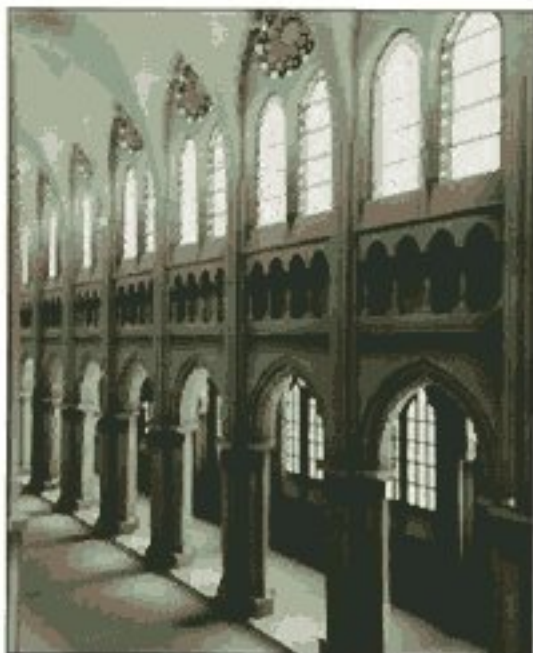


COMPUTER GRAPHICS

C VERSION



DONALD HEARN ■ M. PAULINE BAKER

SECOND EDITION

Contents

PREFACE	xvii	Stereoscopic and Virtual-Reality Systems	50
1 A Survey of Computer Graphics	2	2-2 Raster-Scan Systems	53
		Video Controller	53
		Raster-Scan Display Processor	55
1-1 Computer-Aided Design	4	2-3 Random-Scan Systems	56
1-2 Presentation Graphics	11	2-4 Graphics Monitors and Workstations	57
1-3 Computer Art	13	2-5 Input Devices	60
1-4 Entertainment	18	Keyboards	61
1-5 Education and Training	21	Mouse	61
1-6 Visualization	25	Trackball and Spaceball	63
1-7 Image Processing	32	Joysticks	63
1-8 Graphical User Interfaces	34	Data Glove	64
		Digitizers	64
		Image Scanners	67
		Touch Panels	68
		Light Pens	70
		Voice Systems	70
		2-6 Hard-Copy Devices	72
		2-7 Graphics Software	75
		Coordinate Representations	76
		Graphics Functions	77
		Software Standards	78
		PHIGS Workstations	79
		Summary	79
		References	81
		Exercises	81
2 Overview of Graphics Systems	35		
2-1 Video Display Devices	36		
Refresh Cathode-Ray Tubes	37		
Raster-Scan Displays	40		
Random-Scan Displays	41		
Color CRT Monitors	42		
Direct-View Storage Tubes	45		
Flat-Panel Displays	45		
Three-Dimensional Viewing Devices	49		

3 Output Primitives 83

3-1	Points and Lines	84
3-2	Line-Drawing Algorithms	86
	DDA Algorithm	87
	Bresenham's Line Algorithm	88
	Parallel Line Algorithms	92
3-3	Loading the Frame Buffer	94
3-4	Line Function	95
3-5	Circle-Generating Algorithms	97
	Properties of Circles	97
	Midpoint Circle Algorithm	98
3-6	Ellipse-Generating Algorithms	102
	Properties of Ellipses	102
	Midpoint Ellipse Algorithm	103
3-7	Other Curves	110
	Conic Sections	110
	Polynomials and Spline Curves	112
3-8	Parallel Curve Algorithms	112
3-9	Curve Functions	113
3-10	Pixel Addressing	
	and Object Geometry	114
	Screen Grid Coordinates	114
	Maintaining Geometric Properties	
	of Displayed Objects	114
3-11	Filled-Area Primitives	117
	Scan-Line Polygon Fill Algorithm	117
	Inside-Outside Tests	125
	Scan-Line Fill of Curved Boundary	
	Areas	126
	Boundary-Fill Algorithm	127
	Flood-Fill Algorithm	130
3-12	Fill-Area Functions	131
3-13	Cell Array	131
3-14	Character Generation	131

Summary	134
Applications	136
References	140
Exercises	140

4 Attributes of Output Primitives 143

4-1	Line Attributes	144
	Line Type	144
	Line Width	146
	Pen and Brush Options	149
	Line Color	149
4-2	Curve Attributes	152
4-3	Color and Grayscale Levels	154
	Color Tables	155
	Grayscale	157
4-4	Area-Fill Attributes	158
	Fill Styles	158
	Pattern Fill	159
	Soft Fill	162
4-5	Character Attributes	163
	Text Attributes	163
	Marker Attributes	167
4-6	Bundled Attributes	168
	Bundled Line Attributes	168
	Bundled Area-Fill Attributes	169
	Bundled Text Attributes	169
	Bundled Marker Attributes	170
4-7	Inquiry Functions	170
4-8	Antialiasing	171
	Supersampling Straight Line	
	Segments	172
	Pixel-Weighting Masks	174

Area Sampling Straight Line Segments	174
Filtering Techniques	174
Pixel Phasing	175
Compensating for Line Intensity Differences	175
Antialiasing Area Boundaries	176
Summary	178
References	180
Exercises	180

5-6 Affine Transformations	208
5-7 Transformation Functions	208
5-8 Raster Methods for Transformations	210
Summary	212
References	213
Exercises	213

5 Two-Dimensional Geometric Transformations 183

5-1 Basic Transformations	184
Translation	184
Rotation	186
Scaling	187
5-2 Matrix Representations and Homogeneous Coordinates	188
5-3 Composite Transformations	191
Translations	191
Rotations	191
Scalings	192
General Pivot-Point Rotation	192
General Fixed-Point Scaling	193
General Scaling Directions	193
Concatenation Properties	194
General Composite Transformations and Computational Efficiency	195
5-4 Other Transformations	201
Reflection	201
Shear	203
5-5 Transformations Between Coordinate Systems	205

6 Two-Dimensional Viewing 216

6-1 The Viewing Pipeline	217
6-2 Viewing Coordinate Reference Frame	219
6-3 Window-to-Viewport Coordinate Transformation	220
6-4 Two-Dimensional Viewing Functions	222
6-5 Clipping Operations	224
6-6 Point Clipping	225
6-7 Line Clipping	225
Cohen-Sutherland Line Clipping	226
Liang-Barsky Line Clipping	230
Nicholl-Lee-Nicholl Line Clipping	233
Line Clipping Using Nonrectangular Clip Windows	235
Splitting Concave Polygons	235
6-8 Polygon Clipping	237
Sutherland-Hodgeman Polygon Clipping	238
Weiler-Atherton Polygon Clipping	242
Other Polygon-Clipping Algorithms	243
6-9 Curve Clipping	244
6-10 Text Clipping	244
6-11 Exterior Clipping	245
Summary	245
References	248
Exercises	248

7 Structures and Hierarchical Modeling 250

7-1	Structure Concepts	250
	Basic Structure Functions	250
	Setting Structure Attributes	253
7-2	Editing Structures	254
	Structure Lists and the Element Pointer	255
	Setting the Edit Mode	256
	Inserting Structure Elements	256
	Replacing Structure Elements	257
	Deleting Structure Elements	257
	Labeling Structure Elements	258
	Copying Elements from One Structure to Another	260
7-3	Basic Modeling Concepts	260
	Model Representations	261
	Symbol Hierarchies	262
	Modeling Packages	263
7-4	Hierarchical Modeling with Structures	265
	Local Coordinates and Modeling Transformations	265
	Modeling Transformations	266
	Structure Hierarchies	266
	Summary	268
	References	269
	Exercises	269

8 Graphical User Interfaces and Interactive Input Methods 271

8-1	The User Dialogue	272
	Windows and Icons	273

	Accommodating Multiple Skill Levels	273
	Consistency	274
	Minimizing Memorization	274
	Backup and Error Handling	274
	Feedback	275
8-2	Input of Graphical Data	276
	Logical Classification of Input Devices	276
	Locator Devices	277
	Stroke Devices	277
	String Devices	277
	Valuator Devices	277
	Choice Devices	279
	Pick Devices	279
8-3	Input Functions	281
	Input Modes	281
	Request Mode	282
	Locator and Stroke Input in Request Mode	282
	String Input in Request Mode	283
	Valuator Input in Request Mode	284
	Choice Input in Request Mode	284
	Pick Input in Request Mode	284
	Sample Mode	285
	Event Mode	285
	Concurrent Use of Input Modes	287
8-4	Initial Values for Input-Device Parameters	287
8-5	Interactive Picture-Construction Techniques	288
	Basic Positioning Methods	288
	Constraints	288
	Grids	289
	Gravity Field	290
	Rubber-Band Methods	290
	Dragging	291
	Painting and Drawing	291

8-6	Virtual-Reality Environments	292	10-4	Superquadrics	312
	Summary	293		Superellipse	312
	References	294		Superellipsoid	313
	Exercises	294	10-5	Blobby Objects	314
			10-6	Spline Representations	315
9	Three-Dimensional Concepts	296		Interpolation and Approximation Splines	316
9-1	Three-Dimensional Display Methods	297		Parametric Continuity Conditions	317
	Parallel Projection	298		Geometric Continuity Conditions	318
	Perspective Projection	299	10-7	Spline Specifications	319
	Depth Cueing	299		Cubic Spline Interpolation Methods	320
	Visible Line and Surface Identification	300		Natural Cubic Splines	321
	Surface Rendering	300		Hermite Interpolation	322
	Exploded and Cutaway Views	300		Cardinal Splines	323
	Three-Dimensional and Stereoscopic Views	300		Kochanek-Bartels Splines	325
9-2	Three-Dimensional Graphics Packages	302	10-8	Bézier Curves and Surfaces	327
				Bézier Curves	327
				Properties of Bézier Curves	329
				Design Techniques Using Bézier Curves	330
				Cubic Bézier Curves	331
				Bézier Surfaces	333
			10-9	B-Spline Curves and Surfaces	334
				B-Spline Curves	335
				Uniform, Periodic B-Splines	336
				Cubic, Periodic B-Splines	339
				Open, Uniform B-Splines	341
				Nonuniform B-Splines	344
				B-Spline Surfaces	344
			10-10	Beta-Splines	345
				Beta-Spline Continuity Conditions	345
				Cubic, Periodic Beta-Spline Matrix Representation	346
			10-11	Rational Splines	347
10	Three-Dimensional Object Representations	304			
10-1	Polygon Surfaces	305			
	Polygon Tables	306			
	Plane Equations	307			
	Polygon Meshes	309			
10-2	Curved Lines and Surfaces	310			
10-3	Quadric Surfaces	310			
	Sphere	311			
	Ellipsoid	311			
	Torus	311			

12-1	Viewing Pipeline	432
12-2	Viewing Coordinates	433
	Specifying the View Plane	433
	Transformation from World to Viewing Coordinates	437

12-3	Projections	438
	Parallel Projections	439
	Perspective Projections	443
12-4	View Volumes and General Projection Transformations	447
	General Parallel-Projection Transformations	452
	General Perspective-Projection Transformations	454
12-5	Clipping	456
	Normalized View Volumes	458
	Viewport Clipping	460
	Clipping in Homogeneous Coordinates	461
12-6	Hardware Implementations	463
12-7	Three-Dimensional Viewing Functions	464
	Summary	467
	References	468
	Exercises	468

13 Visible-Surface Detection Methods 469

13-1	Classification of Visible-Surface Detection Algorithms	470
13-2	Back-Face Detection	471
13-3	Depth-Buffer Method	472
13-4	A-Buffer Method	475
13-5	Scan-Line Method	476
13-6	Depth-Sorting Method	478
13-7	BSP-Tree Method	481
13-8	Area-Subdivision Method	482
13-9	Octree Methods	485
13-10	Ray-Casting Method	487
13-11	Curved Surfaces	488
	Curved-Surface Representations	488
	Surface Contour Plots	489

13-12	Wireframe Methods	490
13-13	Visibility-Detection Functions	490
	Summary	491
	References	492
	Exercises	492

14 Illumination Models and Surface-Rendering Methods 494

14-1	Light Sources	496
14-2	Basic Illumination Models	497
	Ambient Light	497
	Diffuse Reflection	497
	Specular Reflection and the Phong Model	500
	Combined Diffuse and Specular Reflections with Multiple Light Sources	504
	Warn Model	504
	Intensity Attenuation	505
	Color Considerations	507
	Transparency	508
	Shadows	511
14-3	Displaying Light Intensities	511
	Assigning Intensity Levels	512
	Gamma Correction and Video Lookup Tables	513
	Displaying Continuous-Tone Images	515
14-4	Halftone Patterns and Dithering Techniques	516
	Halftone Approximations	516
	Dithering Techniques	519
14-5	Polygon-Rendering Methods	522
	Constant-Intensity Shading	522
	Gouraud Shading	523
	Phong Shading	525

	Fast Phong Shading	526	15-6	CMY Color Model	574
14-6	Ray-Tracing Methods	527	15-7	HSV Color Model	575
	Basic Ray-Tracing Algorithm	528	15-8	Conversion Between HSV and RGB Models	578
	Ray-Surface Intersection Calculations	531	15-9	HLS Color Model	579
	Reducing Object-Intersection Calculations	535	15-10	Color Selection and Applications	580
	Space-Subdivision Methods	535		Summary	581
	Antialiased Ray Tracing	538		References	581
	Distributed Ray Tracing	540		Exercises	582
14-7	Radiosity Lighting Model	544			
	Basic Radiosity Model	544			
	Progressive Refinement Radiosity Method	549			
14-8	Environment Mapping	552			
14-9	Adding Surface Detail	553			
	Modeling Surface Detail with Polygons	553			
	Texture Mapping	554			
	Procedural Texturing Methods	556			
	Bump Mapping	558			
	Frame Mapping	559			
	Summary	560			
	References	561			
	Exercises	562			

15 Color Models and Color Applications 564

15-1	Properties of Light	565
15-2	Standard Primaries and the Chromaticity Diagram	568
	XYZ Color Model	569
	CIE Chromaticity Diagram	569
15-3	Intuitive Color Concepts	571
15-4	RGB Color Model	572
15-5	YIQ Color Model	574

16 Computer Animation 583

16-1	Design of Animation Sequences	584
16-2	General Computer-Animation Functions	586
16-3	Raster Animations	586
16-4	Computer-Animation Languages	587
16-5	Key-Frame Systems	588
	Morphing	588
	Simulating Accelerations	591
16-6	Motion Specifications	594
	Direct Motion Specification	594
	Goal-Directed Systems	595
	Kinematics and Dynamics	595
	Summary	596
	References	597
	Exercises	597

A Mathematics for Computer Graphics 599

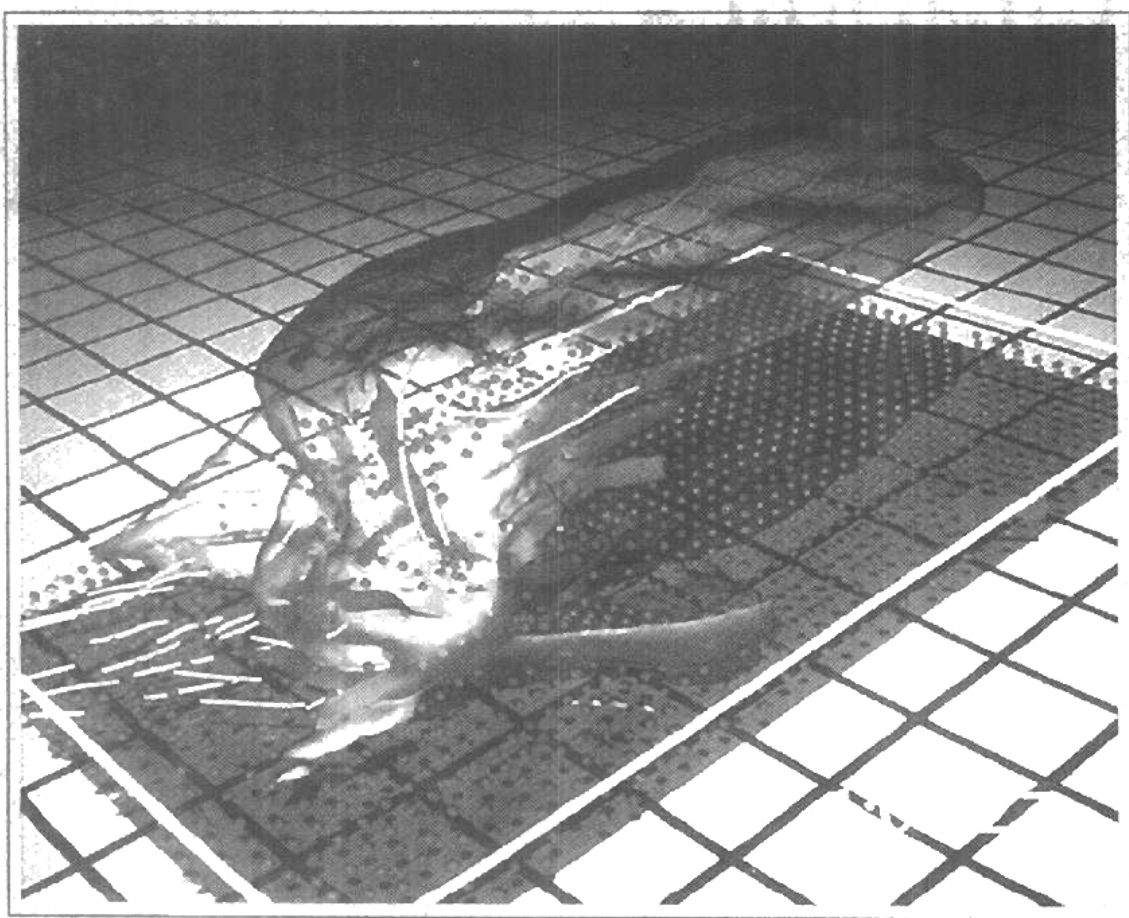
A-1	Coordinate-Reference Frames	600
	Two-Dimensional Cartesian Reference Frames	600
	Polar Coordinates in the xy Plane	601

	Three-Dimensional Cartesian Reference Frames	602		Matrix Transpose	613
	Three-Dimensional Curvilinear Coordinate Systems	602		Determinant of a Matrix	613
	Solid Angle	604	A-5	Matrix Inverse	614
A-2	Points and Vectors	605		Complex Numbers	615
	Vector Addition and Scalar Multiplication	607	A-6	Quaternions	617
	Scalar Product of Two Vectors	607	A-7	Nonparametric Representations	618
	Vector Product of Two Vectors	608	A-8	Parametric Representations	619
A-3	Basis Vectors and the Metric Tensor	609	A-9	Numerical Methods	620
	Orthonormal Basis	609		Solving Sets of Linear Equations	620
	Metric Tensor	610		Finding Roots of Nonlinear Equations	621
A-4	Matrices	611		Evaluating Integrals	622
	Scalar Multiplication and Matrix Addition	612		Fitting Curves to Data Sets	625
	Matrix Multiplication	612		BIBLIOGRAPHY	626
				INDEX	639

Computer Graphics C Version

1

A Survey of Computer Graphics



Computers have become a powerful tool for the rapid and economical production of pictures. There is virtually no area in which graphical displays cannot be used to some advantage, and so it is not surprising to find the use of computer graphics so widespread. Although early applications in engineering and science had to rely on expensive and cumbersome equipment, advances in computer technology have made interactive computer graphics a practical tool. Today, we find computer graphics used routinely in such diverse areas as science, engineering, medicine, business, industry, government, art, entertainment, advertising, education, and training. Figure 1-1 summarizes the many applications of graphics in simulations, education, and graph presentations. Before we get into the details of how to do computer graphics, we first take a short tour through a gallery of graphics applications.

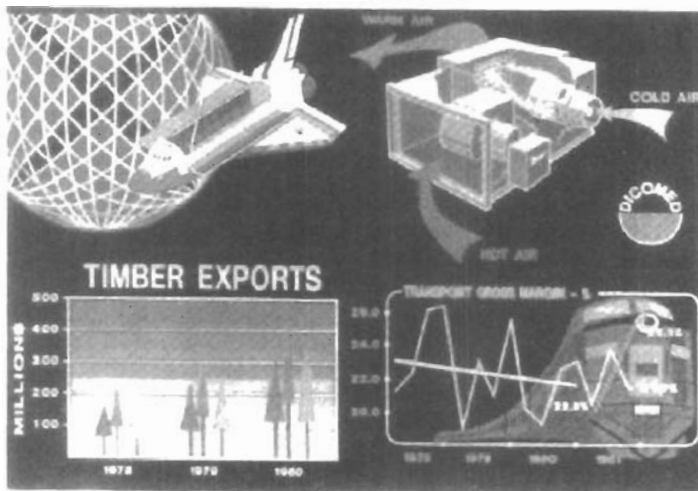


Figure 1-1
Examples of computer graphics applications. (Courtesy of DICOMED Corporation.)

COMPUTER-AIDED DESIGN

A major use of computer graphics is in design processes, particularly for engineering and architectural systems, but almost all products are now computer designed. Generally referred to as CAD, **computer-aided design** methods are now routinely used in the design of buildings, automobiles, aircraft, watercraft, spacecraft, computers, textiles, and many, many other products.

For some design applications, objects are first displayed in a wireframe outline form that shows the overall shape and internal features of objects. Wireframe displays also allow designers to quickly see the effects of interactive adjustments to design shapes. Figures 1-2 and 1-3 give examples of wireframe displays in design applications.

Software packages for CAD applications typically provide the designer with a multi-window environment, as in Figs. 1-4 and 1-5. The various displayed windows can show enlarged sections or different views of objects.

Circuits such as the one shown in Fig. 1-5 and networks for communications, water supply, or other utilities are constructed with repeated placement of a few graphical shapes. The shapes used in a design represent the different network or circuit components. Standard shapes for electrical, electronic, and logic circuits are often supplied by the design package. For other applications, a designer can create personalized symbols that are to be used to construct the network or circuit. The system is then designed by successively placing components into the layout, with the graphics package automatically providing the connections between components. This allows the designer to quickly try out alternate circuit schematics for minimizing the number of components or the space required for the system.

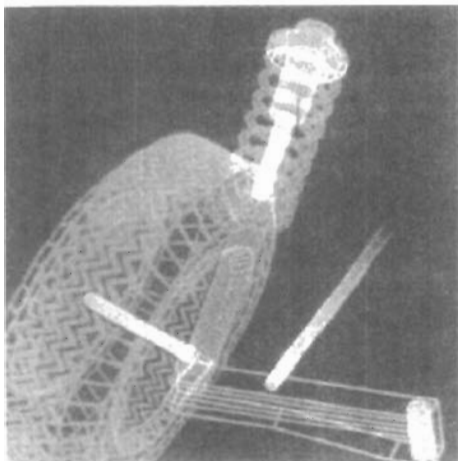
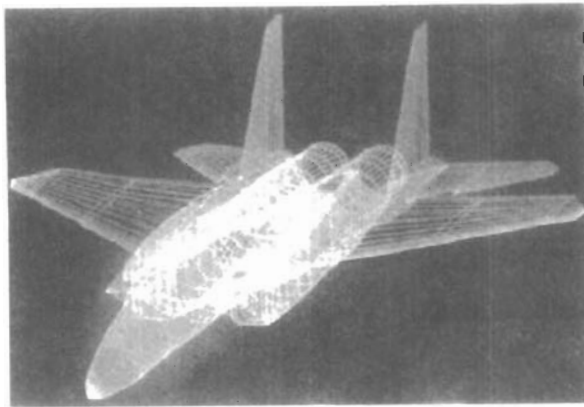


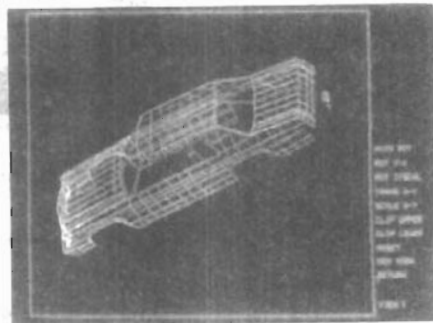
Figure 1-2

Color-coded wireframe display for an automobile wheel assembly.

(Courtesy of Evans & Sutherland.)



(a)



(b)

Figure 1-3

Color-coded wireframe displays of body designs for an aircraft and an automobile.

(Courtesy of (a) Evans & Sutherland and (b) Megatek Corporation.)

Animations are often used in CAD applications. Real-time animations using wireframe displays on a video monitor are useful for testing performance of a vehicle or system, as demonstrated in Fig. 1-6. When we do not display objects with rendered surfaces, the calculations for each segment of the animation can be performed quickly to produce a smooth real-time motion on the screen. Also, wireframe displays allow the designer to see into the interior of the vehicle and to watch the behavior of inner components during motion. Animations in *virtual-reality environments* are used to determine how vehicle operators are affected by

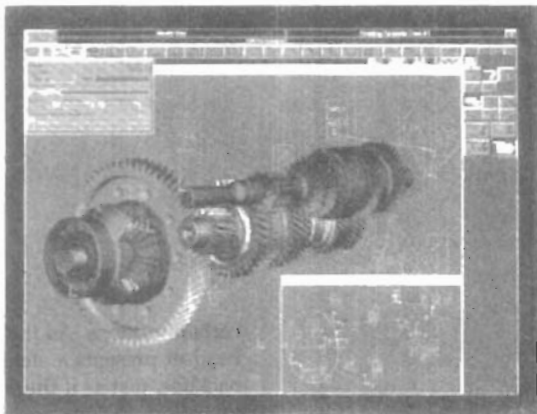
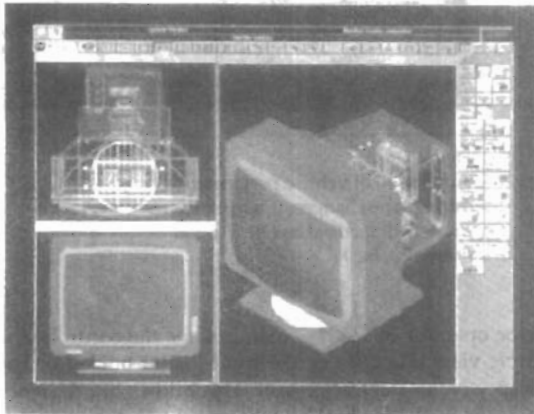


Figure 1-4

Multiple-window, color-coded CAD workstation displays. (Courtesy of Intergraph Corporation.)

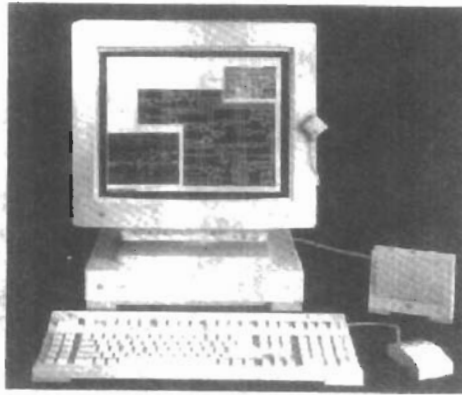


Figure 1-5
A circuit-design application, using multiple windows and color-coded logic components, displayed on a Sun workstation with attached speaker and microphone. (Courtesy of Sun Microsystems.)

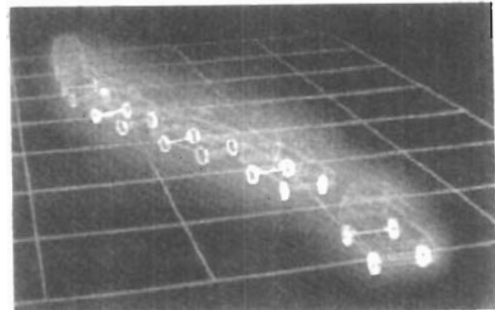


Figure 1-6
Simulation of vehicle performance during lane changes. (Courtesy of Evans & Sutherland and Mechanical Dynamics, Inc.)

certain motions. As the tractor operator in Fig. 1-7 manipulates the controls, the headset presents a stereoscopic view (Fig. 1-8) of the front-loader bucket or the backhoe, just as if the operator were in the tractor seat. This allows the designer to explore various positions of the bucket or backhoe that might obstruct the operator's view, which can then be taken into account in the overall tractor design. Figure 1-9 shows a composite, wide-angle view from the tractor seat, displayed on a standard video monitor instead of in a virtual three-dimensional scene. And Fig. 1-10 shows a view of the tractor that can be displayed in a separate window or on another monitor.

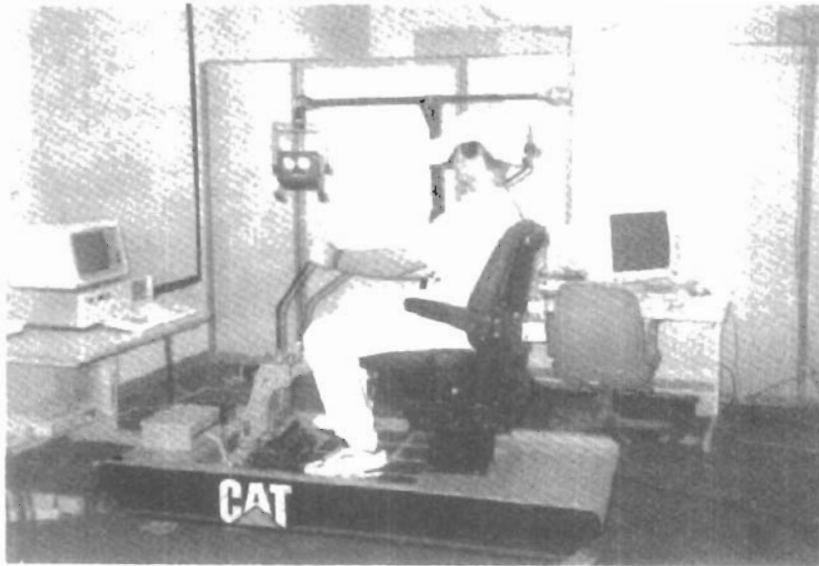


Figure 1-7
Operating a tractor in a virtual-reality environment. As the controls are moved, the operator views the front loader, backhoe, and surroundings through the headset. (Courtesy of the National Center for Supercomputing Applications, University of Illinois at Urbana-Champaign, and Caterpillar, Inc.)

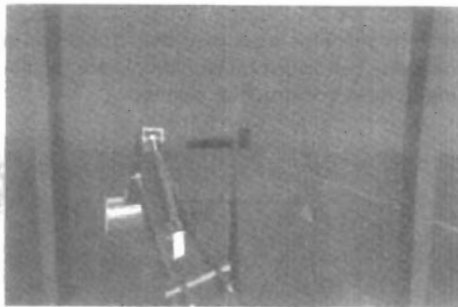


Figure 1-8
A headset view of the backhoe presented to the tractor operator. (Courtesy of the National Center for Supercomputing Applications, University of Illinois at Urbana-Champaign, and Caterpillar, Inc.)

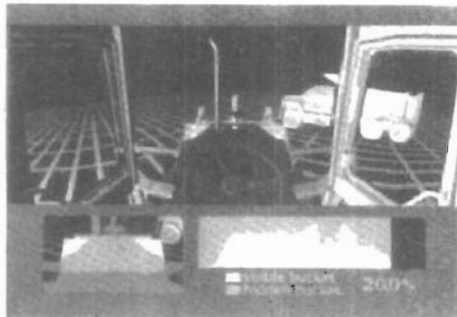


Figure 1-9
Operator's view of the tractor bucket, composited in several sections to form a wide-angle view on a standard monitor. (Courtesy of the National Center for Supercomputing Applications, University of Illinois at Urbana-Champaign, and Caterpillar, Inc.)



Figure 1-10

View of the tractor displayed on a standard monitor. (Courtesy of the National Center for Supercomputing Applications, University of Illinois at Urbana-Champaign, and Caterpillar, Inc.)

When object designs are complete, or nearly complete, realistic lighting models and surface rendering are applied to produce displays that will show the appearance of the final product. Examples of this are given in Fig. 1-11. Realistic displays are also generated for advertising of automobiles and other vehicles using special lighting effects and background scenes (Fig. 1-12).

The manufacturing process is also tied in to the computer description of designed objects to automate the construction of the product. A circuit board layout, for example, can be transformed into a description of the individual processes needed to construct the layout. Some mechanical parts are manufactured by describing how the surfaces are to be formed with machine tools. Figure 1-13 shows the path to be taken by machine tools over the surfaces of an object during its construction. Numerically controlled machine tools are then set up to manufacture the part according to these construction layouts.



(a)



(b)

Figure 1-11

Realistic renderings of design products. (Courtesy of (a) Intergraph Corporation and (b) Evans & Sutherland.)



