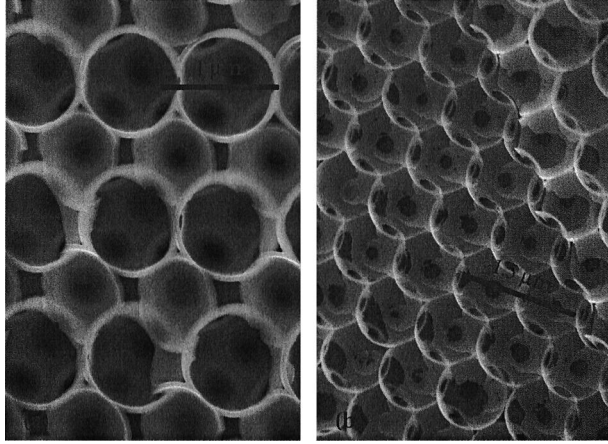


FIGURE 10 Scanning electron microscope (SEM) pictures of cross section along (a) the cubic (110) direction and (b) the cubic (111) direction of a silicon inverse opal with a complete 5% PBG around $1.5 \mu\text{m}$. The structure has been obtained through the infiltration of an artificial opal with silicon (light-shaded regions) and subsequent removal of the SiO_2 spheres of the opal. The air sphere diameter is 870 nm . Clearly visible is the incomplete infiltration (diamond-shaped voids between spheres) and the effect of sintering the artificial opal prior to infiltration (small holes connecting adjacent spheres).

phenomena. The PBG affects the light–matter interaction in a fundamental way. For instance, if the transition frequency of an excited atom embedded in such a material lies within the complete PBG, spontaneous emission may be completely suppressed and a bound photon–atom state is formed instead. Based on these principles numerous applications for active devices have been suggested. Two illustrative examples are given below.

A. Low-Threshold Resonant Nonlinear Optics

The ability to achieve ultrafast nonlinear optical response in a nonabsorbing material is crucial in applications such as all-optical switches and other nonlinear devices for integrated optical circuits. In a conventional Fabry–Perot

response requires relatively high intensity laser fields. The situation may be dramatically different in the context of a three-dimensional PBG material, where the inhibition of spontaneous emission from atoms and molecules is essentially complete.

Consider the injection of a classical monochromatic electromagnetic wave of frequency ω into a 3D PBG material by means of a single-mode waveguide channel. Suppose this waveguide channel contains a small active region of optically excitable two-level systems (confined electron–hole pair excitations or atoms) with a radiative transition at frequency ω_0 such that the detuning

$\Delta = \omega - \omega_0 \ll \omega_0$. For weak fields, the response of the two-level atom is that of a simple harmonic oscillator since the atom spends the majority of time in its ground state. Whenever the external field excites the atom, it quickly returns to its ground state due to the rapid rate of spontaneous emission. As the external field intensity is increased, the upper-state population increases and eventually saturates. This is associated with nonlinear response. The appropriate mechanical analogy for the quantum two-level system is no longer a classical harmonic oscillator, but rather a simple pendulum (atomic Bloch vector) whose coordinates are given by the different components of the 2×2 atomic density matrix. Small-angle oscillation of the atomic Bloch vector describes linear response, whereas large-angle oscillations probe the nonlinear susceptibility of the atom. The threshold external field required to probe nonlinear response is called the line center saturation field strength E_s^0 , which is related to the rate of spontaneous emission $1/T_1$ by the formula (Boyd, 1992)

$$|E_s^0|^2 = \frac{\hbar^2}{4\mu_{ba}T_1T_2}. \quad (3)$$

Here, $1/T_2$ is the rate of dipolar dephasing and μ_{ba} is the electric dipole transition matrix element for the atom. In ordinary vacuum (as opposed to our 3D PBG material), the textbook formula (Boyd, 1992) for the nonlinear susceptibility is

$$\chi = \chi_0 \frac{\Delta T_2 - i}{1 + \Delta^2 T_2^2 + |E|^2/|E_s^0|^2}. \quad (4)$$

