# Image
# Formation

<div style="text-align: right;">2</div>

## 2.1 IMAGES

Image formation occurs when a *sensor* registers *radiation* that has interacted with *physical objects*. Section 2.2 deals with mathematical models of images and image formation. Section 2.3 describes several specific image formation technologies.

The mathematical model of imaging has several different components.

1. An *image function* is the fundamental abstraction of an image.
2. A *geometrical model* describes how three dimensions are projected into two.
3. A *radiometrical model* shows how the imaging geometry, light sources, and reflectance properties of objects affect the light measurement at the sensor.
4. A *spatial* frequency model describes how spatial variations of the image may be characterized in a transform domain.
5. A *color model* describes how different spectral measurements are related to image colors.
6. A *digitizing model* describes the process of obtaining discrete samples.

This material forms the basis of much image-processing work and is developed in much more detail elsewhere, e.g., [Rosenfeld and Kak 1976; Pratt 1978]. Our goals are not those of image processing, so we limit our discussion to a summary of the essentials.

The wide range of possible sources of samples and the resulting different implications for later processing motivate our overview of specific imaging techniques. Our goal is not to provide an exhaustive catalog, but rather to give an idea of the range of techniques available. Very different analysis techniques may be needed depending on how the image was formed. Two examples illustrate this

point. If the image is formed by reflected light intensity, as in a photograph, the image records both light from primary light sources and (more usually) the light reflected off physical surfaces. We show in Chapter 3 that in certain cases we can use these kinds of images together with knowledge about physics to derive the orientation of the surfaces. If, on the other hand, the image is a computed tomogram of the human body (discussed in Section 2.3.4), the image represents tissue density of internal organs. Here orientation calculations are irrelevant, but general segmentation techniques of Chapters 4 and 5 (the agglomeration of neighboring samples of similar density into units representing organs) are appropriate.

## 2.2 IMAGE MODEL

Sophisticated image models of a statistical flavor are useful in image processing [Jain 1981]. Here we are concerned with more geometrical considerations.

### 2.2.1 Image Functions

An *image function* is a mathematical representation of an image. Generally, an image function is a vector-valued function of a small number of arguments. A special case of the image function is the *digital (discrete) image function*, where the arguments to and value of the function are all integers. Different image functions may be used to represent the same image, depending on which of its characteristics are important. For instance, a camera produces an image on black-and-white film which is usually thought of as a real-valued function (whose value could be the density of the photographic negative) of two real-valued arguments, one for each of two spatial dimensions. However, at a very small scale (the order of the film grain) the negative basically has only two densities, "opaque" and "transparent."

Most images are presented by functions of two *spatial* variables $f(\mathbf{x}) = f(x, y)$, where $f(x, y)$ is the brightness of the gray level of the image at a spatial coordinate $(x, y)$. A multispectral image $\mathbf{f}$ is a vector-valued function with components $(f_1 \ldots f_n)$. One special multispectral image is a color image in which, for example, the components measure the brightness values of each of three wavelengths, that is,

$$f(\mathbf{x}) = \left\{ f_{\text{red}}(\mathbf{x}), f_{\text{blue}}(\mathbf{x}), f_{\text{green}}(\mathbf{x}) \right\}$$

Time-varying images $f(\mathbf{x}, t)$ have an added temporal argument. For special three-dimensional images, $\mathbf{x} = (x, y, z)$. Usually, both the domain and range of $f$ are bounded.

An important part of the formation process is the conversion of the image representation from a continuous function to a discrete function; we need some way of describing the images as samples at discrete points. The mathematical tool we shall use is the *delta function.*

Formally, the delta function may be defined by

*Ch. 2   Image Formation*

$$\delta(x) = \begin{cases} 0 & \text{when } x \neq 0 \\ \infty & \text{when } x = 0 \end{cases} \tag{2.1}$$

$$\int_{-\infty}^{\infty} \delta(x)\, dx = 1$$

If some care is exercised, the delta function may be interpreted as the limit of a set of functions:

$$\delta(x) = \lim_{n \to \infty} \delta_n(x)$$

wh∩re

$$\delta_n(x) = \begin{cases} n & \text{if } |x| < \dfrac{1}{2n} \\ 0 & \text{otherwise} \end{cases} \tag{2.2}$$

A useful property of the delta function is the *sifting property:*

$$\int_{-\infty}^{\infty} f(x)\, \delta(x - a)\, dx = f(a) \tag{2.3}$$

A continuous image may be multipled by a two-dimensional "comb," or array of delta functions, to extract a finite number of discrete *samples* (one for each delta function). This mathematical model of the sampling process will be useful later.

### 2.2.2 Imaging Geometry

*Monocular Imaging*

    *Point projection* is the fundamental model for the transformation wrought by our eye, by cameras, or by numerous other imaging devices. To a first-order approximation, these devices act like a pinhole camera in that the image results from projecting scene points through a single point onto an *image plane* (see Fig. 2.1). In Fig. 2.1, the image plane is behind the point of projection, and the image is reversed. However, it is more intuitive to recompose the geometry so that the point of projection corresponds to a *viewpoint* behind the image plane, and the image occurs right side up (Fig. 2.2). The mathematics is the same, but now the viewpoint is $+f$ on the $z$ axis, with $z = 0$ plane being the image plane upon which the image is projected. ($f$ is sometimes called the *focal length* in this context. The use of $f$ in this section should not be confused with the use of $f$ for image function.) As the imaged object approaches the viewpoint, its projection gets bigger (try moving your hand toward your eye). To specify how its imaged size changes, one needs only the geometry of similar triangles. In Fig. 2.2b $y'$, the projected height of the object, is related to its real height $y$, its position $z$, and the focal length $f$ by

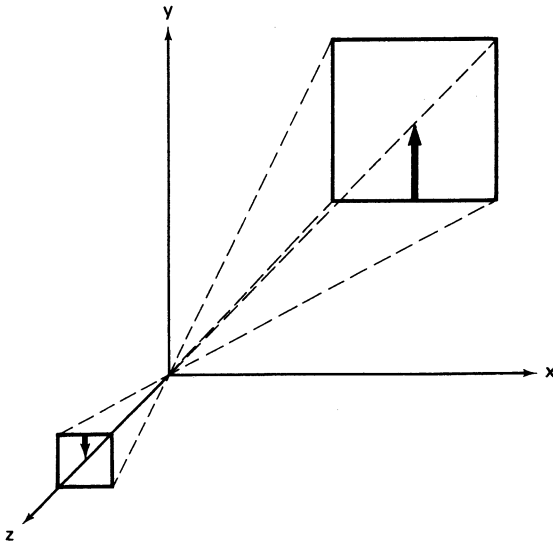$$\frac{y}{f - z} = \frac{y'}{f} \tag{2.4}$$

**Fig. 2.1** A geometric camera model.

The case for $x'$ is treated similarly:

$$\frac{x}{f - z} = \frac{x'}{f} \tag{2.5}$$

The projected image has $\dot{z} = 0$ everywhere. However, projecting away the $z$ component is best considered a separate transformation; the projective transform is usually thought to distort the $z$ component just as it does the $x$ and $y$. *Perspective distortion* thus maps $(x, y, z)$ to

$$(x', y', z') = \left[ \frac{fx}{f - z}, \; \frac{fy}{f - z}, \; \frac{fz}{f - z} \right] \tag{2.6}$$

The perspective transformation yields *orthographic projection* as a special case when the viewpoint is the *point at infinity* in the $z$ direction. Then all objects are projected onto the viewing plane with no distortion of their $x$ and $y$ coordinates.

The perspective distortion yields a three-dimensional object that has been "pushed out of shape"; it is more shrunken the farther it is from the viewpoint. The $z$ component is not available directly from a two-dimensional image, being identically equal to zero. In our model, however, the distorted $z$ component has information about the distance of imaged points from the viewpoint. When this distorted object is projected orthographically onto the image plane, the result is a perspective picture. Thus, to achieve the effect of railroad tracks appearing to come together in the distance, the perspective distortion transforms the tracks so that they *do* come together (at a point at infinity)! The simple orthographic projection that projects away the $z$ component unsurprisingly preserves this distortion. Several properties of the perspective transform are of interest and are investigated further in Appendix 1.

*Binocular Imaging*

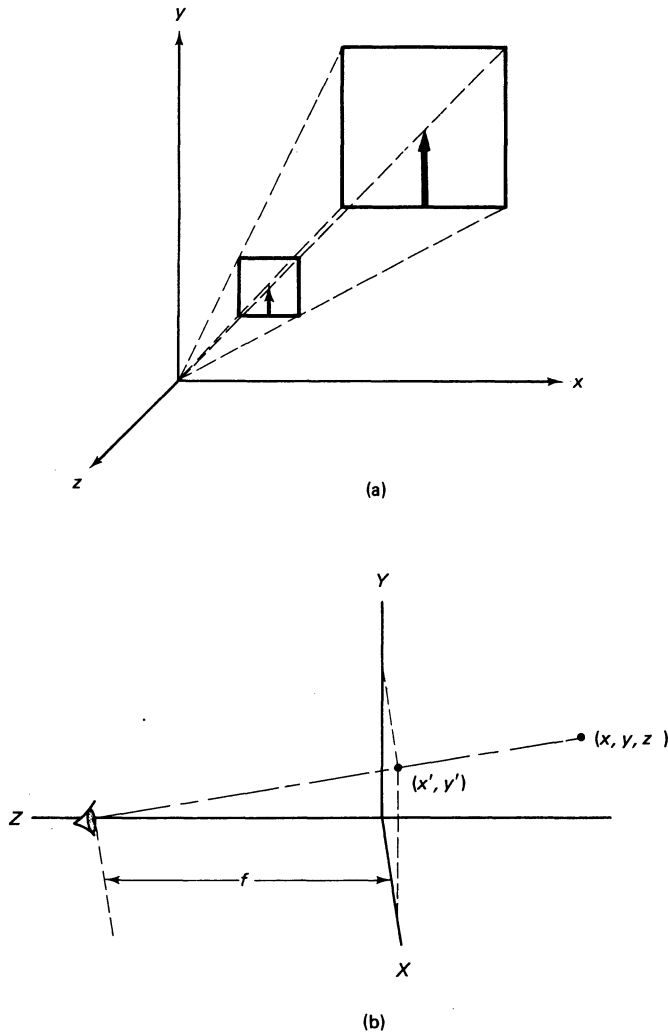Basic binocular imaging geometry is shown in Fig. 2.3a. For simplicity, we

(a)



(b)

Fig. 2.2 (a) Camera model equivalent to that of Fig. 2.1; (b) definition of terms.

use a system with two viewpoints. In this model the eyes do not *converge*; they are aimed in parallel at the point at infinity in the $-z$ direction. The depth information about a point is then encoded only by its different positions (*disparity*) in the two image planes.

With the stereo arrangement of Fig. 2.3,

$$x' = \frac{(x - d)f}{f - z}$$

$$x'' = \frac{(x + d)f}{f - z}$$

where $(x', y')$ and $(x'', y'')$ are the retinal coordinates for the world point imaged
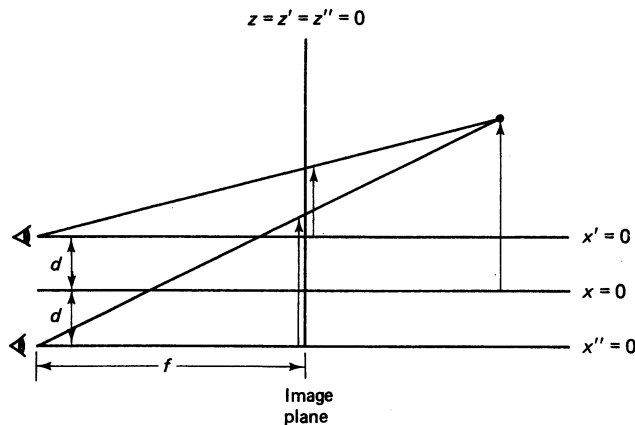
$$z = z' = z'' = 0$$

$$x' = 0$$
$$x = 0$$
$$x'' = 0$$

Image plane

Fig. 2.3 A nonconvergent binocular imaging system.

through each eye. The *baseline* of the binocular system is $2d$. Thus

$$(f - z)x' = (x - d)f \tag{2.7}$$

$$(f - z)x'' = (x + d)f \tag{2.8}$$

Subtracting (2.7) from (2.8) gives

$$(f - z)(x'' - x') = 2df$$

or

$$z = f - \frac{2df}{x'' - x'} \tag{2.9}$$

Thus if points can be matched to determine the disparity $(x'' - x')$ and the baseline and focal length are known, the $z$ coordinate is simple to calculate.

If the system can converge its directions of view to a finite distance, convergence angle may also be used to compute depth. The hardest part of extracting depth information from stereo is the *matching* of points for disparity calculations. "Light striping" is a way to maintain geometric simplicity and also simplify matching (Section 2.3.2).

### 2.2.3 Reflectance

*Terminology*

A basic aspect of the imaging process is the physics of the reflectance of objects, which determines how their "brightness" in an image depends on their inherent characteristics and the geometry of the imaging situation. A clear presentation of the mathematics of reflectance is given in [Horn and Sjoberg 1978; Horn 1977]. Light *energy flux* $\Phi$ is measured in watts; "brightness" is measured with respect to area and solid angle. The *radiant intensity* $I$ of a source is the exitant flux per unit solid angle:

$$I = \frac{d\Phi}{d\omega} \quad \text{watts/steradian} \tag{2.10}$$

*Ch. 2 Image Formation*

Here $d\omega$ is an incremental solid angle. The solid angle of a small area $dA$ measured perpendicular to a radius $r$ is given by

$$d\omega = \frac{dA}{r^2} \qquad (2.11)$$

in units of steradians. (The total solid angle of a sphere is $4\pi$.)

The *irradiance* is flux incident on a surface element $dA$:

$$E = \frac{d\Phi}{dA} \qquad \text{watts/meter}^2 \qquad (2.12)$$

and the flux exitant from the surface is defined in terms of the *radiance L*, which is the flux emitted per unit foreshortened surface area per unit solid angle:

$$L = \frac{d^2\Phi}{dA \, \cos\theta \, d\omega} \qquad \text{watts/(meter}^2 \text{ steradian)} \qquad (2.13)$$

where $\theta$ is the angle between the surface normal and the direction of emission.

*Image irradiance f* is the "brightness" of the image at a point, and is proportional to scene radiance. A "gray-level" is a quantized measurement of image irradiance. Image irradiance depends on the reflective properties of the imaged surfaces as well as on the illumination characteristics. How a surface reflects light depends on its micro-structure and physical properties. Surfaces may be *matte* (dull, flat), *specular* (mirrorlike), or have more complicated reflectivity characteristics (Section 3.5.1). The *reflectance r* of a surface is given quite generally by its Bidirectional Reflectance Distribution Function (BRDF) [Nicodemus et al. 1977]. The BRDF is the ratio of reflected radiance in the direction towards the viewer to the irradiance in the direction towards a small area of the source.

### Effects of Geometry on an Imaging System

Let us now analyze a simple image-forming system shown in Fig. 2.4 with the objective of showing how the gray levels are related to the radiance of imaged objects. Following [Horn and Sjoberg 1978], assume that the imaging device is properly focused; rays originating in the infinitesimal area $dA_o$ on the object's surface are projected into some area $dA_p$ in the image plane and no rays from other portions of the object's surface reach this area of the image. The system is assumed to be an ideal one, obeying the laws of simple geometrical optics.