Figure 1-12
Studio lighting effects and realistic surface-rendering techniques are applied to produce advertising pieces for finished products. The data for this rendering of a Chrysler Laser was supplied by Chrysler Corporation. (Courtesy of Eric Haines, 3D/EYE Inc.)

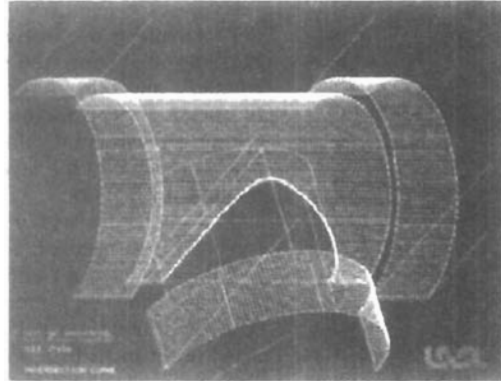


Figure 1-13
A CAD layout for describing the numerically controlled machining of a part. The part surface is displayed in one color and the tool path in another color. (Courtesy of Los Alamos National Laboratory.)

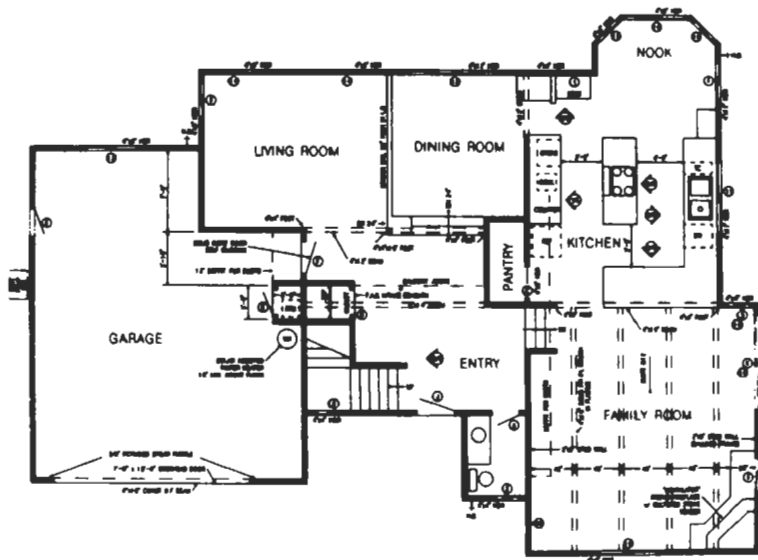


Figure 1-14
Architectural CAD layout for a building design. (Courtesy of Precision Visuals, Inc., Boulder, Colorado.)

Architects use interactive graphics methods to lay out floor plans, such as Fig. 1-14, that show the positioning of rooms, doors, windows, stairs, shelves, counters, and other building features. Working from the display of a building layout on a video monitor, an electrical designer can try out arrangements for wiring, electrical outlets, and fire warning systems. Also, facility-layout packages can be applied to the layout to determine space utilization in an office or on a manufacturing floor.

Realistic displays of architectural designs, as in Fig. 1-15, permit both architects and their clients to study the appearance of a single building or a group of buildings, such as a campus or industrial complex. With virtual-reality systems, designers can even go for a simulated "walk" through the rooms or around the outsides of buildings to better appreciate the overall effect of a particular design. In addition to realistic exterior building displays, architectural CAD packages also provide facilities for experimenting with three-dimensional interior layouts and lighting (Fig. 1-16).

Many other kinds of systems and products are designed using either general CAD packages or specially developed CAD software. Figure 1-17, for example, shows a rug pattern designed with a CAD system.

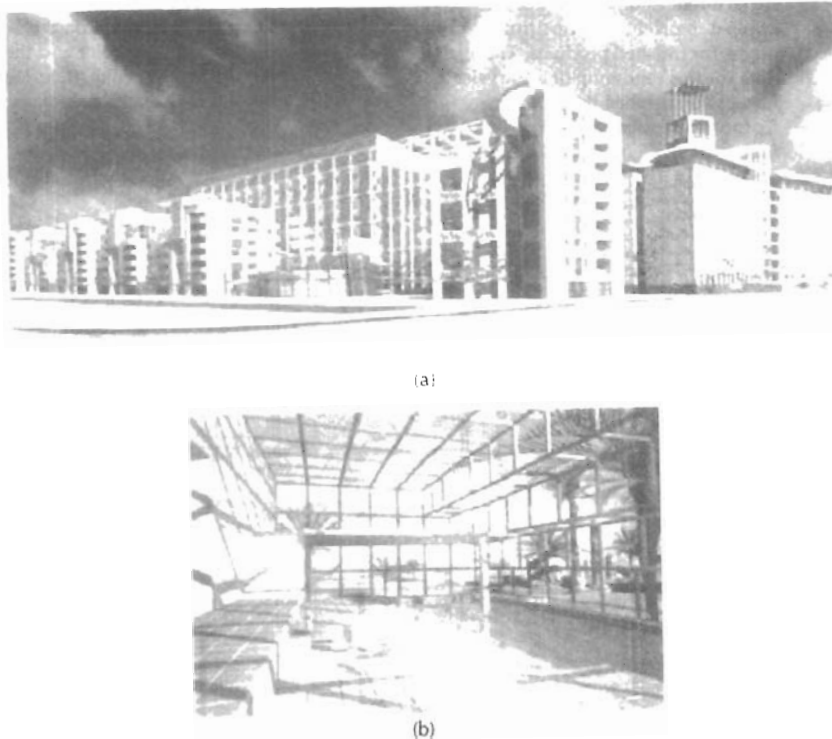


Figure 1-15

Realistic, three-dimensional renderings of building designs. (a) A street-level perspective for the World Trade Center project. (Courtesy of Skidmore, Owings & Merrill.)

(b) Architectural visualization of an atrium, created for a computer animation by Marialine Prieur, Lyon, France. (Courtesy of Thomson Digital Image, Inc.)



Figure 1-16
A hotel corridor providing a sense of movement by placing light fixtures along an undulating path and creating a sense of entry by using light towers at each hotel room. (Courtesy of Skidmore, Owings & Merrill.)

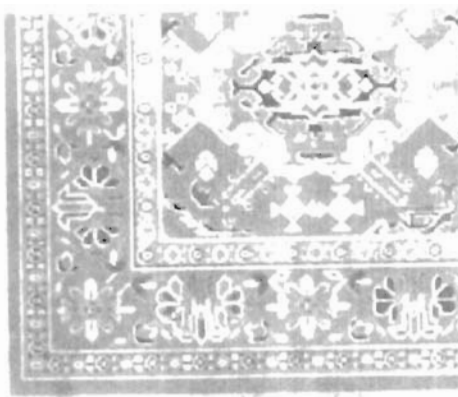


Figure 1-17
Oriental rug pattern created with computer graphics design methods. (Courtesy of Lexidata Corporation.)

1-2

PRESENTATION GRAPHICS

Another major application area is **presentation graphics**, used to produce illustrations for reports or to generate 35-mm slides or transparencies for use with projectors. Presentation graphics is commonly used to summarize financial, statistical, mathematical, scientific, and economic data for research reports, managerial reports, consumer information bulletins, and other types of reports. Workstation devices and service bureaus exist for converting screen displays into 35-mm slides or overhead transparencies for use in presentations. Typical examples of presentation graphics are bar charts, line graphs, surface graphs, pie charts, and other displays showing relationships between multiple parameters.

Figure 1-18 gives examples of two-dimensional graphics combined with geographical information. This illustration shows three color-coded bar charts combined onto one graph and a pie chart with three sections. Similar graphs and charts can be displayed in three dimensions to provide additional information. Three-dimensional graphs are sometimes used simply for effect; they can provide a more dramatic or more attractive presentation of data relationships. The charts in Fig. 1-19 include a three-dimensional bar graph and an exploded pie chart.

Additional examples of three-dimensional graphs are shown in Figs. 1-20 and 1-21. Figure 1-20 shows one kind of surface plot, and Fig. 1-21 shows a two-dimensional contour plot with a height surface.

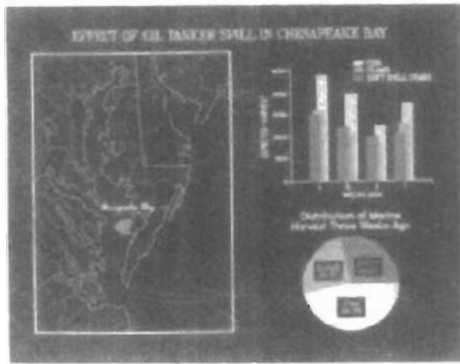


Figure 1-18
Two-dimensional bar chart and pie chart linked to a geographical chart. (Courtesy of Computer Associates, copyright © 1992. All rights reserved.)

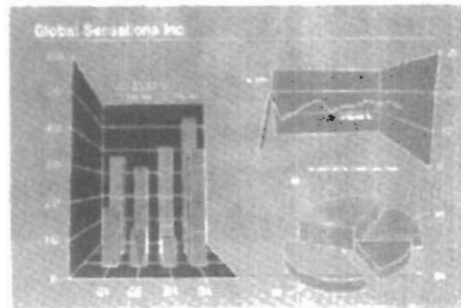


Figure 1-19
Three-dimensional bar chart, exploded pie chart, and line graph. (Courtesy of Computer Associates, copyright © 1992. All rights reserved.)

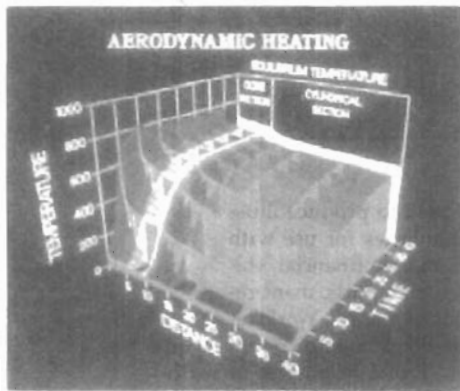


Figure 1-20
Showing relationships with a surface chart. (Courtesy of Computer Associates, copyright © 1992. All rights reserved.)

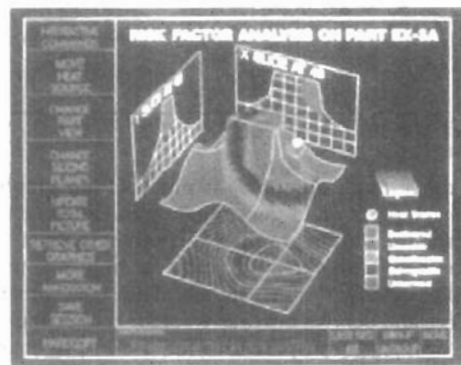


Figure 1-21
Plotting two-dimensional contours in the ground plane, with a height field plotted as a surface above the ground plane. (Courtesy of Computer Associates, copyright © 1992. All rights reserved.)

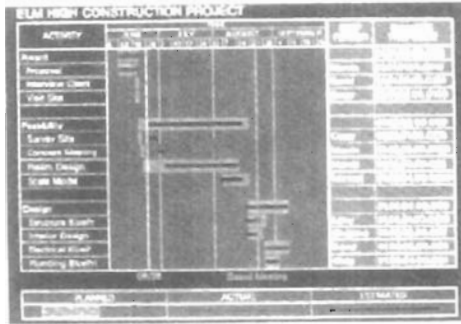


Figure 1-22
Time chart displaying relevant
information about project tasks.
(Courtesy of Computer Associates,
copyright © 1992. All rights reserved.)

Figure 1-22 illustrates a time chart used in task planning. Time charts and task network layouts are used in project management to schedule and monitor the progress of projects.

1-3

COMPUTER ART

Computer graphics methods are widely used in both fine art and commercial art applications. Artists use a variety of computer methods, including special-purpose hardware, artist's paintbrush programs (such as Lumena), other paint packages (such as PixelPaint and SuperPaint), specially developed software, symbolic mathematics packages (such as Mathematica), CAD packages, desktop publishing software, and animation packages that provide facilities for designing object shapes and specifying object motions.

Figure 1-23 illustrates the basic idea behind a *paintbrush* program that allows artists to "paint" pictures on the screen of a video monitor. Actually, the picture is usually painted electronically on a graphics tablet (digitizer) using a stylus, which can simulate different brush strokes, brush widths, and colors. A paintbrush program was used to create the characters in Fig. 1-24, who seem to be busy on a creation of their own.

A paintbrush system, with a Wacom cordless, pressure-sensitive stylus, was used to produce the electronic painting in Fig. 1-25 that simulates the brush strokes of Van Gogh. The stylus translates changing hand pressure into variable line widths, brush sizes, and color gradations. Figure 1-26 shows a watercolor painting produced with this stylus and with software that allows the artist to create watercolor, pastel, or oil brush effects that simulate different drying out times, wetness, and footprint. Figure 1-27 gives an example of paintbrush methods combined with scanned images.

Fine artists use a variety of other computer technologies to produce images. To create pictures such as the one shown in Fig. 1-28, the artist uses a combination of three-dimensional modeling packages, texture mapping, drawing programs, and CAD software. In Fig. 1-29, we have a painting produced on a pen

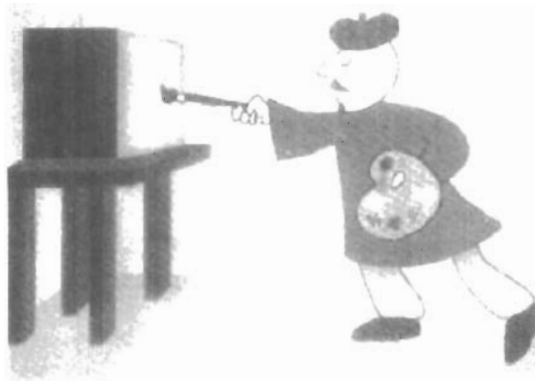


Figure 1-23
Cartoon drawing produced with a paintbrush program, symbolically illustrating an artist at work on a video monitor. (Courtesy of Gould Inc., Imaging & Graphics Division and Aurora Imaging.)

plotter with specially designed software that can create "automatic art" without intervention from the artist.

Figure 1-30 shows an example of "mathematical" art. This artist uses a combination of mathematical functions, fractal procedures, Mathematica software, ink-jet printers, and other systems to create a variety of three-dimensional and two-dimensional shapes and stereoscopic image pairs. Another example of elec-

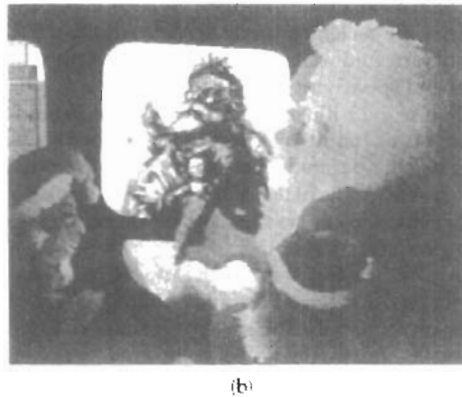
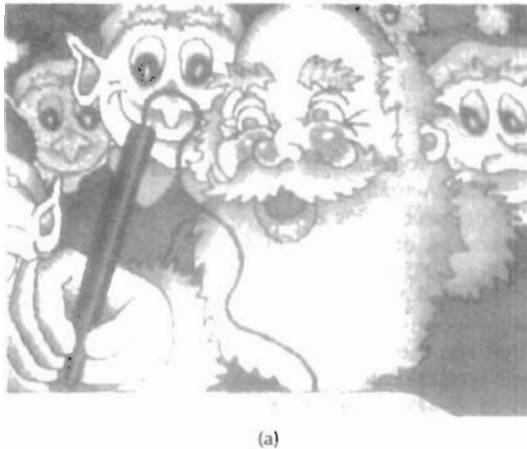


Figure 1-24
Cartoon demonstrations of an "artist" creating a picture with a paintbrush system. The picture, drawn on a graphics tablet, is displayed on the video monitor as the elves look on. In (b), the cartoon is superimposed on the famous Thomas Nast drawing of Saint Nicholas, which was input to the system with a video camera, then scaled and positioned. (Courtesy Gould Inc., Imaging & Graphics Division and Aurora Imaging.)



Figure 1-25

A Van Gogh look-alike created by graphics artist Elizabeth O'Rourke with a cordless, pressure-sensitive stylus. (Courtesy of Wacom Technology Corporation.)



Figure 1-26

An electronic watercolor, painted by John Derry of Time Arts, Inc. using a cordless, pressure-sensitive stylus and Lumena gouache-brush software. (Courtesy of Wacom Technology Corporation.)



Figure 1-27

The artist of this picture, called *Electronic Avalanche*, makes a statement about our entanglement with technology using a personal computer with a graphics tablet and Lumena software to combine renderings of leaves, flower petals, and electronics components with scanned images. (Courtesy of the Williams Gallery. Copyright © 1991 by Joan Truckenbrod, The School of the Art Institute of Chicago.)



Figure 1-28

From a series called *Spheres of Influence*, this electronic painting (entitled, *Whigmalarée*) was created with a combination of methods using a graphics tablet, three-dimensional modeling, texture mapping, and a series of transformations. (Courtesy of the Williams Gallery. Copyright © 1992 by Wynne Ragland, Jr.)

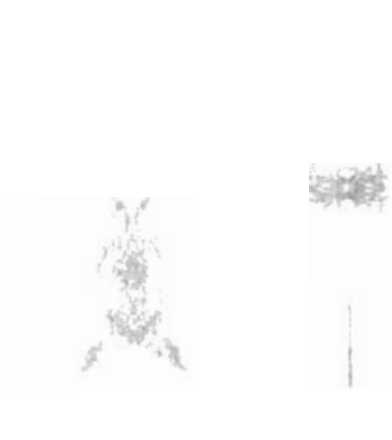


Figure 1-29

Electronic art output to a pen plotter from software specially designed by the artist to emulate his style. The pen plotter includes multiple pens and painting instruments, including Chinese brushes. (Courtesy of the Williams Gallery. Copyright © by Roman Verostko, Minneapolis College of Art & Design.)



Figure 1-30

This creation is based on a visualization of Fermat's Last Theorem, $x^n + y^n = z^n$, with $n = 5$, by Andrew Hanson, Department of Computer Science, Indiana University. The image was rendered using Mathematica and Wavefront software. (Courtesy of the Williams Gallery. Copyright © 1991 by Stewart Dickson.)



Figure 1-31

Using mathematical functions, fractal procedures, and supercomputers, this artist-composer experiments with various designs to synthesize form and color with musical composition. (Courtesy of Brian Evans, Vanderbilt University.)

tronic art created with the aid of mathematical relationships is shown in Fig. 1-31. The artwork of this composer is often designed in relation to frequency variations and other parameters in a musical composition to produce a video that integrates visual and aural patterns.

Although we have spent some time discussing current techniques for generating electronic images in the fine arts, these methods are also applied in commercial art for logos and other designs, page layouts combining text and graphics, TV advertising spots, and other areas. A workstation for producing page layouts that combine text and graphics is illustrated in Fig. 1-32.

For many applications of commercial art (and in motion pictures and other applications), photorealistic techniques are used to render images of a product. Figure 1-33 shows an example of logo design, and Fig. 1-34 gives three computer graphics images for product advertising. Animations are also used frequently in advertising, and television commercials are produced frame by frame, where



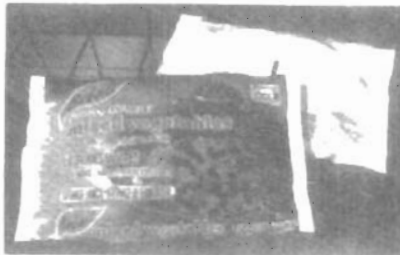
Figure 1-32
Page-layout workstation. (Courtesy
of Visual Technology.)



Figure 1-33
Three-dimensional rendering for a
logo. (Courtesy of Vertigo Technology,
Inc.)



(a)



(b)



Figure 1-34
Product advertising. (Courtesy of (a) Audrey Fleisher and (b) and (c) SOFTIMAGE, Inc.)

each frame of the motion is rendered and saved as an image file. In each successive frame, the motion is simulated by moving object positions slightly from their positions in the previous frame. When all frames in the animation sequence have been rendered, the frames are transferred to film or stored in a video buffer for playback. Film animations require 24 frames for each second in the animation sequence. If the animation is to be played back on a video monitor, 30 frames per second are required.

A common graphics method employed in many commercials is *morphing*, where one object is transformed (metamorphosed) into another. This method has been used in TV commercials to turn an oil can into an automobile engine, an automobile into a tiger, a puddle of water into a tire, and one person's face into another face. An example of morphing is given in Fig. 1-40.

1-4

ENTERTAINMENT

Computer graphics methods are now commonly used in making motion pictures, music videos, and television shows. Sometimes the graphics scenes are displayed by themselves, and sometimes graphics objects are combined with the actors and live scenes.

A graphics scene generated for the movie *Star Trek—The Wrath of Khan* is shown in Fig. 1-35. The planet and spaceship are drawn in wireframe form and will be shaded with rendering methods to produce solid surfaces. Figure 1-36 shows scenes generated with advanced modeling and surface-rendering methods for two award-winning short films.

Many TV series regularly employ computer graphics methods. Figure 1-37 shows a scene produced for the series *Deep Space Nine*. And Fig. 1-38 shows a wireframe person combined with actors in a live scene for the series *Stay Tuned*.

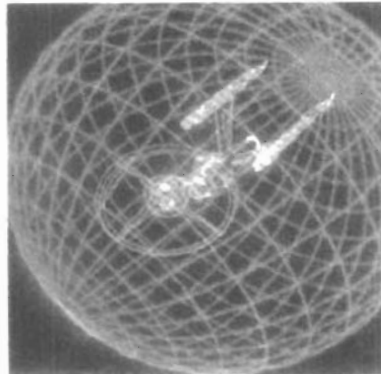


Figure 1-35
Graphics developed for the
Paramount Pictures movie *Star
Trek—The Wrath of Khan*. (Courtesy of
Evans & Sutherland.)

In Fig. 1-39, we have a highly realistic image taken from a reconstruction of thirteenth-century Dadu (now Beijing) for a Japanese broadcast.

Music videos use graphics in several ways. Graphics objects can be combined with the live action, as in Fig.1-38, or graphics and image processing techniques can be used to produce a transformation of one person or object into another (morphing). An example of morphing is shown in the sequence of scenes in Fig. 1-40, produced for the David Byrne video *She's Mad*.



(a)



(b)

Figure 1-36

(a) A computer-generated scene from the film *Red's Dream*, copyright © Pixar 1987. (b) A computer-generated scene from the film *Knickknack*, copyright © Pixar 1989. (Courtesy of Pixar.)

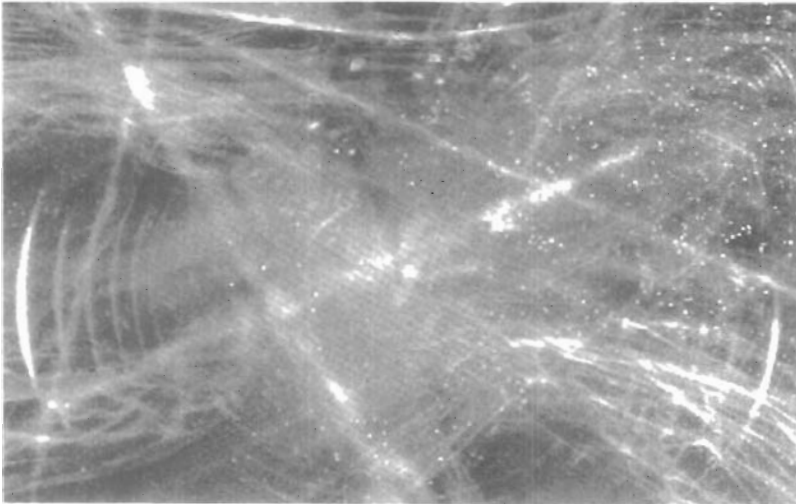


Figure 1-37

A graphics scene in the TV series *Deep Space Nine*. (Courtesy of Rhythm & Hues Studios.)

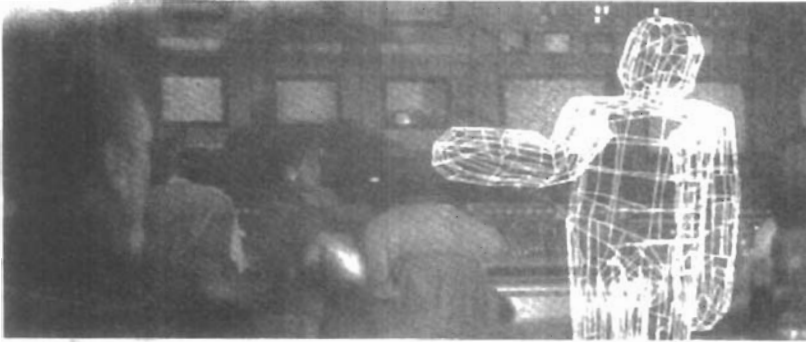


Figure 1-38

Graphics combined with a live scene in the TV series *Stay Tuned*.
(Courtesy of Rhythm & Hues Studios.)



Figure 1-39

An image from a reconstruction of thirteenth-century Dadu (Beijing today), created by Taisei Corporation (Tokyo) and rendered with TDI software. (Courtesy of Thompson Digital Image, Inc.)

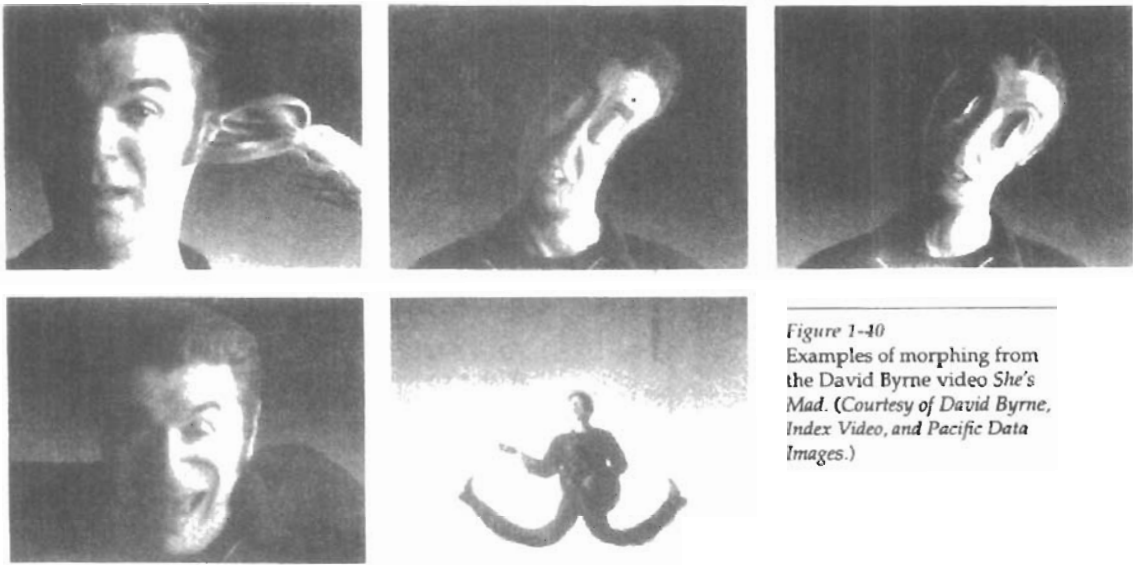


Figure 1-40
Examples of morphing from
the David Byrne video *She's
Mad*. (Courtesy of David Byrne,
Index Video, and Pacific Data
Images.)

1-5 EDUCATION AND TRAINING

Computer-generated models of physical, financial, and economic systems are often used as educational aids. Models of physical systems, physiological systems, population trends, or equipment, such as the color-coded diagram in Fig. 1-41, can help trainees to understand the operation of the system.

For some training applications, special systems are designed. Examples of such specialized systems are the simulators for practice sessions or training of ship captains, aircraft pilots, heavy-equipment operators, and air traffic-control personnel. Some simulators have no video screens; for example, a flight simulator with only a control panel for instrument flying. But most simulators provide graphics screens for visual operation. Two examples of large simulators with internal viewing systems are shown in Figs. 1-42 and 1-43. Another type of viewing system is shown in Fig. 1-44. Here a viewing screen with multiple panels is mounted in front of the simulator, and color projectors display the flight scene on the screen panels. Similar viewing systems are used in simulators for training aircraft control-tower personnel. Figure 1-45 gives an example of the instructor's area in a flight simulator. The keyboard is used to input parameters affecting the airplane performance or the environment, and the pen plotter is used to chart the path of the aircraft during a training session.

Scenes generated for various simulators are shown in Figs. 1-46 through 1-48. An output from an automobile-driving simulator is given in Fig. 1-49. This simulator is used to investigate the behavior of drivers in critical situations. The drivers' reactions are then used as a basis for optimizing vehicle design to maximize traffic safety.

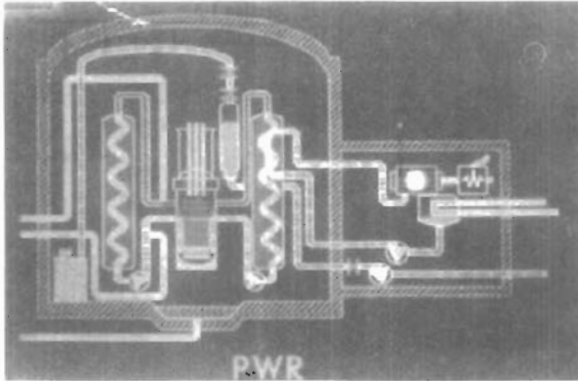


Figure 1-41
Color-coded diagram used to explain the operation of a nuclear reactor. (Courtesy of Los Alamos National Laboratory.)



Figure 1-42
A large, enclosed flight simulator with a full-color visual system and six degrees of freedom in its motion. (Courtesy of Frasca International.)

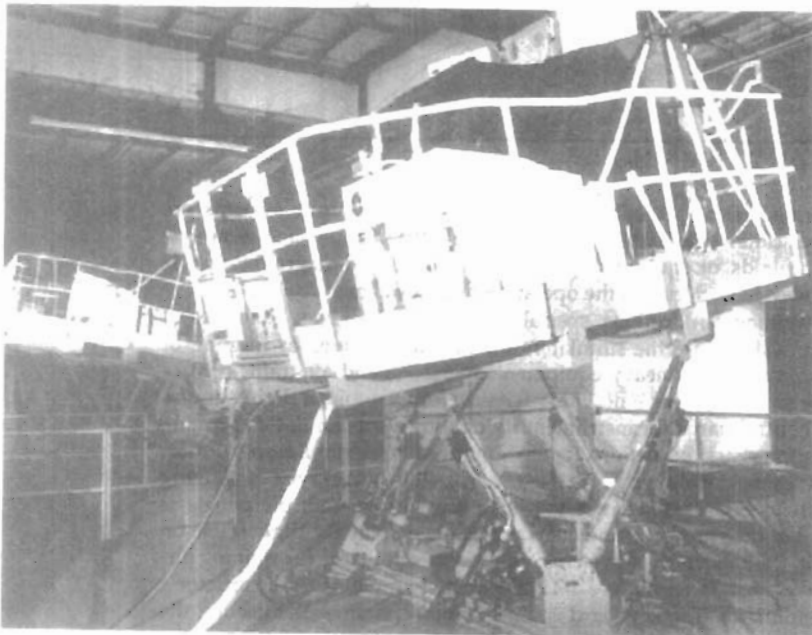


Figure 1-43
A military tank simulator with a visual imagery system. (Courtesy of Mediatech and GE Aerospace.)

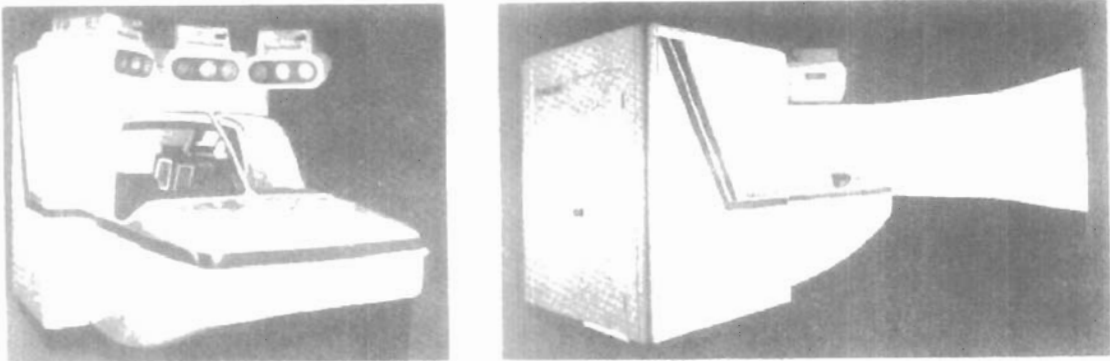


Figure 1-44
A flight simulator with an external full-color viewing system. (Courtesy of Frasca International.)