Here, χ_0 is a frequency independent constant and $\Delta = \omega - \omega_0$ is the detuning of the external field of amplitude E and frequency ω . Clearly, the susceptibility χ is independent of E for $E \ll E_s^0$, whereas it has a highly nonlinear dependence on E for $E > E_s^0$. This widely accepted picture of nonlinear optical response in ordinary materials is no longer applicable in a three-dimensional PBG material in which two-level atoms have a radiative transition at frequency ω_0 that lies within the PBG. Inside a PBG, the rate $1/T_1$ and accordingly the threshold intensity $|E_s^0|^2$ formally vanish. While this is suggestive of low-threshold nonlinear optical response, it is in fact an indication that the entire derivation of the conventional susceptibility, Eq. (4), must be carefully re-examined in the context of the PBG. In particular, Eq. (4) is based on retaining only the leading order photon–atom interaction, namely spontaneous emission. In the PBG, this leading process is almost absent. Therefore, it is necessary to consider the next process, namely resonance dipole–dipole interaction (RDDI) between a pair of two-level atoms. A straightforward derivation (John and Quang, 1996) of the nonlinear susceptibility in this context leads

to some remarkable effects. For example, it is possible (at relatively low field intensities) to completely saturate the imaginary (absorptive) part of χ while retaining a large real part for χ . This arises in a situation where a collection of atoms inside a PBG interact randomly by RDDI and an external field enters the material through a small number of defect (waveguide) modes. Due to the inhibition of spontaneous emission, the single-atom absorptive transition is saturated at nearly the one-photon level. However, due to the random nature of RDDI, the induced atomic dipoles are randomly oriented and remain highly susceptible to alignment (macroscopic polarization) by the external field.

B. Collective Switching and Transistor Effects

A second illustration of novel radiation–matter interaction in a PBG material arises if a collection of two-level atoms is selected such that their radiative transition lies very close to a photonic band edge. Near the photonic band edge, the electromagnetic density of states varies rapidly with frequency and the spontaneous emission processes cannot be adequately described using Fermi’s golden rule. We refer to the case in which the density of states exhibits singularities or abrupt variation as a “colored vacuum.” In the colored vacuum of a photonic band edge it is possible to achieve population inversion for a collection of two-level atoms driven by a weak coherent pump laser field. As the number of atoms in a cubic wavelength increases, this switching from a passive medium (most atoms in ground state) to an active medium (population inverted) occurs sharply as the external pump laser intensity is increased. In the region of this collective jump in atomic population, there is a large differential optical gain. If the pump laser intensity is chosen slightly below threshold for population inversion, a second control laser field can be introduced to act as a gate which determines whether the pump laser field is either attenuated or amplified by the medium (see Fig. 11). Since all of these processes involve coherent radiation–atom interactions, this system may form the basis of a low-threshold optical switch or all-optical transistor (John and Quang, 1997).

When a coherent laser field with average incident energy density W and frequency ω interacts with a collection of N two-level atoms in ordinary vacuum, the steady state behavior of the system is governed by the well-known Einstein rate equations. These equations implicitly make use of the smooth nature of the vacuum density of states $N(\omega) = \omega^2/(\pi^2 c^3)$ in the vicinity of the atomic transition frequency $\omega \approx \omega_0$. In steady state equilibrium, the ratio of the number of excited atoms N_2 to the total number of atoms is given by (Laudon, 1983)