The energy flux/unit area that impinges on the sensor is defined to be $E_p$. To show how $E_p$ is related to the scene radiance $L$, first consider the flux arriving at the lens from a small surface area $dA_o$. From (2.13) this is given as

$$d\Phi = dA_o \int L \cos\theta \, d\omega \qquad (2.14)$$

This flux is assumed to arrive at an area $dA_p$ in the imaging plane. Hence the irradiance is given by [using Eq. (2.12)]

$$E_p = \frac{d\Phi}{dA_p} \qquad (2.15)$$

Now relate $dA_o$ to $dA_p$ by equating the respective solid angles as seen from the lens; that is [making use of Eq. (2.12)],
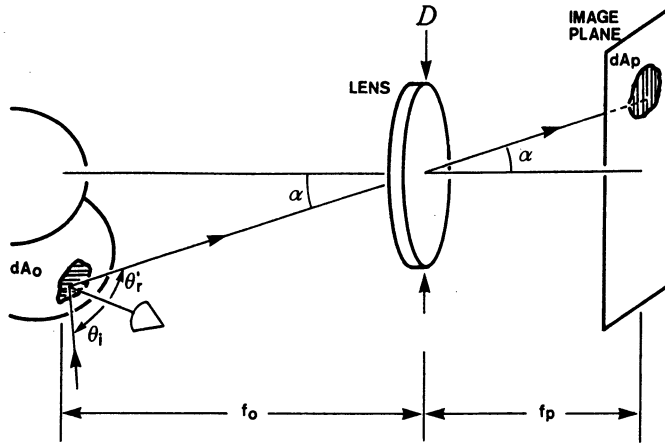
Fig. 2.4 Geometry of an image forming system.

$$dA_o \frac{\cos\theta}{f_o^2} = dA_p \frac{\cos\alpha}{f_p^2} \qquad (2.16)$$

Substituting Eqs. (2.16) and (2.14) into (2.15) gives

$$E = \cos\alpha \left(\frac{f_o}{f_p}\right)^2 \int L\, d\omega \qquad (2.17)$$

The integral is over the solid angle seen by the lens. In most instances we can assume that $L$ is constant over this angle and hence can be removed from the integral. Finally, approximate $d\omega$ by the area of the lens foreshortened by $\cos\alpha$, that is, $(\pi/4)D^2 \cos\alpha$ divided by the distance $f_o/\cos\alpha$ squared:

$$d\omega = \frac{\pi}{4} D^2 \frac{\cos^3\alpha}{f_o^2} \qquad (2.18)$$

so that finally

$$E = \frac{1}{4} \left(\frac{D}{f_p}\right)^2 \cos^4\alpha\, \pi L \qquad (2.19)$$

The interesting results here are that (1) the image irradiance is proportional to the scene radiance $L$, and (2) the factor of proportionality includes the fourth power of the off-axis angle $\alpha$. Ideally, an imaging device should be calibrated so that the variation in sensitivity as a function of $\alpha$ is removed.

### 2.2.4 Spatial Properties

*The Fourier Transform*

An image is a spatially varying function. One way to analyze spatial variations is the decomposition of an image function into a set of orthogonal functions, one such set being the Fourier (sinusoidal) functions. The Fourier transform may be used to transform the intensity image into the domain of *spatial frequency*. For no-

tational convenience and intuition, we shall generally use as an example the continuous one-dimensional Fourier transform. The results can readily be extended to the discrete case and also to higher dimensions [Rosenfeld and Kak 1976]. In two dimensions we shall denote transform domain coordinates by $(u, v)$. The one-dimensional Fourier transform, denoted $\mathcal{F}$, is defined by

$$\mathcal{F}[f(x)] = F(u)$$

where

$$F(u) = \int_{-\infty}^{+\infty} f(x)\exp(-j2\pi ux)\,dx \qquad (2.20)$$

where $j = \sqrt{(-1)}$. Intuitively, Fourier analysis expresses a function as a sum of sine waves of different frequency and phase. The Fourier transform has an *inverse* $\mathcal{F}^{-1}[F(u)] = f(x)$. This inverse is given by

$$f(x) = \int_{-\infty}^{\infty} F(u)\exp(j2\pi ux)\,du \qquad (2.21)$$

The transform has many useful properties, some of which are summarized in Table 2.1. Common one-dimensional Fourier transform pairs are shown in Table 2.2.

The transform $F(u)$ is simply another representation of the image function. Its meaning can be understood by interpreting Eq. (2.21) for a specific value of $x$, say $x_0$:

$$f(x_0) = \int F(u)\exp(j2\pi ux_0)\,du \qquad (2.22)$$

This equation states that a particular point in the image can be represented by a weighted sum of complex exponentials (sinusoidal patterns) at different spatial frequencies $u$. $F(u)$ is thus a *weighting function* for the different frequencies. Low-spatial frequencies account for the "slowly" varying gray levels in an image, such as the variation of intensity over a continuous surface. High-frequency components are associated with "quickly varying" information, such as edges. Figure 2.5 shows the Fourier transform of an image of rectangles, together with the effects of removing low- and high-frequency components.

The Fourier transform is defined above to be a continuous transform. Although it may be performed instantly by optics, a discrete version of it, the "fast Fourier transform," is almost universally used in image processing and computer vision. This is because of the relative versatility of manipulating the transform in the digital domain as compared to the optical domain. Image-processing texts, e.g., [Pratt 1978; Gonzalez and Wintz 1977] discuss the FFT in some detail; we content ourselves with an algorithm for it (Appendix 1).

*The Convolution Theorem*

*Convolution* is a very important image-processing operation, and is a basic operation of linear systems theory. The convolution of two functions $f$ and $g$ is a function $h$ of a displacement $y$ defined as

$$h(y) = f*g = \int_{-\infty}^{\infty} f(x)g(y-x)\,dx \qquad (2.23)$$

**Table 2.1**

## PROPERTIES OF THE FOURIER TRANSFORM

| Spatial Domain | Frequency Domain |
|---|---|
| $f(x)$ | $F(u) = \mathcal{F}[f(x)]$ |
| $g(x)$ | $G(u) = \mathcal{F}[g(x)]$ |

| | | |
|---|---|---|
| (1) | Linearity<br>$c_1 f(x) + c_2 g(x)$<br>$c_1, c_2$ scalars | $c_1 F(u) + c_2 G(u)$ |
| (2) | Scaling<br>$f(ax)$ | $\dfrac{1}{|a|} F\left(\dfrac{u}{a}\right)$ |
| (3) | Shifting<br>$f(x - x_0)$ | $e^{-j2\pi x_0 u} F(u)$ |
| (4) | Symmetry<br>$F(x)$ | $f(-u)$ |
| (5) | Conjugation<br>$f^*(x)$ | $F^*(-u)$ |
| (6) | Convolution<br>$h(x) = f * g = \displaystyle\int_{-\infty}^{\infty} f(x') g(x - x')\, dx'$ | $F(u) G(u)$ |
| (7) | Differentiation<br>$\dfrac{d^n f(x)}{dx^n}$ | $(2\pi j u)^n F(u)$ |

Parseval's theorem:

$$\int_{-\infty}^{\infty} |f(x)|^2 dx = \int_{-\infty}^{\infty} |F(\xi)|^2 d\xi$$

$$\int_{-\infty}^{\infty} f(x) g^*(x)\, dx = \int_{-\infty}^{\infty} F(\xi) G^*(\xi)\, d\xi$$

| $f(x)$ | $F(\xi)$ |
|---|---|
| Real($R$) | Real part even (RE)<br>Imaginary part odd (IO) |
| Imaginary (I) | RO,IE |
| RE,IO | R |
| RE,IE | I |
| RE | RE |
| RO | IO |
| IE | IE |
| IO | RO |
| Complex even (CE) | CE |
| CO | CO |

**Table 2.2**

## FOURIER TRANSFORM PAIRS

| $f(x)$ | $F(u)$ |
|---|---|
| **Rectangle function** <br> Rect $(x)$, value 1 from $-\frac{1}{2}$ to $\frac{1}{2}$ | **Sinc function** <br> $\text{Sinc}(u) = \dfrac{\sin \pi u}{\pi u}$ |
| **Triangle function** <br> value 1, from $\frac{1}{2}$ to $\frac{1}{2}$ | $\text{Sinc}^2(u)$ |
| **Exponential** <br> $e^{-\alpha\lvert x\rvert}$ | $\dfrac{2\alpha}{\alpha^2 + (2\pi u)^2}$ |
| **Gaussian** <br> $e^{-\alpha x^2}$ | $\dfrac{\pi}{\alpha} e^{-\frac{\pi u^2}{\alpha}}$ |
| **Unit impulse** $\delta(x)$ | $1$ |
| **Unit step** | $\frac{1}{2}\delta(u) + \dfrac{1}{2\pi j u}$ |

**Table 2.2**   (cont.)

Comb function

$$\sum_{n=-\infty}^{\infty} \delta(x - nx_0)$$

$$-2x_0 \quad -x_0 \quad x_0 \quad 2x_0$$

$$\frac{1}{x_0} \sum_{n=-\infty}^{\infty} \delta\left(\xi - \frac{n}{x_0}\right)$$

$$\frac{-2}{x_0} \quad \frac{-1}{x_0} \quad \frac{1}{x_0} \quad \frac{2}{x_0}$$

$$\cos 2\pi\omega_0 x$$

$$\tfrac{1}{2}[\,\delta(\xi - \omega_0) + \delta(\xi + \omega_0)\,]$$

$$\sin 2\pi\omega_0 x$$

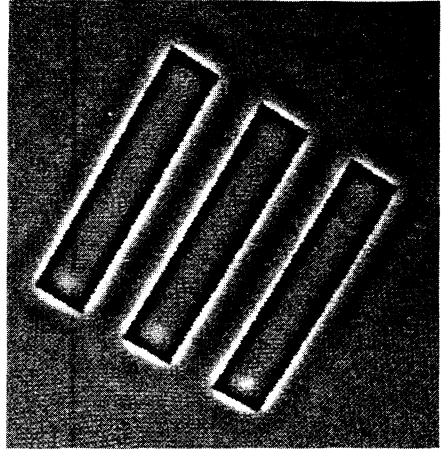$$\tfrac{1}{2}j\,[-\delta(\xi - \omega_0) + \delta(\xi + \omega_0)\,]$$

Im $F$

Intuitively, one function is "swept past" (in one dimension) or "rubbed over" (in two dimensions) the other. The value of the convolution at any displacement is the integral of the product of the (relatively displaced) function values. One common phenomenon that is well expressed by a convolution is the formation of an image by an optical system. The system (say a camera) has a "point-spread function," which is the image of a single point. (In linear systems theory, this is the "impulse response," or response to a delta-function input.) The ideal point-spread function is, of course, a point. A typical point-spread function is a two-dimensional Gaussian spatial distribution of intensities, but may include such phenomena as diffraction rings. In any event, if the camera is modeled as a linear system (ignor-
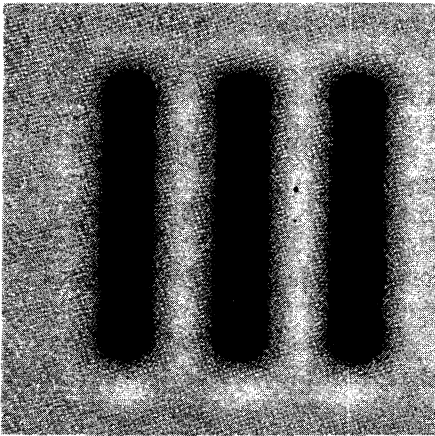
---

**Fig. 2.5**  (on facing page) (a) An image, $f(x, y)$. (b) A rotated version of (a), filtered to enhance high spatial frequencies. (c) Similar to (b), but filtered to enhance low spatial frequencies. (d), (e), and (f) show the logarithm of the power spectrum of (a), (b), and (c). The power spectrum is the log square modulus of the Fourier transform $F(u, v)$. Considered in polar coordinates $(\rho, \theta)$, points of small $\rho$ correspond to low spatial frequencies ("slowly-varying" intensities), large $\rho$ to high spatial frequencies contributed by "fast" variations such as step edges. The power at $(\rho, \theta)$ is determined by the amount of intensity variation at the frequency $\rho$ occurring at the angle $\theta$.
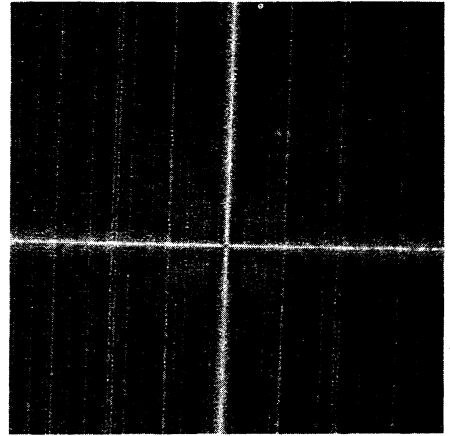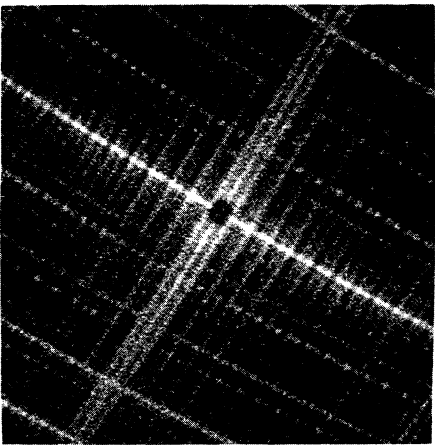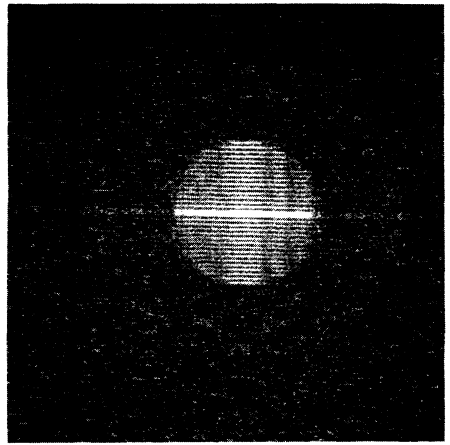
(a)

(b)

(c)

(d)

(e)

(f)

ing the added complexity that the point-spread function usually varies over the field of view), the image is the convolution of the point-spread function and the input signal. The point-spread function is rubbed over the perfect input image, thus blurring it.

Convolution is also a good model for the application of many other linear operators, such as line-detecting templates. It can be used in another guise (called correlation) to perform matching operations (Chapter 3) which detect instances of subimages or features in an image.

In the spatial domain, the obvious implementation of the convolution operation involves a shift–multiply–integrate operation which is hard to do efficiently. However, multiplication and convolution are "transform pairs," so that the calculation of the convolution in one domain (say the spatial) is simplified by first Fourier transforming to the other (the frequency) domain, performing a multiplication, and then transforming back.