Figure 1-45
An instructor's area in a flight simulator. The equipment allows the instructor to monitor flight conditions and to set airplane and environment parameters. (Courtesy of Frasca International.)



Figure 1-46
Flight-simulator imagery. (Courtesy of Evans & Sutherland.)



Figure 1-47
Imagery generated for a naval simulator. (Courtesy of Evans & Sutherland.)

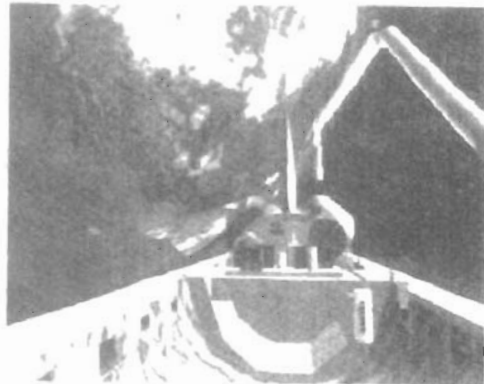


Figure 1-48
Space shuttle imagery. (Courtesy of Mediatech and GE Aerospace.)



Figure 1-49
Imagery from an automobile simulator used to test driver reaction. (Courtesy of Evans & Sutherland.)

1-6

VISUALIZATION

Scientists, engineers, medical personnel, business analysts, and others often need to analyze large amounts of information or to study the behavior of certain processes. Numerical simulations carried out on supercomputers frequently produce data files containing thousands and even millions of data values. Similarly, satellite cameras and other sources are amassing large data files faster than they can be interpreted. Scanning these large sets of numbers to determine trends and relationships is a tedious and ineffective process. But if the data are converted to a visual form, the trends and patterns are often immediately apparent. Figure 1-50 shows an example of a large data set that has been converted to a color-coded display of relative heights above a ground plane. Once we have plotted the density values in this way, we can see easily the overall pattern of the data. Producing graphical representations for scientific, engineering, and medical data sets and processes is generally referred to as *scientific visualization*. And the term *business visualization* is used in connection with data sets related to commerce, industry, and other nonscientific areas.

There are many different kinds of data sets, and effective visualization schemes depend on the characteristics of the data. A collection of data can contain scalar values, vectors, higher-order tensors, or any combination of these data types. And data sets can be two-dimensional or three-dimensional. Color coding is just one way to visualize a data set. Additional techniques include contour plots, graphs and charts, surface renderings, and visualizations of volume interiors. In addition, image processing techniques are combined with computer graphics to produce many of the data visualizations.

Mathematicians, physical scientists, and others use visual techniques to analyze mathematical functions and processes or simply to produce interesting graphical representations. A color plot of mathematical curve functions is shown in Fig. 1-51, and a surface plot of a function is shown in Fig. 1-52. Fractal proce-

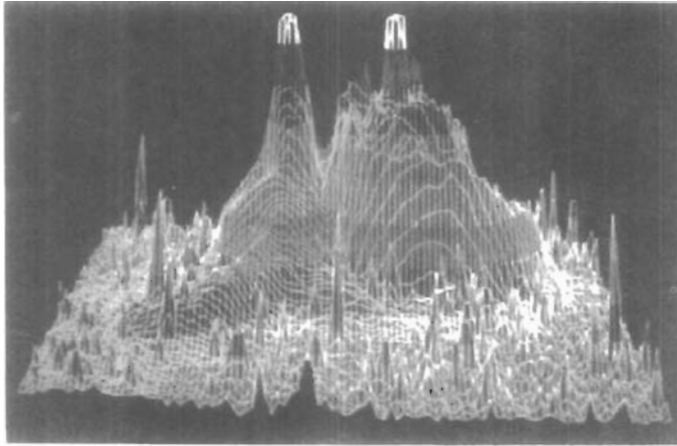


Figure 1-50
A color-coded plot with 16 million density points of relative brightness observed for the Whirlpool Nebula reveals two distinct galaxies.
(Courtesy of Los Alamos National Laboratory.)

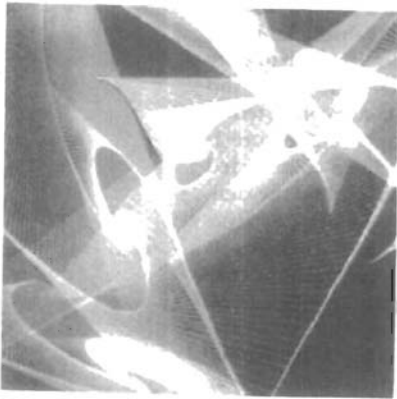


Figure 1-51
Mathematical curve functions plotted in various color combinations. (Courtesy of Melvin L. Prueitt, Los Alamos National Laboratory.)

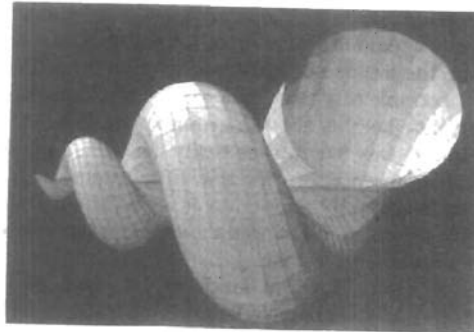


Figure 1-52
Lighting effects and surface-rendering techniques were applied to produce this surface representation for a three-dimensional function. (Courtesy of Wolfram Research, Inc, The Maker of Mathematica.)

dures using quaternions generated the object shown in Fig. 1-53, and a topological structure is displayed in Fig. 1-54. Scientists are also developing methods for visualizing general classes of data. Figure 1-55 shows a general technique for graphing and modeling data distributed over a spherical surface.

A few of the many other visualization applications are shown in Figs. 1-56 through 1-69. These figures show airflow over the surface of a space shuttle, numerical modeling of thunderstorms, study of crack propagation in metals, a color-coded plot of fluid density over an airfoil, a cross-sectional slicer for data sets, protein modeling, stereoscopic viewing of molecular structure, a model of the ocean floor, a Kuwaiti oil-fire simulation, an air-pollution study, a corn-growing study, reconstruction of Arizona's Chaco Canyon ruins, and a graph of automobile accident statistics.

Section 1-6

Visualization



Figure 1-53

A four-dimensional object projected into three-dimensional space, then projected to a video monitor, and color coded. The object was generated using quaternions and fractal squaring procedures, with an octant subtracted to show the complex Julia set. (Courtesy of John C. Hart, School of Electrical Engineering and Computer Science, Washington State University.)



Figure 1-54

Four views from a real-time, interactive computer-animation study of minimal surfaces ("snails") in the 3-sphere projected to three-dimensional Euclidean space.

(Courtesy of George Francis, Department of Mathematics and the National Center for Supercomputing Applications, University of Illinois at Urbana-Champaign. Copyright © 1993.)

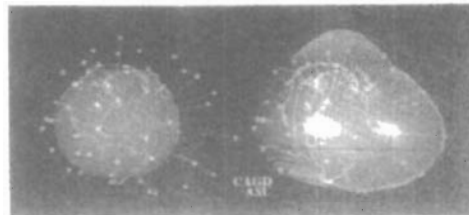


Figure 1-55

A method for graphing and modeling data distributed over a spherical surface. (Courtesy of Greg Nielson, Computer Science Department, Arizona State University.)

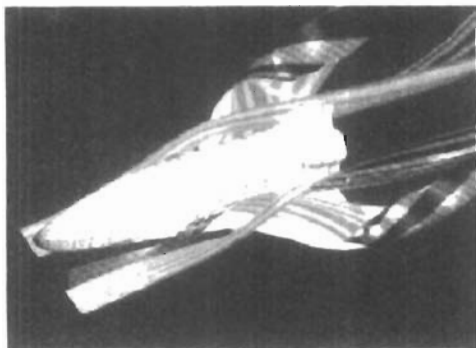


Figure 1-56
A visualization of stream surfaces
flowing past a space shuttle by Jeff
Hultquist and Eric Raible, NASA
Ames. (Courtesy of Sam Useton,
NASA Ames Research Center.)

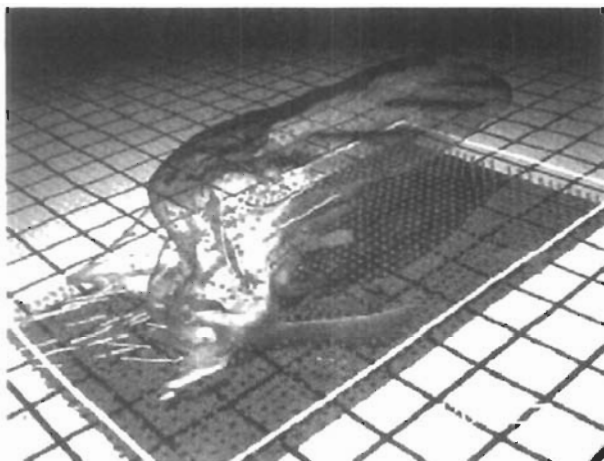


Figure 1-57
Numerical model of airflow inside
a thunderstorm. (Courtesy of Bob
Wilhelmson, Department of
Atmospheric Sciences and the National
Center for Supercomputing
Applications, University of Illinois at
Urbana-Champaign.)

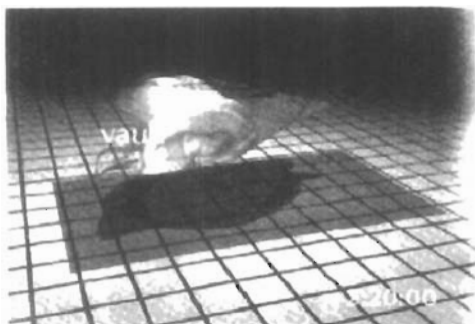


Figure 1-58
Numerical model of the surface of a
thunderstorm. (Courtesy of Bob
Wilhelmson, Department of
Atmospheric Sciences and the National
Center for Supercomputing
Applications, University of Illinois at
Urbana-Champaign.)

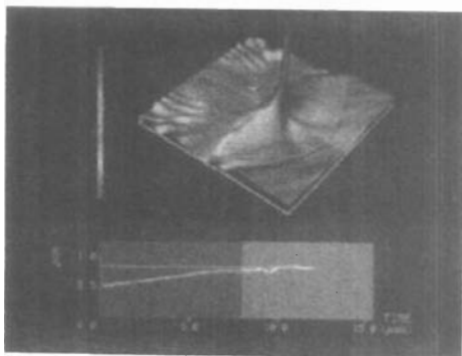


Figure 1-59
Color-coded visualization of stress energy density in a crack-propagation study for metal plates, modeled by Bob Haber. (Courtesy of the National Center for Supercomputing Applications, University of Illinois at Urbana-Champaign.)

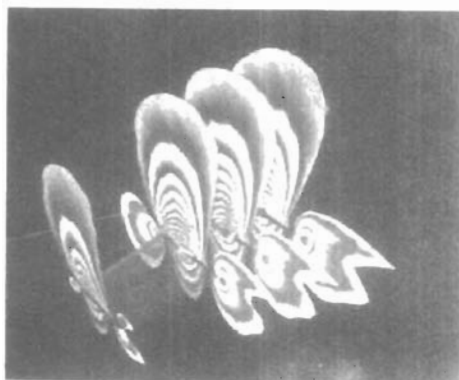


Figure 1-60
A fluid dynamic simulation, showing a color-coded plot of fluid density over a span of grid planes around an aircraft wing, developed by Lee-Hian Quek, John Eickemeyer, and Jeffery Tan. (Courtesy of the Information Technology Institute, Republic of Singapore.)

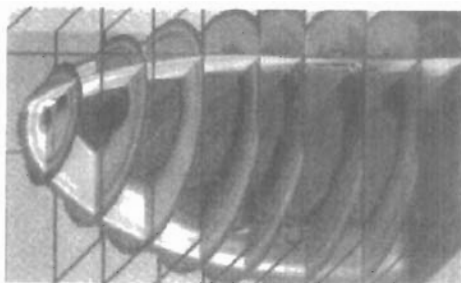


Figure 1-61
Commercial slicer-dicer software, showing color-coded data values over cross-sectional slices of a data set. (Courtesy of Spyglass, Inc.)

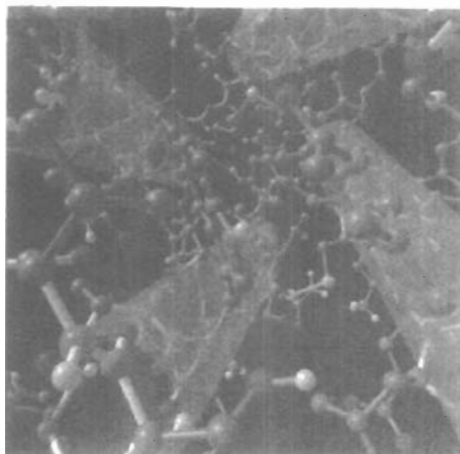


Figure 1-62
Visualization of a protein structure by Jay Siegel and Kim Baldrige, SDSC. (Courtesy of Stephanie Sides, San Diego Supercomputer Center.)

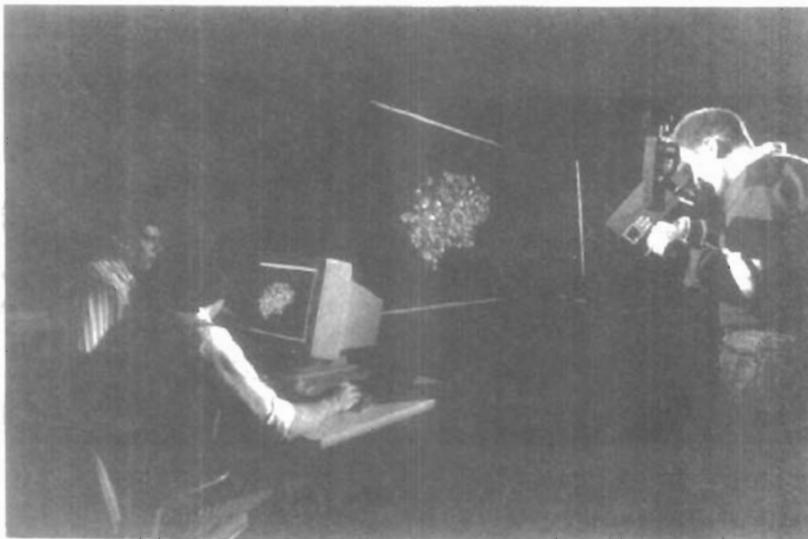


Figure 1-63
Stereoscopic viewing of a molecular structure using a "boom" device.
(Courtesy of the National Center for Supercomputing Applications, University of Illinois at Urbana-Champaign.)

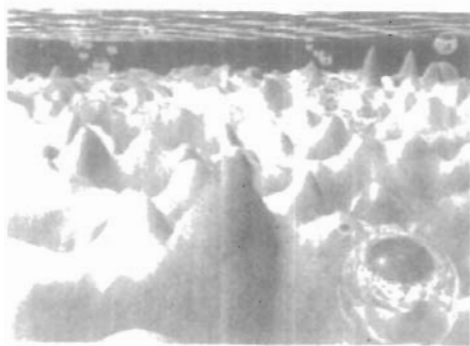


Figure 1-64
One image from a stereoscopic pair, showing a visualization of the ocean floor obtained from satellite data, by David Sandwell and Chris Small, Scripps Institution of Oceanography, and Jim McLeod, SDSC.
(Courtesy of Stephanie Sides, San Diego Supercomputer Center.)

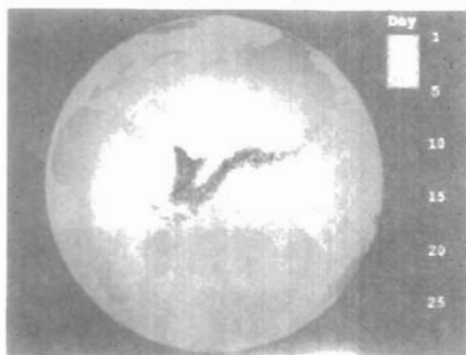


Figure 1-65
A simulation of the effects of the Kuwaiti oil fire, by Gary Glatzmeier, Chuck Hanson, and Paul Hinker. (Courtesy of Mike Krogh, Advanced Computing Laboratory at Los Alamos National Laboratory.)

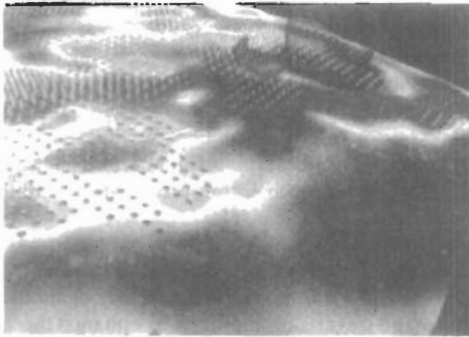


Figure 1-66
A visualization of pollution over the earth's surface by Tom Palmer, Cray Research Inc./NCSC; Chris Landreth, NCSC; and Dave Bock, NCSC. Pollutant SO_4 is plotted as a blue surface, acid-rain deposition is a color plane on the map surface, and rain concentration is shown as clear cylinders. (Courtesy of the North Carolina Supercomputing Center/MCNC.)

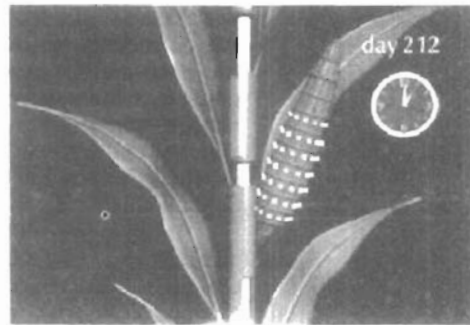


Figure 1-67
One frame of an animation sequence showing the development of a corn ear. (Courtesy of the National Center for Supercomputing Applications, University of Illinois at Urbana-Champaign.)



Figure 1-68
A visualization of the reconstruction of the ruins at Chaco Canyon, Arizona. (Courtesy of Melvin L. Prueitt, Los Alamos National Laboratory. Data supplied by Stephen H. Lekson.)

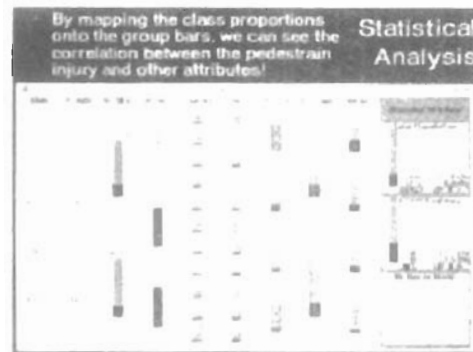


Figure 1-69
A prototype technique, called WinViz, for visualizing tabular multidimensional data is used here to correlate statistical information on pedestrians involved in automobile accidents, developed by a visualization team at ITT. (Courtesy of Lee-Hian Quek, Information Technology Institute, Republic of Singapore.)

IMAGE PROCESSING

Although methods used in computer graphics and image processing overlap, the two areas are concerned with fundamentally different operations. In computer graphics, a computer is used to create a picture. **Image processing**, on the other hand, applies techniques to modify or interpret existing pictures, such as photographs and TV scans. Two principal applications of image processing are (1) improving picture quality and (2) machine perception of visual information, as used in robotics.

To apply image-processing methods, we first digitize a photograph or other picture into an image file. Then digital methods can be applied to rearrange picture parts, to enhance color separations, or to improve the quality of shading. An example of the application of image-processing methods to enhance the quality of a picture is shown in Fig. 1-70. These techniques are used extensively in commercial art applications that involve the retouching and rearranging of sections of photographs and other artwork. Similar methods are used to analyze satellite photos of the earth and photos of galaxies.

Medical applications also make extensive use of image-processing techniques for picture enhancements, in tomography and in simulations of operations. Tomography is a technique of X-ray photography that allows cross-sectional views of physiological systems to be displayed. Both *computed X-ray tomography* (CT) and *position emission tomography* (PET) use projection methods to reconstruct cross sections from digital data. These techniques are also used to

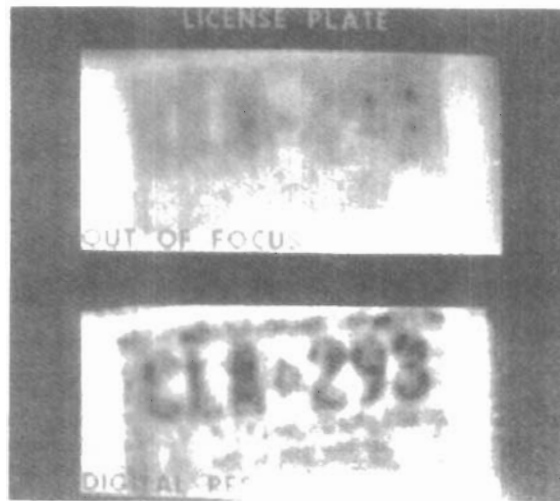


Figure 1-70

A blurred photograph of a license plate becomes legible after the application of image-processing techniques. (Courtesy of Los Alamos National Laboratory.)

monitor internal functions and show cross sections during surgery. Other medical imaging techniques include ultrasonics and nuclear medicine scanners. With ultrasonics, high-frequency sound waves, instead of X-rays, are used to generate digital data. Nuclear medicine scanners collect digital data from radiation emitted from ingested radionuclides and plot color-coded images.

Image processing and computer graphics are typically combined in many applications. Medicine, for example, uses these techniques to model and study physical functions, to design artificial limbs, and to plan and practice surgery. The last application is generally referred to as *computer-aided surgery*. Two-dimensional cross sections of the body are obtained using imaging techniques. Then the slices are viewed and manipulated using graphics methods to simulate actual surgical procedures and to try out different surgical cuts. Examples of these medical applications are shown in Figs. 1-71 and 1-72.

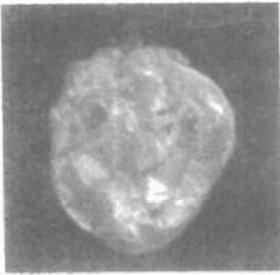


Figure 1-71

One frame from a computer animation visualizing cardiac activation levels within regions of a semitransparent volume-rendered dog heart. Medical data provided by William Smith, Ed Simpson, and G. Allan Johnson, Duke University. Image-rendering software by Tom Palmer, Cray Research, Inc./NCSC. (Courtesy of Dave Bock, North Carolina Supercomputing Center/MCNC.)

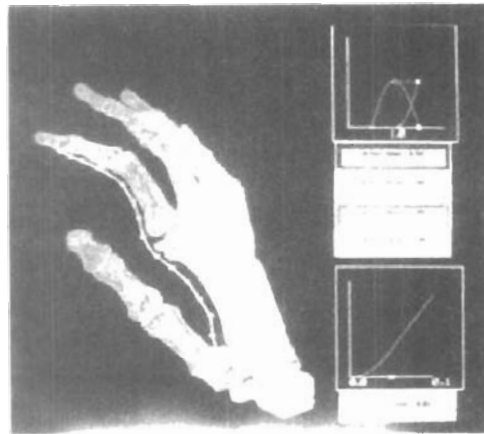


Figure 1-72

One image from a stereoscopic pair showing the bones of a human hand. The images were rendered by Inmo Yoon, D. E. Thompson, and W. N. Waggenspack, Jr., LSU, from a data set obtained with CT scans by Rehabilitation Research, GWLNHDC. These images show a possible tendon path for reconstructive surgery. (Courtesy of IMRLAB, Mechanical Engineering, Louisiana State University.)

GRAPHICAL USER INTERFACES

It is common now for software packages to provide a **graphical interface**. A major component of a graphical interface is a window manager that allows a user to display multiple-window areas. Each window can contain a different process that can contain graphical or nongraphical displays. To make a particular window active, we simply click in that window using an interactive pointing device.

Interfaces also display menus and icons for fast selection of processing options or parameter values. An **icon** is a graphical symbol that is designed to look like the processing option it represents. The advantages of icons are that they take up less screen space than corresponding textual descriptions and they can be understood more quickly if well designed. Menus contain lists of textual descriptions and icons.

Figure 1-73 illustrates a typical graphical interface, containing a window manager, menu displays, and icons. In this example, the menus allow selection of processing options, color values, and graphics parameters. The icons represent options for painting, drawing, zooming, typing text strings, and other operations connected with picture construction.

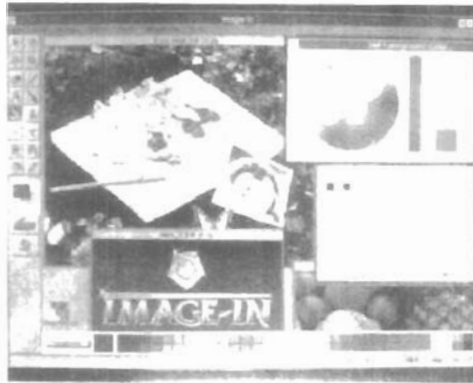


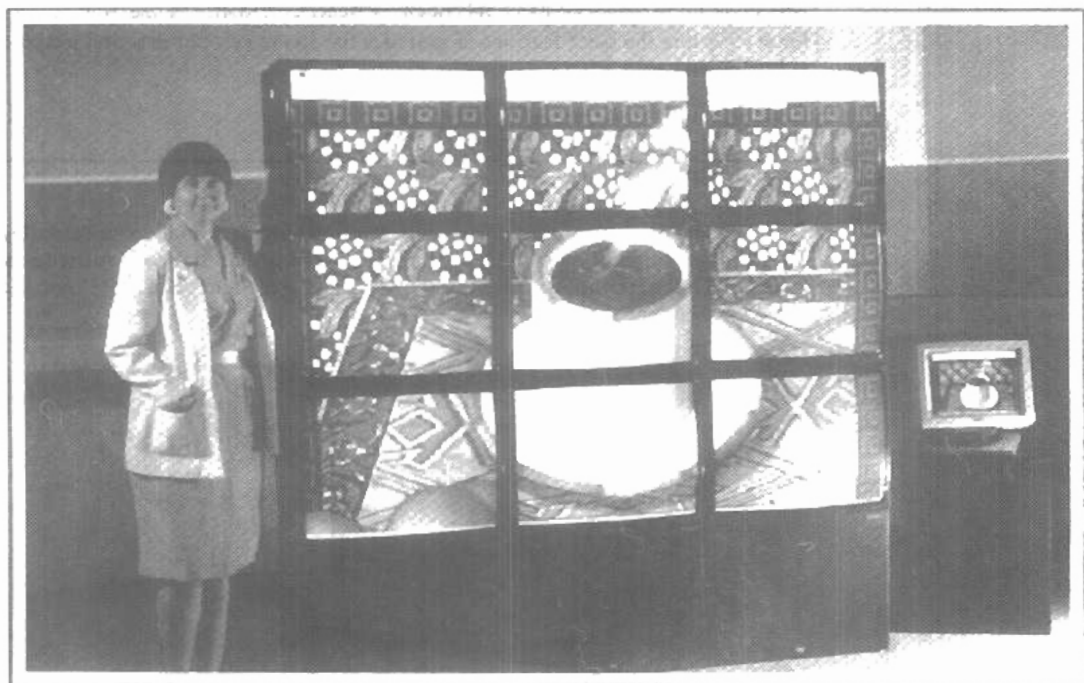
Figure 1-73

A graphical user interface, showing multiple window areas, menus, and icons. (Courtesy of Image-In Corporation.)

CHAPTER

2

Overview of Graphics Systems



Due to the widespread recognition of the power and utility of computer graphics in virtually all fields, a broad range of graphics hardware and software systems is now available. Graphics capabilities for both two-dimensional and three-dimensional applications are now common on general-purpose computers, including many hand-held calculators. With personal computers, we can use a wide variety of interactive input devices and graphics software packages. For higher-quality applications, we can choose from a number of sophisticated special-purpose graphics hardware systems and technologies. In this chapter, we explore the basic features of graphics hardware components and graphics software packages.

2-1

VIDEO DISPLAY DEVICES

Typically, the primary output device in a graphics system is a video monitor (Fig. 2-1). The operation of most video monitors is based on the standard **cathode-ray tube (CRT)** design, but several other technologies exist and solid-state monitors may eventually predominate.



Figure 2-1
A computer graphics workstation. (Courtesy of Tektronix, Inc.)

Refresh Cathode-Ray Tubes

Section 2-1

Video Display Devices

Figure 2-2 illustrates the basic operation of a CRT. A beam of electrons (*cathode rays*), emitted by an electron gun, passes through focusing and deflection systems that direct the beam toward specified positions on the phosphor-coated screen. The phosphor then emits a small spot of light at each position contacted by the electron beam. Because the light emitted by the phosphor fades very rapidly, some method is needed for maintaining the screen picture. One way to keep the phosphor glowing is to redraw the picture repeatedly by quickly directing the electron beam back over the same points. This type of display is called a **refresh CRT**.

The primary components of an electron gun in a CRT are the heated metal cathode and a control grid (Fig. 2-3). Heat is supplied to the cathode by directing a current through a coil of wire, called the filament, inside the cylindrical cathode structure. This causes electrons to be "boiled off" the hot cathode surface. In the vacuum inside the CRT envelope, the free, negatively charged electrons are then accelerated toward the phosphor coating by a high positive voltage. The acceler-

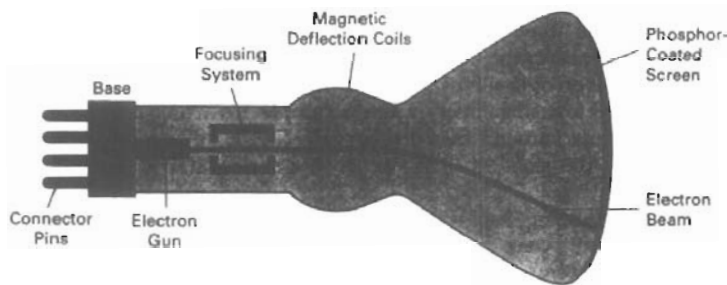


Figure 2-2
Basic design of a magnetic-deflection CRT.

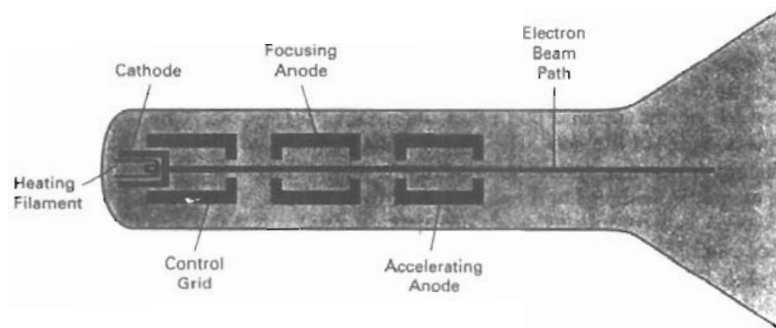


Figure 2-3
Operation of an electron gun with an accelerating anode.

ating voltage can be generated with a positively charged metal coating on the inside of the CRT envelope near the phosphor screen, or an accelerating anode can be used, as in Fig. 2-3. Sometimes the electron gun is built to contain the accelerating anode and focusing system within the same unit.

Intensity of the electron beam is controlled by setting voltage levels on the control grid, which is a metal cylinder that fits over the cathode. A high negative voltage applied to the control grid will shut off the beam by repelling electrons and stopping them from passing through the small hole at the end of the control grid structure. A smaller negative voltage on the control grid simply decreases the number of electrons passing through. Since the amount of light emitted by the phosphor coating depends on the number of electrons striking the screen, we control the brightness of a display by varying the voltage on the control grid. We specify the intensity level for individual screen positions with graphics software commands, as discussed in Chapter 3.

The focusing system in a CRT is needed to force the electron beam to converge into a small spot as it strikes the phosphor. Otherwise, the electrons would repel each other, and the beam would spread out as it approaches the screen. Focusing is accomplished with either electric or magnetic fields. Electrostatic focusing is commonly used in television and computer graphics monitors. With electrostatic focusing, the electron beam passes through a positively charged metal cylinder that forms an electrostatic lens, as shown in Fig. 2-3. The action of the electrostatic lens focuses the electron beam at the center of the screen, in exactly the same way that an optical lens focuses a beam of light at a particular focal distance. Similar lens focusing effects can be accomplished with a magnetic field set up by a coil mounted around the outside of the CRT envelope. Magnetic lens focusing produces the smallest spot size on the screen and is used in special-purpose devices.

Additional focusing hardware is used in high-precision systems to keep the beam in focus at all screen positions. The distance that the electron beam must travel to different points on the screen varies because the radius of curvature for most CRTs is greater than the distance from the focusing system to the screen center. Therefore, the electron beam will be focused properly only at the center of the screen. As the beam moves to the outer edges of the screen, displayed images become blurred. To compensate for this, the system can adjust the focusing according to the screen position of the beam.