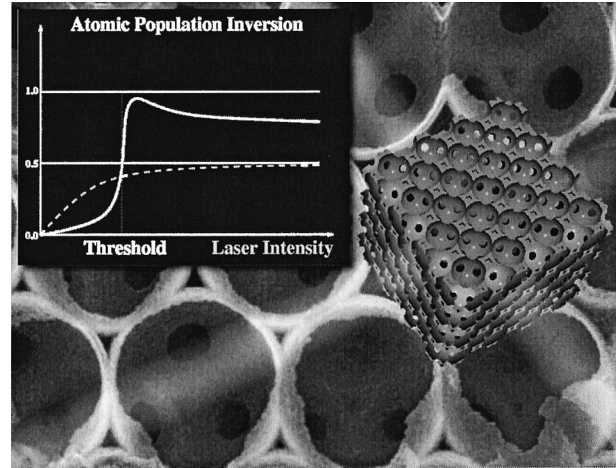


FIGURE 11 Microphotonic all-optical transistor may consist of an active region buried in the intersection of two waveguide channels in a 3D PBG material. The two-level systems (“atoms”) in the active region are coherently pumped and controlled by laser beams passing through the waveguides. In addition, the 3D PBG material is chosen to exhibit an abrupt variation in the photon density of states near the transition frequency of the atoms. This leads to atomic “population inversion” through coherent pumping, an effect which is forbidden in ordinary vacuum. The inversion threshold is characterized by a narrow region of large differential optical gain (solid curve in the inset). A second, “control laser” allows the device to pass through this threshold region leading to strong amplification of the output signal. In ordinary vacuum, population inversion is unattainable (dashed curve in the inset).

$$\frac{N_2}{N} = \frac{W}{\hbar\omega N(\omega) + 2W}. \quad (5)$$

Clearly, as W increases, the maximum value of N_2/N is less than $1/2$. In other words, it is not possible to invert a collection of two-level atoms with a coherent laser field. From a more quantum mechanical point of view, the external laser field may be regarded as consisting of a large collection of n photons. The atom–radiation field interaction Hamiltonian H_{int} breaks the degeneracy between a state consisting of a given atom in its excited state and $n - 1$ photons in the radiation field (which we denote by the ket $|2, n - 1\rangle$) and a state consisting of the given atom in its ground state and n photons in the radiation field (which we denote by the ket $|1, n\rangle$). The matrix element $\langle 2, n - 1 | H_{\text{int}} | 1, n \rangle$ is nonzero and the true eigenstates of the system are called dressed atomic states. These “dressed states” are linear combinations of the “bare” kets listed above. Accordingly, the eigenenergies of the dressed states are shifted from their bare values by an amount $\Delta \sim \mu_{ba}|E|/\hbar$ (Rabi frequency), where μ_{ba} is the atomic dipole matrix element and $|E|$ is the laser field amplitude. This leads to the well-known Mollow fluorescence spectrum in ordinary vacuum (Laudon, 1983). Rather than a

single peak in the atomic spectrum centered at the bare atomic transition frequency ω_0 , the fluorescence spectrum exhibits three peaks centered at ω_0 , $\omega_0 + \Delta$, and $\omega_0 - \Delta$.

The Einstein rate equation picture of the steady state atomic inversion, Eq. (5), relies on the fact that the vacuum density of electromagnetic modes $N(\omega)$ is relatively smooth on the scale of Δ . That is to say, the Einstein picture assumes that the rate of spontaneous emission in the Mollow sidebands at $\omega_0 + \Delta$ and $\omega_0 - \Delta$ is roughly the same. In ordinary vacuum [$N(\omega) = \omega^2/(\pi^2 c^3)$] this assumption is easily satisfied. Moreover, in ordinary vacuum, very high intensity fields are required to observe any Mollow splitting whatsoever. The situation is dramatically different in a PBG material, where the density of states itself exhibits rapid variation from frequency $\omega_0 - \Delta$ to $\omega_0 + \Delta$. Another striking property of the photonic band edge is that atomic line splitting may be achieved with very low intensity fields. In particular, vacuum Rabi splitting of an atomic transition placed directly at the band edge has been predicted (John and Wang, 1990, 1991; John and Quang, 1994). In other words, at a band edge, significant splitting can be expected in an atomic line even in the presence of a single photon! This leads to a dramatic modification of the Einstein picture.

In a weak applied laser field, atoms with a bare transition frequency ω_0 which coincides with a photonic band edge will exhibit a pair of dressed states that straddle the band edge (John and Wang, 1990, 1991; John and Quang, 1994). These two spectral sidebands will experience vastly different rates of spontaneous emission. The spectral component that is pulled into the gap will exhibit negligible decay. This component corresponds to a photon-atom bound state (John and Wang, 1990). The spectral component that is pushed away from the gap can exhibit emission into the allowed modes of the photonic band. Population inversion for a collection of such atoms can readily be achieved by an external laser field due to trapping of excitation energy in the photon-atom bound-state component (John and Quang, 1997).