The convolution of $f$ and $g$ in the spatial domain is equivalent to the point-wise product of $F$ and $G$ in the frequency domain,

$$\mathcal{F}(f*g) = FG \tag{2.24}$$

We shall show this in a manner similar to [Duda and Hart 1973]. First we prove the *shift theorem*. If the Fourier transform of $f(x)$ is $F(u)$, defined as

$$F(u) = \int_x f(x) \exp\left[-j2\pi(ux)\right]dx \tag{2.25}$$

then

$$\mathcal{F}\left[f(x-a)\right] = \int_x f(x-a) \exp\left[-j2\pi(ux)\right]dx \tag{2.26}$$

changing variables so that $x' = x - a$ and $dx = dx'$

$$= \int_{x'} f(x') \exp\left\{-j2\pi[u(x'+a)]\right\}dx' \tag{2.27}$$

Now $\exp[-j2\pi u(x'+a)] = \exp(-j2\pi ua) \exp(-j2\pi ux')$, where the first term is a constant. This means that

$$\mathcal{F}\left[f(x-a)\right] = \exp(-j2\pi ua)F(u) \qquad \text{(shift theorem)}$$

Now we are ready to show that $\mathcal{F}[f(x)*g(x)] = F(u)G(u)$.

$$\mathcal{F}(f*g) = \int_y \left\{\int_x f(x)g(y-x)\right\} \exp(-j2\pi uy)\ dx\ dy \tag{2.28}$$

$$= \int_x f(x)\left\{\int_y g(y-x) \exp(-j2\pi uy)\ dy\right\} dx \tag{2.29}$$

Recognizing that the terms in braces represent $\mathcal{F}[g(y-x)]$ and applying the shift theorem, we obtain

$$\mathcal{F}(f*g) = \int_x f(x)\exp(-j2\pi ux)G(u)\ dx \tag{2.30}$$

$$= F(u)G(u) \tag{2.31}$$

### 2.2.5 Color

Not all images are monochromatic; in fact, applications using multispectral images are becoming increasingly common (Section 2.3.2). Further, human beings intuitively feel that color is an important part of their visual experience, and is useful or even necessary for powerful visual processing in the real world. Color vision provides a host of research issues, both for psychology and computer vision. We briefly discuss two aspects of color vision: color spaces and color perception. Several models of the human visual system not only include color but have proven useful in applications [Granrath 1981].

*Color Spaces*

*Color spaces* are a way of organizing the colors perceived by human beings. It happens that weighted combinations of stimuli at three principal wavelengths are sufficient to define almost all the colors we perceive. These wavelengths form a natural basis or coordinate system from which the color measurement process can be described. Color perception is not related in a simple way to color measurement, however.

Color is a perceptual phenomenon related to human response to different wavelengths in the visible *electromagnetic spectrum* [400 (blue) to 700 nanometers (red); a nanometer (nm) is $10^{-9}$ meter]. The sensation of color arises from the sensitivities of three types of neurochemical sensors in the retina to the visible spectrum. The relative response of these sensors is shown in Fig. 2.6. Note that each sensor responds to a range of wavelengths. The illumination source has its own spectral composition $f(\lambda)$ which is modified by the reflecting surface. Let $r(\lambda)$ be this reflectance function. Then the measurement $R$ produced by the "red" sensor is given by

$$R = \int f(\lambda) r(\lambda) h_R(\lambda) \, d\lambda \qquad (2.32)$$

So the sensor output is actually the integral of three different wavelength-dependent components: the source $f$, the surface reflectance $r$, and the sensor $h_R$.

Surprisingly, only weighted combinations of three delta-function approximations to the different $f(\lambda) h(\lambda)$, that is, $\delta(\lambda_R)$, $\delta(\lambda_G)$, and $\delta(\lambda_B)$, are necessary to
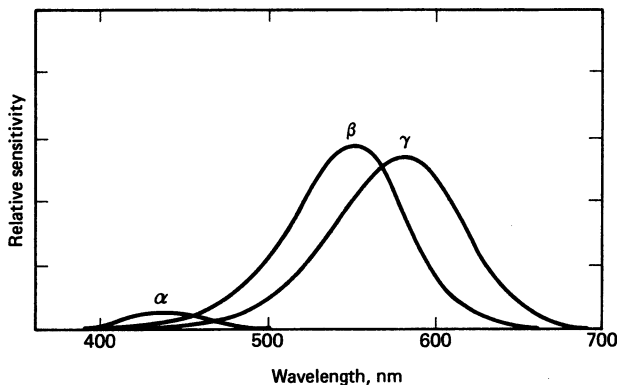


Fig. 2.6 Spectral response of human color sensors.

produce the sensation of nearly all the colors. This result is displayed on a *chromaticity diagram*. Such a diagram is obtained by first normalizing the three sensor measurements:
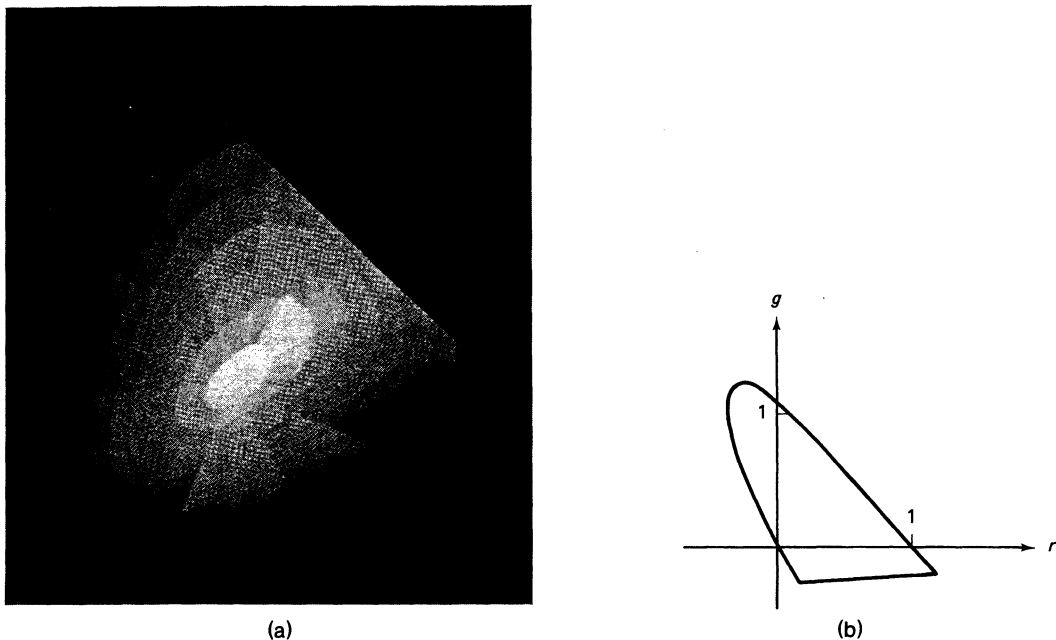
$$r = \frac{R}{R + G + B}$$
$$g = \frac{G}{R + G + B} \qquad (2.33)$$
$$b = \frac{B}{R + G + B}$$

and then plotting perceived color as a function of any two (usually red and green). Chromaticity explicitly ignores intensity or brightness; it is a section through the three-dimensional color space (Fig. 2.7). The choice of $(\lambda_R, \lambda_G, \lambda_B) = (410, 530, 650) \, nm$ maximizes the realizable colors, but some colors still cannot be realized since they would require negative values for some of $r$, $g$, and $b$.

Another more intuitive way of visualizing the possible colors from the $RGB$ space is to view these measurements as Euclidean coordinates. Here any color can be visualized as a point in the unit cube. Other coordinate systems are useful for different applications; computer graphics has proved a strong stimulus for investigation of different color space bases.

### Color Perception

Color perception is complex, but the essential step is a transformation of three input intensity measurements into another basis. The coordinates of the new



(a)            (b)

**Fig. 2.7** (a) An artist's conception of the chromaticity diagram—*see color insert*; (b) a more useful depiction. Spectral colors range along the curved boundary; the straight boundary is the line of purples.

*Ch. 2   Image Formation*

basis are more directly related to human color judgments.

Although the *RGB* basis is good for the acquisition or display of color information, it is not a particularly good basis to explain the perception of colors. Human vision systems can make good judgments about the relative surface reflectance $r(\lambda)$ despite different illuminating wavelengths; this reflectance seems to be what we mean by surface color.

Another important feature of the color basis is revealed by an ability to perceive in "black and white," effectively deriving intensity information from the color measurements. From an evolutionary point of view, we might expect that color perception in animals would be compatible with preexisting noncolor perceptual mechanisms.

These two needs—the need to make good color judgments and the need to retain and use intensity information—imply that we use a transformed, non-*RGB* basis for color space. Of the different bases in use for color vision, all are variations on this theme: Intensity forms one dimension and color is a two-dimensional subspace. The differences arise in how the color subspace is described. We categorize such bases into two groups.

1. *Intensity/Saturation/Hue (IHS)*. In this basis, we compute intensity as

$$\text{intensity:} = R + G + B \qquad (2.34)$$

The saturation measures the lack of whiteness in the color. Colors such as "fire engine" red and "grass" green are saturated; pastels (e.g., pinks and pale blues) are desaturated. Saturation can be computed from *RGB* coordinates by the formula [Tenenbaum and Weyl 1975]

$$\text{saturation:} = 1 - \frac{3 \min (R, G, B)}{\text{intensity}} \qquad (2.35)$$

Hue is roughly proportional to the average wavelength of the color. It can be defined using *RGB* by the following program fragment:

$$\text{hue:} = \cos^{-1} \left\{ \frac{\{\frac{1}{2}[(R - G) + (R - B)]\}}{\sqrt{(R - G)^2 + (R - B)(G - B)}} \right\} \qquad (2.36)$$

$$\text{If } B > G \text{ then hue:} = 2pi - \text{hue}$$

The IHS basis transforms the *RGB* basis in the following way. Thinking of the color cube, the diagonal from the origin to (1, 1, 1) becomes the intensity axis. Saturation is the distance of a point from that axis and hue is the angle with regard to the point about that axis from some reference (Fig. 2.8).

This basis is essentially that used by artists [Munsell 1939], who term saturation *chroma*. Also, this basis has been used in graphics [Smith 1978; Joblove and Greenberg 1978].

One problem with the IHS basis, particularly as defined by (2.34) through (2.36), is that it contains essential singularities where it is impossible to define the color in a consistent manner [Kender 1976]. For example, hue has an essential singularity for all values of $(R, G, B)$, where $R = G = B$. This means that special care must be taken in algorithms that use hue.

2. *Opponent processes*. The opponent process basis uses Cartesian rather than
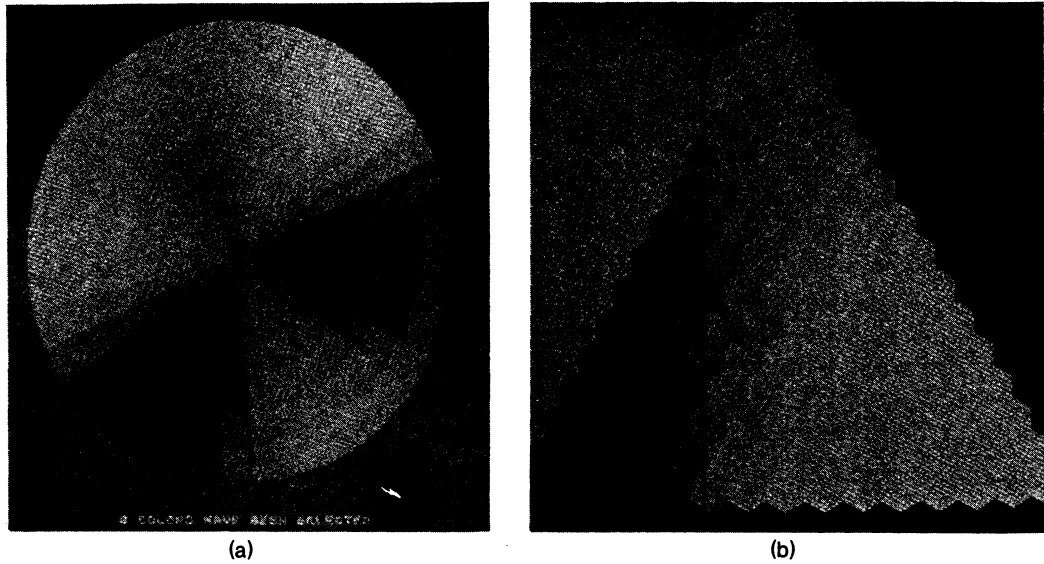
**Fig. 2.8** An IHS Color Space. (a) Cross section at one intensity; (b) cross section at one hue—*see color inserts.*

cylindrical coordinates for the color subspace, and was first proposed by Hering [Teevan and Birney 1961]. The simplest form of basis is a linear transformation from $R$, $G$, $B$ coordinates. The new coordinates are termed "$R - G$", "$Bl - Y$", and "$W - Bk$":

$$\begin{bmatrix} R - G \\ Bl - Y \\ W - Bk \end{bmatrix} = \begin{bmatrix} 1 & -2 & 1 \\ -1 & -1 & 2 \\ 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} R \\ G \\ B \end{bmatrix}$$

The advocates of this representation, such as [Hurvich and Jameson 1957], theorize that this basis has neurological correlates and is in fact the way human beings represent ("name") colors. For example, in this basis it makes sense to talk about a "reddish blue" but not a "reddish green." Practical opponent process models usually have more complex weights in the transform matrix to account for psychophysical data. Some startling experiments [Land 1977] show our ability to make correct color judgments even when the illumination consists of only two principal wavelengths. The opponent process, at the level at which we have developed it, does not demonstrate how such judgments are made, but does show how stimulus at only two wavelengths will project into the color subspace. Readers interested in the details of the theory should consult the references.

Commercial television transmission needs an intensity, or "$W - Bk$" component for black-and-white television sets while still spanning the color space. The National Television Systems Committee (NTSC) uses a "YIQ" basis extracted from *RGB* via

$$\begin{bmatrix} I \\ Q \\ Y \end{bmatrix} = \begin{bmatrix} 0.60 & -0.28 & -0.32 \\ 0.21 & -0.52 & 0.31 \\ 0.30 & 0.59 & 0.11 \end{bmatrix} \begin{bmatrix} R \\ G \\ B \end{bmatrix}$$

This basis is a weighted form of

$$(I, \ Q, \ Y) = (\text{``}R-\text{cyan, '' ``magenta}-\text{green, '' ``}W-Bk\text{''})$$
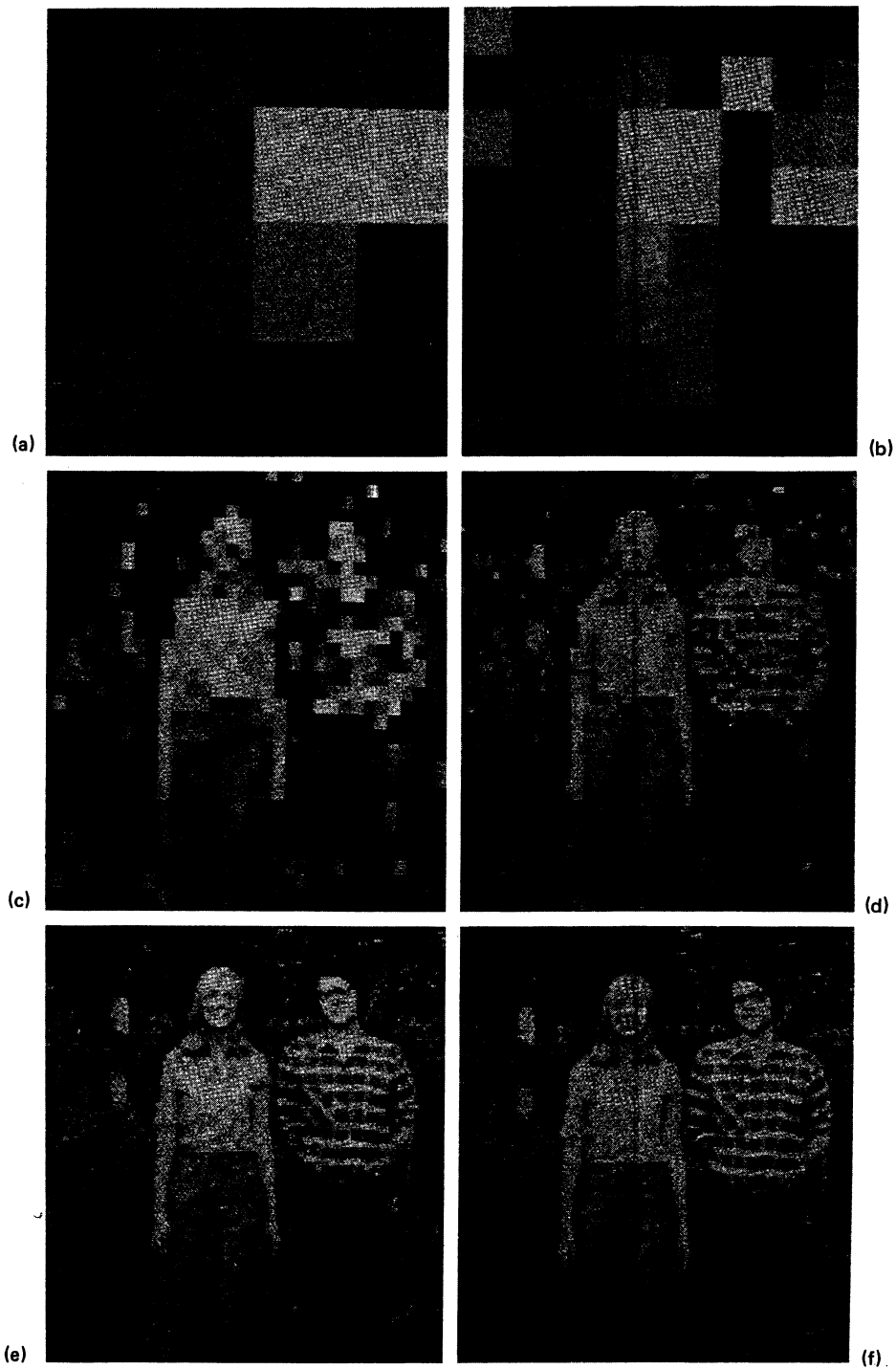
### 2.2.6 Digital Images

The *digital images* with which computer vision deals are represented by $m$-vector discrete-valued image functions $f(\mathbf{x})$, usually of one, two, three, or four dimensions.