As with focusing, deflection of the electron beam can be controlled either with electric fields or with magnetic fields. Cathode-ray tubes are now commonly constructed with magnetic deflection coils mounted on the outside of the CRT envelope, as illustrated in Fig. 2-2. Two pairs of coils are used, with the coils in each pair mounted on opposite sides of the neck of the CRT envelope. One pair is mounted on the top and bottom of the neck, and the other pair is mounted on opposite sides of the neck. The magnetic field produced by each pair of coils results in a transverse deflection force that is perpendicular both to the direction of the magnetic field and to the direction of travel of the electron beam. Horizontal deflection is accomplished with one pair of coils, and vertical deflection by the other pair. The proper deflection amounts are attained by adjusting the current through the coils. When electrostatic deflection is used, two pairs of parallel plates are mounted inside the CRT envelope. One pair of plates is mounted horizontally to control the vertical deflection, and the other pair is mounted vertically to control horizontal deflection (Fig. 2-4).

Spots of light are produced on the screen by the transfer of the CRT beam energy to the phosphor. When the electrons in the beam collide with the phos-

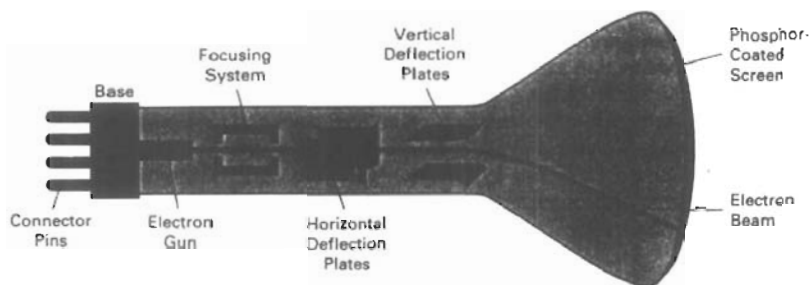


Figure 2-4
Electrostatic deflection of the electron beam in a CRT.

phor coating, they are stopped and their kinetic energy is absorbed by the phosphor. Part of the beam energy is converted by friction into heat energy, and the remainder causes electrons in the phosphor atoms to move up to higher quantum-energy levels. After a short time, the “excited” phosphor electrons begin dropping back to their stable ground state, giving up their extra energy as small quanta of light energy. What we see on the screen is the combined effect of all the electron light emissions: a glowing spot that quickly fades after all the excited phosphor electrons have returned to their ground energy level. The frequency (or color) of the light emitted by the phosphor is proportional to the energy difference between the excited quantum state and the ground state.

Different kinds of phosphors are available for use in a CRT. Besides color, a major difference between phosphors is their **persistence**: how long they continue to emit light (that is, have excited electrons returning to the ground state) after the CRT beam is removed. Persistence is defined as the time it takes the emitted light from the screen to decay to one-tenth of its original intensity. Lower-persistence phosphors require higher refresh rates to maintain a picture on the screen without flicker. A phosphor with low persistence is useful for animation; a high-persistence phosphor is useful for displaying highly complex, static pictures. Although some phosphors have a persistence greater than 1 second, graphics monitors are usually constructed with a persistence in the range from 10 to 60 microseconds.

Figure 2-5 shows the intensity distribution of a spot on the screen. The intensity is greatest at the center of the spot, and decreases with a Gaussian distribution out to the edges of the spot. This distribution corresponds to the cross-sectional electron density distribution of the CRT beam.

The maximum number of points that can be displayed without overlap on a CRT is referred to as the **resolution**. A more precise definition of resolution is the number of points per centimeter that can be plotted horizontally and vertically, although it is often simply stated as the total number of points in each direction. Spot intensity has a Gaussian distribution (Fig. 2-5), so two adjacent spots will appear distinct as long as their separation is greater than the diameter at which each spot has an intensity of about 60 percent of that at the center of the spot. This overlap position is illustrated in Fig. 2-6. Spot size also depends on intensity. As more electrons are accelerated toward the phosphor per second, the CRT beam diameter and the illuminated spot increase. In addition, the increased excitation energy tends to spread to neighboring phosphor atoms not directly in the

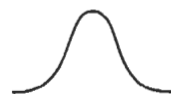


Figure 2-5
Intensity distribution of an illuminated phosphor spot on a CRT screen.



Figure 2-6

Two illuminated phosphor spots are distinguishable when their separation is greater than the diameter at which a spot intensity has fallen to 60 percent of maximum.

path of the beam, which further increases the spot diameter. Thus, resolution of a CRT is dependent on the type of phosphor, the intensity to be displayed, and the focusing and deflection systems. Typical resolution on high-quality systems is 1280 by 1024, with higher resolutions available on many systems. High-resolution systems are often referred to as *high-definition systems*. The physical size of a graphics monitor is given as the length of the screen diagonal, with sizes varying from about 12 inches to 27 inches or more. A CRT monitor can be attached to a variety of computer systems, so the number of screen points that can actually be plotted depends on the capabilities of the system to which it is attached.

Another property of video monitors is **aspect ratio**. This number gives the ratio of vertical points to horizontal points necessary to produce equal-length lines in both directions on the screen. (Sometimes aspect ratio is stated in terms of the ratio of horizontal to vertical points.) An aspect ratio of 3/4 means that a vertical line plotted with three points has the same length as a horizontal line plotted with four points.

Raster-Scan Displays

The most common type of graphics monitor employing a CRT is the **raster-scan display**, based on television technology. In a raster-scan system, the electron beam is swept across the screen, one row at a time from top to bottom. As the electron beam moves across each row, the beam intensity is turned on and off to create a pattern of illuminated spots. Picture definition is stored in a memory area called the **refresh buffer** or **frame buffer**. This memory area holds the set of intensity values for all the screen points. Stored intensity values are then retrieved from the refresh buffer and "painted" on the screen one row (**scan line**) at a time (Fig. 2-7). Each screen point is referred to as a **pixel** or **pel** (shortened forms of **picture element**). The capability of a raster-scan system to store intensity information for each screen point makes it well suited for the realistic display of scenes containing subtle shading and color patterns. Home television sets and printers are examples of other systems using raster-scan methods.

Intensity range for pixel positions depends on the capability of the raster system. In a simple black-and-white system, each screen point is either on or off, so only one bit per pixel is needed to control the intensity of screen positions. For a bilevel system, a bit value of 1 indicates that the electron beam is to be turned on at that position, and a value of 0 indicates that the beam intensity is to be off. Additional bits are needed when color and intensity variations can be displayed. Up to 24 bits per pixel are included in high-quality systems, which can require several megabytes of storage for the frame buffer, depending on the resolution of the system. A system with 24 bits per pixel and a screen resolution of 1024 by 1024 requires 3 megabytes of storage for the frame buffer. On a black-and-white system with one bit per pixel, the frame buffer is commonly called a **bitmap**. For systems with multiple bits per pixel, the frame buffer is often referred to as a **pixmap**.

Refreshing on raster-scan displays is carried out at the rate of 60 to 80 frames per second, although some systems are designed for higher refresh rates. Sometimes, refresh rates are described in units of cycles per second, or Hertz (Hz), where a cycle corresponds to one frame. Using these units, we would describe a refresh rate of 60 frames per second as simply 60 Hz. At the end of each scan line, the electron beam returns to the left side of the screen to begin displaying the next scan line. The return to the left of the screen, after refreshing each

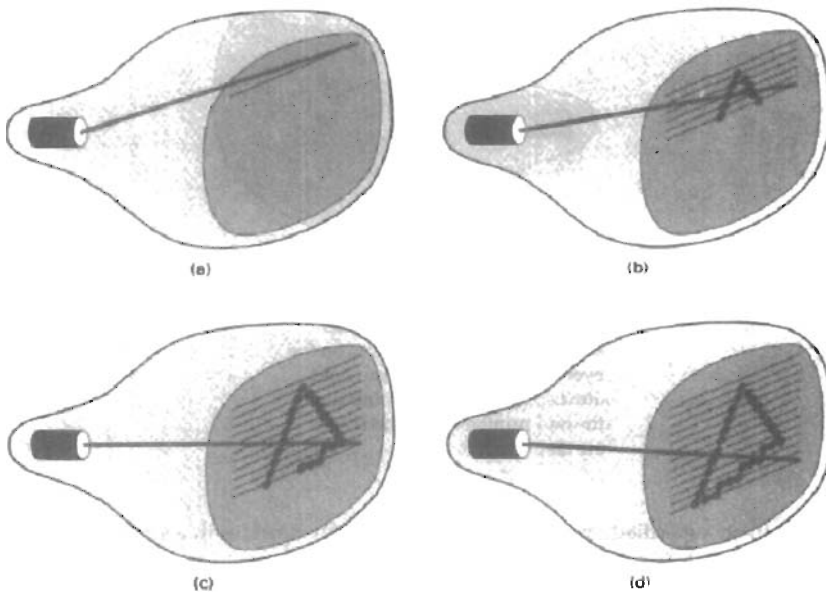


Figure 2-7

A raster-scan system displays an object as a set of discrete points across each scan line.

scan line, is called the **horizontal retrace** of the electron beam. And at the end of each frame (displayed in 1/80th to 1/60th of a second), the electron beam returns (**vertical retrace**) to the top left corner of the screen to begin the next frame.

On some raster-scan systems (and in TV sets), each frame is displayed in two passes using an *interlaced* refresh procedure. In the first pass, the beam sweeps across every other scan line from top to bottom. Then after the vertical retrace, the beam sweeps out the remaining scan lines (Fig. 2-8). Interlacing of the scan lines in this way allows us to see the entire screen displayed in one-half the time it would have taken to sweep across all the lines at once from top to bottom. Interlacing is primarily used with slower refreshing rates. On an older, 30 frame-per-second, noninterlaced display, for instance, some flicker is noticeable. But with interlacing, each of the two passes can be accomplished in 1/60th of a second, which brings the refresh rate nearer to 60 frames per second. This is an effective technique for avoiding flicker, providing that adjacent scan lines contain similar display information.

Random-Scan Displays

When operated as a **random-scan** display unit, a CRT has the electron beam directed only to the parts of the screen where a picture is to be drawn. Random-scan monitors draw a picture one line at a time and for this reason are also referred to as **vector** displays (or **stroke-writing** or **calligraphic** displays). The component lines of a picture can be drawn and refreshed by a random-scan sys-

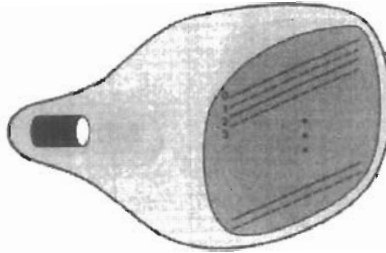


Figure 2-8
Interlacing scan lines on a raster-scan display. First, all points on the even-numbered (solid) scan lines are displayed; then all points along the odd-numbered (dashed) lines are displayed.

tem in any specified order (Fig. 2-9). A pen plotter operates in a similar way and is an example of a random-scan, hard-copy device.

Refresh rate on a random-scan system depends on the number of lines to be displayed. Picture definition is now stored as a set of line-drawing commands in an area of memory referred to as the **refresh display file**. Sometimes the refresh display file is called the **display list**, **display program**, or simply the **refresh buffer**. To display a specified picture, the system cycles through the set of commands in the display file, drawing each component line in turn. After all line-drawing commands have been processed, the system cycles back to the first line command in the list. Random-scan displays are designed to draw all the component lines of a picture 30 to 60 times each second. High-quality vector systems are capable of handling approximately 100,000 "short" lines at this refresh rate. When a small set of lines is to be displayed, each refresh cycle is delayed to avoid refresh rates greater than 60 frames per second. Otherwise, faster refreshing of the set of lines could burn out the phosphor.

Random-scan systems are designed for line-drawing applications and cannot display realistic shaded scenes. Since picture definition is stored as a set of line-drawing instructions and not as a set of intensity values for all screen points, vector displays generally have higher resolution than raster systems. Also, vector displays produce smooth line drawings because the CRT beam directly follows the line path. A raster system, in contrast, produces jagged lines that are plotted as discrete point sets.

Color CRT Monitors

A CRT monitor displays color pictures by using a combination of phosphors that emit different-colored light. By combining the emitted light from the different phosphors, a range of colors can be generated. The two basic techniques for producing color displays with a CRT are the beam-penetration method and the shadow-mask method.

The **beam-penetration** method for displaying color pictures has been used with random-scan monitors. Two layers of phosphor, usually red and green, are

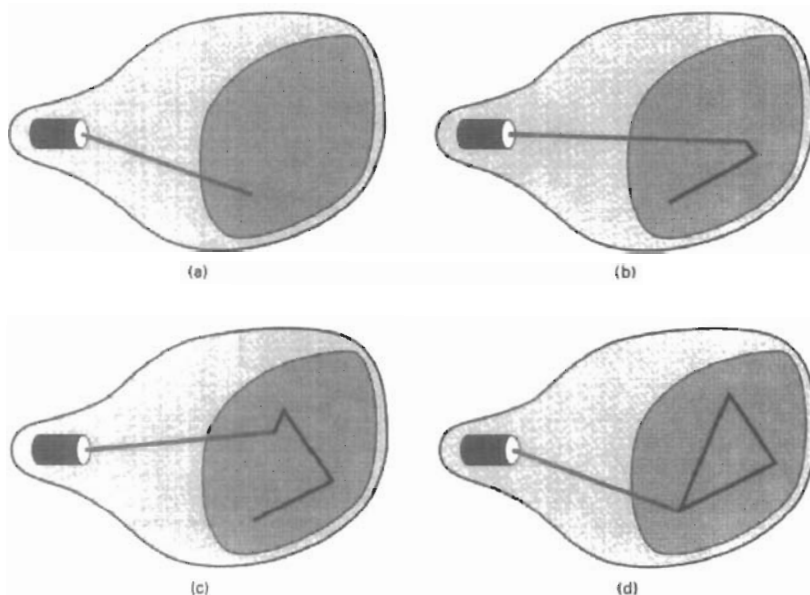


Figure 2-9
A random-scan system draws the component lines of an object in any order specified.

coated onto the inside of the CRT screen, and the displayed color depends on how far the electron beam penetrates into the phosphor layers. A beam of slow electrons excites only the outer red layer. A beam of very fast electrons penetrates through the red layer and excites the inner green layer. At intermediate beam speeds, combinations of red and green light are emitted to show two additional colors, orange and yellow. The speed of the electrons, and hence the screen color at any point, is controlled by the beam-acceleration voltage. Beam penetration has been an inexpensive way to produce color in random-scan monitors, but only four colors are possible, and the quality of pictures is not as good as with other methods.

Shadow-mask methods are commonly used in raster-scan systems (including color TV) because they produce a much wider range of colors than the beam-penetration method. A shadow-mask CRT has three phosphor color dots at each pixel position. One phosphor dot emits a red light, another emits a green light, and the third emits a blue light. This type of CRT has three electron guns, one for each color dot, and a shadow-mask grid just behind the phosphor-coated screen. Figure 2-10 illustrates the *delta-delta* shadow-mask method, commonly used in color CRT systems. The three electron beams are deflected and focused as a group onto the shadow mask, which contains a series of holes aligned with the phosphor-dot patterns. When the three beams pass through a hole in the shadow mask, they activate a dot triangle, which appears as a small color spot on the screen. The phosphor dots in the triangles are arranged so that each electron beam can activate only its corresponding color dot when it passes through the

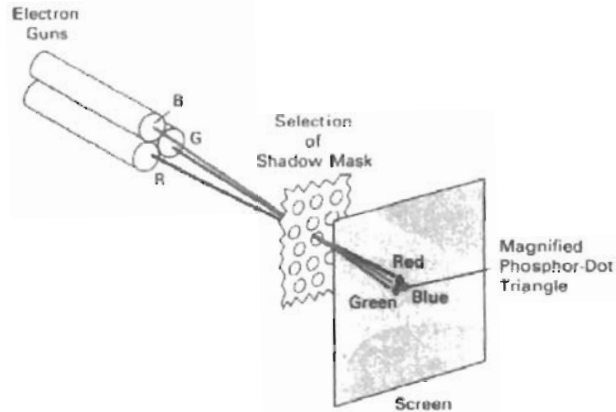


Figure 2-10
Operation of a delta-delta, shadow-mask CRT. Three electron guns, aligned with the triangular color-dot patterns on the screen, are directed to each dot triangle by a shadow mask.

shadow mask. Another configuration for the three electron guns is an *in-line* arrangement in which the three electron guns, and the corresponding red-green-blue color dots on the screen, are aligned along one scan line instead of in a triangular pattern. This in-line arrangement of electron guns is easier to keep in alignment and is commonly used in high-resolution color CRTs.

We obtain color variations in a shadow-mask CRT by varying the intensity levels of the three electron beams. By turning off the red and green guns, we get only the color coming from the blue phosphor. Other combinations of beam intensities produce a small light spot for each pixel position, since our eyes tend to merge the three colors into one composite. The color we see depends on the amount of excitation of the red, green, and blue phosphors. A white (or gray) area is the result of activating all three dots with equal intensity. Yellow is produced with the green and red dots only, magenta is produced with the blue and red dots, and cyan shows up when blue and green are activated equally. In some low-cost systems, the electron beam can only be set to on or off, limiting displays to eight colors. More sophisticated systems can set intermediate intensity levels for the electron beams, allowing several million different colors to be generated.

Color graphics systems can be designed to be used with several types of CRT display devices. Some inexpensive home-computer systems and video games are designed for use with a color TV set and an RF (radio-frequency) modulator. The purpose of the RF modulator is to simulate the signal from a broadcast TV station. This means that the color and intensity information of the picture must be combined and superimposed on the broadcast-frequency carrier signal that the TV needs to have as input. Then the circuitry in the TV takes this signal from the RF modulator, extracts the picture information, and paints it on the screen. As we might expect, this extra handling of the picture information by the RF modulator and TV circuitry decreases the quality of displayed images.

Composite monitors are adaptations of TV sets that allow bypass of the broadcast circuitry. These display devices still require that the picture informa-

tion be combined, but no carrier signal is needed. Picture information is combined into a composite signal and then separated by the monitor, so the resulting picture quality is still not the best attainable.

Color CRTs in graphics systems are designed as **RGB monitors**. These monitors use shadow-mask methods and take the intensity level for each electron gun (red, green, and blue) directly from the computer system without any intermediate processing. High-quality raster-graphics systems have 24 bits per pixel in the frame buffer, allowing 256 voltage settings for each electron gun and nearly 17 million color choices for each pixel. An RGB color system with 24 bits of storage per pixel is generally referred to as a **full-color system** or a **true-color system**.

Direct-View Storage Tubes

An alternative method for maintaining a screen image is to store the picture information inside the CRT instead of refreshing the screen. A **direct-view storage tube (DVST)** stores the picture information as a charge distribution just behind the phosphor-coated screen. Two electron guns are used in a DVST. One, the primary gun, is used to store the picture pattern; the second, the flood gun, maintains the picture display.

A DVST monitor has both disadvantages and advantages compared to the refresh CRT. Because no refreshing is needed, very complex pictures can be displayed at very high resolutions without flicker. Disadvantages of DVST systems are that they ordinarily do not display color and that selected parts of a picture cannot be erased. To eliminate a picture section, the entire screen must be erased and the modified picture redrawn. The erasing and redrawing process can take several seconds for a complex picture. For these reasons, storage displays have been largely replaced by raster systems.

Flat-Panel Displays

Although most graphics monitors are still constructed with CRTs, other technologies are emerging that may soon replace CRT monitors. The term **flat-panel display** refers to a class of video devices that have reduced volume, weight, and power requirements compared to a CRT. A significant feature of flat-panel displays is that they are thinner than CRTs, and we can hang them on walls or wear them on our wrists. Since we can even write on some flat-panel displays, they will soon be available as pocket notepads. Current uses for flat-panel displays include small TV monitors, calculators, pocket video games, laptop computers, armrest viewing of movies on airlines, as advertisement boards in elevators, and as graphics displays in applications requiring rugged, portable monitors.

We can separate flat-panel displays into two categories: **emissive displays** and **nonemissive displays**. The emissive displays (or **emitters**) are devices that convert electrical energy into light. Plasma panels, thin-film electroluminescent displays, and light-emitting diodes are examples of emissive displays. Flat CRTs have also been devised, in which electron beams are accelerated parallel to the screen, then deflected 90° to the screen. But flat CRTs have not proved to be as successful as other emissive devices. Nonemissive displays (or **nonemitters**) use optical effects to convert sunlight or light from some other source into graphics patterns. The most important example of a nonemissive flat-panel display is a liquid-crystal device.

Plasma panels, also called **gas-discharge displays**, are constructed by filling the region between two glass plates with a mixture of gases that usually in-

cludes neon. A series of vertical conducting ribbons is placed on one glass panel, and a set of horizontal ribbons is built into the other glass panel (Fig. 2-11). Firing voltages applied to a pair of horizontal and vertical conductors cause the gas at the intersection of the two conductors to break down into a glowing plasma of electrons and ions. Picture definition is stored in a refresh buffer, and the firing voltages are applied to refresh the pixel positions (at the intersections of the conductors) 60 times per second. Alternating-current methods are used to provide faster application of the firing voltages, and thus brighter displays. Separation between pixels is provided by the electric field of the conductors. Figure 2-12 shows a high-definition plasma panel. One disadvantage of plasma panels has been that they were strictly monochromatic devices, but systems have been developed that are now capable of displaying color and grayscale.

Thin-film electroluminescent displays are similar in construction to a plasma panel. The difference is that the region between the glass plates is filled with a phosphor, such as zinc sulfide doped with manganese, instead of a gas (Fig. 2-13). When a sufficiently high voltage is applied to a pair of crossing electrodes, the phosphor becomes a conductor in the area of the intersection of the two electrodes. Electrical energy is then absorbed by the manganese atoms, which then release the energy as a spot of light similar to the glowing plasma effect in a plasma panel. Electroluminescent displays require more power than plasma panels, and good color and gray scale displays are hard to achieve.

A third type of emissive device is the **light-emitting diode (LED)**. A matrix of diodes is arranged to form the pixel positions in the display, and picture definition is stored in a refresh buffer. As in scan-line refreshing of a CRT, information

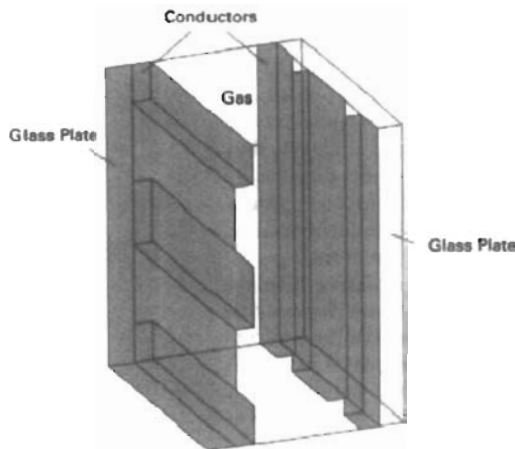


Figure 2-11
Basic design of a plasma-panel display device.

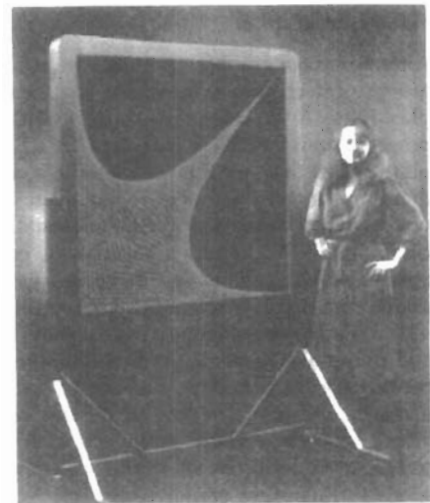


Figure 2-12
A plasma-panel display with a resolution of 2048 by 2048 and a screen diagonal of 1.5 meters.
(Courtesy of Photonics Systems.)

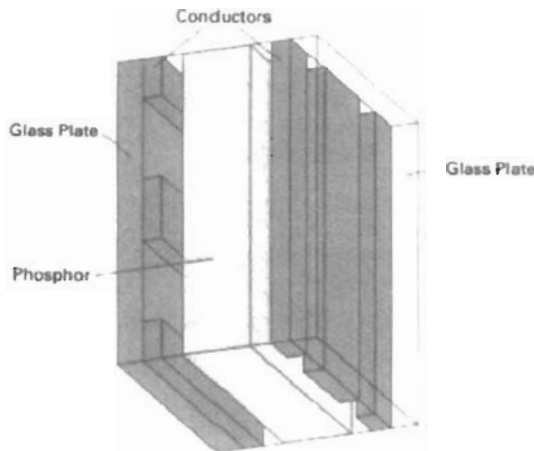


Figure 2-13
Basic design of a thin-film
electroluminescent display device.

is read from the refresh buffer and converted to voltage levels that are applied to the diodes to produce the light patterns in the display.

Liquid-crystal displays (LCDs) are commonly used in small systems, such as calculators (Fig. 2-14) and portable, laptop computers (Fig. 2-15). These non-emissive devices produce a picture by passing polarized light from the surroundings or from an internal light source through a liquid-crystal material that can be aligned to either block or transmit the light.

The term *liquid crystal* refers to the fact that these compounds have a crystalline arrangement of molecules, yet they flow like a liquid. Flat-panel displays commonly use nematic (threadlike) liquid-crystal compounds that tend to keep the long axes of the rod-shaped molecules aligned. A flat-panel display can then be constructed with a nematic liquid crystal, as demonstrated in Fig. 2-16. Two glass plates, each containing a light polarizer at right angles to the other plate, sandwich the liquid-crystal material. Rows of horizontal transparent conductors are built into one glass plate, and columns of vertical conductors are put into the other plate. The intersection of two conductors defines a pixel position. Normally, the molecules are aligned as shown in the "on state" of Fig. 2-16. Polarized light passing through the material is twisted so that it will pass through the opposite polarizer. The light is then reflected back to the viewer. To turn off the pixel, we apply a voltage to the two intersecting conductors to align the molecules so that the light is not twisted. This type of flat-panel device is referred to as a **passive-matrix LCD**. Picture definitions are stored in a refresh buffer, and the screen is refreshed at the rate of 60 frames per second, as in the emissive devices. Back lighting is also commonly applied using solid-state electronic devices, so that the system is not completely dependent on outside light sources. Colors can be displayed by using different materials or dyes and by placing a triad of color pixels at each screen location. Another method for constructing LCDs is to place a transistor at each pixel location, using thin-film transistor technology. The transistors are used to control the voltage at pixel locations and to prevent charge from gradually leaking out of the liquid-crystal cells. These devices are called **active-matrix displays**.



Figure 2-14
A hand calculator with an
LCD screen. (Courtesy of Texas
Instruments.)



Figure 2-15
A backlit, passive-matrix, liquid-crystal display in a laptop computer, featuring 256 colors, a screen resolution of 640 by 400, and a screen diagonal of 9 inches.
(Courtesy of Apple Computer, Inc.)

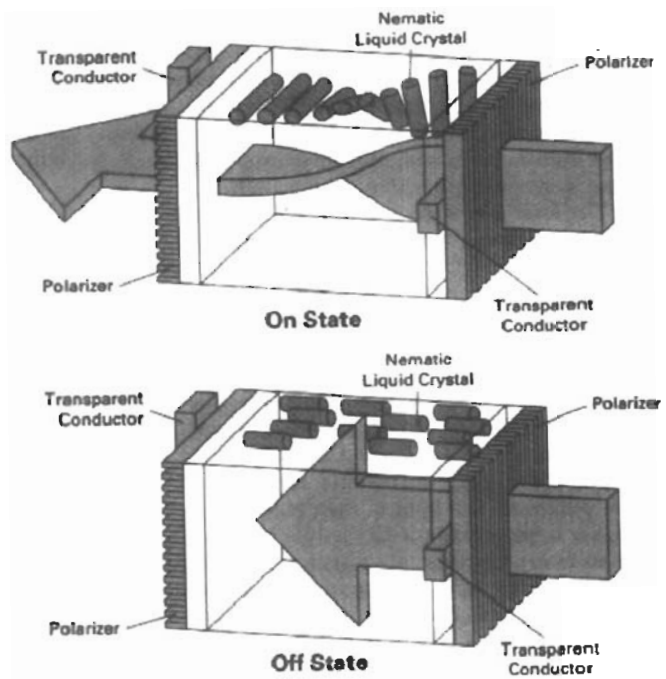


Figure 2-16
The light-twisting, shutter effect used in the design of most liquid-crystal display devices.

Graphics monitors for the display of three-dimensional scenes have been devised using a technique that reflects a CRT image from a vibrating, flexible mirror. The operation of such a system is demonstrated in Fig. 2-17. As the varifocal mirror vibrates, it changes focal length. These vibrations are synchronized with the display of an object on a CRT so that each point on the object is reflected from the mirror into a spatial position corresponding to the distance of that point from a specified viewing position. This allows us to walk around an object or scene and view it from different sides.

Figure 2-18 shows the Genisco SpaceGraph system, which uses a vibrating mirror to project three-dimensional objects into a 25-cm by 25-cm by 25-cm volume. This system is also capable of displaying two-dimensional cross-sectional "slices" of objects selected at different depths. Such systems have been used in medical applications to analyze data from ultrasonography and CAT scan devices, in geological applications to analyze topological and seismic data, in design applications involving solid objects, and in three-dimensional simulations of systems, such as molecules and terrain.

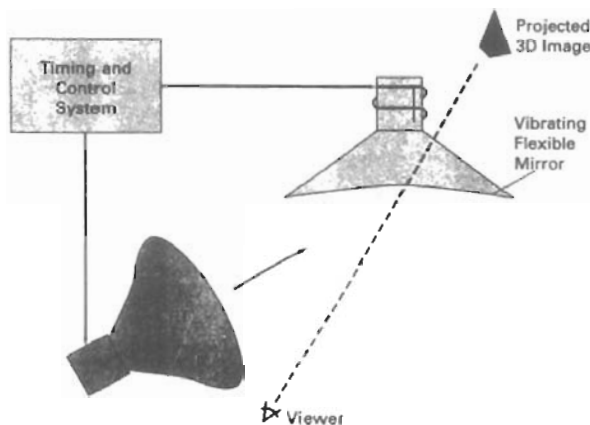


Figure 2-17
Operation of a three-dimensional display system using a vibrating mirror that changes focal length to match the depth of points in a scene.



Figure 2-18
The SpaceGraph interactive graphics system displays objects in three dimensions using a vibrating, flexible mirror. (Courtesy of Genisco Computers Corporation.)

Stereoscopic and Virtual-Reality Systems

Another technique for representing three-dimensional objects is displaying stereoscopic views. This method does not produce true three-dimensional images, but it does provide a three-dimensional effect by presenting a different view to each eye of an observer so that scenes do appear to have depth (Fig. 2-19).

To obtain a stereoscopic projection, we first need to obtain two views of a scene generated from a viewing direction corresponding to each eye (left and right). We can construct the two views as computer-generated scenes with different viewing positions, or we can use a stereo camera pair to photograph some object or scene. When we simultaneously look at the left view with the left eye and the right view with the right eye, the two views merge into a single image and we perceive a scene with depth. Figure 2-20 shows two views of a computer-generated scene for stereographic projection. To increase viewing comfort, the areas at the left and right edges of this scene that are visible to only one eye have been eliminated.



Figure 2-19
Viewing a stereoscopic projection.
(Courtesy of StereoGraphics Corporation.)

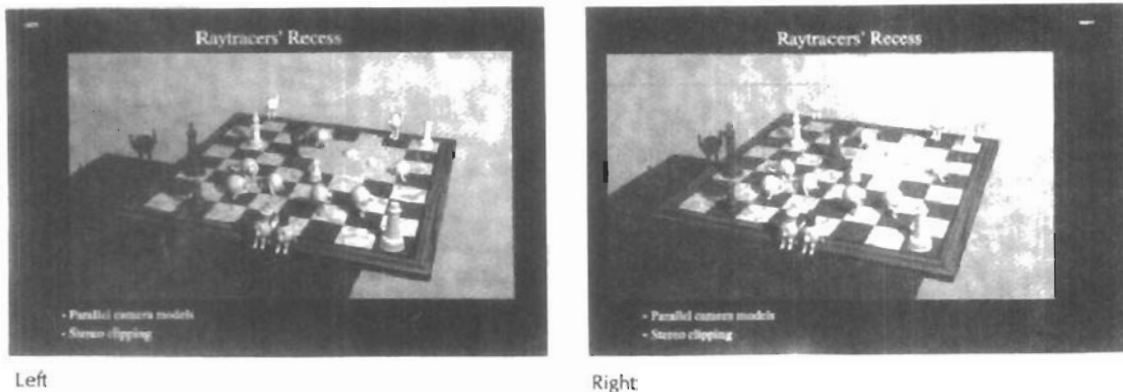


Figure 2-20
A stereoscopic viewing pair. (Courtesy of Jerry Farm.)