The resulting collective switch in atomic population as a function of applied field intensity is depicted schematically within the inset of Fig. 11. Details of this result may be found in John and Quang (1997). This transition becomes increasingly sharp as the number of atoms increases and defines a region of large differential optical gain. It is expected that the collective switch can be achieved for a very low applied field and that the switching effect is robust with respect to variety of dephasing effects due to lattice vibrations in the host PBG material.

VII. TUNABLE PBG MATERIALS

For many applications it is advantageous to obtain some degree of tunability of the photonic band structure through electrooptic effects. This may be useful in routing of signals through an optical communication network and to provide flexibility to reconfigure the routing as conditions in the network change. One of the great advantages of PBG materials is that by volume, most of the material consists of pores. These pores can be infiltrated with other electrooptically active materials, which enables reconfiguration of the underlying photonic band structure either globally or locally. Tunability may be obtained by controlling one or several forms of optical anisotropy of the constituent materials. For example the science of liquid crystals has spawned an entire industry related to these electrooptic effects. Inverse opal structures provide a highly efficient scattering system as illustrated by the complete PBG of silicon inverse opals (Figs. 10 and 11). The nearly 75% empty volume of this structure is ideally suited for infiltration by a low-refractive index liquid crystal with strong optical anisotropy, making it efficacious for electrooptic tuning effects. In particular, a change in the orientation of the nematic director field with respect to the inverse opal backbone by an external electric field can completely open or close the full, three-dimensional PBG (Busch and John, 1999). This clearly demonstrates an electrooptic shutter effect to the PBG which may be realized by external means. The resulting tunability of spontaneous emission, waveguiding effects, and light localization may considerably enhance the technological value of a composite liquid crystal PBG material over and above that of either a bulk liquid crystal or a conventional PBG material by itself. A tunable optical microchip which routes light from a set of optical fibers is shown in Fig. 12.

VIII. OUTLOOK

The synergistic interplay among advanced material science, theoretical analysis, and modern spectroscopy has been and continues to be the driving force for the field of PBG materials. Recent advances in microstructuring technology have allowed the realization and controlled engineering of three-dimensional PBG structures at the near-IR as well as the visible frequency spectrum of electromagnetic radiation. In a parallel development, the theoretical description of PBG materials has matured to the point where it provides a reliable interpretative as well as predictive tool for both material synthesis and spectroscopic analysis of these novel semiconductors for light. The current state of PBG research suggests that this field

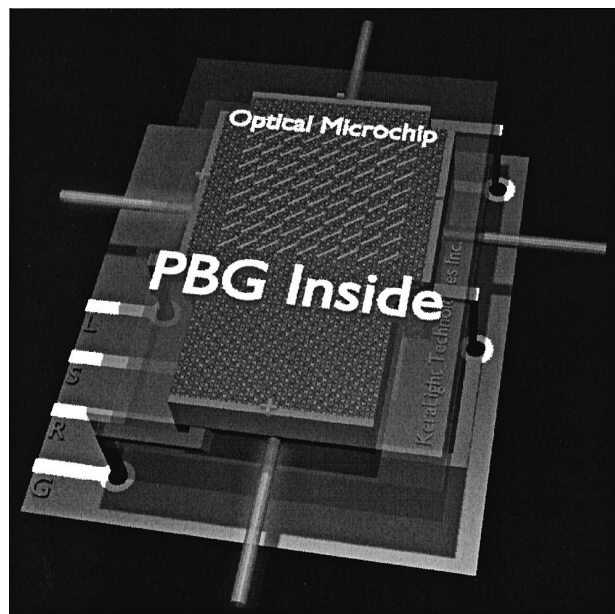


FIGURE 12 Artist's depiction of an electroactively tunable PBG routing device. Here the PBG material has been infiltrated with an optically anisotropic material (such as a liquid crystal) exhibiting a large electrooptic response. When a voltage is applied to the electrooptically tunable PBG, the polarization state (yellow arrows) can be rotated, leading to corresponding shifts in the photonic band structure. This allows light from an optical fiber to be routed into one of several output fibers.

is at a stage comparable to the early years of semiconductor technology shortly before the invention of the solid state electronic transistor by W. Shockley, J. Bardeen, and W. H. Brattain. If this analogy continues to hold, one may find the PBG materials at the heart of a 21st century revolution in optical information technology similar to the revolution in electronics we witnessed over the latter half of the 20th century.