Usually $m = 1$, and both the domain and range of $f(\mathbf{x})$ are discrete. The domain of $f$ is finite, usually a rectangle, and the range of $f$ is positive and bounded: $0 \leqslant f(\mathbf{x}) \leqslant M$ for some integer $M$. For all practical purposes, the image is a continuous function which is represented by measurements or *samples* at regularly spaced intervals. At the time the image is sampled, the intensity is usually *quantized* into a number of different *gray levels*. For a discrete image, $f(\mathbf{x})$ is an integer gray level, and $\mathbf{x} = (x, y)$ is a pair of *integer* coordinates representing a sample point in a two-dimensional image plane. Sampling involves two important choices: (1) the *sampling interval*, which determines in a basic way whether all the information in the image is represented, and (2) the *tesselation* or spatial pattern of sample points, which affects important notions of connectivity and distance. In our presentation, we first show qualitatively the effects of sampling and gray-level quantization. Second, we discuss the simplest kinds of tesselations of the plane. Finally, and most important, we describe the sampling theorem, which specifies how close the image samples must be to represent the image unambiguously.

The choice of integers to represent the gray levels and coordinates is dictated by limitations in sensing. Also, of course, there are hardware limitations in representing images arising from their sheer size. Table 2.3 shows the storage required for an image in 8-bit bytes as a function of m, the number of bits per sample, and N, the linear dimension of a square image.

For reasons of economy (and others discussed in Chapter 3) we often use images of considerably less spatial resolution than that required to preserve fidelity to the human viewer. Figure 2.9 provides a qualitative idea of image degradation with decreasing spatial resolution.

As shown in Table 2.3, another way to save space besides using less spatial resolution is to use fewer bits per gray level sample. Figure 2.10 shows an image represented with different numbers of bits per sample. One striking effect is the "contouring" introduced with small numbers of gray levels. This is, in general, a problem for computer vision algorithms, which cannot easily discount the false contours. The choice of spatial and gray-level resolution for any particular computer vision task is an important one which depends on many factors. It is typical in

**Fig. 2.9** Using different numbers of samples. (a) $N = 16$; (b) $N = 32$; (c) $N = 64$; (d) $N = 128$; (e) $N = 256$; (f) $N = 512$.

**Table 2.3**

**NUMBER OF 8-BIT BYTES OF STORAGE FOR
VARIOUS VALUES OF N AND M**

| $N$ | 32 | 64 | 128 | 256 | 512 |
|---|---|---|---|---|---|
| $m$ | | | | | |
| 1 | 128 | 512 | 2,048 | 8,192 | 32,768 |
| 2 | 256 | 1,024 | 4,096 | 16,384 | 65,536 |
| 3 | 512 | 2,048 | 8,192 | 32,768 | 131,072 |
| 4 | 512 | 2,048 | 8,192 | 32,768 | 131,072 |
| 5 | 1,024 | 4,096 | 16,384 | 65,536 | 262,144 |
| 6 | 1,024 | 4,096 | 16,384 | 65,536 | 262,144 |
| 7 | 1,024 | 4,096 | 16,384 | 65,536 | 262,144 |
| 8 | 1,024 | 4,096 | 16,384 | 65,536 | 262,144 |

computer vision to have to balance the desire for increased resolution (both gray scale and spatial) against its cost. Better data can often make algorithms easier to write, but a small amount of data can make processing more efficient. Of course, the image domain, choice of algorithms, and image characteristics all heavily influence the choice of resolutions.
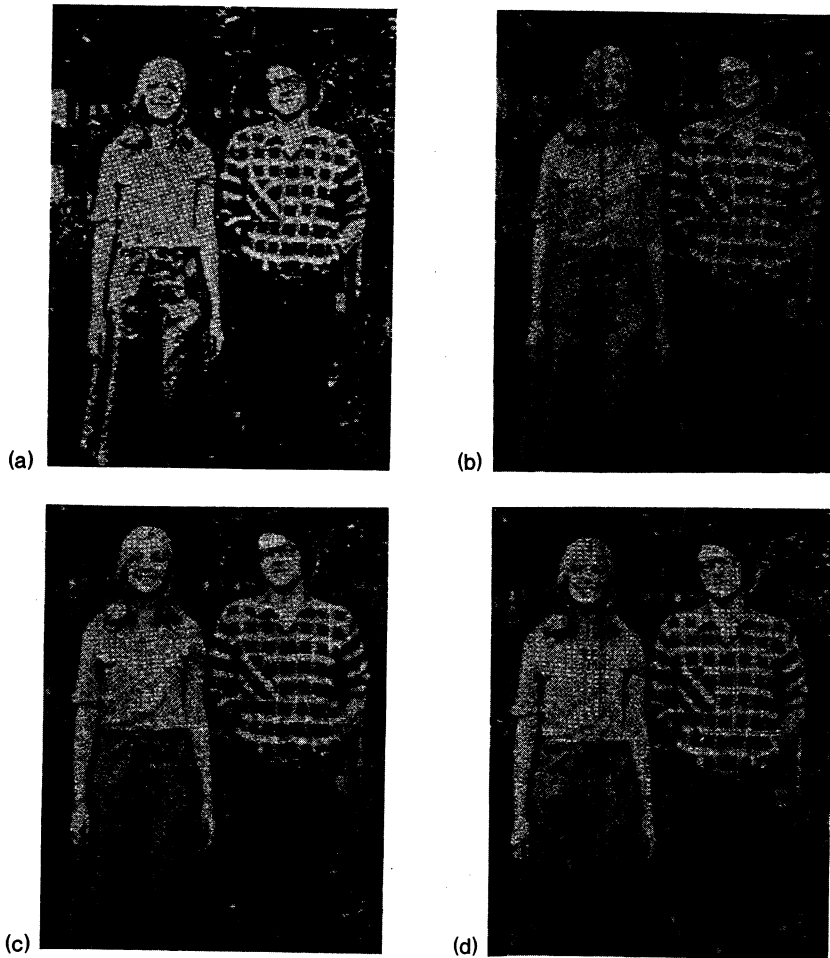
*Tesselations and Distance Metrics*

Although the spatial samples for $f(x)$ can be represented as points, it is more satisfying to the intuition and a closer approximation to the acquisition process to think of these samples as finite-sized cells of constant gray-level partitioning the image. These cells are termed *pixels*, an acronym for *picture elements*. The pattern into which the plane is divided is called its *tesselation*. The most common regular tesselations of the plane are shown in Fig. 2.11.

Although rectangular tesselations are almost universally used in computer vision, they have a structural problem known as the "connectivity paradox." Given a pixel in a rectangular tesselation, how should we define the pixels to which it is connected? Two common ways are *four-connectivity* and *eight-connectivity*, shown in Fig. 2.12.

However, each of these schemes has complications. Consider Fig. 2.12c, consisting of a black object with a hole on a white background. If we use four-connectedness, the figure consists of four disconnected pieces, yet the hole is separated from the "outside" background. Alternatively, if we use eight-connectedness, the figure is one connected piece, yet the hole is now connected to the outside. This paradox poses complications for many geometric algorithms. Triangular and hexagonal tesselations do not suffer from connectivity difficulties (if we use three-connectedness for triangles); however, *distance* can be more difficult to compute on these arrays than for rectangular arrays.

The distance between two pixels in an image is an important measure that is fundamental to many algorithms. In general, a distance $d$ is a *metric*. That is,

**Fig. 2.10** Using different numbers of bits per sample. (a) $m = 1$; (b) $m = 2$; (c) $m = 4$; (d) $m = 8$.

(1)  $d(\mathbf{x}, \mathbf{y}) = 0$ iff $\mathbf{x} = \mathbf{y}$

(2)  $d(\mathbf{x}, \mathbf{y}) = d(\mathbf{y}, \mathbf{x})$

(3)  $d(\mathbf{x}, \mathbf{y}) + d(\mathbf{y}, \mathbf{z}) \geqslant d(\mathbf{x}, \mathbf{z})$
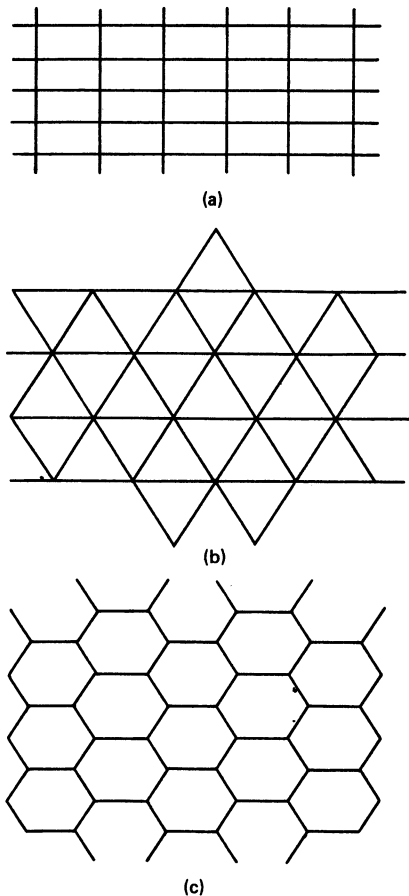
For square arrays with unit spacing between pixels, we can use any of the following common distance metrics (Fig. 2.13) for two pixels $\mathbf{x} = (x_1, y_1)$ and $\mathbf{y} = (x_2, y_2)$.

Euclidean:

$$d_e(\mathbf{x}, \mathbf{y}) = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2} \tag{2.37}$$

City block:

$$d_{cb}(\mathbf{x}, \mathbf{y}) = |x_1 - x_2| + |y_1 - y_2| \tag{2.38}$$

Fig. 2.11 Different tesselations of the image plane. (a) Rectangular; (b) triangular; (c) hexagonal.

Chessboard:

$$d_{ch}(\mathbf{x}, \mathbf{y}) = \max\left\{|x_1 - x_2|, |y_1 - y_2|\right\} \qquad (2.39)$$

Other definitions are possible, and all such measures extend to multiple dimensions. The tesselation of higher-dimensional space into pixels usually is confined to (n-dimensional) cubical pixels.

### The Sampling Theorem

Consider the one-dimensional "image" shown in Fig. 2.14. To digitize this image one must sample the image function. These samples will usually be separated at regular intervals as shown. How far apart should these samples be to allow reconstruction (to a given accuracy) of the underlying continuous image from its samples? This question is answered by the Shannon sampling theorem. An excellent rigorous presentation of the sampling theorem may be found in [Rosenfeld and Kak 1976]. Here we shall present a shorter graphical interpretation using the results of Table 2.2. For simplicity we consider the image to be periodic in order to avoid small edge effects introduced by the finite image domain. A more rigorous
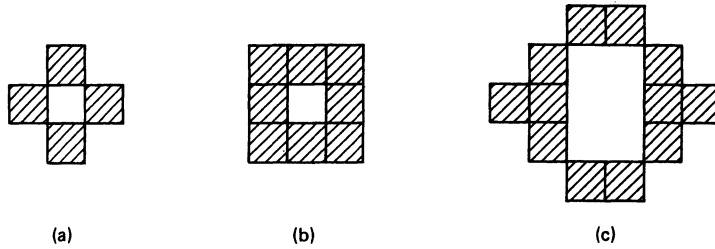
**Fig. 2.12** Connectivity paradox for rectangular tesselations. (a) A central pixel and its 4-connected neighbors; (b) a pixel and its 8-connected neighbors; (c) a figure with ambiguous connectivity.
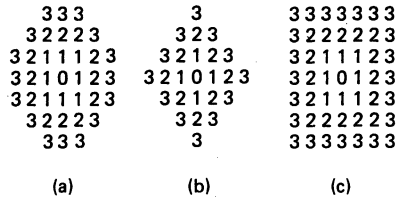
```
     333              3        3333333
   32223             323       3222223
 3211123           32123       3211123
 3210123         3210123       3210123
 3211123           32123       3211123
   32223             323       3222223
     333              3        3333333

    (a)              (b)           (c)
```

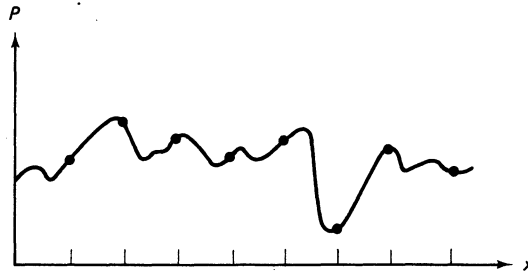**Fig. 2.13** Equidistant contours for different metrics.



**Fig. 2.14** One-dimensional image and its samples.

treatment, which considers these effects, is given in [Andrews and Hunt 1977].

Suppose that the image is sampled with a "comb" function of spacing $x_0$ (see Table 2.2). Then the sampled image can be modeled by

$$f_s(x) = f(x) \sum_n \delta(x - nx_0) \tag{2.40}$$

where the image function modulates the comb function. Equivalently, this can be written as

$$f_s(x) = \sum_n f(nx_0) \delta(x - nx_0) \tag{2.41}$$

The right-hand side of Eq. (2.40) is the product of two functions, so that property

(6) in Table 2.1 is appropriate. The Fourier transform of $f_s(x)$ is equal to the convolution of the transforms of each of the two functions. Using this result yields

$$F_s(u) = F(u) * \frac{1}{x_0} \sum_n \delta(u - \frac{n}{x_0}) \tag{2.42}$$

But from Eq. (2.3),

$$F(u) * \delta(u - \frac{n}{x_0}) = F(u - \frac{n}{x_0}) \tag{2.43}$$

so that

$$F_s(u) = \frac{1}{x_0} \sum_n F(u - \frac{n}{x_0}) \tag{2.44}$$

Therefore, sampling the image function $f(x)$ at intervals of $x_0$ is equivalent in the frequency domain to replicating the transform of $f$ at intervals of $\frac{1}{x_0}$. This limits the recovery of $f(x)$ from its sampled representation, $f_s(x)$. There are two basic situations to consider. If the transform of $f(x)$ is *bandlimited* such that $F(u) = 0$ for $|u| > 1/(2x_0)$, then there is no overlap between successive replications of $F(u)$ in the frequency domain. This is shown for the case of Fig. 2.15a, where we have arbitrarily used a triangular-shaped image transform to illustrate the effects of sampling. Incidentally, note that for this transform $F(u) = F(-u)$ and that it has no imaginary part; from Table 2.2, the one-dimensional image must also be real and even. Now if $F(u)$ is not bandlimited, i.e., there are $u > \frac{1}{2x_0}$ for which $F(u) \neq 0$, then components of different replications of $F(u)$ will interact to produce the composite function $F_s(u)$, as shown in Fig. 2.15b. In the first case $f(x)$ can be recovered from $F_s(u)$ by multiplying $F_s(u)$ by a suitable $G(u)$:

$$G(u) = \begin{cases} 1 & |u| < \frac{1}{2x_0} \\ 0 & \text{otherwise} \end{cases} \tag{2.45}$$

Then

$$f(x) = \mathcal{F}^{-1}[F_s(u)G(u)] \tag{2.46}$$

However, in the second case, $F_s(u)G(u)$ is very different from the original $F(u)$. This is shown in Fig. 2.15c. Sampling a $F(u)$ that is not bandlimited allows information at high spatial frequencies to interfere with that at low frequencies, a phenomenon known as *aliasing*.