One way to produce a stereoscopic effect is to display each of the two views with a raster system on alternate refresh cycles. The screen is viewed through glasses, with each lens designed to act as a rapidly alternating shutter that is synchronized to block out one of the views. Figure 2-21 shows a pair of stereoscopic glasses constructed with liquid-crystal shutters and an infrared emitter that synchronizes the glasses with the views on the screen.

Stereoscopic viewing is also a component in **virtual-reality** systems, where users can step into a scene and interact with the environment. A headset (Fig. 2-22) containing an optical system to generate the stereoscopic views is commonly used in conjunction with interactive input devices to locate and manipulate objects in the scene. A sensing system in the headset keeps track of the viewer's position, so that the front and back of objects can be seen as the viewer

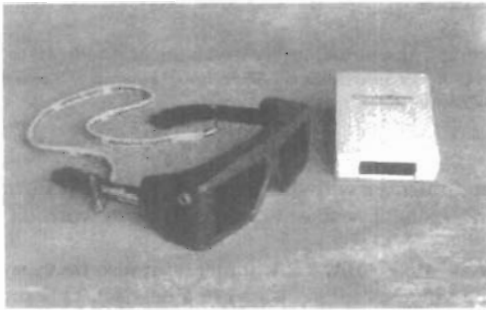


Figure 2-21
Glasses for viewing a
stereoscopic scene and an
infrared synchronizing emitter.
(Courtesy of StereoGraphics Corporation.)



Figure 2-22
A headset used in virtual-reality systems. (Courtesy of Virtual
Research.)



Figure 2-23
Interacting with a virtual-reality environment. (Courtesy of the
National Center for Supercomputing Applications, University of Illinois at
Urbana-Champaign.)

“walks through” and interacts with the display. Figure 2-23 illustrates interaction with a virtual scene, using a headset and a data glove worn on the right hand (Section 2-5).

An interactive virtual-reality environment can also be viewed with stereoscopic glasses and a video monitor, instead of a headset. This provides a means for obtaining a lower-cost virtual-reality system. As an example, Fig. 2-24 shows an ultrasound tracking device with six degrees of freedom. The tracking device is placed on top of the video display and is used to monitor head movements so that the viewing position for a scene can be changed as head position changes.

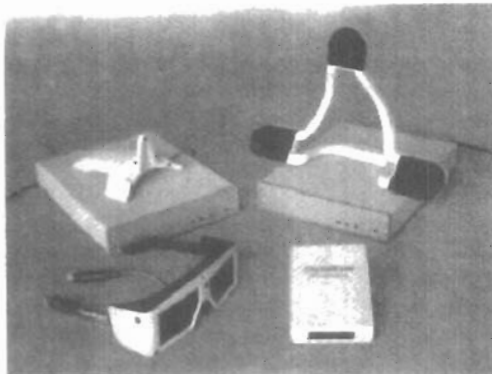


Figure 2-24
An ultrasound tracking device used
with stereoscopic glasses to track
head position. (Courtesy of
StereoGraphics Corporation.)

Interactive raster graphics systems typically employ several processing units. In addition to the central processing unit, or CPU, a special-purpose processor, called the **video controller** or **display controller**, is used to control the operation of the display device. Organization of a simple raster system is shown in Fig. 2-25. Here, the frame buffer can be anywhere in the system memory, and the video controller accesses the frame buffer to refresh the screen. In addition to the video controller, more sophisticated raster systems employ other processors as co-processors and accelerators to implement various graphics operations.

Video Controller

Figure 2-26 shows a commonly used organization for raster systems. A fixed area of the system memory is reserved for the frame buffer, and the video controller is given direct access to the frame-buffer memory.

Frame-buffer locations, and the corresponding screen positions, are referenced in Cartesian coordinates. For many graphics monitors, the coordinate ori-

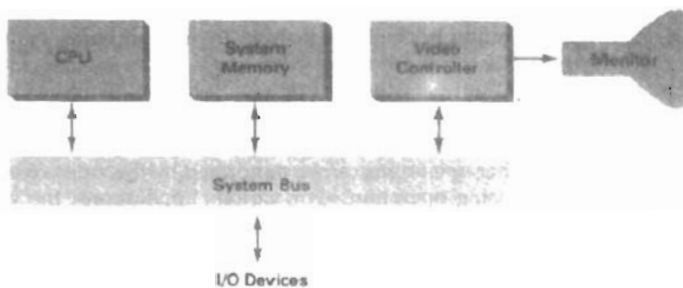


Figure 2-25

Architecture of a simple raster graphics system.

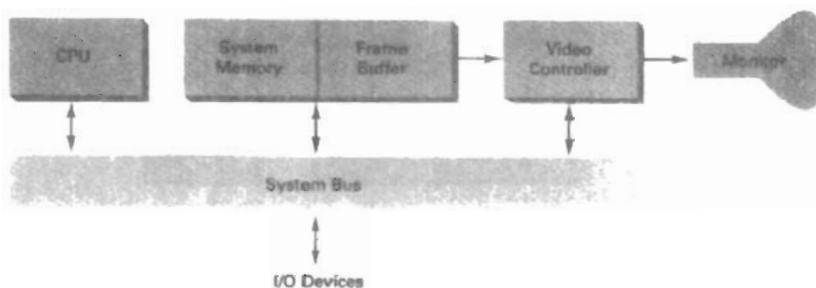


Figure 2-26

Architecture of a raster system with a fixed portion of the system memory reserved for the frame buffer.

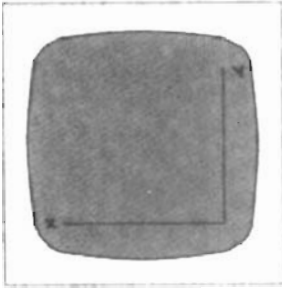


Figure 2-27

The origin of the coordinate system for identifying screen positions is usually specified in the lower-left corner.

gin is defined at the lower left screen corner (Fig. 2-27). The screen surface is then represented as the first quadrant of a two-dimensional system, with positive x values increasing to the right and positive y values increasing from bottom to top. (On some personal computers, the coordinate origin is referenced at the upper left corner of the screen, so the y values are inverted.) Scan lines are then labeled from y_{\max} at the top of the screen to 0 at the bottom. Along each scan line, screen pixel positions are labeled from 0 to x_{\max} .

In Fig. 2-28, the basic refresh operations of the video controller are diagrammed. Two registers are used to store the coordinates of the screen pixels. Initially, the x register is set to 0 and the y register is set to y_{\max} . The value stored in the frame buffer for this pixel position is then retrieved and used to set the intensity of the CRT beam. Then the x register is incremented by 1, and the process repeated for the next pixel on the top scan line. This procedure is repeated for each pixel along the scan line. After the last pixel on the top scan line has been processed, the x register is reset to 0 and the y register is decremented by 1. Pixels along this scan line are then processed in turn, and the procedure is repeated for each successive scan line. After cycling through all pixels along the bottom scan line ($y = 0$), the video controller resets the registers to the first pixel position on the top scan line and the refresh process starts over.

Since the screen must be refreshed at the rate of 60 frames per second, the simple procedure illustrated in Fig. 2-28 cannot be accommodated by typical RAM chips. The cycle time is too slow. To speed up pixel processing, video controllers can retrieve multiple pixel values from the refresh buffer on each pass. The multiple pixel intensities are then stored in a separate register and used to control the CRT beam intensity for a group of adjacent pixels. When that group of pixels has been processed, the next block of pixel values is retrieved from the frame buffer.

A number of other operations can be performed by the video controller, besides the basic refreshing operations. For various applications, the video con-

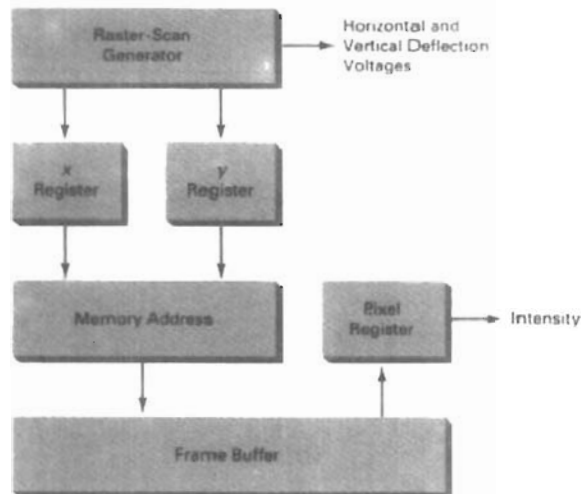


Figure 2-28

Basic video-controller refresh operations.

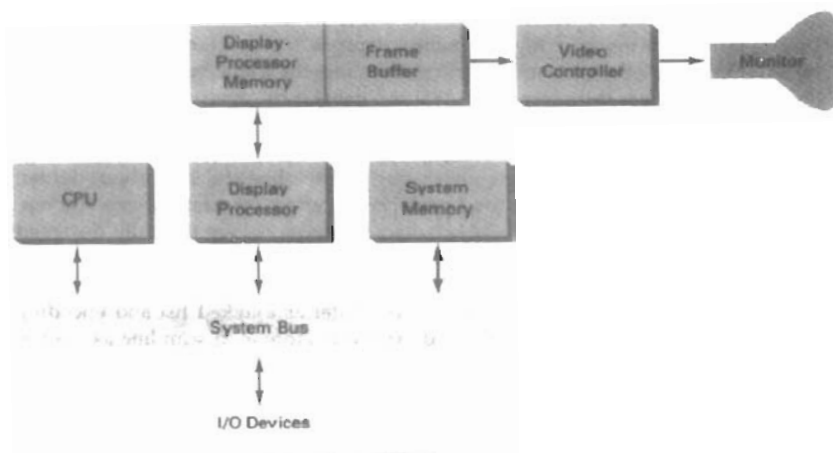


Figure 2-29
Architecture of a raster-graphics system with a display processor.

troller can retrieve pixel intensities from different memory areas on different refresh cycles. In high-quality systems, for example, two frame buffers are often provided so that one buffer can be used for refreshing while the other is being filled with intensity values. Then the two buffers can switch roles. This provides a fast mechanism for generating real-time animations, since different views of moving objects can be successively loaded into the refresh buffers. Also, some transformations can be accomplished by the video controller. Areas of the screen can be enlarged, reduced, or moved from one location to another during the refresh cycles. In addition, the video controller often contains a lookup table, so that pixel values in the frame buffer are used to access the lookup table instead of controlling the CRT beam intensity directly. This provides a fast method for changing screen intensity values, and we discuss lookup tables in more detail in Chapter 4. Finally, some systems are designed to allow the video controller to mix the frame-buffer image with an input image from a television camera or other input device.

Raster-Scan Display Processor

Figure 2-29 shows one way to set up the organization of a raster system containing a separate display processor, sometimes referred to as a graphics controller or a display coprocessor. The purpose of the display processor is to free the CPU from the graphics chores. In addition to the system memory, a separate display-processor memory area can also be provided.

A major task of the display processor is digitizing a picture definition given in an application program into a set of pixel-intensity values for storage in the frame buffer. This digitization process is called **scan conversion**. Graphics commands specifying straight lines and other geometric objects are scan converted into a set of discrete intensity points. Scan converting a straight-line segment, for example, means that we have to locate the pixel positions closest to the line path and store the intensity for each position in the frame buffer. Similar methods are used for scan converting curved lines and polygon outlines. Characters can be defined with rectangular grids, as in Fig. 2-30, or they can be defined with curved

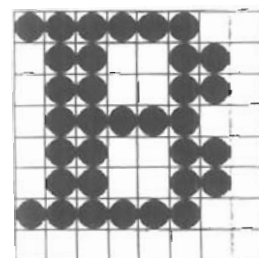


Figure 2-30
A character defined as a rectangular grid of pixel positions.

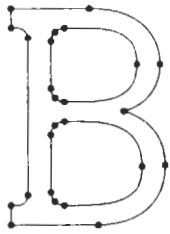


Figure 2-31
A character defined as a
curve outline.

outlines, as in Fig. 2-31. The array size for character grids can vary from about 5 by 7 to 9 by 12 or more for higher-quality displays. A character grid is displayed by superimposing the rectangular grid pattern into the frame buffer at a specified coordinate position. With characters that are defined as curve outlines, character shapes are scan converted into the frame buffer.

Display processors are also designed to perform a number of additional operations. These functions include generating various line styles (dashed, dotted, or solid), displaying color areas, and performing certain transformations and manipulations on displayed objects. Also, display processors are typically designed to interface with interactive input devices, such as a mouse.

In an effort to reduce memory requirements in raster systems, methods have been devised for organizing the frame buffer as a linked list and encoding the intensity information. One way to do this is to store each scan line as a set of integer pairs. One number of each pair indicates an intensity value, and the second number specifies the number of adjacent pixels on the scan line that are to have that intensity. This technique, called **run-length encoding**, can result in a considerable saving in storage space if a picture is to be constructed mostly with long runs of a single color each. A similar approach can be taken when pixel intensities change linearly. Another approach is to encode the raster as a set of rectangular areas (**cell encoding**). The disadvantages of encoding runs are that intensity changes are difficult to make and storage requirements actually increase as the length of the runs decreases. In addition, it is difficult for the display controller to process the raster when many short runs are involved.

2-3

RANDOM-SCAN SYSTEMS

The organization of a simple random-scan (vector) system is shown in Fig. 2-32. An application program is input and stored in the system memory along with a graphics package. Graphics commands in the application program are translated by the graphics package into a display file stored in the system memory. This display file is then accessed by the display processor to refresh the screen. The display processor cycles through each command in the display file program once during every refresh cycle. Sometimes the display processor in a random-scan system is referred to as a **display processing unit** or a **graphics controller**.

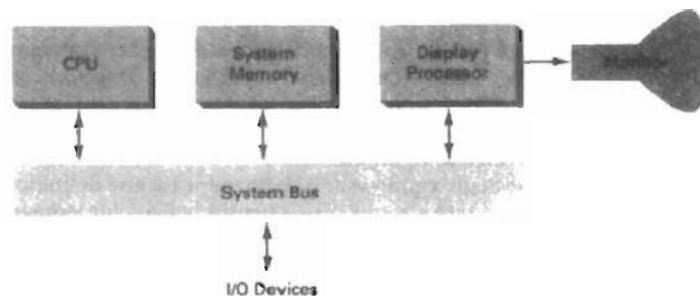


Figure 2-32
Architecture of a simple random-scan system.

Graphics patterns are drawn on a random-scan system by directing the electron beam along the component lines of the picture. Lines are defined by the values for their coordinate endpoints, and these input coordinate values are converted to x and y deflection voltages. A scene is then drawn one line at a time by positioning the beam to fill in the line between specified endpoints.

2-4

GRAPHICS MONITORS AND WORKSTATIONS

Most graphics monitors today operate as raster-scan displays, and here we survey a few of the many graphics hardware configurations available. Graphics systems range from small general-purpose computer systems with graphics capabilities (Fig. 2-33) to sophisticated full-color systems that are designed specifically for graphics applications (Fig. 2-34). A typical screen resolution for personal com-



Figure 2-33
A desktop general-purpose
computer system that can be used
for graphics applications. (Courtesy of
Apple Computer, Inc.)

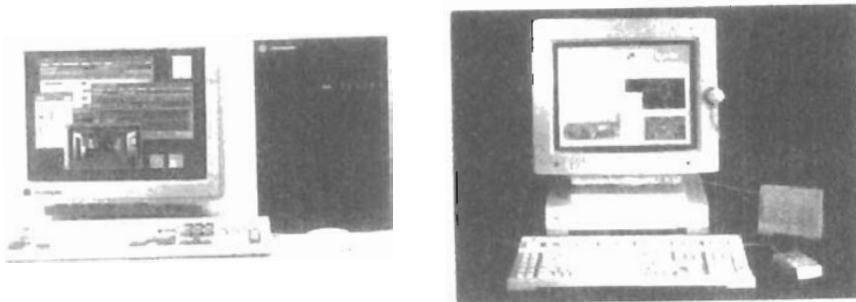


Figure 2-34
Computer graphics workstations with keyboard and mouse input devices. (a) The Iris
Indigo. (Courtesy of Silicon Graphics Corporation.) (b) SPARCstation 10. (Courtesy of Sun Microsystems.)

puter systems, such as the Apple Quadra shown in Fig. 2-33, is 640 by 480, although screen resolution and other system capabilities vary depending on the size and cost of the system. Diagonal screen dimensions for general-purpose personal computer systems can range from 12 to 21 inches, and allowable color selections range from 16 to over 32,000. For workstations specifically designed for graphics applications, such as the systems shown in Fig. 2-34, typical screen resolution is 1280 by 1024, with a screen diagonal of 16 inches or more. Graphics workstations can be configured with from 8 to 24 bits per pixel (full-color systems), with higher screen resolutions, faster processors, and other options available in high-end systems.

Figure 2-35 shows a high-definition graphics monitor used in applications such as air traffic control, simulation, medical imaging, and CAD. This system has a diagonal screen size of 27 inches, resolutions ranging from 2048 by 1536 to 2560 by 2048, with refresh rates of 80 Hz or 60 Hz noninterlaced.

A multiscreen system called the MediaWall, shown in Fig. 2-36, provides a large "wall-sized" display area. This system is designed for applications that require large area displays in brightly lighted environments, such as at trade shows, conventions, retail stores, museums, or passenger terminals. MediaWall operates by splitting images into a number of sections and distributing the sections over an array of monitors or projectors using a graphics adapter and satellite control units. An array of up to 5 by 5 monitors, each with a resolution of 640 by 480, can be used in the MediaWall to provide an overall resolution of 3200 by 2400 for either static scenes or animations. Scenes can be displayed behind mullions, as in Fig. 2-36, or the mullions can be eliminated to display a continuous picture with no breaks between the various sections.

Many graphics workstations, such as some of those shown in Fig. 2-37, are configured with two monitors. One monitor can be used to show all features of an object or scene, while the second monitor displays the detail in some part of the picture. Another use for dual-monitor systems is to view a picture on one monitor and display graphics options (menus) for manipulating the picture components on the other monitor.



Figure 2-35
A very high-resolution (2560 by 2048) color monitor. (Courtesy of BARCO Chromatics.)

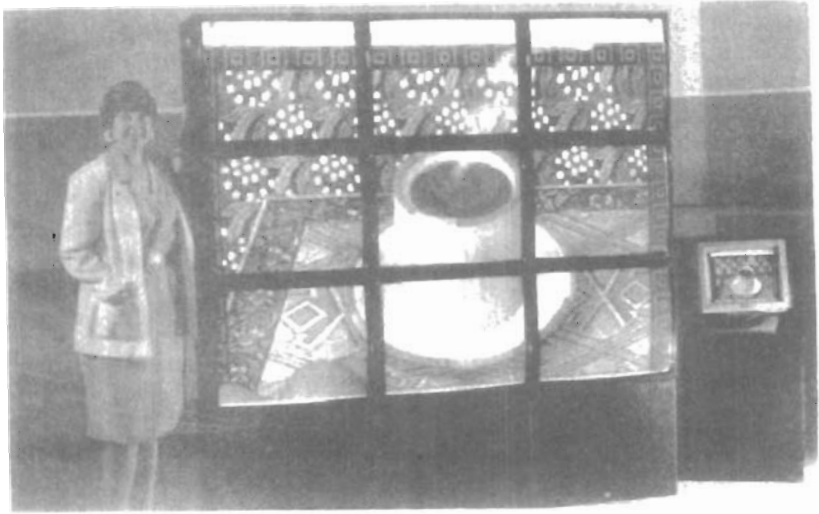


Figure 2-36
The MediaWall: A multiscreen display system. The image displayed on this 3-by-3 array of monitors was created by Deneba Software. (Courtesy of RGB Spectrum.)



Figure 2-37
Single- and dual-monitor graphics workstations. (Courtesy of Intergraph Corporation.)

Figures 2-38 and 2-39 illustrate examples of interactive graphics workstations containing multiple input and other devices. A typical setup for CAD applications is shown in Fig. 2-38. Various keyboards, button boxes, tablets, and mice are attached to the video monitors for use in the design process. Figure 2-39 shows features of some types of artist's workstations.



Figure 2-38

Multiple workstations for a CAD group. (Courtesy of Hewlett-Packard Company.)



Figure 2-39

An artist's workstation, featuring a color raster monitor, keyboard, graphics tablet with hand cursor, and a light table, in addition to data storage and telecommunications devices. (Courtesy of DICOMED Corporation.)

2-5

INPUT DEVICES

Various devices are available for data input on graphics workstations. Most systems have a keyboard and one or more additional devices specially designed for interactive input. These include a mouse, trackball, spaceball, joystick, digitizers,

dials, and button boxes. Some other input devices used in particular applications are data gloves, touch panels, image scanners, and voice systems.

Keyboards

An alphanumeric keyboard on a graphics system is used primarily as a device for entering text strings. The keyboard is an efficient device for inputting such nongraphic data as picture labels associated with a graphics display. Keyboards can also be provided with features to facilitate entry of screen coordinates, menu selections, or graphics functions.

Cursor-control keys and function keys are common features on general-purpose keyboards. Function keys allow users to enter frequently used operations in a single keystroke, and cursor-control keys can be used to select displayed objects or coordinate positions by positioning the screen cursor. Other types of cursor-positioning devices, such as a trackball or joystick, are included on some keyboards. Additionally, a numeric keypad is often included on the keyboard for fast entry of numeric data. Typical examples of general-purpose keyboards are given in Figs. 2-1, 2-33, and 2-34. Fig. 2-40 shows an ergonomic keyboard design.

For specialized applications, input to a graphics application may come from a set of buttons, dials, or switches that select data values or customized graphics operations. Figure 2-41 gives an example of a button box and a set of input dials. Buttons and switches are often used to input predefined functions, and dials are common devices for entering scalar values. Real numbers within some defined range are selected for input with dial rotations. Potentiometers are used to measure dial rotations, which are then converted to deflection voltages for cursor movement.

Mouse

A **mouse** is small hand-held box used to position the screen cursor. Wheels or rollers on the bottom of the mouse can be used to record the amount and direc-



Figure 2-40
Ergonomically designed keyboard with removable palm rests. The slope of each half of the keyboard can be adjusted separately. (Courtesy of Apple Computer, Inc.)

tion of movement. Another method for detecting mouse motion is with an optical sensor. For these systems, the mouse is moved over a special mouse pad that has a grid of horizontal and vertical lines. The optical sensor detects movement across the lines in the grid.

Since a mouse can be picked up and put down at another position without change in cursor movement, it is used for making relative changes in the position of the screen cursor. One, two, or three buttons are usually included on the top of the mouse for signaling the execution of some operation, such as recording cursor position or invoking a function. Most general-purpose graphics systems now include a mouse and a keyboard as the major input devices, as in Figs. 2-1, 2-33, and 2-34.

Additional devices can be included in the basic mouse design to increase the number of allowable input parameters. The Z mouse in Fig. 2-42 includes

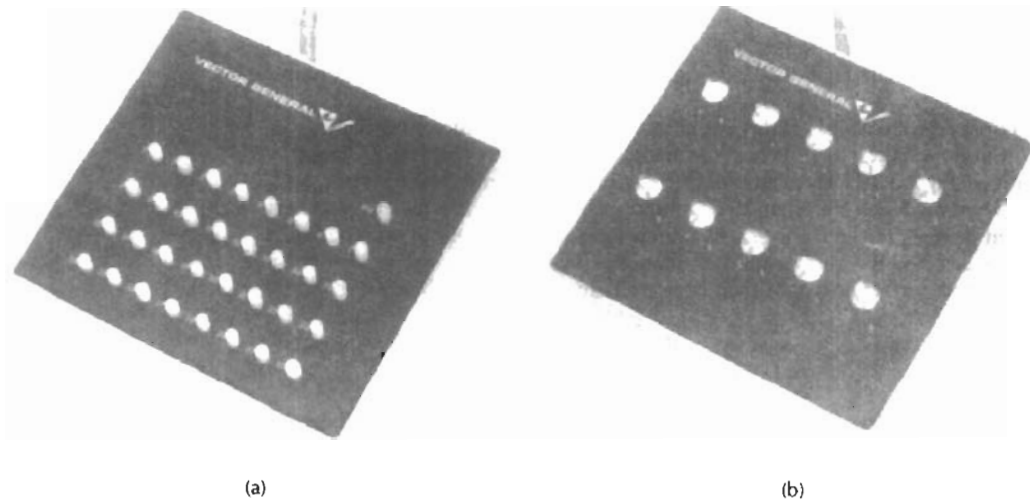


Figure 2-41
A button box (a) and a set of input dials (b). (Courtesy of Vector General.)



Figure 2-42
The Z mouse features three buttons, a mouse ball underneath, a thumbwheel on the side, and a trackball on top. (Courtesy of Multipoint Technology Corporation.)

three buttons, a thumbwheel on the side, a trackball on the top, and a standard mouse ball underneath. This design provides six degrees of freedom to select spatial positions, rotations, and other parameters. With the Z mouse, we can pick up an object, rotate it, and move it in any direction, or we can navigate our viewing position and orientation through a three-dimensional scene. Applications of the Z mouse include virtual reality, CAD, and animation.

Trackball and Spaceball

As the name implies, a **trackball** is a ball that can be rotated with the fingers or palm of the hand, as in Fig. 2-43, to produce screen-cursor movement. Potentiometers, attached to the ball, measure the amount and direction of rotation. Trackballs are often mounted on keyboards (Fig. 2-15) or other devices such as the Z mouse (Fig. 2-42).

While a trackball is a two-dimensional positioning device, a **spaceball** (Fig. 2-45) provides six degrees of freedom. Unlike the trackball, a spaceball does not actually move. Strain gauges measure the amount of pressure applied to the spaceball to provide input for spatial positioning and orientation as the ball is pushed or pulled in various directions. Spaceballs are used for three-dimensional positioning and selection operations in virtual-reality systems, modeling, animation, CAD, and other applications.

Joysticks

A **joystick** consists of a small, vertical lever (called the stick) mounted on a base that is used to steer the screen cursor around. Most joysticks select screen positions with actual stick movement; others respond to pressure on the stick. Figure 2-44 shows a movable joystick. Some joysticks are mounted on a keyboard; others function as stand-alone units.

The distance that the stick is moved in any direction from its center position corresponds to screen-cursor movement in that direction. Potentiometers mounted at the base of the joystick measure the amount of movement, and springs return the stick to the center position when it is released. One or more buttons can be programmed to act as input switches to signal certain actions once a screen position has been selected.



Figure 2-43
A three-button track ball. (Courtesy of Measurement Systems Inc., Norwalk, Connecticut.)

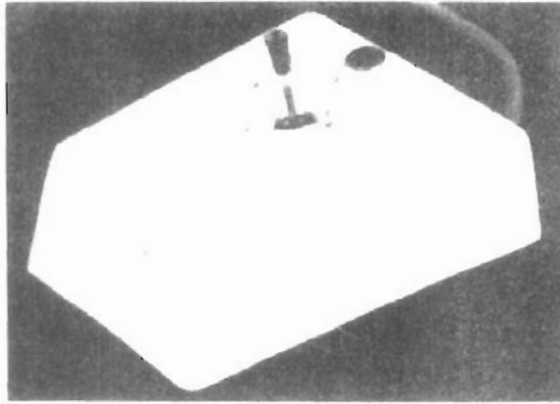


Figure 2-44
A moveable joystick. (Courtesy of CalComp Group; Sanders Associates, Inc.)

In another type of movable joystick, the stick is used to activate switches that cause the screen cursor to move at a constant rate in the direction selected. Eight switches, arranged in a circle, are sometimes provided, so that the stick can select any one of eight directions for cursor movement. Pressure-sensitive joysticks, also called isometric joysticks, have a nonmovable stick. Pressure on the stick is measured with strain gauges and converted to movement of the cursor in the direction specified.

Data Glove

Figure 2-45 shows a **data glove** that can be used to grasp a “virtual” object. The glove is constructed with a series of sensors that detect hand and finger motions. Electromagnetic coupling between transmitting antennas and receiving antennas is used to provide information about the position and orientation of the hand. The transmitting and receiving antennas can each be structured as a set of three mutually perpendicular coils, forming a three-dimensional Cartesian coordinate system. Input from the glove can be used to position or manipulate objects in a virtual scene. A two-dimensional projection of the scene can be viewed on a video monitor, or a three-dimensional projection can be viewed with a headset.

Digitizers

A common device for drawing, painting, or interactively selecting coordinate positions on an object is a **digitizer**. These devices can be used to input coordinate values in either a two-dimensional or a three-dimensional space. Typically, a digitizer is used to scan over a drawing or object and to input a set of discrete coordinate positions, which can be joined with straight-line segments to approximate the curve or surface shapes.

One type of digitizer is the **graphics tablet** (also referred to as a data tablet), which is used to input two-dimensional coordinates by activating a hand cursor or stylus at selected positions on a flat surface. A hand cursor contains cross hairs for sighting positions, while a stylus is a pencil-shaped device that is pointed at



Figure 2-45
A virtual-reality scene, displayed on a two-dimensional video monitor, with input from a data glove and a spaceball. (Courtesy of The Computer Graphics Center, Darmstadt, Germany.)

positions on the tablet. Figures 2-46 and 2-47 show examples of desktop and floor-model tablets, using hand cursors that are available with 2, 4, or 16 buttons. Examples of stylus input with a tablet are shown in Figs. 2-48 and 2-49. The artist's digitizing system in Fig. 2-49 uses electromagnetic resonance to detect the three-dimensional position of the stylus. This allows an artist to produce different brush strokes with different pressures on the tablet surface. Tablet size varies from 12 by 12 inches for desktop models to 44 by 60 inches or larger for floor models. Graphics tablets provide a highly accurate method for selecting coordinate positions, with an accuracy that varies from about 0.2 mm on desktop models to about 0.05 mm or less on larger models.

Many graphics tablets are constructed with a rectangular grid of wires embedded in the tablet surface. Electromagnetic pulses are generated in sequence

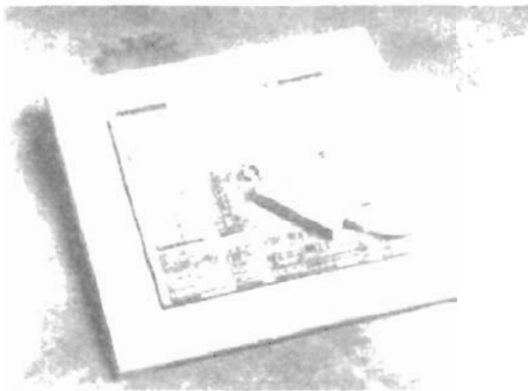


Figure 2-46
The SummaSketch III desktop tablet with a 16-button hand cursor. (Courtesy of Summagraphics Corporation.)



Figure 2-47
The Microgrid III tablet with a 16-button hand cursor, designed for digitizing larger drawings. (Courtesy of Summagraphics Corporation.)



Figure 2-48
The NotePad desktop tablet with stylus. (Courtesy of CalComp Digitizer Division, a part of CalComp, Inc.)

along the wires, and an electric signal is induced in a wire coil in an activated stylus or hand cursor to record a tablet position. Depending on the technology, either signal strength, coded pulses, or phase shifts can be used to determine the position on the tablet.

Acoustic (or sonic) tablets use sound waves to detect a stylus position. Either strip microphones or point microphones can be used to detect the sound emitted by an electrical spark from a stylus tip. The position of the stylus is calcu-



Figure 2-49
An artist's digitizer system, with a pressure-sensitive, cordless stylus. (Courtesy of Wacom Technology Corporation.)

lated by timing the arrival of the generated sound at the different microphone positions. An advantage of two-dimensional acoustic tablets is that the microphones can be placed on any surface to form the "tablet" work area. This can be convenient for various applications, such as digitizing drawings in a book.

Three-dimensional digitizers use sonic or electromagnetic transmissions to record positions. One electromagnetic transmission method is similar to that used in the data glove: A coupling between the transmitter and receiver is used to compute the location of a stylus as it moves over the surface of an object. Figure 2-50 shows a three-dimensional digitizer designed for Apple Macintosh computers. As the points are selected on a nonmetallic object, a wireframe outline of the surface is displayed on the computer screen. Once the surface outline is constructed, it can be shaded with lighting effects to produce a realistic display of the object. Resolution of this system is from 0.8 mm to 0.08 mm, depending on the model.