SEE ALSO THE FOLLOWING ARTICLES

ELECTROMAGNETISM • LIGHT EMITTING DIODES (LEDs)
• MULTIPHOTON SPECTROSCOPY • OPTICAL AMPLIFIERS

(SEMICONDUCTOR) • OPTICAL FIBERS, FABRICATION AND APPLICATIONS • PHOTOCHEMISTRY BY VUV PHOTONS
• RADIOMETRY AND PHOTOMETRY • SEMICONDUCTOR ALLOYS

BIBLIOGRAPHY

- Anderson, P. W. (1985). *Phil. Mag.* **B 52**, 505.
 Blanco, A., *et al.* (2000). *Nature* **405**, 437.
 Boyd, R. (1992). "Nonlinear Optics," Academic Press.
 Busch, K., and John, S. (1999). *Phys. Rev. Lett.* **83**, 967.
 Bykov, V. P. (1975). *Sov. J. Quant. Electron.* **4**, 861.
 Fan, S., Villeneuve, P. R., Joannopoulos, J. D., and Haus, H. A. (1998). *Phys. Rev. Lett.* **80**, 960.
 Grüning, U., Lehman, V., Ottow, S., and Busch, K. (1996). *Appl. Phys. Lett.* **68**, 747.
 Ho, K.-M., Chan, C. T., Soukoulis, C. M., Biswas, R., and Sigalas, M. (1993). *State. Commun.* **89**, 413.
 Imada, M., Noda, S., Chutinan, A., Tokuda, T., Murata, M., and Sasaki, G. (1999). *Appl. Phys. Lett.* **75**, 316.
 John, S. (1984). *Phys. Rev. Lett.* **53**, 2169.
 John, S. (1987). *Phys. Rev. Lett.* **58**, 2486.
 John, S., and Quang, T. (1994). *Phys. Rev. A* **50**, 1764.
 John, S., and Quang, T. (1996). *Phys. Rev. Lett.* **76**, 2484.
 John, S., and Quang, T. (1997). *Phys. Rev. Lett.* **78**, 1888.
 John, S., and Wang, J. (1990). *Phys. Rev. Lett.* **64**, 2418.
 John, S., and Wang, J. (1991). *Phys. Rev. B* **43**, 12772.
 Lin, S.-Y., and Fleming, J. G. (1999). *J. Lightwave Technol.* **17**, 1944.
 Loudon, R. (1983). "The Quantum Theory of Light," Clarendon, Oxford.
 Noda, S., Yamamoto, N., Imada, M., Kobayashi, H., and Okano, M. (1999). *J. Lightwave Technol.* **17**, 1948.
 Noda, S., Chutinan, A., and Imada, M. (2000). *Nature* **407**, 608.
 Painter, O., *et al.* (1999). *Science* **284**, 1819.
 Quang, T., Woldeyohannes, M., John, S., and Agarwal, G. S. (1997). *Phys. Rev. Lett.* **79**, 5238.
 Robbie, K., and Brett, M. J. (1996). *Nature* **384**, 616.
 Robbie, K., Sit, J. C., and Brett, M. (1998). *J. Vac. Sci. Technol. B* **16**, 1115.
 Schilling, J., Birner, A., Müller, F., Wehrspohn, R. B., Hillebrand, R., Gösele, U., Busch, K., John, S., Leonard, S. W., and van Driel, H. M. (2001). *Opt. Materials*, to appear.
 Toader, O., and John, S. (2001). *Science* **292**, 1133.
 Vlasov, *et al.*, and Yu. A. (1997). *Phys. Rev. B* **55**, R13357.
 Vos, W. L., Sprik, R., van Blaaderen, A., Imhof, A., Lagendijk, A., and Weydam, G. H. (1996). *Phys. Rev. B* **53**, 16231.
 Yablonovitch, E. (1987). *Phys. Rev. Lett.* **58**, 2059.
 Yablonovitch, E., Gmitter, T. J., and Leung, K. M. (1991). *Phys. Rev. Lett.* **67**, 2295.