Thus the sampling theorem has this very important result: As long as the image contains no spatial frequencies greater than one-half the sampling frequency, the underlying continuous image is unambiguously represented by its samples. However, lest one be tempted to insist on images that have been so sampled, note that it may be useful to sample at lower frequencies than would be required for total reconstruction. Such sampling is usually preceded by some form of blurring of
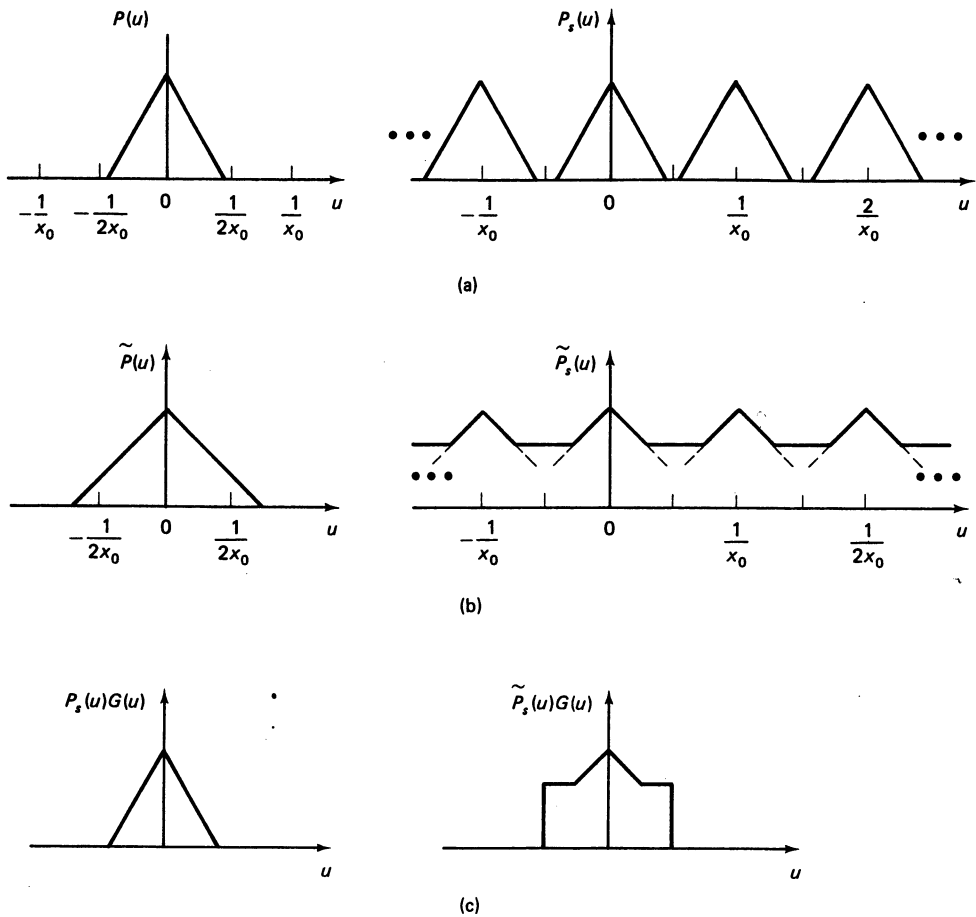
**Fig. 2.15** (a) $F(u)$ bandlimited so that $F(u) = 0$ for $|u| > 1/2x_0$. (b) $F(u)$ not bandlimited as in (a). (c) reconstructed transform.

the image, or can be incorporated with such blurring (by integrating the image intensity over a finite area for each sample). Image blurring can bury irrelevant details, reduce certain forms of noise, and also reduce the effects of aliasing.

## 2.3 IMAGING DEVICES FOR COMPUTER VISION

There is a vast array of methods for obtaining a digital image in a computer. In this section we have in mind only "traditional" images produced by various forms of radiation impinging on a sensor after having been affected by physical objects.

Many sensors are best modeled as an *analog* device whose response must be *digitized* for computer representation. The types of imaging devices possible are limited only by the technical ingenuity of their developers; attempting a definitive
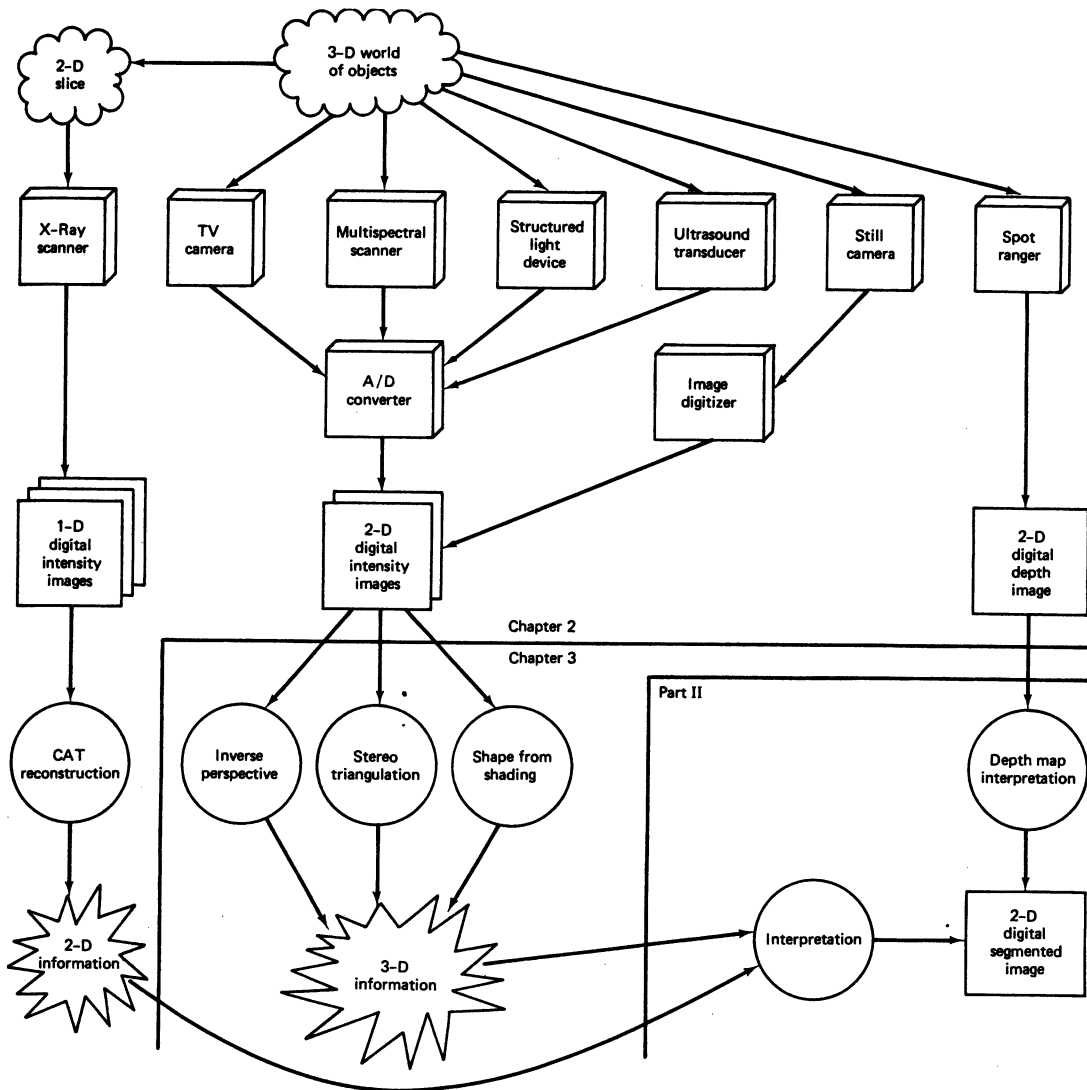
**Fig. 2.16** Imaging devices (boxes), information structurés (rectangles), and processes (circles).

taxonomy is probably unwise. Figure 2.16 is a flowchart of devices, information structures, and processes addressed in this and succeeding sections.

When the image already exists in some form, or physical considerations limit choice of imaging technology, the choice of digitizing technology may still be open. Most images are carried on a permanent medium, such as film, or at least are available in (essentially) analog form to a digitizing device. Generally, the relevant technical characteristics of imaging or digitizing devices should be foremost in mind when a technique is being selected. Such considerations as the signal-to-noise ratio of the device, its resolution, the speed at which it works, and its expense are important issues.

### 2.3.1 Photographic Imaging

The camera is the most familiar producer of optical images on a permanent medium. We shall not address here the multitudes of still- and movie-camera options; rather, we briefly treat the characteristics of the photographic film and of the digitizing devices that convert the image to machine-readable form. More on these topics is well presented in the References.

Photographic (black-and-white) film consists of an emulsion of silver halide crystals on a film base. (Several other layers are identifiable, but are not essential to an understanding of the relevant properties of film.) Upon exposure to light, the silver halide crystals form *development centers,* which are small grains of metallic silver. The photographic development process extends the formation of metallic silver to the entire silver halide crystal, which thus becomes a binary ("light" or "no light") detector. Subsequent processing removes undeveloped silver halide. The resulting film *negative* is dark where many crystals were developed and light where few were. The resolution of the film is determined by the *grain* size, which depends on the original halide crystals and on development techniques. Generally, the *faster* the film (the less light needed to expose it), the coarser the grain. Film exists that is sensitive to infrared radiation; x-ray film typically has two emulsion layers, giving it more gray-level range than that of normal film.

A repetition of the negative-forming process is used to obtain a photographic *print.* The negative is projected onto photographic paper, which responds roughly in the same way as the negative. Most photographic print paper cannot capture in one print the range of densities that can be present in a negative. Positive films do exist that do not require printing; the most common example is color slide film.

The response of film to light is not completely linear. The photographic *density* obtained by a negative is defined as the logarithm (base 10) of the ratio of incident light to transmitted light.

$$D = \log_{10}\left(\frac{I_i}{I_t}\right)$$

The *exposure* of a negative dictates (approximately) its response. Exposure is defined as the energy per unit area that exposed the film (in its sensitive spectral range). Thus exposure is the product of the *intensity* and the time of exposure. This mathematical model of the behavior of the photographic exposure process is correct for a wide operating range of the film, but *reciprocity failure* effects in the film keep one from being able always to trade light level for exposure time. At very low light levels, longer exposure times are needed than are predicted by the product rule.

The response of film to light is usually plotted in an "H&D curve" (named for Hurter and Driffield), which plots density versus exposure. The H&D curve of film displays many of its important characteristics. Figure 2.17 exhibits a typical H&D curve for a black and white film.

The *toe* of the curve is the lower region of low slope. It expresses reciprocity failure and the fact that the film has a certain bias, or *fog* response, which dominates its behavior at the lowest exposure levels. As one would expect, there is an upper limit to the density of the film, attained when a maximum number of silver
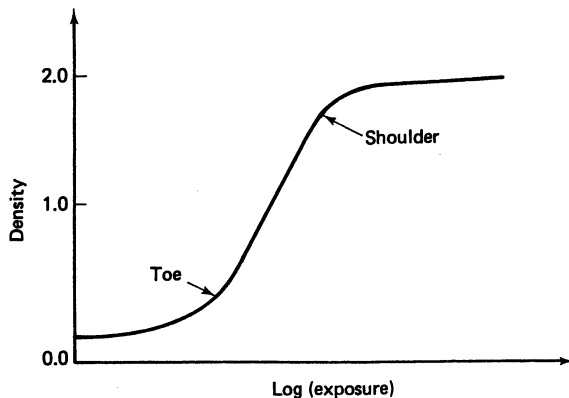
Fig. 2.17 Typical H & D curve.

halide crystals are rendered developable. Increasing exposure beyond this maximum level has little effect, accounting for the *shoulder* in the H&D curve, or its flattened upper end.

In between the toe and shoulder, there is typically a linear operating region of the curve. High-contrast films are those with high slope (traditionally called *gamma*); they respond dramatically to small changes in exposure. A high-contrast film may have a gamma between about 1.5 and 10. Films with gammas of approximately 10 are used in graphics arts to copy line drawings. General-purpose films have gammas of about 0.5 to 1.0.

The resolution of a general film is about 40 lines/mm, which means that a 1400 × 1400 image may be digitized from a 35mm slide. At any greater sampling frequency, the individual film grains will occupy more than a pixel, and the resolution will thus be grain-limited.

### Image Digitizers (Scanners)

Accuracy and speed are the main considerations in converting an image on film into digital form. Accuracy has two aspects: spatial resolution, loosely the level of image spatial detail to which the digitizer can respond, and gray-level resolution, defined generally as the range of densities or reflectances to which the digitizer responds and how finely it divides the range. Speed is also important because usually many data are involved; images of 1 million samples are commonplace.

Digitizers broadly take two forms: mechanical and "flying spot." In a mechanical digitizer, the film and a sensing assembly are mechanically transported past one another while readings are made. In a flying-spot digitizer, the film and sensor are static. What moves is the "flying spot," which is a point of light on the face of a cathode-ray tube, or a laser beam directed by mirrors. In all digitizers a very narrow beam of light is directed through the film or onto the print at a known coordinate point. The light transmittance or reflectance is measured, transformed from analog to digital form, and made available to the computer through interfacing electronics. The location on the medium where density is being measured may also be transmitted with each reading, but it is usually determined by relative offset from positions transmitted less frequently. For example, a "new scan line" impulse is transmitted for TV output; the position along the current scan line yields an $x$ position, and the number of scan lines yields a $y$ position.

The mechanical scanners are mostly of two types, *flat-bed* and *drum*. In a flat-bed digitizer, the film is laid flat on a surface over which the light source and the sensor (usually a very accurate photoelectric cell) are transported in a raster fashion. In a drum digitizer, the film is fastened to a circular drum which revolves as the sensor and light source are transported down the drum parallel to its axis of rotation.

Color mechanical digitizers also exist; they work by using colored filters, effectively extracting in three scans three "color overlays" which when superimposed would yield the original color image. Extracting some "composite" color signal with one reading presents technical problems and would be difficult to do as accurately.