Image Scanners

Drawings, graphs, color and black-and-white photos, or text can be stored for computer processing with an **image scanner** by passing an optical scanning mechanism over the information to be stored. The gradations of gray scale or color are then recorded and stored in an array. Once we have the internal representation of a picture, we can apply transformations to rotate, scale, or crop the picture to a particular screen area. We can also apply various image-processing methods to modify the array representation of the picture. For scanned text input, various editing operations can be performed on the stored documents. Some scanners are able to scan either graphical representations or text, and they come in a variety of sizes and capabilities. A small hand-model scanner is shown in Fig. 2-51, while Figs 2-52 and 2-53 show larger models.



Figure 2-50
A three-dimensional digitizing
system for use with Apple
Macintosh computers. (Courtesy of
Mira Imaging.)

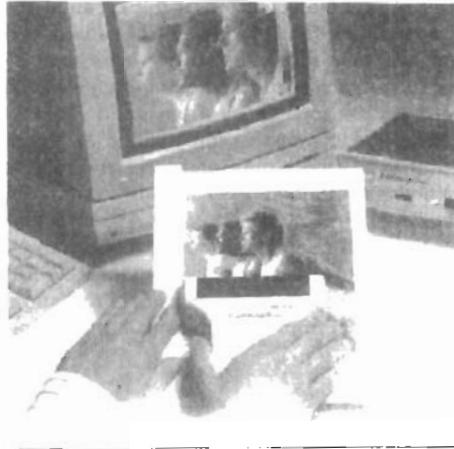


Figure 2-51
A hand-held scanner that can be used to input either text or graphics images. (Courtesy of Thunderware, Inc.)

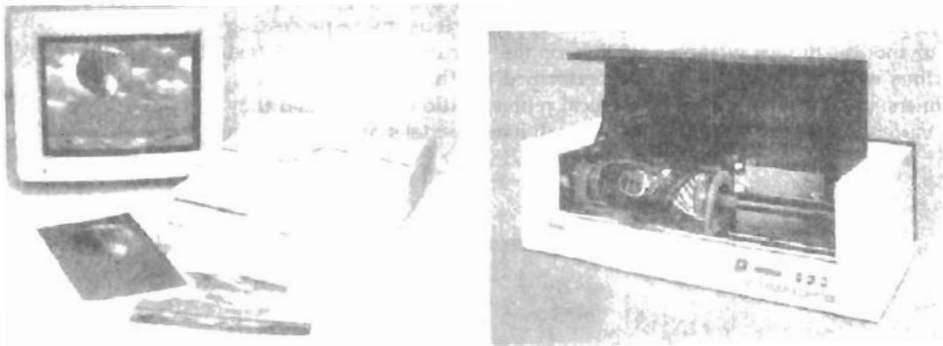


Figure 2-52
Desktop full-color scanners: (a) Flatbed scanner with a resolution of 600 dots per inch. (Courtesy of Sharp Electronics Corporation.) (b) Drum scanner with a selectable resolution from 50 to 4000 dots per inch. (Courtesy of Howtek, Inc.)

Touch Panels

As the name implies, **touch panels** allow displayed objects or screen positions to be selected with the touch of a finger. A typical application of touch panels is for the selection of processing options that are represented with graphical icons. Some systems, such as the plasma panels shown in Fig. 2-54, are designed with touch screens. Other systems can be adapted for touch input by fitting a transparent device with a touch-sensing mechanism over the video monitor screen. Touch input can be recorded using optical, electrical, or acoustical methods.

Optical touch panels employ a line of infrared light-emitting diodes (LEDs) along one vertical edge and along one horizontal edge of the frame. The opposite vertical and horizontal edges contain light detectors. These detectors are used to record which beams are interrupted when the panel is touched. The two crossing

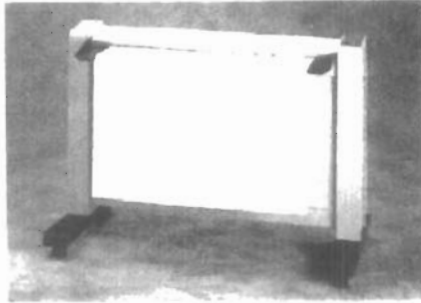


Figure 2-53
A large floor-model scanner used to scan architectural and engineering drawings up to 40 inches wide and 100 feet long. (Courtesy of Summagraphics Corporation.)

beams that are interrupted identify the horizontal and vertical coordinates of the screen position selected. Positions can be selected with an accuracy of about 1/4 inch. With closely spaced LEDs, it is possible to break two horizontal or two vertical beams simultaneously. In this case, an average position between the two interrupted beams is recorded. The LEDs operate at infrared frequencies, so that the light is not visible to a user. Figure 2-55 illustrates the arrangement of LEDs in an optical touch panel that is designed to match the color and contours of the system to which it is to be fitted.

An electrical touch panel is constructed with two transparent plates separated by a small distance. One of the plates is coated with a conducting material, and the other plate is coated with a resistive material. When the outer plate is touched, it is forced into contact with the inner plate. This contact creates a voltage drop across the resistive plate that is converted to the coordinate values of the selected screen position.

In acoustical touch panels, high-frequency sound waves are generated in the horizontal and vertical directions across a glass plate. Touching the screen causes part of each wave to be reflected from the finger to the emitters. The screen position at the point of contact is calculated from a measurement of the time interval between the transmission of each wave and its reflection to the emitter.



Figure 2-54
Plasma panels with touch screens. (Courtesy of Photonics Systems.)

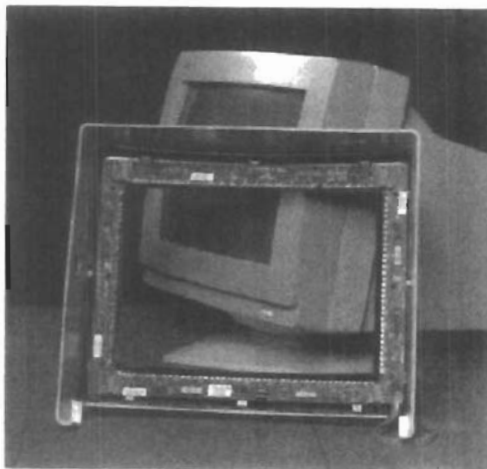


Figure 2-55
An optical touch panel, showing the arrangement of infrared LED units and detectors around the edges of the frame. (Courtesy of Carroll Touch, Inc.)

Light Pens

Figure 2-56 shows the design of one type of **light pen**. Such pencil-shaped devices are used to select screen positions by detecting the light coming from points on the CRT screen. They are sensitive to the short burst of light emitted from the phosphor coating at the instant the electron beam strikes a particular point. Other light sources, such as the background light in the room, are usually not detected by a light pen. An activated light pen, pointed at a spot on the screen as the electron beam lights up that spot, generates an electrical pulse that causes the coordinate position of the electron beam to be recorded. As with cursor-positioning devices, recorded light-pen coordinates can be used to position an object or to select a processing option.

Although light pens are still with us, they are not as popular as they once were since they have several disadvantages compared to other input devices that have been developed. For one, when a light pen is pointed at the screen, part of the screen image is obscured by the hand and pen. And prolonged use of the light pen can cause arm fatigue. Also, light pens require special implementations for some applications because they cannot detect positions within black areas. To be able to select positions in any screen area with a light pen, we must have some nonzero intensity assigned to each screen pixel. In addition, light pens sometimes give false readings due to background lighting in a room.

Voice Systems

Speech recognizers are used in some graphics workstations as input devices to accept voice commands. The voice-system input can be used to initiate graphics



Figure 2-56
A light pen activated with a button switch. (Courtesy of Interactive Computer Products.)

operations or to enter data. These systems operate by matching an input against a predefined dictionary of words and phrases.

A dictionary is set up for a particular operator by having the operator speak the command words to be used into the system. Each word is spoken several times, and the system analyzes the word and establishes a frequency pattern for that word in the dictionary along with the corresponding function to be performed. Later, when a voice command is given, the system searches the dictionary for a frequency-pattern match. Voice input is typically spoken into a microphone mounted on a headset, as in Fig. 2-57. The microphone is designed to minimize input of other background sounds. If a different operator is to use the system, the dictionary must be reestablished with that operator's voice patterns. Voice systems have some advantage over other input devices, since the attention of the operator does not have to be switched from one device to another to enter a command.



Figure 2-57
A speech-recognition system. (Courtesy of Threshold Technology, Inc.)

HARD-COPY DEVICES

We can obtain hard-copy output for our images in several formats. For presentations or archiving, we can send image files to devices or service bureaus that will produce 35-mm slides or overhead transparencies. To put images on film, we can simply photograph a scene displayed on a video monitor. And we can put our pictures on paper by directing graphics output to a printer or plotter.

The quality of the pictures obtained from a device depends on dot size and the number of dots per inch, or lines per inch, that can be displayed. To produce smooth characters in printed text strings, higher-quality printers shift dot positions so that adjacent dots overlap.

Printers produce output by either impact or nonimpact methods. *Impact* printers press formed character faces against an inked ribbon onto the paper. A line printer is an example of an impact device, with the typefaces mounted on bands, chains, drums, or wheels. *Nonimpact* printers and plotters use laser techniques, ink-jet sprays, xerographic processes (as used in photocopying machines), electrostatic methods, and electrothermal methods to get images onto paper.

Character impact printers often have a *dot-matrix* print head containing a rectangular array of protruding wire pins, with the number of pins depending on the quality of the printer. Individual characters or graphics patterns are obtained by retracting certain pins so that the remaining pins form the pattern to be printed. Figure 2-58 shows a picture printed on a dot-matrix printer.

In a *laser* device, a laser beam creates a charge distribution on a rotating drum coated with a photoelectric material, such as selenium. Toner is applied to the drum and then transferred to paper. Figure 2-59 shows examples of desktop laser printers with a resolution of 360 dots per inch.

Ink-jet methods produce output by squirting ink in horizontal rows across a roll of paper wrapped on a drum. The electrically charged ink stream is deflected by an electric field to produce dot-matrix patterns. A desktop ink-jet plotter with



Figure 2-58

A picture generated on a dot-matrix printer showing how the density of the dot patterns can be varied to produce light and dark areas. (Courtesy of Apple Computer, Inc.)

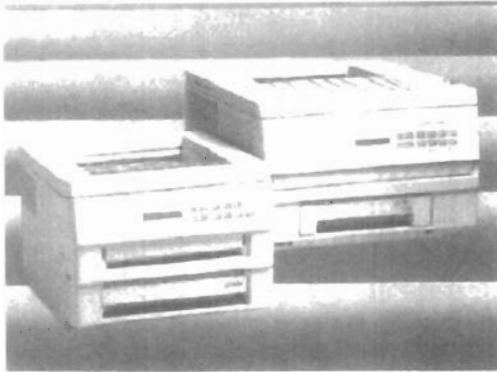


Figure 2-59
Small-footprint laser printers.
(Courtesy of Texas Instruments.)

a resolution of 360 dots per inch is shown in Fig. 2-60, and examples of larger high-resolution ink-jet printer/plotters are shown in Fig. 2-61.

An *electrostatic* device places a negative charge on the paper, one complete row at a time along the length of the paper. Then the paper is exposed to a toner. The toner is positively charged and so is attracted to the negatively charged areas, where it adheres to produce the specified output. A color electrostatic printer/plotter is shown in Fig. 2-62. *Electrothermal* methods use heat in a dot-matrix print head to output patterns on heat-sensitive paper.

We can get limited color output on an impact printer by using different-colored ribbons. Nonimpact devices use various techniques to combine three color pigments (cyan, magenta, and yellow) to produce a range of color patterns. Laser and xerographic devices deposit the three pigments on separate passes; ink-jet methods shoot the three colors simultaneously on a single pass along each print line on the paper.

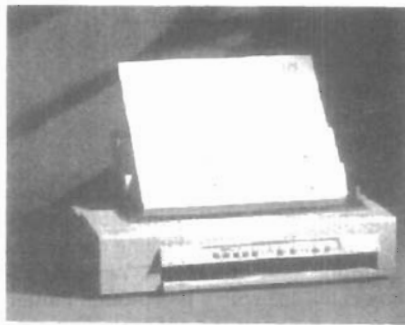


Figure 2-60
A 360-dot-per-inch desktop ink-jet
plotter. (Courtesy of Summagraphics
Corporation.)

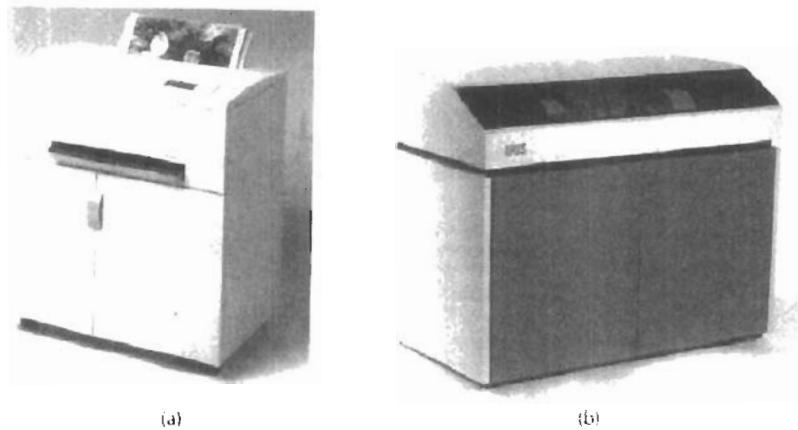


Figure 2-61
Floor-model, ink-jet color printers that use variable dot size to achieve an equivalent resolution of 1500 to 1800 dots per inch. (Courtesy of IRIS Graphics Inc., Bedford, Massachusetts.)

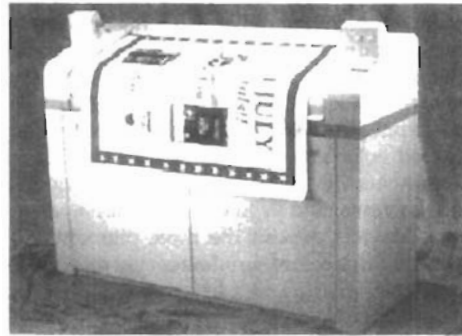


Figure 2-62
An electrostatic printer that can display 400 dots per inch. (Courtesy of CalComp Digitizer Division, a part of CalComp, Inc.)

Drafting layouts and other drawings are typically generated with ink-jet or pen plotters. A pen plotter has one or more pens mounted on a carriage, or cross-bar, that spans a sheet of paper. Pens with varying colors and widths are used to produce a variety of shadings and line styles. Wet-ink, ball-point, and felt-tip pens are all possible choices for use with a pen plotter. Plotter paper can lie flat or be rolled onto a drum or belt. Crossbars can be either moveable or stationary, while the pen moves back and forth along the bar. Either clamps, a vacuum, or an electrostatic charge hold the paper in position. An example of a table-top flatbed pen plotter is given in Figure 2-63, and a larger, rollfeed pen plotter is shown in Fig. 2-64.

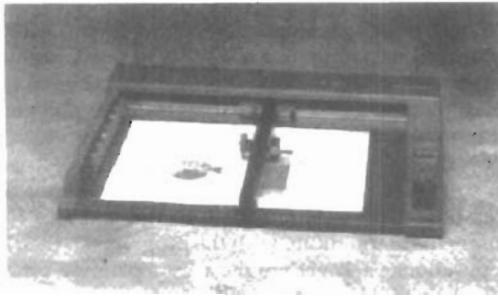


Figure 2-63
A desktop pen plotter with a
resolution of 0.025 mm. (Courtesy of
Summagraphics Corporation.)

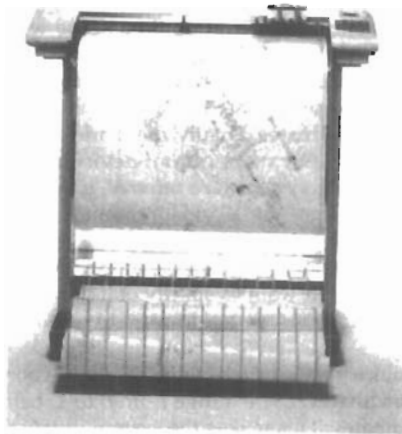


Figure 2-64
A large, rollfeed pen plotter with
automatic multicolor 8-pen changer
and a resolution of 0.0127 mm.
(Courtesy of Summagraphics Corporation.)

2-7

GRAPHICS SOFTWARE

There are two general classifications for graphics software: general programming packages and special-purpose applications packages. A general graphics programming package provides an extensive set of graphics functions that can be

used in a high-level programming language, such as C or FORTRAN. An example of a general graphics programming package is the GL (Graphics Library) system on Silicon Graphics equipment. Basic functions in a general package include those for generating picture components (straight lines, polygons, circles, and other figures), setting color and intensity values, selecting views, and applying transformations. By contrast, application graphics packages are designed for nonprogrammers, so that users can generate displays without worrying about how graphics operations work. The interface to the graphics routines in such packages allows users to communicate with the programs in their own terms. Examples of such applications packages are the artist's painting programs and various business, medical, and CAD systems.

Coordinate Representations

With few exceptions, general graphics packages are designed to be used with Cartesian coordinate specifications. If coordinate values for a picture are specified in some other reference frame (spherical, hyperbolic, etc.), they must be converted to Cartesian coordinates before they can be input to the graphics package. Special-purpose packages may allow use of other coordinate frames that are appropriate to the application. In general, several different Cartesian reference frames are used to construct and display a scene. We can construct the shape of individual objects, such as trees or furniture, in a scene within separate coordinate reference frames called **modeling coordinates**, or sometimes **local coordinates** or **master coordinates**. Once individual object shapes have been specified, we can place the objects into appropriate positions within the scene using a reference frame called **world coordinates**. Finally, the world-coordinate description of the scene is transferred to one or more output-device reference frames for display. These display coordinate systems are referred to as **device coordinates**, or **screen coordinates** in the case of a video monitor. Modeling and world-coordinate definitions allow us to set any convenient floating-point or integer dimensions without being hampered by the constraints of a particular output device. For some scenes, we might want to specify object dimensions in fractions of a foot, while for other applications we might want to use millimeters, kilometers, or light-years.

Generally, a graphics system first converts world-coordinate positions to **normalized device coordinates**, in the range from 0 to 1, before final conversion to specific device coordinates. This makes the system independent of the various devices that might be used at a particular workstation. Figure 2-65 illustrates the sequence of coordinate transformations from modeling coordinates to device coordinates for a two-dimensional application. An initial modeling-coordinate position (x_{mc}, y_{mc}) in this illustration is transferred to a device coordinate position (x_{dc}, y_{dc}) with the sequence:

$$(x_{mc}, y_{mc}) \rightarrow (x_{wc}, y_{wc}) \rightarrow (x_{nc}, y_{nc}) \rightarrow (x_{dc}, y_{dc})$$

The modeling and world-coordinate positions in this transformation can be any floating-point values; normalized coordinates satisfy the inequalities: $0 \leq x_{nc} \leq 1$, $0 \leq y_{nc} \leq 1$; and the device coordinates x_{dc} and y_{dc} are integers within the range $(0, 0)$ to (x_{max}, y_{max}) for a particular output device. To accommodate differences in scales and aspect ratios, normalized coordinates are mapped into a square area of the output device so that proper proportions are maintained.

A general-purpose graphics package provides users with a variety of functions for creating and manipulating pictures. These routines can be categorized according to whether they deal with output, input, attributes, transformations, viewing, or general control.

The basic building blocks for pictures are referred to as **output primitives**. They include character strings and geometric entities, such as points, straight lines, curved lines, filled areas (polygons, circles, etc.), and shapes defined with arrays of color points. Routines for generating output primitives provide the basic tools for constructing pictures.

Attributes are the properties of the output primitives; that is, an attribute describes how a particular primitive is to be displayed. They include intensity and color specifications, line styles, text styles, and area-filling patterns. Functions within this category can be used to set attributes for an individual primitive class or for groups of output primitives.

We can change the size, position, or orientation of an object within a scene using **geometric transformations**. Similar **modeling transformations** are used to construct a scene using object descriptions given in modeling coordinates.

Given the primitive and attribute definition of a picture in world coordinates, a graphics package projects a selected view of the picture on an output device. **Viewing transformations** are used to specify the view that is to be presented and the portion of the output display area that is to be used.

Pictures can be subdivided into component parts, called **structures** or **segments** or **objects**, depending on the software package in use. Each structure defines one logical unit of the picture. A scene with several objects could reference each individual object in a separate named structure. Routines for processing

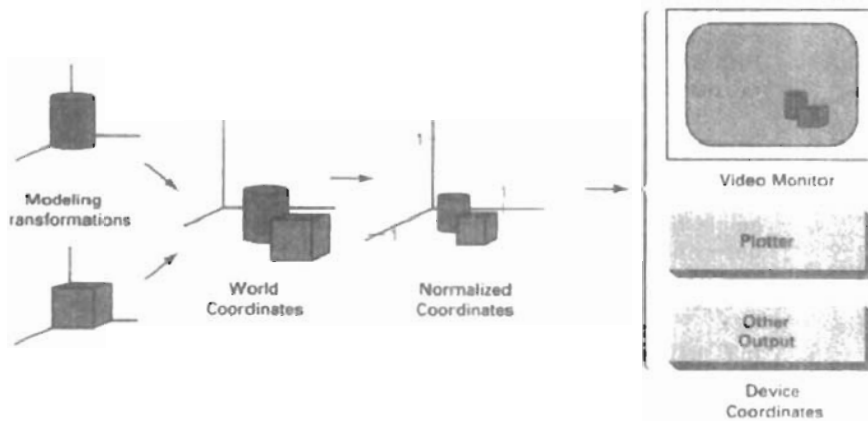


Figure 2-65

The transformation sequence from modeling coordinates to device coordinates for a two-dimensional scene. Object shapes are defined in local modeling-coordinate systems, then positioned within the overall world-coordinate scene. World-coordinate specifications are then transformed into normalized coordinates. At the final step, individual device drivers transfer the normalized-coordinate representation of the scene to the output devices for display.

structures carry out operations such as the creation, modification, and transformation of structures.

Interactive graphics applications use various kinds of input devices, such as a mouse, a tablet, or a joystick. **Input functions** are used to control and process the data flow from these interactive devices.

Finally, a graphics package contains a number of housekeeping tasks, such as clearing a display screen and initializing parameters. We can lump the functions for carrying out these chores under the heading **control operations**.

Software Standards

The primary goal of standardized graphics software is portability. When packages are designed with standard graphics functions, software can be moved easily from one hardware system to another and used in different implementations and applications. Without standards, programs designed for one hardware system often cannot be transferred to another system without extensive rewriting of the programs.

International and national standards planning organizations in many countries have cooperated in an effort to develop a generally accepted standard for computer graphics. After considerable effort, this work on standards led to the development of the **Graphical Kernel System (GKS)**. This system was adopted as the first graphics software standard by the International Standards Organization (ISO) and by various national standards organizations, including the American National Standards Institute (ANSI). Although GKS was originally designed as a two-dimensional graphics package, a three-dimensional GKS extension was subsequently developed. The second software standard to be developed and approved by the standards organizations was **PHIGS (Programmer's Hierarchical Interactive Graphics Standard)**, which is an extension of GKS. Increased capabilities for object modeling, color specifications, surface rendering, and picture manipulations are provided in PHIGS. Subsequently, an extension of PHIGS, called PHIGS+, was developed to provide three-dimensional surface-shading capabilities not available in PHIGS.

Standard graphics functions are defined as a set of specifications that is independent of any programming language. A **language binding** is then defined for a particular high-level programming language. This binding gives the syntax for accessing the various standard graphics functions from this language. For example, the general form of the PHIGS (and GKS) function for specifying a sequence of $n - 1$ connected two-dimensional straight line segments is

`polyline(n, x, y)`

In FORTRAN, this procedure is implemented as a subroutine with the name *GPL*. A graphics programmer, using FORTRAN, would invoke this procedure with the subroutine call statement `CALL GPL(N, X, Y)`, where *X* and *Y* are one-dimensional arrays of coordinate values for the line endpoints. In C, the procedure would be invoked with `ppolyline(n, pts)`, where *pts* is the list of coordinate endpoint positions. Each language binding is defined to make best use of the corresponding language capabilities and to handle various syntax issues, such as data types, parameter passing, and errors.

In the following chapters, we use the standard functions defined in PHIGS as a framework for discussing basic graphics concepts and the design and application of graphics packages. Example programs are presented in Pascal to illus-

trate the algorithms for implementation of the graphics functions and to illustrate also some applications of the functions. Descriptive names for functions, based on the PHIGS definitions, are used whenever a graphics function is referenced in a program.

Although PHIGS presents a specification for basic graphics functions, it does not provide a standard methodology for a graphics interface to output devices. Nor does it specify methods for storing and transmitting pictures. Separate standards have been developed for these areas. Standardization for device interface methods is given in the **Computer Graphics Interface (CGI)** system. And the **Computer Graphics Metafile (CGM)** system specifies standards for archiving and transporting pictures.

PHIGS Workstations

Generally, the term *workstation* refers to a computer system with a combination of input and output devices that is designed for a single user. In PHIGS and GKS, however, the term **workstation** is used to identify various combinations of graphics hardware and software. A PHIGS workstation can be a single output device, a single input device, a combination of input and output devices, a file, or even a window displayed on a video monitor.

To define and use various “workstations” within an applications program, we need to specify a *workstation identifier* and the workstation type. The following statements give the general structure of a PHIGS program:

```
openPhigs (errorFile, memorySize)
openWorkstation (ws, connection, type)
    { create and display picture }
closeWorkstation (ws)
closePhigs
```

where parameter *errorFile* is to contain any error messages that are generated, and parameter *memorySize* specifies the size of an internal storage area. The workstation identifier (an integer) is given in parameter *ws*, and parameter *connection* states the access mechanism for the workstation. Parameter *type* specifies the particular category for the workstation, such as an *input* device, an *output* device, a combination *outin* device, or an input or output metafile.

Any number of workstations can be open in a particular application, with input coming from the various open input devices and output directed to all the open output devices. We discuss input and output methods in applications programs in Chapter 8, after we have explored the basic procedures for creating and manipulating pictures.

SUMMARY

In this chapter, we have surveyed the major hardware and software features of computer graphics systems. Hardware components include video monitors, hard-copy devices, keyboards, and other devices for graphics input or output. Graphics software includes special applications packages and general programming packages.

The predominant graphics display device is the raster refresh monitor, based on television technology. A raster system uses a frame buffer to store intensity information for each screen position (pixel). Pictures are then painted on the

screen by retrieving this information from the frame buffer as the electron beam in the CRT sweeps across each scan line, from top to bottom. Older vector displays construct pictures by drawing lines between specified line endpoints. Picture information is then stored as a set of line-drawing instructions.

Many other video display devices are available. In particular, flat-panel display technology is developing at a rapid rate, and these devices may largely replace raster displays in the near future. At present, flat-panel displays are commonly used in small systems and in special-purpose systems. Flat-panel displays include plasma panels and liquid-crystal devices. Although vector monitors can be used to display high-quality line drawings, improvements in raster display technology have caused vector monitors to be largely replaced with raster systems.

Other display technologies include three-dimensional and stereoscopic viewing systems. Virtual-reality systems can include either a stereoscopic headset or a standard video monitor.

For graphical input, we have a range of devices to choose from. Keyboards, button boxes, and dials are used to input text, data values, or programming options. The most popular "pointing" device is the mouse, but trackballs, spaceballs, joysticks, cursor-control keys, and thumbwheels are also used to position the screen cursor. In virtual-reality environments, data gloves are commonly used. Other input devices include image scanners, digitizers, touch panels, light pens, and voice systems.

Hard-copy devices for graphics workstations include standard printers and plotters, in addition to devices for producing slides, transparencies, and film output. Printing methods include dot matrix, laser, ink jet, electrostatic, and electrophoretic. Plotter methods include pen plotting and combination printer-plotter devices.

Graphics software can be roughly classified as applications packages or programming packages. Applications graphics software include CAD packages, drawing and painting programs, graphing packages, and visualization programs. Common graphics programming packages include PHIGS, PHIGS+, GKS, 3D GKS, and GL. Software standards, such as PHIGS, GKS, CGI, and CGM, are evolving and are becoming widely available on a variety of machines.

Normally, graphics packages require coordinate specifications to be given with respect to Cartesian reference frames. Each object for a scene can be defined in a separate modeling Cartesian coordinate system, which is then mapped to world coordinates to construct the scene. From world coordinates, objects are transferred to normalized device coordinates, then to the final display device coordinates. The transformations from modeling coordinates to normalized device coordinates are independent of particular devices that might be used in an application. Device drivers are then used to convert normalized coordinates to integer device coordinates.

Functions in graphics programming packages can be divided into the following categories: output primitives, attributes, geometric and modeling transformations, viewing transformations, structure operations, input functions, and control operations.

Some graphics systems, such as PHIGS and GKS, use the concept of a "workstation" to specify devices or software that are to be used for input or output in a particular application. A workstation identifier in these systems can refer to a file; a single device, such as a raster monitor; or a combination of devices, such as a monitor, keyboard, and a mouse. Multiple workstations can be open to provide input or to receive output in a graphics application.

REFERENCES

Exercises

A general treatment of electronic displays, including flat-panel devices, is available in Sherr (1993). Flat-panel devices are discussed in Depp and Howard (1993). Tannas (1985) provides a reference for both flat-panel displays and CRTs. Additional information on raster-graphics architecture can be found in Foley, et al. (1990). Three-dimensional terminals are discussed in Fuchs et al. (1982), Johnson (1982), and Ikeda (1984). Head-mounted displays and virtual-reality environments are discussed in Chung et al. (1989).

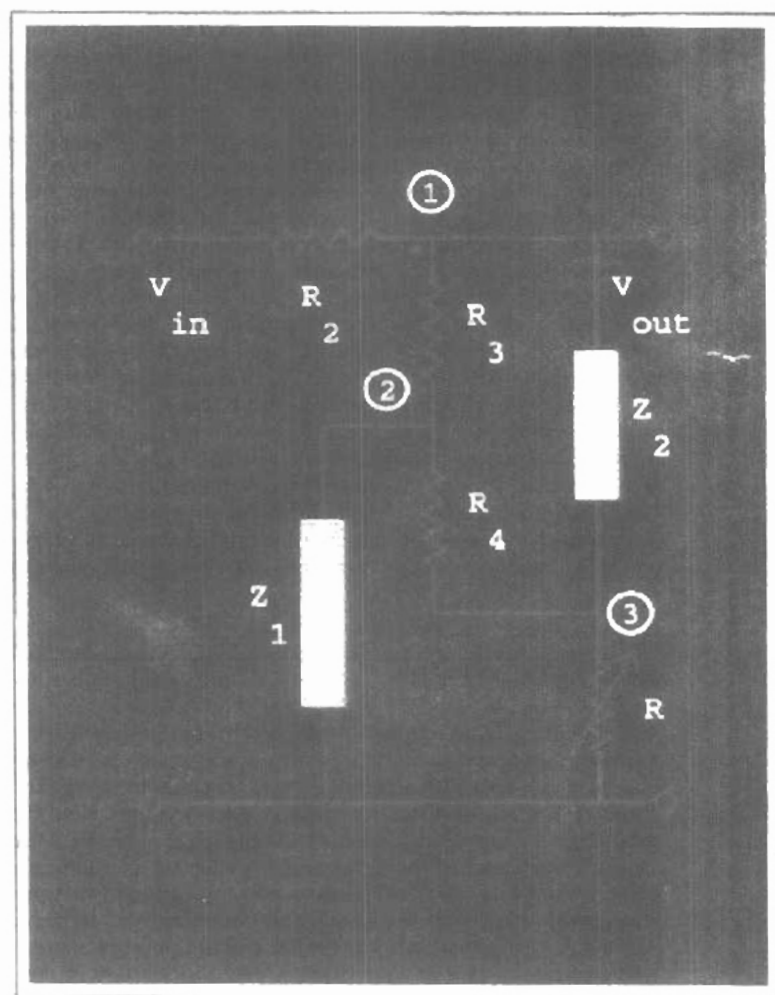
For information on PHIGS and PHIGS+, see Hopgood and Duce (1991), Howard et al. (1991), Gaskins (1992), and Blake (1993). Information on the two-dimensional GKS standard and on the evolution of graphics standards is available in Hopgood et al. (1983). An additional reference for GKS is Enderle, Kansy, and Plaff (1984).

EXERCISES

- 2-1. List the operating characteristics for the following display technologies: raster refresh systems, vector refresh systems, plasma panels, and LCDs.
- 2-2. List some applications appropriate for each of the display technologies in Exercise 2-1.
- 2-3. Determine the resolution (pixels per centimeter) in the x and y directions for the video monitor in use on your system. Determine the aspect ratio, and explain how relative proportions of objects can be maintained on your system.
- 2-4. Consider three different raster systems with resolutions of 640 by 480, 1280 by 1024, and 2560 by 2048. What size frame buffer (in bytes) is needed for each of these systems to store 12 bits per pixel? How much storage is required for each system if 24 bits per pixel are to be stored?
- 2-5. Suppose an RGB raster system is to be designed using an 8-inch by 10-inch screen with a resolution of 100 pixels per inch in each direction. If we want to store 6 bits per pixel in the frame buffer, how much storage (in bytes) do we need for the frame buffer?
- 2-6. How long would it take to load a 640 by 480 frame buffer with 12 bits per pixel, if 10^5 bits can be transferred per second? How long would it take to load a 24-bit per pixel frame buffer with a resolution of 1280 by 1024 using this same transfer rate?
- 2-7. Suppose we have a computer with 32 bits per word and a transfer rate of 1 mip (one million instructions per second). How long would it take to fill the frame buffer of a 300-dpi (dot per inch) laser printer with a page size of 8 1/2 inches by 11 inches?
- 2-8. Consider two raster systems with resolutions of 640 by 480 and 1280 by 1024. How many pixels could be accessed per second in each of these systems by a display controller that refreshes the screen at a rate of 60 frames per second? What is the access time per pixel in each system?
- 2-9. Suppose we have a video monitor with a display area that measures 12 inches across and 9.6 inches high. If the resolution is 1280 by 1024 and the aspect ratio is 1, what is the diameter of each screen point?
- 2-10. How much time is spent scanning across each row of pixels during screen refresh on a raster system with a resolution of 1280 by 1024 and a refresh rate of 60 frames per second?
- 2-11. Consider a noninterlaced raster monitor with a resolution of n by m (m scan lines and n pixels per scan line), a refresh rate of r frames per second, a horizontal retrace time of t_{horiz} , and a vertical retrace time of t_{vert} . What is the fraction of the total refresh time per frame spent in retrace of the electron beam?
- 2-12. What is the fraction of the total refresh time per frame spent in retrace of the electron beam for a noninterlaced raster system with a resolution of 1280 by 1024, a refresh rate of 60 Hz, a horizontal retrace time of 5 microseconds, and a vertical retrace time of 500 microseconds?

- 2-13. Assuming that a certain full-color (24-bit per pixel) RGB raster system has a 512-by-512 frame buffer, how many distinct color choices (intensity levels) would we have available? How many different colors could we display at any one time?
- 2-14. Compare the advantages and disadvantages of a three-dimensional monitor using a varifocal mirror with a stereoscopic system.
- 2-15. List the different input and output components that are typically used with virtual-reality systems. Also explain how users interact with a virtual scene displayed with different output devices, such as two-dimensional and stereoscopic monitors.
- 2-16. Explain how virtual-reality systems can be used in design applications. What are some other applications for virtual-reality systems?
- 2-17. List some applications for large-screen displays.
- 2-18. Explain the differences between a general graphics system designed for a programmer and one designed for a specific application, such as architectural design?

Output Primitives



A picture can be described in several ways. Assuming we have a raster display, a picture is completely specified by the set of intensities for the pixel positions in the display. At the other extreme, we can describe a picture as a set of complex objects, such as trees and terrain or furniture and walls, positioned at specified coordinate locations within the scene. Shapes and colors of the objects can be described internally with pixel arrays or with sets of basic geometric structures, such as straight line segments and polygon color areas. The scene is then displayed either by loading the pixel arrays into the frame buffer or by scan converting the basic geometric-structure specifications into pixel patterns. Typically, graphics programming packages provide functions to describe a scene in terms of these basic geometric structures, referred to as **output primitives**, and to group sets of output primitives into more complex structures. Each output primitive is specified with input coordinate data and other information about the way that object is to be displayed. Points and straight line segments are the simplest geometric components of pictures. Additional output primitives that can be used to construct a picture include circles and other conic sections, quadric surfaces, spline curves and surfaces, polygon color areas, and character strings. We begin our discussion of picture-generation procedures by examining device-level algorithms for displaying two-dimensional output primitives, with particular emphasis on scan-conversion methods for raster graphics systems. In this chapter, we also consider how output functions can be provided in graphics packages, and we take a look at the output functions available in the PHIGS language.

3-1

POINTS AND LINES

Point plotting is accomplished by converting a single coordinate position furnished by an application program into appropriate operations for the output device in use. With a CRT monitor, for example, the electron beam is turned on to illuminate the screen phosphor at the selected location. How the electron beam is positioned depends on the display technology. A random-scan (vector) system stores point-plotting instructions in the display list, and coordinate values in these instructions are converted to deflection voltages that position the electron beam at the screen locations to be plotted during each refresh cycle. For a black-and-white raster system, on the other hand, a point is plotted by setting the bit value corresponding to a specified screen position within the frame buffer to 1. Then, as the electron beam sweeps across each horizontal scan line, it emits a

burst of electrons (plots a point) whenever a value of 1 is encountered in the frame buffer. With an RGB system, the frame buffer is loaded with the color codes for the intensities that are to be displayed at the screen pixel positions.

Line drawing is accomplished by calculating intermediate positions along the line path between two specified endpoint positions. An output device is then directed to fill in these positions between the endpoints. For analog devices, such as a vector pen plotter or a random-scan display, a straight line can be drawn smoothly from one endpoint to the other. Linearly varying horizontal and vertical deflection voltages are generated that are proportional to the required changes in the x and y directions to produce the smooth line.

Digital devices display a straight line segment by plotting discrete points between the two endpoints. Discrete coordinate positions along the line path are calculated from the equation of the line. For a raster video display, the line color (intensity) is then loaded into the frame buffer at the corresponding pixel coordinates. Reading from the frame buffer, the video controller then "plots" the screen pixels. Screen locations are referenced with integer values, so plotted positions may only approximate actual line positions between two specified endpoints. A computed line position of (10.48, 20.51), for example, would be converted to pixel position (10, 21). This rounding of coordinate values to integers causes lines to be displayed with a stairstep appearance ("the jaggies"), as represented in Fig 3-1. The characteristic stairstep shape of raster lines is particularly noticeable on systems with low resolution, and we can improve their appearance somewhat by displaying them on high-resolution systems. More effective techniques for smoothing raster lines are based on adjusting pixel intensities along the line paths.

For the raster-graphics device-level algorithms discussed in this chapter, object positions are specified directly in integer device coordinates. For the time being, we will assume that pixel positions are referenced according to scan-line number and column number (pixel position across a scan line). This addressing scheme is illustrated in Fig. 3-2. Scan lines are numbered consecutively from 0, starting at the bottom of the screen; and pixel columns are numbered from 0, left to right across each scan line. In Section 3-10, we consider alternative pixel addressing schemes.

To load a specified color into the frame buffer at a position corresponding to column x along scan line y , we will assume we have available a low-level procedure of the form

```
setPixel (x, y)
```



Figure 3-1
Stairstep effect (jaggies) produced
when a line is generated as a series
of pixel positions.

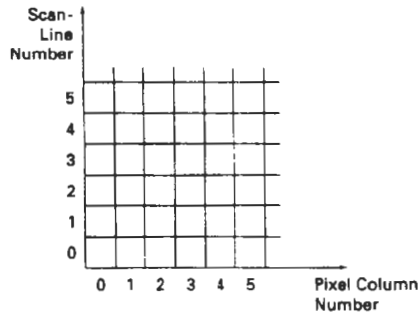


Figure 3-2
Pixel positions referenced by scan-line number and column number.

We sometimes will also want to be able to retrieve the current frame-buffer intensity setting for a specified location. We accomplish this with the low-level function

```
getPixel (x, y)
```

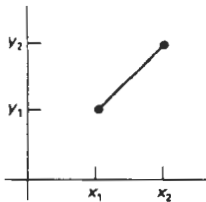
3-2

LINE-DRAWING ALGORITHMS

The Cartesian *slope-intercept equation* for a straight line is

$$y = m \cdot x + b \quad (3-1)$$

with m representing the slope of the line and b as the y intercept. Given that the two endpoints of a line segment are specified at positions (x_1, y_1) and (x_2, y_2) , as shown in Fig. 3-3, we can determine values for the slope m and y intercept b with the following calculations:



$$m = \frac{y_2 - y_1}{x_2 - x_1} \quad (3-2)$$

$$b = y_1 - m \cdot x_1 \quad (3-3)$$

Algorithms for displaying straight lines are based on the line equation 3-1 and the calculations given in Eqs. 3-2 and 3-3.

For any given x interval Δx along a line, we can compute the corresponding y interval Δy from Eq. 3-2 as

$$\Delta y = m \Delta x \quad (3-4)$$

Similarly, we can obtain the x interval Δx corresponding to a specified Δy as

$$\Delta x = \frac{\Delta y}{m} \quad (3-5)$$

These equations form the basis for determining deflection voltages in analog de-

Figure 3-3
Line path between endpoint positions (x_1, y_1) and (x_2, y_2) .

vices. For lines with slope magnitudes $|m| < 1$, Δx can be set proportional to a small horizontal deflection voltage and the corresponding vertical deflection is then set proportional to Δy as calculated from Eq. 3-4. For lines whose slopes have magnitudes $|m| > 1$, Δy can be set proportional to a small vertical deflection voltage with the corresponding horizontal deflection voltage set proportional to Δx , calculated from Eq. 3-5. For lines with $m = 1$, $\Delta x = \Delta y$ and the horizontal and vertical deflections voltages are equal. In each case, a smooth line with slope m is generated between the specified endpoints.

On raster systems, lines are plotted with pixels, and step sizes in the horizontal and vertical directions are constrained by pixel separations. That is, we must "sample" a line at discrete positions and determine the nearest pixel to the line at each sampled position. This scan-conversion process for straight lines is illustrated in Fig. 3-4, for a near horizontal line with discrete sample positions along the x axis.

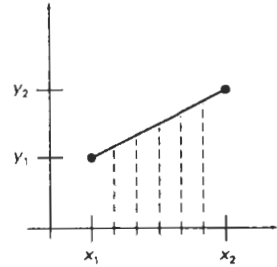


Figure 3-4
Straight line segment with five sampling positions along the x axis between x_1 and x_2 .

DDA Algorithm

The *digital differential analyzer* (DDA) is a scan-conversion line algorithm based on calculating either Δy or Δx , using Eq. 3-4 or Eq. 3-5. We sample the line at unit intervals in one coordinate and determine corresponding integer values nearest the line path for the other coordinate.

Consider first a line with positive slope, as shown in Fig. 3-3. If the slope is less than or equal to 1, we sample at unit x intervals ($\Delta x = 1$) and compute each successive y value as

$$y_{k+1} = y_k + m \quad (3-6)$$

Subscript k takes integer values starting from 1, for the first point, and increases by 1 until the final endpoint is reached. Since m can be any real number between 0 and 1, the calculated y values must be rounded to the nearest integer.

For lines with a positive slope greater than 1, we reverse the roles of x and y . That is, we sample at unit y intervals ($\Delta y = 1$) and calculate each succeeding x value as

$$x_{k+1} = x_k + \frac{1}{m} \quad (3-7)$$

Equations 3-6 and 3-7 are based on the assumption that lines are to be processed from the left endpoint to the right endpoint (Fig. 3-3). If this processing is reversed, so that the starting endpoint is at the right, then either we have $\Delta x = -1$ and

$$y_{k+1} = y_k - m \quad (3-8)$$

or (when the slope is greater than 1) we have $\Delta y = -1$ with

$$x_{k+1} = x_k - \frac{1}{m} \quad (3-9)$$

Equations 3-6 through 3-9 can also be used to calculate pixel positions along a line with negative slope. If the absolute value of the slope is less than 1 and the start endpoint is at the left, we set $\Delta x = 1$ and calculate y values with Eq. 3-6.

When the start endpoint is at the right (for the same slope), we set $\Delta x = -1$ and obtain y positions from Eq. 3-8. Similarly, when the absolute value of a negative slope is greater than 1, we use $\Delta y = -1$ and Eq. 3-9 or we use $\Delta y = 1$ and Eq. 3-7.

This algorithm is summarized in the following procedure, which accepts as input the two endpoint pixel positions. Horizontal and vertical differences between the endpoint positions are assigned to parameters dx and dy . The difference with the greater magnitude determines the value of parameter $steps$. Starting with pixel position (x_a, y_a) , we determine the offset needed at each step to generate the next pixel position along the line path. We loop through this process $steps$ times. If the magnitude of dx is greater than the magnitude of dy and x_a is less than x_b , the values of the increments in the x and y directions are 1 and m , respectively. If the greater change is in the x direction, but x_a is greater than x_b , then the decrements -1 and $-m$ are used to generate each new point on the line. Otherwise, we use a unit increment (or decrement) in the y direction and an x increment (or decrement) of $1/m$.

```
#include "device.h"

#define ROUND(a) ((int)(a+0.5))

void lineDDA (int xa, int ya, int xb, int yb)
{
    int dx = xb - xa, dy = yb - ya, steps, k;
    float xIncrement, yIncrement, x = xa, y = ya;

    if (abs (dx) > abs (dy)) steps = abs (dx);
    else steps = abs (dy);
    xIncrement = dx / (float) steps;
    yIncrement = dy / (float) steps;

    setPixel (ROUND(x), ROUND(y));
    for (k=0; k<steps; k++) {
        x += xIncrement;
        y += yIncrement;
        setPixel (ROUND(x), ROUND(y));
    }
}
```

The DDA algorithm is a faster method for calculating pixel positions than the direct use of Eq. 3-1. It eliminates the multiplication in Eq. 3-1 by making use of raster characteristics, so that appropriate increments are applied in the x or y direction to step to pixel positions along the line path. The accumulation of roundoff error in successive additions of the floating-point increment, however, can cause the calculated pixel positions to drift away from the true line path for long line segments. Furthermore, the rounding operations and floating-point arithmetic in procedure `lineDDA` are still time-consuming. We can improve the performance of the DDA algorithm by separating the increments m and $1/m$ into integer and fractional parts so that all calculations are reduced to integer operations. A method for calculating $1/m$ increments in integer steps is discussed in Section 3-11. In the following sections, we consider more general scan-line procedures that can be applied to both lines and curves.

Bresenham's Line Algorithm

An accurate and efficient raster line-generating algorithm, developed by Bresen-

ham, scan converts lines using only incremental integer calculations that can be adapted to display circles and other curves. Figures 3-5 and 3-6 illustrate sections of a display screen where straight line segments are to be drawn. The vertical axes show scan-line positions, and the horizontal axes identify pixel columns. Sampling at unit x intervals in these examples, we need to decide which of two possible pixel positions is closer to the line path at each sample step. Starting from the left endpoint shown in Fig. 3-5, we need to determine at the next sample position whether to plot the pixel at position (11, 11) or the one at (11, 12). Similarly, Fig. 3-6 shows a negative slope line path starting from the left endpoint at pixel position (50, 50). In this one, do we select the next pixel position as (51, 50) or as (51, 49)? These questions are answered with Bresenham's line algorithm by testing the sign of an integer parameter, whose value is proportional to the difference between the separations of the two pixel positions from the actual line path.

To illustrate Bresenham's approach, we first consider the scan-conversion process for lines with positive slope less than 1. Pixel positions along a line path are then determined by sampling at unit x intervals. Starting from the left endpoint (x_0, y_0) of a given line, we step to each successive column (x position) and plot the pixel whose scan-line y value is closest to the line path. Figure 3-7 demonstrates the k th step in this process. Assuming we have determined that the pixel at (x_k, y_k) is to be displayed, we next need to decide which pixel to plot in column x_{k+1} . Our choices are the pixels at positions $(x_k + 1, y_k)$ and $(x_k + 1, y_k + 1)$.

At sampling position $x_k + 1$, we label vertical pixel separations from the mathematical line path as d_1 and d_2 (Fig. 3-8). The y coordinate on the mathematical line at pixel column position $x_k + 1$ is calculated as

$$y = m(x_k + 1) + b \quad (3-10)$$

Then

$$\begin{aligned} d_1 &= y - y_k \\ &= m(x_k + 1) + b - y_k \end{aligned}$$

and

$$\begin{aligned} d_2 &= (y_k + 1) - y \\ &= y_k + 1 - m(x_k + 1) - b \end{aligned}$$

The difference between these two separations is

$$d_1 - d_2 = 2m(x_k + 1) - 2y_k + 2b - 1 \quad (3-11)$$

A decision parameter p_k for the k th step in the line algorithm can be obtained by rearranging Eq. 3-11 so that it involves only integer calculations. We accomplish this by substituting $m = \Delta y / \Delta x$, where Δy and Δx are the vertical and horizontal separations of the endpoint positions, and defining:

$$\begin{aligned} p_k &= \Delta x(d_1 - d_2) \\ &= 2\Delta y \cdot x_k - 2\Delta x \cdot y_k + c \end{aligned} \quad (3-12)$$

The sign of p_k is the same as the sign of $d_1 - d_2$, since $\Delta x > 0$ for our example. Parameter c is constant and has the value $2\Delta y + \Delta x(2b - 1)$, which is independent

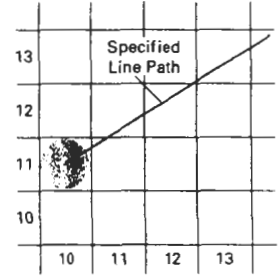


Figure 3-5

Section of a display screen where a straight line segment is to be plotted, starting from the pixel at column 10 on scan line 11.

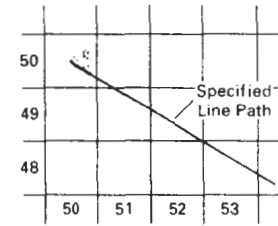


Figure 3-6

Section of a display screen where a negative slope line segment is to be plotted, starting from the pixel at column 50 on scan line 50.

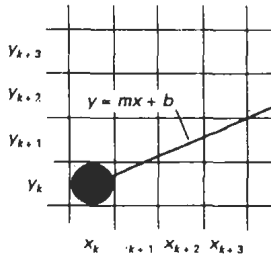


Figure 3-7

Section of the screen grid showing a pixel in column x_k on scan line y_k that is to be plotted along the path of a line segment with slope $0 < m < 1$.

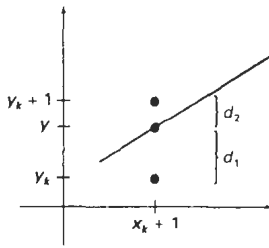


Figure 3-8

Distances between pixel positions and the line y coordinate at sampling position $x_k + 1$.

of pixel position and will be eliminated in the recursive calculations for p_k . If the pixel at y_k is closer to the line path than the pixel at $y_k + 1$ (that is, $d_1 < d_2$), then decision parameter p_k is negative. In that case, we plot the lower pixel; otherwise, we plot the upper pixel.

Coordinate changes along the line occur in unit steps in either the x or y directions. Therefore, we can obtain the values of successive decision parameters using incremental integer calculations. At step $k + 1$, the decision parameter is evaluated from Eq. 3-12 as

$$p_{k+1} = 2\Delta y \cdot x_{k+1} - 2\Delta x \cdot y_{k+1} + c$$

Subtracting Eq. 3-12 from the preceding equation, we have

$$p_{k+1} - p_k = 2\Delta y(x_{k+1} - x_k) - 2\Delta x(y_{k+1} - y_k)$$

But $x_{k+1} = x_k + 1$, so that

$$p_{k+1} = p_k + 2\Delta y - 2\Delta x(y_{k+1} - y_k) \quad (3-13)$$

where the term $y_{k+1} - y_k$ is either 0 or 1, depending on the sign of parameter p_k .

This recursive calculation of decision parameters is performed at each integer x position, starting at the left coordinate endpoint of the line. The first parameter, p_0 , is evaluated from Eq. 3-12 at the starting pixel position (x_0, y_0) and with m evaluated as $\Delta y / \Delta x$:

$$p_0 = 2\Delta y - \Delta x \quad (3-14)$$

We can summarize Bresenham line drawing for a line with a positive slope less than 1 in the following listed steps. The constants $2\Delta y$ and $2\Delta y - 2\Delta x$ are calculated once for each line to be scan converted, so the arithmetic involves only integer addition and subtraction of these two constants.

Bresenham's Line-Drawing Algorithm for $|m| < 1$

1. Input the two line endpoints and store the left endpoint in (x_0, y_0) .
2. Load (x_0, y_0) into the frame buffer; that is, plot the first point.
3. Calculate constants Δx , Δy , $2\Delta y$, and $2\Delta y - 2\Delta x$, and obtain the starting value for the decision parameter as

$$p_0 = 2\Delta y - \Delta x$$

4. At each x_k along the line, starting at $k = 0$, perform the following test:
If $p_k < 0$, the next point to plot is $(x_k + 1, y_k)$ and

$$p_{k+1} = p_k + 2\Delta y$$

Otherwise, the next point to plot is $(x_k + 1, y_k + 1)$ and

$$p_{k+1} = p_k + 2\Delta y - 2\Delta x$$

5. Repeat step 4 Δx times.

Example 3-1 Bresenham Line Drawing

To illustrate the algorithm, we digitize the line with endpoints (20, 10) and (30, 18). This line has a slope of 0.8, with

$$\Delta x = 10, \quad \Delta y = 8$$

The initial decision parameter has the value

$$\begin{aligned} p_0 &= 2\Delta y - \Delta x \\ &= 6 \end{aligned}$$

and the increments for calculating successive decision parameters are

$$2\Delta y = 16, \quad 2\Delta y - 2\Delta x = -4$$

We plot the initial point $(x_0, y_0) = (20, 10)$, and determine successive pixel positions along the line path from the decision parameter as

k	p_k	(x_{k+1}, y_{k+1})	k	p_k	(x_{k+1}, y_{k+1})
0	6	(21, 11)	5	6	(26, 15)
1	2	(22, 12)	6	2	(27, 16)
2	-2	(23, 12)	7	-2	(28, 16)
3	14	(24, 13)	8	14	(29, 17)
4	10	(25, 14)	9	10	(30, 18)

A plot of the pixels generated along this line path is shown in Fig. 3-9.

An implementation of Bresenham line drawing for slopes in the range $0 < m < 1$ is given in the following procedure. Endpoint pixel positions for the line are passed to this procedure, and pixels are plotted from the left endpoint to the right endpoint. The call to setPixel loads a preset color value into the frame buffer at the specified (x, y) pixel position.

```
#include "device.h"

void lineBres (int xa, int ya, int xb, int yb)
{
    int dx = abs (xa - xb), dy = abs (ya - yb);
    int p = 2 * dy - dx;
    int twoDy = 2 * dy, twoDyDx = 2 * (dy - dx);
    int x, y, xEnd;

    /* Determine which point to use as start, which as end */
    if (xa > xb) {
        x = xb;
        y = yb;
        xEnd = xa;
    }
    else {
```

```
x = xa;
y = ya;
xEnd = xb;
}
setPixel (x, y);

while (x < xEnd) {
    x++;
    if (p < 0)
        p += twoDy;
    else {
        y++;
        p += twoDyDx;
    }
    setPixel (x, y);
}
```

Bresenham's algorithm is generalized to lines with arbitrary slope by considering the symmetry between the various octants and quadrants of the xy plane. For a line with positive slope greater than 1, we interchange the roles of the x and y directions. That is, we step along the y direction in unit steps and calculate successive x values nearest the line path. Also, we could revise the program to plot pixels starting from either endpoint. If the initial position for a line with positive slope is the right endpoint, both x and y decrease as we step from right to left. To ensure that the same pixels are plotted regardless of the starting endpoint, we always choose the upper (or the lower) of the two candidate pixels whenever the two vertical separations from the line path are equal ($d_1 = d_2$). For negative slopes, the procedures are similar, except that now one coordinate decreases as the other increases. Finally, special cases can be handled separately: Horizontal lines ($\Delta y = 0$), vertical lines ($\Delta x = 0$), and diagonal lines with $|\Delta x| = |\Delta y|$ each can be loaded directly into the frame buffer without processing them through the line-plotting algorithm.

Parallel Line Algorithms

The line-generating algorithms we have discussed so far determine pixel positions sequentially. With a parallel computer, we can calculate pixel positions

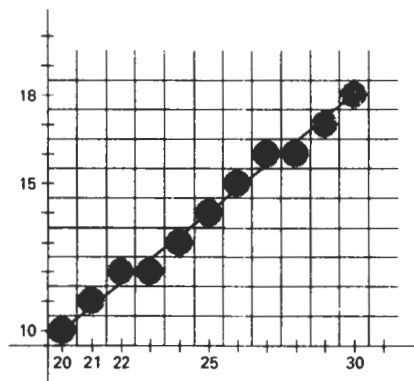


Figure 3-9
Pixel positions along the line path
between endpoints (20, 10) and
(30, 18), plotted with Bresenham's
line algorithm.

along a line path simultaneously by partitioning the computations among the various processors available. One approach to the partitioning problem is to adapt an existing sequential algorithm to take advantage of multiple processors. Alternatively, we can look for other ways to set up the processing so that pixel positions can be calculated efficiently in parallel. An important consideration in devising a parallel algorithm is to balance the processing load among the available processors.

Given n_p processors, we can set up a parallel Bresenham line algorithm by subdividing the line path into n_p partitions and simultaneously generating line segments in each of the subintervals. For a line with slope $0 < m < 1$ and left endpoint coordinate position (x_0, y_0) , we partition the line along the positive x direction. The distance between beginning x positions of adjacent partitions can be calculated as

$$\Delta x_p = \frac{\Delta x + n_p - 1}{n_p} \quad (3-15)$$

where Δx is the width of the line, and the value for partition width Δx_p is computed using integer division. Numbering the partitions, and the processors, as 0, 1, 2, up to $n_p - 1$, we calculate the starting x coordinate for the k th partition as

$$x_k = x_0 + k\Delta x_p \quad (3-16)$$

As an example, suppose $\Delta x = 15$ and we have $n_p = 4$ processors. Then the width of the partitions is 4 and the starting x values for the partitions are x_0 , $x_0 + 4$, $x_0 + 8$, and $x_0 + 12$. With this partitioning scheme, the width of the last (rightmost) subinterval will be smaller than the others in some cases. In addition, if the line endpoints are not integers, truncation errors can result in variable width partitions along the length of the line.

To apply Bresenham's algorithm over the partitions, we need the initial value for the y coordinate and the initial value for the decision parameter in each partition. The change Δy_p in the y direction over each partition is calculated from the line slope m and partition width Δx_p :

$$\Delta y_p = m\Delta x_p \quad (3-17)$$

At the k th partition, the starting y coordinate is then

$$y_k = y_0 + \text{round}(k\Delta y_p) \quad (3-18)$$

The initial decision parameter for Bresenham's algorithm at the start of the k th subinterval is obtained from Eq. 3-12:

$$p_k = (k\Delta x_p)(2\Delta y) - \text{round}(k\Delta y_p)(2\Delta x) + 2\Delta y - \Delta x \quad (3-19)$$

Each processor then calculates pixel positions over its assigned subinterval using the starting decision parameter value for that subinterval and the starting coordinates (x_k, y_k) . We can also reduce the floating-point calculations to integer arithmetic in the computations for starting values y_k and p_k by substituting $m = \Delta y / \Delta x$ and rearranging terms. The extension of the parallel Bresenham algorithm to a line with slope greater than 1 is achieved by partitioning the line in the y di-

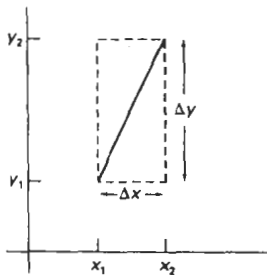


Figure 3-10
Bounding box for a line with
coordinate extents Δx and Δy .

rection and calculating beginning x values for the partitions. For negative slopes, we increment coordinate values in one direction and decrement in the other.

Another way to set up parallel algorithms on raster systems is to assign each processor to a particular group of screen pixels. With a sufficient number of processors (such as a Connection Machine CM-2 with over 65,000 processors), we can assign each processor to one pixel within some screen region. This approach can be adapted to line display by assigning one processor to each of the pixels within the limits of the line coordinate extents (*bounding rectangle*) and calculating pixel distances from the line path. The number of pixels within the bounding box of a line is $\Delta x \cdot \Delta y$ (Fig. 3-10). Perpendicular distance d from the line in Fig. 3-10 to a pixel with coordinates (x, y) is obtained with the calculation

$$d = Ax + By + C \quad (3-20)$$

where

$$\begin{aligned} A &= \frac{-\Delta y}{\text{linelength}} \\ B &= \frac{\Delta x}{\text{linelength}} \\ C &= \frac{x_0 \Delta y - y_0 \Delta x}{\text{linelength}} \end{aligned}$$

with

$$\text{linelength} = \sqrt{\Delta x^2 + \Delta y^2}$$

Once the constants A , B , and C have been evaluated for the line, each processor needs to perform two multiplications and two additions to compute the pixel distance d . A pixel is plotted if d is less than a specified line-thickness parameter.

Instead of partitioning the screen into single pixels, we can assign to each processor either a scan line or a column of pixels depending on the line slope. Each processor then calculates the intersection of the line with the horizontal row or vertical column of pixels assigned that processor. For a line with slope $|m| < 1$, each processor simply solves the line equation for y , given an x column value. For a line with slope magnitude greater than 1, the line equation is solved for x by each processor, given a scan-line y value. Such direct methods, although slow on sequential machines, can be performed very efficiently using multiple processors.

3-3

LOADING THE FRAME BUFFER

When straight line segments and other objects are scan converted for display with a raster system, frame-buffer positions must be calculated. We have assumed that this is accomplished with the `setPixel` procedure, which stores intensity values for the pixels at corresponding addresses within the frame-buffer array. Scan-conversion algorithms generate pixel positions at successive unit in-

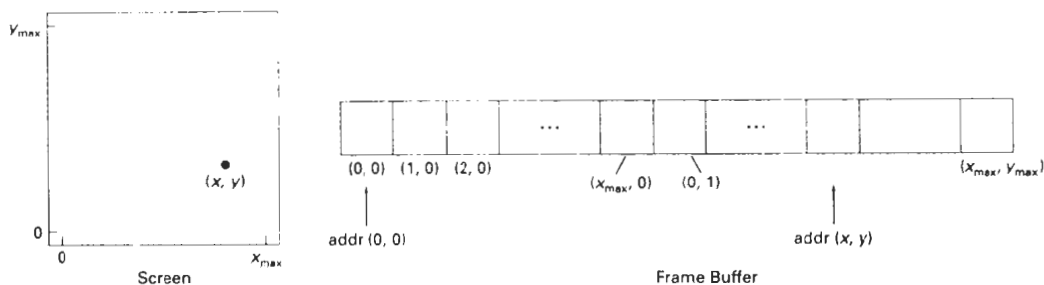


Figure 3-11
Pixel screen positions stored linearly in row-major order within the frame buffer.

tervals. This allows us to use incremental methods to calculate frame-buffer addresses.

As a specific example, suppose the frame-buffer array is addressed in row-major order and that pixel positions vary from (0, 0) at the lower left screen corner to (x_{\max}, y_{\max}) at the top right corner (Fig. 3-11). For a bilevel system (1 bit per pixel), the frame-buffer bit address for pixel position (x, y) is calculated as

$$\text{addr}(x, y) = \text{addr}(0, 0) + y(x_{\max} + 1) + x \quad (3-21)$$

Moving across a scan line, we can calculate the frame-buffer address for the pixel at $(x + 1, y)$ as the following offset from the address for position (x, y) :

$$\text{addr}(x + 1, y) = \text{addr}(x, y) + 1 \quad (3-22)$$

Stepping diagonally up to the next scan line from (x, y) , we get to the frame-buffer address of $(x + 1, y + 1)$ with the calculation

$$\text{addr}(x + 1, y + 1) = \text{addr}(x, y) + x_{\max} + 2 \quad (3-23)$$

where the constant $x_{\max} + 2$ is precomputed once for all line segments. Similar incremental calculations can be obtained from Eq. 3-21 for unit steps in the negative x and y screen directions. Each of these address calculations involves only a single integer addition.

Methods for implementing the `setPixel` procedure to store pixel intensity values depend on the capabilities of a particular system and the design requirements of the software package. With systems that can display a range of intensity values for each pixel, frame-buffer address calculations would include pixel width (number of bits), as well as the pixel screen location.