### Satellite Imagery

LANDSAT and ERTS (Earth Resources Technology Satellites) have similar scanners which produce images of 2340 x 3380 7-bit pixels in four spectral bands, covering an area of 100 × 100 nautical miles. The scanner is mechanical, scanning six horizontal scan lines at a time; the rotation of the earth accounts for the advancement of the scan in the vertical direction.
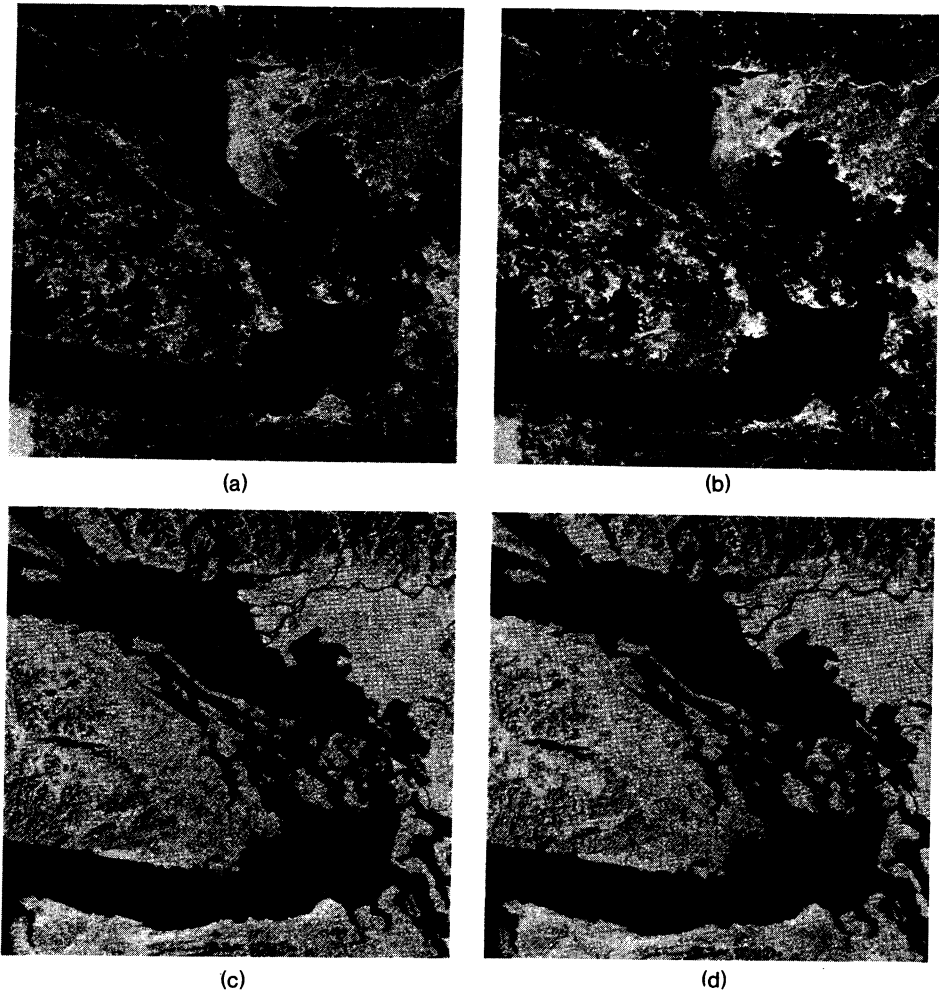
A set of four images is shown in Fig. 2.18. The four spectral bands are numbered 4, 5, 6, and 7. Band 4 [0.5 to 0.6 $\mu$m (green)] accentuates sediment-laden water and shallow water, band 5 [0.6 to 0.7 $\mu$m (red)] emphasizes cultural features such as roads and cities; band 6 [0.7 to 0.8 $\mu$m (near infrared)] emphasizes vegetation and accentuates the contrast between land and water, band 7 [0.8 to 1.1 $\mu$m (near infrared)] is like band 6 except that it is better at penetrating atmospheric haze.

The LANDSAT images are available at nominal cost from the U.S. government (The EROS Data Center, Sioux Falls, South Dakota 57198). They are furnished on tape, and cover the entire surface of the earth (often the buyer has a choice of the amount of cloud cover). These images form a huge data base of multispectral imagery, useful for land-use and geological studies; they furnish something of an image analysis challenge, since one satellite can produce some 6 billion bits of image data per day.

### Television Imaging

Television cameras are appealing devices for computer vision applications for several reasons. For one thing, the image is immediate; the camera can show events as they happen. For another, the image is already in electrical, if not digital form. "Television camera" is basically a nontechnical term, because many different technologies produce video signals conforming to the standards set by the FCC and NTSC. Cameras exist with a wide variety of technical specifications.

Usually, TV cameras have associated electronics which scan an entire "picture" at a time. This operation is closely related to broadcast and receiver standards, and is more oriented to human viewing than to computer vision. An entire image (of some 525 scan lines in the United States) is called a *frame*, and consists of two *fields*, each made up of alternate scan lines from the frame. These fields are generated and transmitted sequentially by the camera electronics. The transmitted image is thus *interlaced*, with all odd-numbered scan lines being "painted" on the

**Fig. 2.18** The straits of Juan de Fuca as seen by the LANDSAT multispectral scanner. (a) Band 4; (b) band 5; (c) band 6; (d) band 7.

screen alternating with all even-numbered scan lines. In the United States, each field takes $\frac{1}{60}$ sec to scan, so a whole frame is scanned every $\frac{1}{30}$ sec. The interlacing is largely to prevent flickering of the image, which would become noticeable if the frame were painted from top to bottom only once in $\frac{1}{30}$ sec. These automatic scanning electronics may be replaced or overridden in many cameras, allowing "random access" to the image. In some technologies, such as the image dissector, the longer the signal is collected from any location, the better the signal-to-noise performance.

There are a number of different systems used to generate television images. We discuss five main methods below.

*Image orthicon tube.* This is one of the two main methods in use today (in addition to the vidicon). It offers very stable performance at all incident light levels

and is widely used in commercial television. It is a storage-type tube, since it depends on the neutralization of positive charges by a scanning electron beam.

The image orthicon (Fig. 2.19) is divided into an imaging and readout section. In the imaging section, light from the scene is focused onto a semitransparent photocathode. This photocathode operates the same way as the cathode in a phototube. It emits electrons which are magnetically focused by a coil and are accelerated toward a positively charged target. The target is a thin glass disk with a fine-wire-mesh screen facing the photocathode. When electrons strike it, secondary emission from the glass takes place. As electrons are emitted from the photocathode side of the disk, positive charges build up on the scanning side. These charges correspond to the pattern of light intensity in the scene being viewed.

In the readout section, the back of the target is scanned by a low velocity electron beam from an electron gun at the rear of the tube. Electrons in this beam are absorbed by the target in varying amounts, depending on the charge on the target. The image is represented by the amplitude-modulated intensity of the returned beam.

*Vidicon tube.* The vidicon is smaller, lighter, and more rugged than the image orthicon, making it ideal for portable use. Here the target (the inner surface of the face plate) is coated with a transparent conducting film which forms a video signal electrode (Fig. 2.20). A thin photosensitive layer is deposited on the film, consisting of a large number of tiny resistive globules whose resistance decreases on illumination. This layer is scanned in raster fashion by a low velocity electron beam from the electron gun at the rear of the tube. The beam deposits electrons on the layer, thus reducing its surface potential. The two surfaces of the target essentially form a capacitor, and the scanning action of the beam produces a capacitive current at the video signal electrode which represents the video signal.

The plumbicon is essentially a vidicon with a lead oxide photosensitive layer. It offers the following advantages over the vidicon: higher sensitivity, lower dark current, and negligible persistence or lag.
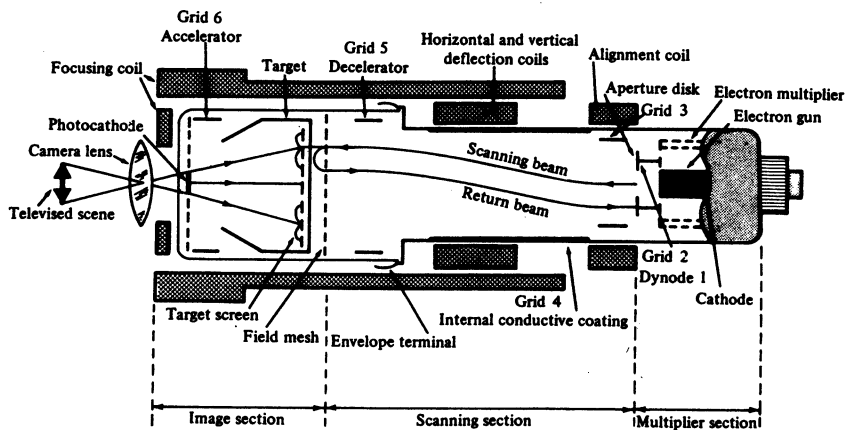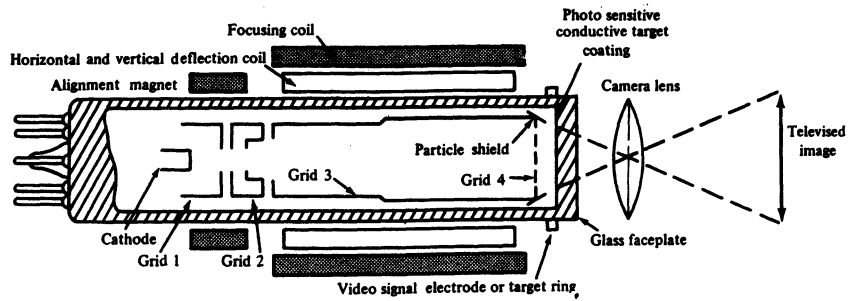


Fig. 2.19 The image orthicon.

Fig. 2.20 The vidicon.

*Iconoscope tube.* The iconoscope is now largely of historical interest. In it, an electron beam scans a target consisting of a thin mica sheet or mosaic coated with a photosensitive layer. In contrast to the vidicon and orthicon, the electron beam and the light both strike the same side of the target surface. The back of the mosaic is covered with a conductive film connected to an output load. The arrangement is equivalent to a matrix of small capacitors which discharge through a common lead.

*Image dissector. tube.* The image dissector tube operates on instantaneous scanning rather than by neutralizing positive charges. Light from the scene is focused on a cathode coated with a photosensitive layer (Fig. 2.21). The cathode emits electrons in proportion to the amount of light striking it. These electrons are accelerated toward a target by the anode. The target is an electron multiplier covered by a small aperture which allows only a small part of the "electron image" emitted by the cathode to reach the target. The electron image is focused by a focusing coil that produces an axial magnetic field. The deflection coils then scan the electron image past the target aperture, where the electron multiplier produces a varying voltage representing the video signal. The image is thus "dissected" as it is scanned past the target, in an electronic version of a flat-bed digitizing process.

*Charge transfer devices.* A more recent development in image formation is that of solid-state image sensors, known as charge transfer devices (CTDs). There are two main classes of CTDs: charge-coupled devices (CCDs) and charge-injection devices (CIDs).

CCDs resemble MOSFETs (metal-oxide semiconductor field-effect transistor) in that they contain a "source" region and a "drain" region coupled by a depletion-region channel (Fig. 2.22). For imaging purposes, they can be considered as a monolithic array of closely spaced MOS capacitors forming a shift register (Fig. 2.23). Charges in the depletion region are transferred to the output by applying a series of clocking pulses to a row of electrodes between the source and the drain.

Photons incident on the semiconductor generate a series of charges on the CCD array. They are transferred to an output register either directly one line at a time (line transfer) or via a temporary storage area (frame transfer). The storage
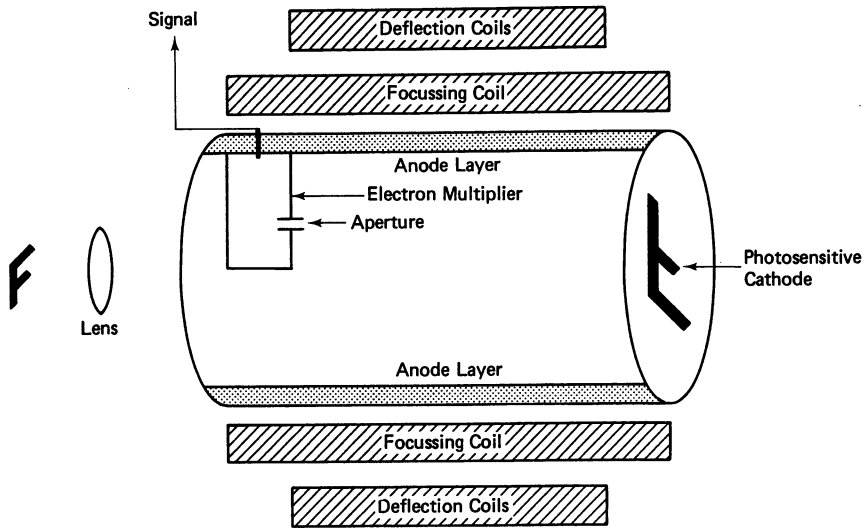
Fig. 2.21   Image dissector.

area is needed in frame transfer because the CCD array is scanned more rapidly than the output can be directly accommodated.

Charge injection devices (CIDs) resemble CCDs except that during sensing the charge is confined to the image site where it was generated (Fig. 2.24). The charges are read using an *X-Y* addressing technique similar to that used in computer memories. Basically, the stored charge is "injected" into the substrate and the resulting displacement current is detected to create the video signal.

CTD technology offers a number of advantages over conventional-tube-type cameras: light weight, small size, low power consumption, resistance to burn-in, low blooming, low dark current, high sensitivity, wide spectral and dynamic range, and lack of persistence. CIDs have the further advantages over CCDs of tolerance to processing defects, simple mechanization, avoidance of charge transfer losses, and minimized blooming. CTD cameras are now available commercially.

*Analog-to-Digital Conversion*

With current technology, the representation of an image as an analog electrical waveform is usually an unavoidable precursor to further processing. Thus the operation of deriving a digital representation of an analog voltage is basic to computer vision input devices.
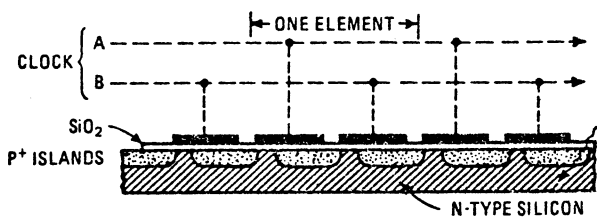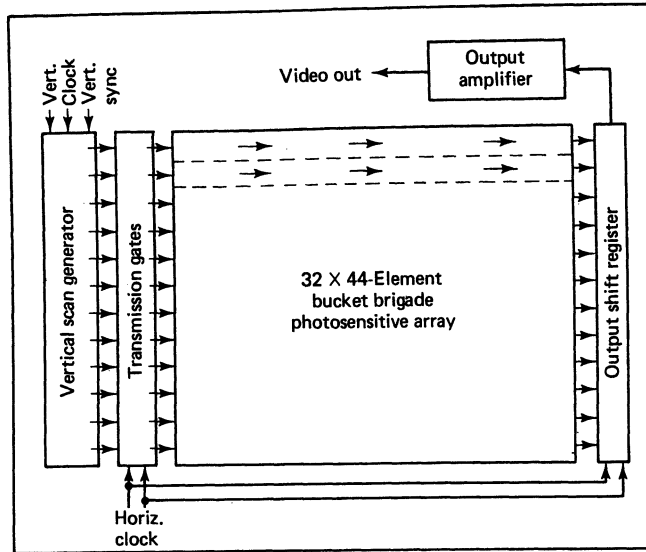


Fig. 2.22   Charge coupled device.

**Fig. 2.23** A CCD array (line transfer).

The function of an analog-to-digital (A/D) converter is to take as input a voltage such as a video signal and to produce as output a representation of the voltage in digital memory, suitable for reading by an interface to a digital computer. The quality of an A/D converter is measured by its temporal resolution (the speed at which it can perform conversions) and the accuracy of its digital output. Analog-
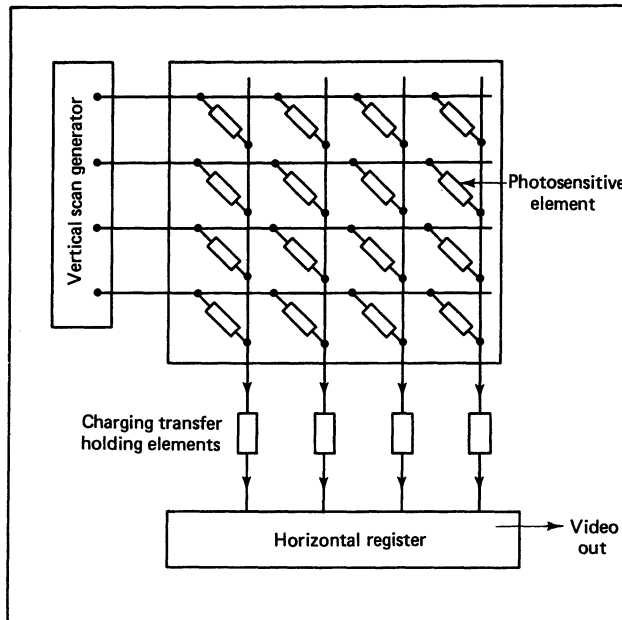


**Fig. 2.24** A CID array.

to-digital converters are being produced as integrated circuit chips, but high-quality models are still expensive. The output precision is usually in the 8- to 12-bit range.

It is quite possible to digitize an entire frame of a TV camera (i.e., approximately 525 scan lines by 300 or so samples along a scan line) in a single frame time (1/30 sec in the United States). Several commercial systems can provide such fast digitization into a "frame buffer" memory, along with raster graphics display capabilities from the same frame buffer, and "video rate processing" of the digital data. The latter term refers to any of various low-level operations (such as averaging, convolution with small templates, image subtraction) which may be performed as fast as the images are acquired.

One inexpensive alternative to digitizing entire TV frames at once is to use an interface that acquires the TV signal for a particular point when the scan passes the requested location. With efficient programming, this point-by-point digitization can acquire an entire frame in a few seconds.

### 2.3.2 Sensing Range

The third dimension may be derived from binocular images by triangulation, as we saw earlier, or inferred from single monocular visual input by a variety of "depth cues," such as size and occlusion. Specialized technology exists to acquire "depth images" directly and reliably. Here we outline two such techniques: "light striping," which is based on triangulation, and "spot ranging," which is based on different principles.

#### Light Striping

Light striping is a particularly simple case of the use of *structured light* [Will and Pennington 1971]. The basic idea is to use geometric information in the illumination to help extract geometric information from the scene. The spatial frequencies and angles of bars of light falling on a scene may be clustered to find faces; randomly structured light may allow blank, featureless surfaces to be matched in stereo views; and so forth.

Many researchers [Popplestone et al. 1975; Agin 1972; Sugihara 1977] have used striping to derive three dimensions. In light striping, a single plane of light is projected onto a scene, which causes a stripe of light to appear on the scene (Fig. 2.25). Only the part of the scene illuminated by the plane is sensed by the vision system. This restricts the "image" to be an essentially one-dimensional entity, and simplifies matching corresponding points. The plane itself has a known position (equation in world coordinates), determinable by any number of methods involving either the measurement of the projecting device or the measurement of the final resulting plane of light. Every image point determines a single "line of sight" in three-space upon which the world point that produces the image point must lie. This line is determined by the focal point of the imaging system and the image point upon which the world point projects. In a light-striping system, any point that is sensed in the image is also guaranteed to lie on the light plane in three-space. But the light plane and the line of sight intersect in just one point (as long as
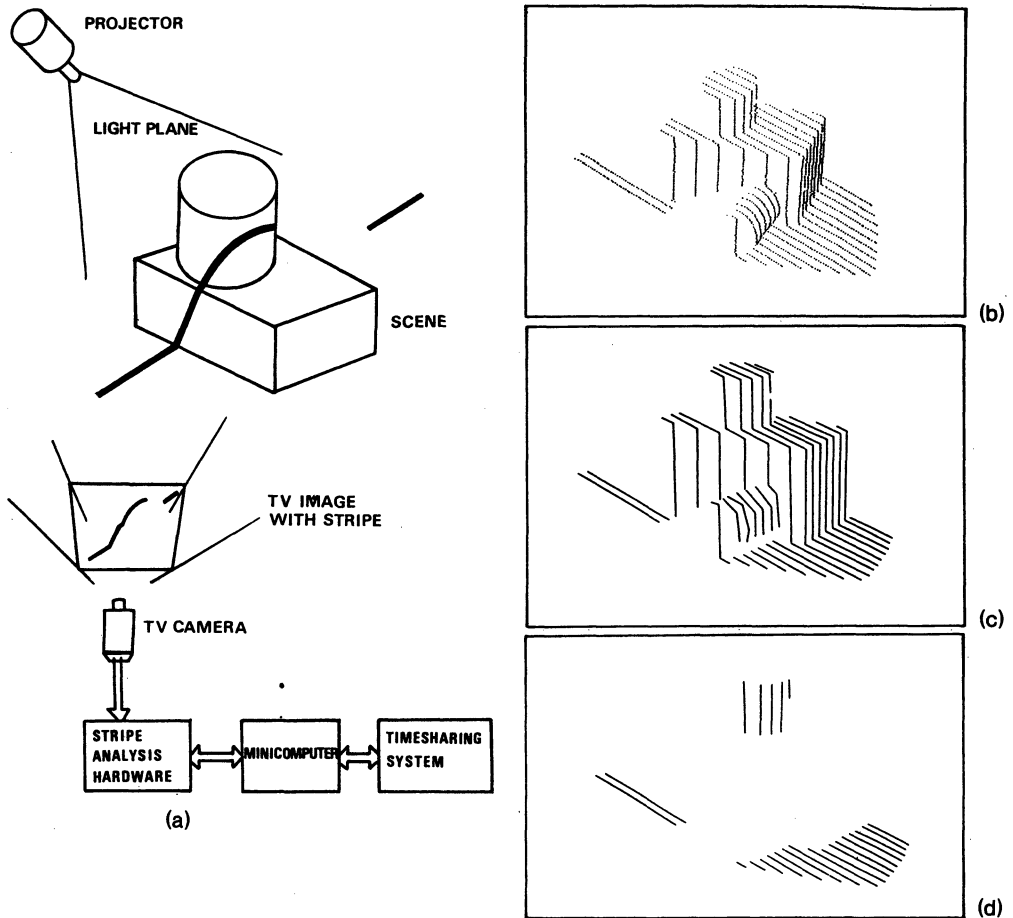
**Fig. 2.25** Light striping. (a) A typical arrangement; (b) raw data; (c) data segmented into strips; (d) strips segmented into two surfaces.

the camera's focal point is not in the light plane). Thus by computation of the intersection of the line of sight with the plane of light, we derive the three-dimensional point that corresponds to any image point visible as part of a stripe.

The plane of light may result from a laser or from the projection of a slit. Only the light stripe should be visible to the imaging device; unless a laser is used, this implies a darkened room. If a camera is fitted with the proper filter, a laser-based system can be operated in normal light. Another advantage of the laser is that it can be focused into a narrower plane than can a slit image.

The only points whose three-dimensional coordinates can be computed are those that can be "seen" by both the light-stripe source and the camera at once. Since there must be a nonzero baseline if triangulation is to derive three-dimensional information, the camera cannot be too close to the projector, and thus concavities in the scene are potential trouble spots, since both the striper and the

camera may not be able to "see" into them. Surfaces in the scene that are nearly parallel with the light plane will have a relatively small number of stripes projected onto them by any uniform stripe placement strategy. This problem is ameliorated by striping with two sets of parallel planes at right angles to each other [Agin 1972]. A major advantage of light striping over spot ranging is that (barring shadows) its continuity and discontinuity indicate similar conditions on the surface. It is easy to "segment" stripe images (Part II): Stripes falling on the same surface may easily be gathered together. This set of related stripes may be used in a number of ways to derive further information on the characteristics of the surface (Fig. 2.25b).
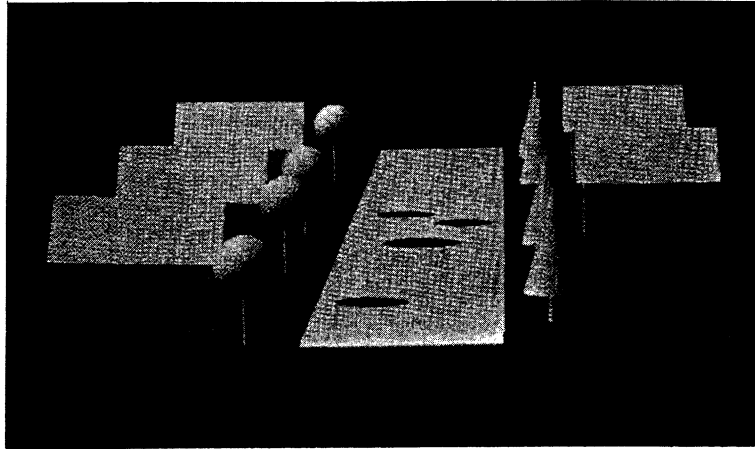
### Spot Ranging

Civil engineers have used laser-based "spot range finders" for some time. In laboratory-size environments, they are a relatively new development. There are two basic techniques. First, one can emit a very sharp pulse and time its return ("lidar," the light equivalent of radar). This requires a sophisticated laser and electronics, since light moves 1 ft every billionth of a second, approximately. The second technique is to modulate the laser light in amplitude and upon its return compare the phase of the returning light with that of the modulator. The phase differences are related to the distance traveled [Nitzan et al. 1977]. A representative image is shown in Fig. 2.26.

Both these techniques produce results that are accurate to within about 1% of the range. Both of them allow the laser to be placed close to a camera, and thus "intensity maps" (images) and range maps may be produced from single viewpoints. The laser beam can easily poke into holes, and the return beam may be sensed close to the emitted one, so concavities do not present a serious problem. Since the laser beam is attenuated by absorption, it can yield intensity information as well. If the laser produces light of several wavelengths, it is possible to use filters and obtain multispectral reflectance information as well as depth information from the same device [Garvey 1976; Nitzan et al. 1977].
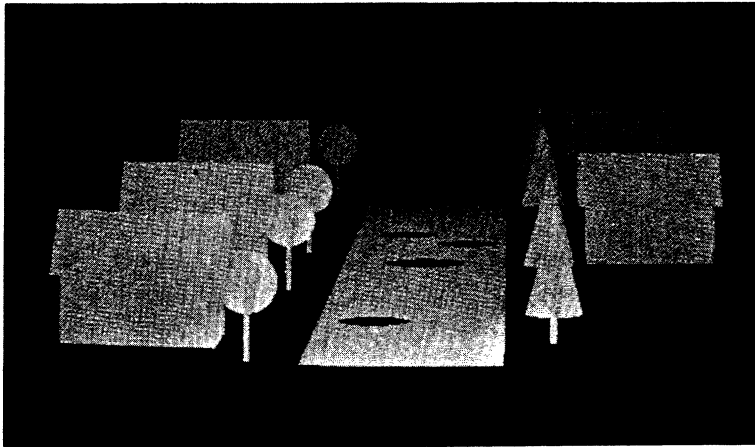
The usual mode of use of a spot ranging device is to produce a range map that corresponds to an intensity map. This has its advantages in that the correspondence may be close. The structural properties of light stripes are lost: It can be hard to "segment" the image into surfaces (to tell which "range pixels" are associated with the same surface). Range maps are amenable to the same sorts of segmentation techniques that are used for intensity images: Hough techniques, region growing, or differentiation-based methods of edge finding (Part II).

### Ultrasonic Ranging

Just as light can be pulsed to determine range, so can sound and ultrasound (frequencies much higher than the audible range). Ultrasound has been used extensively in medicine to produce images of human organs (e.g., [Waag and Gramiak 1976]). The time between the transmitted and received signal determines range; the sound signal travels much slower than light, making the problem of timing the returning signal rather easier than it is in pulsed laser devices. However, the signal is severely attenuated as it travels through biological tissue, so that the detection apparatus must be very sensitive.

(a)



(b)

**Fig. 2.26** Intensity and range images. (a) A (synthesized) intensity image of a
street scene with potholes. The roofs all have the same intensity, which is different
from the walls; (b) a corresponding range image. The wall and roof of each house
have similar ranges, but the ranges differ from house to house.

One basic difference between sound and visible light ranging is that a light
beam is usually reflected off just one surface, but that a sound beam is generally
partially transmitted and partially reflected by "surfaces." The returning sound
pulse has structure determined by the discontinuities in impedence to sound found
in the medium through which it has passed. Roughly, a light beam returns infor-
mation about a spot, whereas a sound beam can return information about the
medium in the entire column of material. Thus, although sound itself travels rela-
tively slowly, the data rate implicit in the returning structured sound pulse is quite
high. Figure 2.27 shows an image made using the range data from ultrasound. The
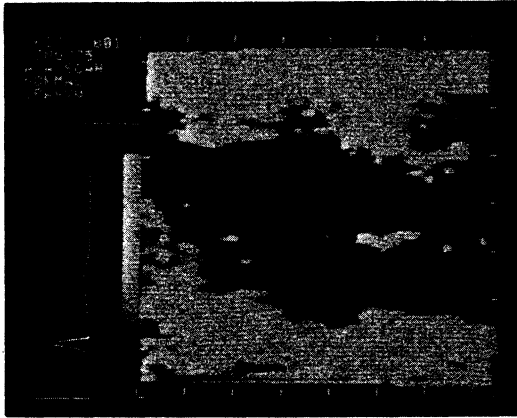
**Fig. 2.27** Image made from ultrasound ranging.

sound pulses emanate from the top of the image and proceed toward the bottom, being partially reflected and transmitted along the way. In the figure, it is as if we were looking perpendicular to the beams, which are being displayed as brighter where strong reflectance is taking place. A single "scan line" of sound thus produces an image of an entire planar slice of medium.

### 2.3.3 Reconstruction Imaging

Two-dimensional reconstruction has been the focus of much research attention because of its important medical applications. High-quality images such as that shown in Fig. 1.2b can be formed by multiple images of x-ray projection data. This section contains the principles behind the most important reconstruction algorithms. These techniques are discussed in more detail with an expanded list of references in [Gordon and Herman 1974]. For a view of the many applications of two-dimensional reconstruction other than transmission scanning, the reader is referred to [Gordon et al. 1975].

Figure 2.28 shows the basic geometry to collect one-dimensional projections of two-dimensional data. (Most systems construct the image in a plane and repeat this technique for other planes; there are few true three-dimensional reconstruction systems that use planes of projection data simultaneously to construct volumes.)

In many applications sensors can measure the one-dimensional *projection* of two-dimensional image data. The projection $g(x')$ of an ideal image $f(x, y)$ in the direction $\theta$ is given by $\int f(x', y') \, dy'$ where $x' = R_\theta x$. If enough different projections are obtained, a good approximation to the image can be obtained with two-dimensional reconstruction techniques.

From Fig. 2.28, with the source at the first position along line $AA'$, we can obtain the first projection datum from the detector at the first position along $BB'$. The line $AB$ is termed a ray and the measurement at $B$ a ray sum. Moving the source
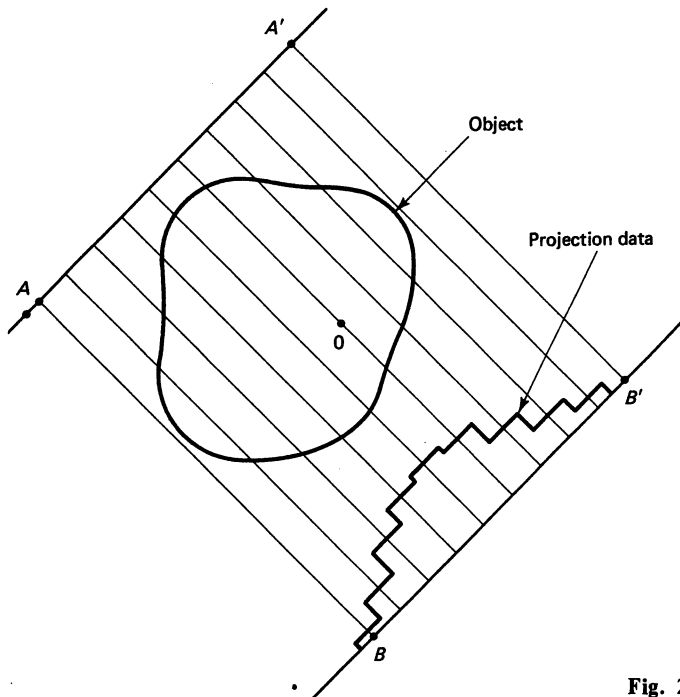
**Fig. 2.28** Projection geometry.

and detector along lines $AA'$ and $BB'$ in synchrony allows us to obtain the entire data for projection 1. Now the lines $AA'$ and $BB'$ are rotated by a small angle $d\theta$ about 0 and the process is repeated. In the original x-ray systems $d\theta$ was $1°$ of angle, and 180 projections were taken. Each projection comprised 160 transmission measurements. The reconstruction problem is simply this: Given the projection data $g_k(x')$, $k = 0, \ldots, N - 1$, construct the original image $f(\mathbf{x})$.

Systems in use today use a fan beam rather than the parallel rays shown. However, the mathematics is simpler for parallel rays and illustrates the fundamental ideas. We describe three related techniques: summation, Fourier interpolation, and convolution.

### The Summation Method

The summation method is simple: Distribute every ray sum $g_k(x')$ over the image cells along the ray. Where there are $N$ cells along a ray, each such cell is incremented by $\frac{1}{N} g(x')$. This step is termed *back projection*. Repeating this process for every ray results in an approximate version of the original [DeRosier 1971]. This technique is equivalent (within a scale factor) to blurring the image, or convolving it with a certain point-spread function. In the continuous case of infinitely many projections, this function is simply the radically symmetric $h(r) = 1/r$.
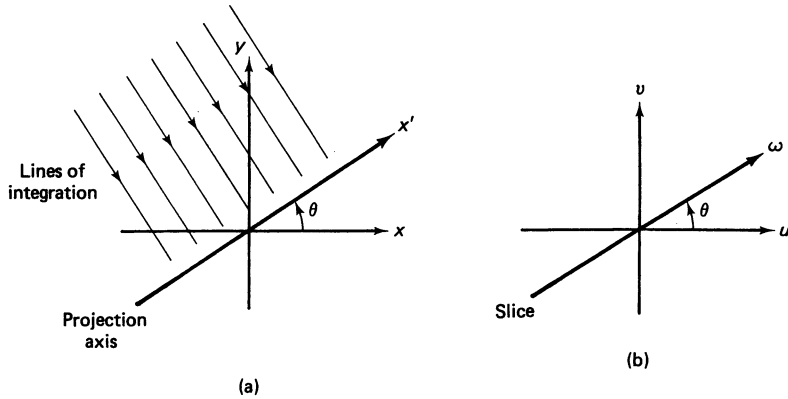
**Fig. 2.29** Basis of Fourier techniques. (a) Projection axis x'; (b) corresponding axis in Fourier Space.

*Fourier Algorithms*

If a projection is Fourier-transformed, it defines a line through the origin in frequency space (Fig. 2.29). To show this formally, consider the expression for the two-dimensional transform

$$F(\mathbf{u}) = \int\int f(x, y) \exp\left[-j2\pi(ux + vy)\right] dx\, dy \qquad (2.47)$$

Now consider $y = 0$ (projection onto the $x$ axis): $x' = x$ and

$$g_0(x') = \int f(x, y)\, dy \qquad (2.48)$$

The Fourier transform of this equation is

$$\mathcal{F}[g_0(x')] = \int\int [f(x, y)\, dy] \exp{-j2\pi ux}\, dx \qquad (2.49)$$

$$= \int\int f(x, y) \exp{-j2\pi ux}\, dy\, dx$$

which, by comparison with (2.47), is

$$\mathcal{F}[g_0(x')] = F(u, 0) \qquad (2.50)$$

Generalizing to any $\theta$, the transform of an arbitrary $g(x')$ defines a line in the Fourier space representation of the cross section. Where $S_k(\omega)$ is the cross section of the Fourier transform along this line,

$$S_k(\omega) = F(u\cos\theta,\ u\sin\theta) \qquad (2.51)$$

$$= \int g_k(x') \exp\left[j2\pi u(x')\right] dx'$$

Thus one way of reconstructing the original image is to use the Fourier transform of the projections to define points in the transform of $f(x)$, interpolate the undefined points of the transform from the known points, and finally take the inverse transform to obtain the reconstructed image.
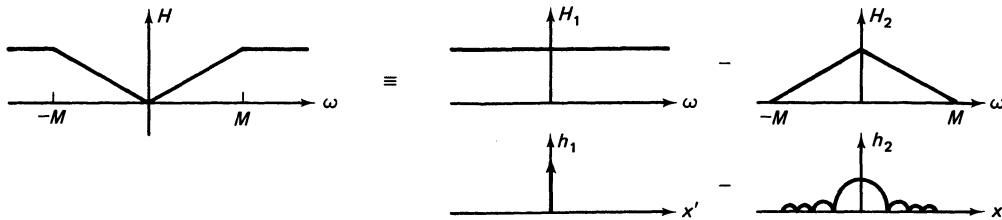
*Ch. 2   Image Formation*

Fig. 2.30   Convolution method.

This technique can be applied with transforms other than the Fourier transform, and such methods are discussed in [DeRosier 1971; Crowther and Klug 1971].

### The Convolution Method

The convolution method is the natural extension of the summation method. Since the summation method produces an image degraded from its convolution with some function $h$, one can remove the degradation by a "deconvolution." The straightforward way to accomplish this is to Fourier-transform the degraded image, multiply the result by an estimate of the transformed $h^{-1}$, and inverse-Fourier-transform the result. However, since all the operations are linear, a faster approach is to deconvolve the projections before performing the back projection. To show this formally, we use the inverse transform

$$f(\mathbf{x}) = \iint F(u, v) \exp\left[j2\pi(ux + vy)\right] du\, dv \tag{2.52}$$

Changing to cylindrical coordinates $(\omega, \theta)$ yields

$$f(\mathbf{x}) = \iint F_\theta(\omega) \exp\left[j2\pi\omega(x\cos\theta + y\sin\theta)\right] |\omega|\, d\omega\, d\theta \tag{2.53}$$

Since $x' = x\cos\theta + y\sin\theta$, rewrite Eq. (2.53) as

$$f(\mathbf{x}) = \int \mathcal{F}^{-1}\{F_\theta(\omega)H(\omega)\} d\theta \tag{2.54}$$

Since the image is bandlimited at some interval $(-\omega_m, \omega_m)$ one can define $H(\omega)$ arbitrarily outside of this interval. Therefore, $H(\omega)$ can be defined as a constant minus a triangular peak as shown in Fig. 2.30. Finally, the operation inside the integral in Eq. (2.54) is a convolution. Using the transforms shown in Fig. 2.30,
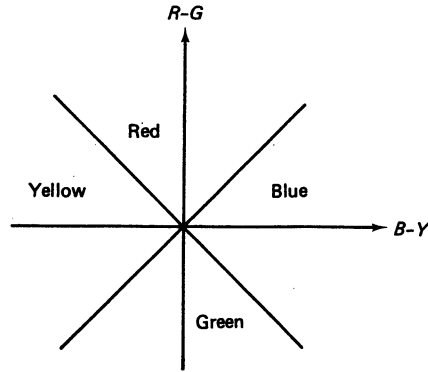
$$f(\mathbf{x}) = \int \left[f_\theta(x') - f_\theta(x')\omega_m \mathrm{sinc}^2(\omega_m x')\right] d\theta \tag{2.55}$$

Owing to its speed and the fact that the deconvolutions can be performed while the data are being acquired, the convolution method is the method employed in the majority of systems.

### EXERCISES

**2.1**  In a binocular animal vision system, assume a focal length $f$ of an eye of 50 mm and a separation distance $d$ of 5 cm. Make a plot of $\Delta x$ vs. $-z$ using Eq. (2.9). If the resolution of each eye is on the order of 50 line pairs/mm, what is the useful range of the binocular system?

**2.2** In an opponent-process color vision system, assume that the following relations hold:



For example, if the $(R-G, B -Y, W-Bk)$ components of the opponent-process system are $(0.5, 3, 4)$, the perceived color will be blue.

Work out the perceived colors for the following (R,G,B) measurements:

(a) $(0.2, 0.3, 0.4)$     (b) $(0.2, 0.3, 0)$     (c) $(7, 4, 1)$

**2.3** Develop an indexing scheme for a hexagonal array and define a Euclidean distance measure between points in the array.

**2.4** Assume that a one-dimensional image has the following form:

$$f(x) = \cos(2\pi u_o x)$$

and is sampled with $u_s = u_o$. Using the graphical method of Section 2.2.6, find an expression for $f(x)$ as given by Eq. (2.49). Is this expression equal to the original image? Explain.

**2.5** A certain image has the following Fourier transform:

$$F(\mathbf{u}) = \begin{cases} \text{nonzero} & \text{inside a hexagonal domain} \\ 0 & \text{otherwise} \end{cases}$$

(a) What are the smallest values for $u$ and $v$ so that $F(\mathbf{u})$ can be reconstructed from $F_x(\mathbf{u})$?

(b) Suppose now that rectangular sampling is *not* used but that now the $u$ and $v$ directions subtend an angle of $\pi/3$. Does this change your answer as to the smallest $u$ and $v$? Explain.

**2.6** Extend the binocular imaging model of Fig. 2.3 to include convergence: Let the two imaging systems pivot in the $y = 0$ plane about the viewpoint. Let the system have a baseline of $2d$ and be converged at some angle $\theta$ such that a point $(x, y, z)$ appears at the origin of each image plane.

(a) Solve for $z$ in terms of $r$ and $\theta$.

(b) Solve for $z$ in this situation for points with nonzero disparity.

**2.7** Compute the convolution of two Rect functions, where

$$\text{Rect}(x) = \begin{cases} 1 & 0 < x < 1 \\ 0 & \text{otherwise} \end{cases}$$

Show the steps in your calculations.

**2.8**

$$\text{Rect}(x) = \begin{cases} b & \text{for } |x| < a \\ 0 & \text{otherwise} \end{cases}$$

(a)   What is $\text{Rect}(x) * \delta(x-a)$?

(b)   What is the Fourier transform of $f(x)$ where $f(x) = \text{Rect}(x+c) + \text{Rect}(x-c)$ and $c > a$?

**2.9**   A digitizer has a sampling interval of $\Delta x = \Delta y = \Delta$. Which of the following images can be represented unambiguously by their samples? (Assume that effects of a finite image domain can be neglected.)

(a)   $(sin\,(\pi x/\Delta))/(\pi x/\Delta)$

(b)   $cos\,(\pi x/2\Delta)\cos(3\pi x/4\Delta)$

(c)   $\text{Rect}(x)$ (see Problem 2.8)

(d)   $e^{-ax^2}$

# REFERENCES

AGIN, G. J. "Representation and description of curved objects" (Ph.D. dissertation). AIM-173, Stanford AI Lab, October 1972.

ANDREWS, H. C. and B. R. HUNT. *Digital Image Restoration.* Englewood Cliffs, NJ: Prentice-Hall, Inc., 1977.

CROWTHER, R. A. and A. KLUG. "ART and science, or, conditions for 3-d reconstruction from electron microscope images." *J. Theoretical Biology 32,* 1971.

DEROSIER, D. J. "The reconstruction of three-dimensional images from electron micrographs." *Contemporary Physics 12,* 1971.

DUDA, R. O. and P. E. HART. *Pattern Recognition and Scene Analysis.* New York: Wiley, 1973.

GARVEY, T. D. "Perceptual strategies for purposive vision." Technical Note 117, AI Center, SRI International, September 1976.

GONZALEZ, R. C. and P. WINTZ. *Digital Image Processing.* Reading, MA: Addison-Wesley, 1977.

GORDON, R. and G. T. HERMAN. "Three-dimensional reconstruction from projections: a review of algorithms." *International Review of Cytology 38,* 1974, 111–151.

GORDON, R., G. T. HERMAN, and S. A. JOHNSON. "Image reconstruction from projections." *Scientific American,* October 1975.

HERING, E. "Principles of a new theory of color sense." *In Color Vision,* R.C. Teevan and R.C. Birney (Eds.). Princeton, NJ: D. Van Nostrand, 1961.

HORN, B. K. P. "Understanding image intensities." *Artificial Intelligence 8,* 2, April 1977, 201–231.

HORN, B. K. P. and R. W. SJOBERG. "Calculating the reflectance map." *Proc.,* DARPA IU Workshop, November 1978, 115–126.

HURVICH, L. M. and D. JAMESON. "An opponent-process theory of color vision." *Psychological Review 64,* 1957, 384–390.

JAIN, A. K. "Advances in mathematical models for image processing." *Proc. IEEE 69,* 5, May 1981, 502-528.

JOBLOVE, G. H. and D. GREENBERG. "Color spaces for computer graphics." *Computer Graphics 12,* 3, August 1978, 20–25.

KENDER, J. R. "Saturation, hue, and normalized color: calculation, digitization effects, and use." Technical Report, Dept. of Computer Science, Carnegie-Mellon Univ., November 1976.

LAND, E. H. "The retinex theory of color vision." *Scientific American,* December 1977, 108–128.

MUNSELL, A. H. *A Color Notation,* 8th ed. Baltimore, MD: Munsell Color Co., 1939.

NICODEMUS, F. E., J. C. RICHMOND, J. J. HSIA, I. W. GINSBERG, and T. LIMPERIS. "Geometrical considerations and nomenclature for reflectance." NBS Monograph 160, National Bureau of Standards, U.S. Department of Commerce, Washington, DC, October 1977.

NITZAN, D., A. BRAIN, and R. DUDA. "The measurement and use of registered reflectance and range data in scene analysis." *Proc. IEEE 65,* 2, February 1977.

POPPLESTONE, R. J., C. M. BROWN, A. P. AMBLER, and G. F. CRAWFORD. "Forming models of plane-and-cylinder faceted bodies from light stripes." *Proc.,* 4th IJCAI, September 1975, 664-668.

PRATT, W. K. *Digital Image Processing.* New York: Wiley-Interscience, 1978.

ROSENFIELD A. and A. C. KAK. *Digital Picture Processing.* New York: Academic Press, 1976.

SMITH, A. R. "Color gamut transform pairs." *Computer Graphics 12,* 3, August 1978, 12–19.

SUGIHARA, K. "Dictionary-guided scene analysis based on depth information." In *Progress Report on 3-D Object Recognition.* Bionics Research Section, ETL, Tokyo, March 1977.

TENENBAUM, J. M. and S. WEYL. "A region-analysis subsystem for interactive scene analysis." *Proc.,* 4th IJCAI, September 1975, 682–687.

WAAG, R. B. and R. GRAMIAK. "Methods for ultrasonic imaging of the heart." *Ultrasound in Medicine and Biology 2,* 1976, 163–170.

WILL, P. M. and K. S. PENNINGTON. "Grid coding: a preprocessing technique for robot and machine vision." *Artificial Intelligence 2,* 3/4, Winter 1971, 319–329.