3-4

LINE FUNCTION

A procedure for specifying straight-line segments can be set up in a number of different forms. In PHIGS, GKS, and some other packages, the two-dimensional line function is

```
polyline (n, wcPoints)
```

where parameter *n* is assigned an integer value equal to the number of coordinate positions to be input, and *wcPoints* is the array of input world-coordinate values for line segment endpoints. This function is used to define a set of $n - 1$ connected straight line segments. Because series of connected line segments occur more often than isolated line segments in graphics applications, *polyline* provides a more general line function. To display a single straight-line segment, we set $n = 2$ and list the *x* and *y* values of the two endpoint coordinates in *wcPoints*.

As an example of the use of *polyline*, the following statements generate two connected line segments, with endpoints at (50, 100), (150, 250), and (250, 100):

```
wcPoints[1].x = 50;
wcPoints[1].y = 100;
wcPoints[2].x = 150;
wcPoints[2].y = 250;
wcPoints[3].x = 250;
wcPoints[3].y = 100;
polyline (3, wcPoints);
```

Coordinate references in the *polyline* function are stated as **absolute coordinate** values. This means that the values specified are the actual point positions in the coordinate system in use.

Some graphics systems employ line (and point) functions with **relative coordinate** specifications. In this case, coordinate values are stated as offsets from the last position referenced (called the **current position**). For example, if location (3, 2) is the last position that has been referenced in an application program, a relative coordinate specification of (2, -1) corresponds to an absolute position of (5, 1). An additional function is also available for setting the current position before the line routine is summoned. With these packages, a user lists only the single pair of offsets in the line command. This signals the system to display a line starting from the current position to a final position determined by the offsets. The current position is then updated to this final line position. A series of connected lines is produced with such packages by a sequence of line commands, one for each line section to be drawn. Some graphics packages provide options allowing the user to specify line endpoints using either relative or absolute coordinates.

Implementation of the *polyline* procedure is accomplished by first performing a series of coordinate transformations, then making a sequence of calls to a device-level line-drawing routine. In PHIGS, the input line endpoints are actually specified in modeling coordinates, which are then converted to world coordinates. Next, world coordinates are converted to normalized coordinates, then to device coordinates. We discuss the details for carrying out these two-dimensional coordinate transformations in Chapter 6. Once in device coordinates, we display the *polyline* by invoking a line routine, such as Bresenham's algorithm, $n - 1$ times to connect the *n* coordinate points. Each successive call passes the coordinate pair needed to plot the next line section, where the first endpoint of each coordinate pair is the last endpoint of the previous section. To avoid setting the intensity of some endpoints twice, we could modify the line algorithm so that the last endpoint of each segment is not plotted. We discuss methods for avoiding overlap of displayed objects in more detail in Section 3-10.

3-5

CIRCLE-GENERATING ALGORITHMS

Since the circle is a frequently used component in pictures and graphs, a procedure for generating either full circles or circular arcs is included in most graphics packages. More generally, a single procedure can be provided to display either circular or elliptical curves.

Properties of Circles

A circle is defined as the set of points that are all at a given distance r from a center position (x_c, y_c) (Fig. 3-12). This distance relationship is expressed by the Pythagorean theorem in Cartesian coordinates as

$$(x - x_c)^2 + (y - y_c)^2 = r^2 \quad (3-24)$$

We could use this equation to calculate the position of points on a circle circumference by stepping along the x axis in unit steps from $x_c - r$ to $x_c + r$ and calculating the corresponding y values at each position as

$$y = y_c \pm \sqrt{r^2 - (x - x_c)^2} \quad (3-25)$$

But this is not the best method for generating a circle. One problem with this approach is that it involves considerable computation at each step. Moreover, the spacing between plotted pixel positions is not uniform, as demonstrated in Fig. 3-13. We could adjust the spacing by interchanging x and y (stepping through y values and calculating x values) whenever the absolute value of the slope of the circle is greater than 1. But this simply increases the computation and processing required by the algorithm.

Another way to eliminate the unequal spacing shown in Fig. 3-13 is to calculate points along the circular boundary using polar coordinates r and θ (Fig. 3-12). Expressing the circle equation in parametric polar form yields the pair of equations

$$\begin{aligned} x &= x_c + r \cos \theta \\ y &= y_c + r \sin \theta \end{aligned} \quad (3-26)$$

When a display is generated with these equations using a fixed angular step size, a circle is plotted with equally spaced points along the circumference. The step size chosen for θ depends on the application and the display device. Larger angular separations along the circumference can be connected with straight line segments to approximate the circular path. For a more continuous boundary on a raster display, we can set the step size at $1/r$. This plots pixel positions that are approximately one unit apart.

Computation can be reduced by considering the symmetry of circles. The shape of the circle is similar in each quadrant. We can generate the circle section in the second quadrant of the xy plane by noting that the two circle sections are symmetric with respect to the y axis. And circle sections in the third and fourth quadrants can be obtained from sections in the first and second quadrants by

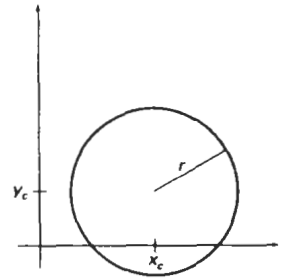


Figure 3-12
Circle with center coordinates (x_c, y_c) and radius r .

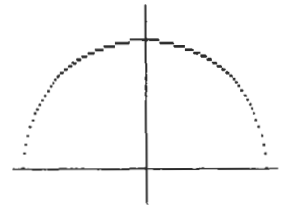


Figure 3-13
Positive half of a circle plotted with Eq. 3-25 and with $(x_c, y_c) = (0, 0)$.

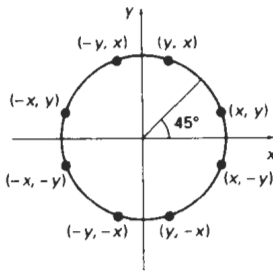


Figure 3-14
Symmetry of a circle.
Calculation of a circle point
(x, y) in one octant yields the
circle points shown for the
other seven octants.

considering symmetry about the x axis. We can take this one step further and note that there is also symmetry between octants. Circle sections in adjacent octants within one quadrant are symmetric with respect to the 45° line dividing the two octants. These symmetry conditions are illustrated in Fig.3-14, where a point at position (x, y) on a one-eighth circle sector is mapped into the seven circle points in the other octants of the xy plane. Taking advantage of the circle symmetry in this way, we can generate all pixel positions around a circle by calculating only the points within the sector from $x = 0$ to $x = y$.

Determining pixel positions along a circle circumference using either Eq. 3-24 or Eq. 3-26 still requires a good deal of computation time. The Cartesian equation 3-24 involves multiplications and square-root calculations, while the parametric equations contain multiplications and trigonometric calculations. More efficient circle algorithms are based on incremental calculation of decision parameters, as in the Bresenham line algorithm, which involves only simple integer operations.

Bresenham's line algorithm for raster displays is adapted to circle generation by setting up decision parameters for finding the closest pixel to the circumference at each sampling step. The circle equation 3-24, however, is nonlinear, so that square-root evaluations would be required to compute pixel distances from a circular path. Bresenham's circle algorithm avoids these square-root calculations by comparing the squares of the pixel separation distances.

A method for direct distance comparison is to test the halfway position between two pixels to determine if this midpoint is inside or outside the circle boundary. This method is more easily applied to other conics; and for an integer circle radius, the midpoint approach generates the same pixel positions as the Bresenham circle algorithm. Also, the error involved in locating pixel positions along any conic section using the midpoint test is limited to one-half the pixel separation.

Midpoint Circle Algorithm

As in the raster line algorithm, we sample at unit intervals and determine the closest pixel position to the specified circle path at each step. For a given radius r and screen center position (x_c, y_c), we can first set up our algorithm to calculate pixel positions around a circle path centered at the coordinate origin (0, 0). Then each calculated position (x, y) is moved to its proper screen position by adding x_c to x and y_c to y . Along the circle section from $x = 0$ to $x = y$ in the first quadrant, the slope of the curve varies from 0 to -1 . Therefore, we can take unit steps in the positive x direction over this octant and use a decision parameter to determine which of the two possible y positions is closer to the circle path at each step. Positions in the other seven octants are then obtained by symmetry.

To apply the midpoint method, we define a circle function:

$$f_{\text{circle}}(x, y) = x^2 + y^2 - r^2 \quad (3-27)$$

Any point (x, y) on the boundary of the circle with radius r satisfies the equation $f_{\text{circle}}(x, y) = 0$. If the point is in the interior of the circle, the circle function is negative. And if the point is outside the circle, the circle function is positive. To summarize, the relative position of any point (x, y) can be determined by checking the sign of the circle function:

$$f_{\text{circle}}(x, y) \begin{cases} < 0, & \text{if } (x, y) \text{ is inside the circle boundary} \\ = 0, & \text{if } (x, y) \text{ is on the circle boundary} \\ > 0, & \text{if } (x, y) \text{ is outside the circle boundary} \end{cases} \quad (3-28)$$

The circle-function tests in 3-28 are performed for the midpositions between pixels near the circle path at each sampling step. Thus, the circle function is the decision parameter in the midpoint algorithm, and we can set up incremental calculations for this function as we did in the line algorithm.

Figure 3-15 shows the midpoint between the two candidate pixels at sampling position $x_k + 1$. Assuming we have just plotted the pixel at (x_k, y_k) , we next need to determine whether the pixel at position $(x_k + 1, y_k)$ or the one at position $(x_k + 1, y_k - 1)$ is closer to the circle. Our decision parameter is the circle function 3-27 evaluated at the midpoint between these two pixels:

$$\begin{aligned} p_k &= f_{\text{circle}}\left(x_k + 1, y_k - \frac{1}{2}\right) \\ &= (x_k + 1)^2 + \left(y_k - \frac{1}{2}\right)^2 - r^2 \end{aligned} \quad (3-29)$$

If $p_k < 0$, this midpoint is inside the circle and the pixel on scan line y_k is closer to the circle boundary. Otherwise, the midpoint is outside or on the circle boundary, and we select the pixel on scanline $y_k - 1$.

Successive decision parameters are obtained using incremental calculations. We obtain a recursive expression for the next decision parameter by evaluating the circle function at sampling position $x_{k+1} + 1 = x_k + 2$:

$$\begin{aligned} p_{k+1} &= f_{\text{circle}}\left(x_{k+1} + 1, y_{k+1} - \frac{1}{2}\right) \\ &= [(x_k + 1) + 1]^2 + \left(y_{k+1} - \frac{1}{2}\right)^2 - r^2 \end{aligned}$$

or

$$p_{k+1} = p_k + 2(x_k + 1) + (y_{k+1}^2 - y_k^2) - (y_{k+1} - y_k) + 1 \quad (3-30)$$

where y_{k+1} is either y_k or y_{k-1} , depending on the sign of p_k .

Increments for obtaining p_{k+1} are either $2x_{k+1} + 1$ (if p_k is negative) or $2x_{k+1} + 1 - 2y_{k+1}$. Evaluation of the terms $2x_{k+1}$ and $2y_{k+1}$ can also be done incrementally as

$$2x_{k+1} = 2x_k + 2$$

$$2y_{k+1} = 2y_k - 2$$

At the start position $(0, r)$, these two terms have the values 0 and $2r$, respectively. Each successive value is obtained by adding 2 to the previous value of $2x$ and subtracting 2 from the previous value of $2y$.

The initial decision parameter is obtained by evaluating the circle function at the start position $(x_0, y_0) = (0, r)$:

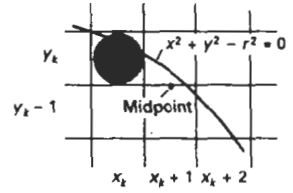


Figure 3-15
Midpoint between candidate pixels at sampling position $x_k + 1$ along a circular path.

$$\begin{aligned} p_0 &= f_{\text{circle}}\left(1, r - \frac{1}{2}\right) \\ &= 1 + \left(r - \frac{1}{2}\right)^2 - r^2 \end{aligned}$$

or

$$p_0 = \frac{5}{4} - r \quad (3-31)$$

If the radius r is specified as an integer, we can simply round p_0 to

$$p_0 = 1 - r \quad (\text{for } r \text{ an integer})$$

since all increments are integers.

As in Bresenham's line algorithm, the midpoint method calculates pixel positions along the circumference of a circle using integer additions and subtractions, assuming that the circle parameters are specified in integer screen coordinates. We can summarize the steps in the midpoint circle algorithm as follows.

Midpoint Circle Algorithm

1. Input radius r and circle center (x_c, y_c) , and obtain the first point on the circumference of a circle centered on the origin as

$$(x_0, y_0) = (0, r)$$

2. Calculate the initial value of the decision parameter as

$$p_0 = \frac{5}{4} - r$$

3. At each x_k position, starting at $k = 0$, perform the following test: If $p_k < 0$, the next point along the circle centered on $(0, 0)$ is (x_{k+1}, y_k) and

$$p_{k+1} = p_k + 2x_{k+1} + 1$$

Otherwise, the next point along the circle is $(x_k + 1, y_k - 1)$ and

$$p_{k+1} = p_k + 2x_{k+1} + 1 - 2y_{k+1}$$

where $2x_{k+1} = 2x_k + 2$ and $2y_{k+1} = 2y_k - 2$.

4. Determine symmetry points in the other seven octants.
5. Move each calculated pixel position (x, y) to the circular path centered on (x_c, y_c) and plot the coordinate values:

$$x = x + x_c \quad y = y + y_c$$

6. Repeat steps 3 through 5 until $x \geq y$.

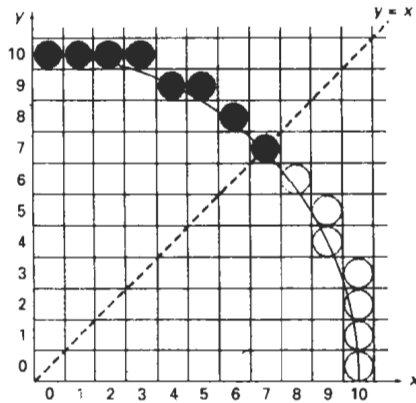


Figure 3-16

Selected pixel positions (solid circles) along a circle path with radius $r = 10$ centered on the origin, using the midpoint circle algorithm. Open circles show the symmetry positions in the first quadrant.

Example 3-2 Midpoint Circle-Drawing

Given a circle radius $r = 10$, we demonstrate the midpoint circle algorithm by determining positions along the circle octant in the first quadrant from $x = 0$ to $x = y$. The initial value of the decision parameter is

$$p_0 = 1 - r = -9$$

For the circle centered on the coordinate origin, the initial point is $(x_0, y_0) = (0, 10)$, and initial increment terms for calculating the decision parameters are

$$2x_0 = 0, \quad 2y_0 = 20$$

Successive decision parameter values and positions along the circle path are calculated using the midpoint method as

k	p_k	(x_{k+1}, y_{k+1})	$2x_{k+1}$	$2y_{k+1}$
0	-9	(1, 10)	2	20
1	-6	(2, 10)	4	20
2	-1	(3, 10)	6	20
3	6	(4, 9)	8	18
4	-3	(5, 9)	10	18
5	8	(6, 8)	12	16
6	5	(7, 7)	14	14

A plot of the generated pixel positions in the first quadrant is shown in Fig. 3-16.

The following procedure displays a raster circle on a bilevel monitor using the midpoint algorithm. Input to the procedure are the coordinates for the circle center and the radius. Intensities for pixel positions along the circle circumference are loaded into the frame-buffer array with calls to the `setPixel` routine.

```

#include "device.h"

void circleMidpoint (int xCenter, int yCenter, int radius)
{
    int x = 0;
    int y = radius;
    int p = 1 - radius;
    void circlePlotPoints (int, int, int, int);

    /* Plot first set of points */
    circlePlotPoints (xCenter, yCenter, x, y);

    while (x < y) {
        x++;
        if (p < 0)
            p += 2 * x + 1;
        else {
            y--;
            p += 2 * (x - y) + 1;
        }
        circlePlotPoints (xCenter, yCenter, x, y);
    }
}

void circlePlotPoints (int xCenter, int yCenter, int x, int y)
{
    setPixel (xCenter + x, yCenter + y);
    setPixel (xCenter - x, yCenter + y);
    setPixel (xCenter + x, yCenter - y);
    setPixel (xCenter - x, yCenter - y);
    setPixel (xCenter + y, yCenter + x);
    setPixel (xCenter - y, yCenter + x);
    setPixel (xCenter + y, yCenter - x);
    setPixel (xCenter - y, yCenter - x);
}

```

3-6

ELLIPSE-GENERATING ALGORITHMS

Loosely stated, an ellipse is an elongated circle. Therefore, elliptical curves can be generated by modifying circle-drawing procedures to take into account the different dimensions of an ellipse along the major and minor axes.

Properties of Ellipses

An ellipse is defined as the set of points such that the sum of the distances from two fixed positions (foci) is the same for all points (Fig. 3-17). If the distances to the two foci from any point $P = (x, y)$ on the ellipse are labeled d_1 and d_2 , then the general equation of an ellipse can be stated as

$$d_1 + d_2 = \text{constant} \quad (3-32)$$

Expressing distances d_1 and d_2 in terms of the focal coordinates $F_1 = (x_1, y_1)$ and $F_2 = (x_2, y_2)$, we have

$$\sqrt{(x - x_1)^2 + (y - y_1)^2} + \sqrt{(x - x_2)^2 + (y - y_2)^2} = \text{constant} \quad (3-33)$$

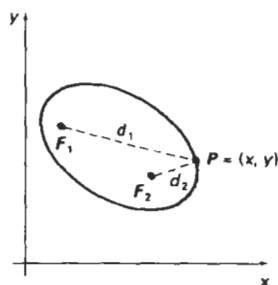


Figure 3-17

Ellipse generated about foci F_1 and F_2 .

By squaring this equation, isolating the remaining radical, and then squaring again, we can rewrite the general ellipse equation in the form

$$Ax^2 + By^2 + Cxy + Dx + Ey + F = 0 \quad (3-34)$$

where the coefficients A, B, C, D, E , and F are evaluated in terms of the focal coordinates and the dimensions of the major and minor axes of the ellipse. The major axis is the straight line segment extending from one side of the ellipse to the other through the foci. The minor axis spans the shorter dimension of the ellipse, bisecting the major axis at the halfway position (ellipse center) between the two foci.

An interactive method for specifying an ellipse in an arbitrary orientation is to input the two foci and a point on the ellipse boundary. With these three coordinate positions, we can evaluate the constant in Eq. 3-33. Then, the coefficients in Eq. 3-34 can be evaluated and used to generate pixels along the elliptical path.

Ellipse equations are greatly simplified if the major and minor axes are oriented to align with the coordinate axes. In Fig. 3-18, we show an ellipse in "standard position" with major and minor axes oriented parallel to the x and y axes. Parameter r_x for this example labels the semimajor axis, and parameter r_y labels the semiminor axis. The equation of the ellipse shown in Fig. 3-18 can be written in terms of the ellipse center coordinates and parameters r_x and r_y as

$$\left(\frac{x - x_c}{r_x}\right)^2 + \left(\frac{y - y_c}{r_y}\right)^2 = 1 \quad (3-35)$$

Using polar coordinates r and θ , we can also describe the ellipse in standard position with the parametric equations:

$$\begin{aligned} x &= x_c + r_x \cos \theta \\ y &= y_c + r_y \sin \theta \end{aligned} \quad (3-36)$$

Symmetry considerations can be used to further reduce computations. An ellipse in standard position is symmetric between quadrants, but unlike a circle, it is not symmetric between the two octants of a quadrant. Thus, we must calculate pixel positions along the elliptical arc throughout one quadrant, then we obtain positions in the remaining three quadrants by symmetry (Fig 3-19).

Midpoint Ellipse Algorithm

Our approach here is similar to that used in displaying a raster circle. Given parameters r_x , r_y , and (x_c, y_c) , we determine points (x, y) for an ellipse in standard position centered on the origin, and then we shift the points so the ellipse is centered at (x_c, y_c) . If we wish also to display the ellipse in nonstandard position, we could then rotate the ellipse about its center coordinates to reorient the major and minor axes. For the present, we consider only the display of ellipses in standard position. We discuss general methods for transforming object orientations and positions in Chapter 5.

The midpoint ellipse method is applied throughout the first quadrant in two parts. Figure 3-20 shows the division of the first quadrant according to the slope of an ellipse with $r_x < r_y$. We process this quadrant by taking unit steps in the x direction where the slope of the curve has a magnitude less than 1, and taking unit steps in the y direction where the slope has a magnitude greater than 1.

Regions 1 and 2 (Fig. 3-20), can be processed in various ways. We can start at position $(0, r_y)$ and step clockwise along the elliptical path in the first quadrant,

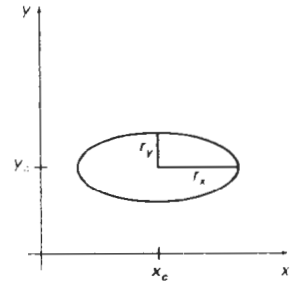


Figure 3-18
Ellipse centered at (x_c, y_c) with semimajor axis r_x and semiminor axis r_y .

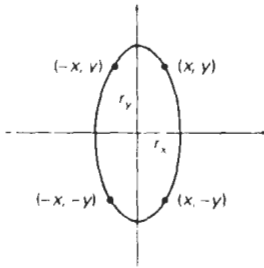


Figure 3-19
Symmetry of an ellipse.
Calculation of a point (x, y)
in one quadrant yields the
ellipse points shown for the
other three quadrants.

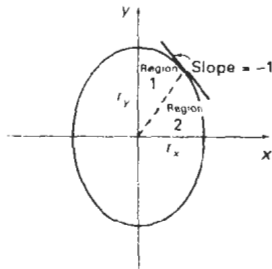


Figure 3-20
Ellipse processing regions.
Over region 1, the magnitude
of the ellipse slope is less
than 1; over region 2, the
magnitude of the slope is
greater than 1.

shifting from unit steps in x to unit steps in y when the slope becomes less than -1 . Alternatively, we could start at $(r_x, 0)$ and select points in a counterclockwise order, shifting from unit steps in y to unit steps in x when the slope becomes greater than -1 . With parallel processors, we could calculate pixel positions in the two regions simultaneously. As an example of a sequential implementation of the midpoint algorithm, we take the start position at $(0, r_y)$ and step along the ellipse path in clockwise order throughout the first quadrant.

We define an ellipse function from Eq. 3-35 with $(x_c, y_c) = (0, 0)$ as

$$f_{\text{ellipse}}(x, y) = r_y^2 x^2 + r_x^2 y^2 - r_x^2 r_y^2 \quad (3-37)$$

which has the following properties:

$$f_{\text{ellipse}}(x, y) \begin{cases} < 0, & \text{if } (x, y) \text{ is inside the ellipse boundary} \\ = 0, & \text{if } (x, y) \text{ is on the ellipse boundary} \\ > 0, & \text{if } (x, y) \text{ is outside the ellipse boundary} \end{cases} \quad (3-38)$$

Thus, the ellipse function $f_{\text{ellipse}}(x, y)$ serves as the decision parameter in the midpoint algorithm. At each sampling position, we select the next pixel along the ellipse path according to the sign of the ellipse function evaluated at the midpoint between the two candidate pixels.

Starting at $(0, r_y)$, we take unit steps in the x direction until we reach the boundary between region 1 and region 2 (Fig. 3-20). Then we switch to unit steps in the y direction over the remainder of the curve in the first quadrant. At each step, we need to test the value of the slope of the curve. The ellipse slope is calculated from Eq. 3-37 as

$$\frac{dy}{dx} = -\frac{2r_y^2 x}{2r_x^2 y} \quad (3-39)$$

At the boundary between region 1 and region 2, $dy/dx = -1$ and

$$2r_y^2 x = 2r_x^2 y$$

Therefore, we move out of region 1 whenever

$$2r_y^2 x \geq 2r_x^2 y \quad (3-40)$$

Figure 3-21 shows the midpoint between the two candidate pixels at sampling position $x_k + 1$ in the first region. Assuming position (x_k, y_k) has been selected at the previous step, we determine the next position along the ellipse path by evaluating the decision parameter (that is, the ellipse function 3-37) at this midpoint:

$$\begin{aligned} p1_k &= f_{\text{ellipse}}\left(x_k + 1, y_k - \frac{1}{2}\right) \\ &= r_y^2 (x_k + 1)^2 + r_x^2 \left(y_k - \frac{1}{2}\right)^2 - r_x^2 r_y^2 \end{aligned} \quad (3-41)$$

If $p1_k < 0$, the midpoint is inside the ellipse and the pixel on scan line y_k is closer to the ellipse boundary. Otherwise, the midpoint is outside or on the ellipse boundary, and we select the pixel on scan line $y_k - 1$.

At the next sampling position $(x_{k+1} + 1 = x_k + 2)$, the decision parameter for region 1 is evaluated as

$$\begin{aligned} p1_{k+1} &= f_{\text{ellipse}}\left(x_{k+1} + 1, y_{k+1} - \frac{1}{2}\right) \\ &= r_y^2[(x_k + 1) + 1]^2 + r_x^2\left(y_{k+1} - \frac{1}{2}\right)^2 - r_x^2 r_y^2 \end{aligned}$$

or

$$p1_{k+1} = p1_k + 2r_y^2(x_k + 1) + r_y^2 + r_x^2\left[\left(y_{k+1} - \frac{1}{2}\right)^2 - \left(y_k - \frac{1}{2}\right)^2\right] \quad (3-42)$$

where y_{k+1} is either y_k or $y_k - 1$, depending on the sign of $p1_k$.

Decision parameters are incremented by the following amounts:

$$\text{increment} = \begin{cases} 2r_y^2 x_{k+1} + r_y^2, & \text{if } p1_k < 0 \\ 2r_y^2 x_{k+1} + r_y^2 - 2r_x^2 y_{k+1}, & \text{if } p1_k \geq 0 \end{cases}$$

As in the circle algorithm, increments for the decision parameters can be calculated using only addition and subtraction, since values for the terms $2r_y^2 x$ and $2r_x^2 y$ can also be obtained incrementally. At the initial position $(0, r_y)$, the two terms evaluate to

$$2r_y^2 x = 0 \quad (3-43)$$

$$2r_x^2 y = 2r_x^2 r_y \quad (3-44)$$

As x and y are incremented, updated values are obtained by adding $2r_y^2$ to 3-43 and subtracting $2r_x^2$ from 3-44. The updated values are compared at each step, and we move from region 1 to region 2 when condition 3-40 is satisfied.

In region 1, the initial value of the decision parameter is obtained by evaluating the ellipse function at the start position $(x_0, y_0) = (0, r_y)$:

$$\begin{aligned} p1_0 &= f_{\text{ellipse}}\left(1, r_y - \frac{1}{2}\right) \\ &= r_y^2 + r_x^2\left(r_y - \frac{1}{2}\right)^2 - r_x^2 r_y^2 \end{aligned}$$

or

$$p1_0 = r_y^2 - r_x^2 r_y + \frac{1}{4} r_x^2 \quad (3-45)$$

Over region 2, we sample at unit steps in the negative y direction, and the midpoint is now taken between horizontal pixels at each step (Fig. 3-22). For this region, the decision parameter is evaluated as

$$\begin{aligned} p2_k &= f_{\text{ellipse}}\left(x_k + \frac{1}{2}, y_k - 1\right) \\ &= r_y^2\left(x_k + \frac{1}{2}\right)^2 + r_x^2(y_k - 1)^2 - r_x^2 r_y^2 \end{aligned} \quad (3-46)$$

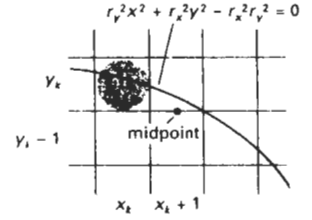


Figure 3-21
Midpoint between candidate pixels at sampling position $x_k + 1$ along an elliptical path.

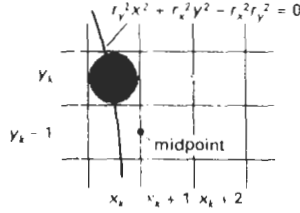


Figure 3-22
Midpoint between candidate pixels at sampling position $y_k - 1$ along an elliptical path.

If $p2_k > 0$, the midposition is outside the ellipse boundary, and we select the pixel at x_k . If $p2_k \leq 0$, the midpoint is inside or on the ellipse boundary, and we select pixel position x_{k+1} .

To determine the relationship between successive decision parameters in region 2, we evaluate the ellipse function at the next sampling step $y_{k+1} - 1 = y_k - 2$:

$$\begin{aligned} p2_{k+1} &= f_{\text{ellipse}}\left(x_{k+1} + \frac{1}{2}, y_{k+1} - 1\right) \\ &= r_y^2 \left(x_{k+1} + \frac{1}{2}\right)^2 + r_x^2 (y_k - 1)^2 - r_x^2 r_y^2 \end{aligned} \quad (3-47)$$

or

$$p2_{k+1} = p2_k - 2r_x^2(y_k - 1) + r_x^2 + r_y^2 \left[\left(x_{k+1} + \frac{1}{2}\right)^2 - \left(x_k + \frac{1}{2}\right)^2 \right] \quad (3-48)$$

with x_{k+1} set either to x_k or to $x_k + 1$, depending on the sign of $p2_k$.

When we enter region 2, the initial position (x_0, y_0) is taken as the last position selected in region 1 and the initial decision parameter in region 2 is then

$$\begin{aligned} p2_0 &= f_{\text{ellipse}}\left(x_0 + \frac{1}{2}, y_0 - 1\right) \\ &= r_y^2 \left(x_0 + \frac{1}{2}\right)^2 + r_x^2 (y_0 - 1)^2 - r_x^2 r_y^2 \end{aligned} \quad (3-49)$$

To simplify the calculation of $p2_0$, we could select pixel positions in counterclockwise order starting at $(r_x, 0)$. Unit steps would then be taken in the positive y direction up to the last position selected in region 1.

The midpoint algorithm can be adapted to generate an ellipse in nonstandard position using the ellipse function Eq. 3-34 and calculating pixel positions over the entire elliptical path. Alternatively, we could reorient the ellipse axes to standard position, using transformation methods discussed in Chapter 5, apply the midpoint algorithm to determine curve positions, then convert calculated pixel positions to path positions along the original ellipse orientation.

Assuming r_x , r_y , and the ellipse center are given in integer screen coordinates, we only need incremental integer calculations to determine values for the decision parameters in the midpoint ellipse algorithm. The increments r_x^2 , r_y^2 , $2r_x^2$, and $2r_y^2$ are evaluated once at the beginning of the procedure. A summary of the midpoint ellipse algorithm is listed in the following steps:

Midpoint Ellipse Algorithm

1. Input r_x , r_y , and ellipse center (x_c, y_c) , and obtain the first point on an ellipse centered on the origin as

$$(x_0, y_0) = (0, r_y)$$

2. Calculate the initial value of the decision parameter in region 1 as

$$p1_0 = r_y^2 - r_x^2 r_y + \frac{1}{4} r_x^2$$

3. At each x_k position in region 1, starting at $k = 0$, perform the following test: If $p1_k < 0$, the next point along the ellipse centered on $(0, 0)$ is (x_{k+1}, y_k) and

$$p1_{k+1} = p1_k + 2r_y^2 x_{k+1} + r_y^2$$

Otherwise, the next point along the circle is $(x_k + 1, y_k - 1)$ and

$$p1_{k+1} = p1_k + 2r_y^2 x_{k+1} - 2r_x^2 y_{k+1} + r_y^2$$

with

$$2r_y^2 x_{k+1} = 2r_y^2 x_k + 2r_y^2, \quad 2r_x^2 y_{k+1} = 2r_x^2 y_k - 2r_x^2$$

and continue until $2r_y^2 x \geq 2r_x^2 y$.

4. Calculate the initial value of the decision parameter in region 2 using the last point (x_0, y_0) calculated in region 1 as

$$p2_0 = r_y^2 \left(x_0 + \frac{1}{2} \right)^2 + r_x^2 (y_0 - 1)^2 - r_x^2 r_y^2$$

5. At each y_k position in region 2, starting at $k = 0$, perform the following test: If $p2_k > 0$, the next point along the ellipse centered on $(0, 0)$ is $(x_k, y_{k+1} - 1)$ and

$$p2_{k+1} = p2_k - 2r_x^2 y_{k+1} + r_x^2$$

Otherwise, the next point along the circle is $(x_k + 1, y_k - 1)$ and

$$p2_{k+1} = p2_k + 2r_y^2 x_{k+1} - 2r_x^2 y_{k+1} + r_x^2$$

using the same incremental calculations for x and y as in region 1.

6. Determine symmetry points in the other three quadrants.
7. Move each calculated pixel position (x, y) onto the elliptical path centered on (x_c, y_c) and plot the coordinate values:

$$x = x + x_c \quad y = y + y_c$$

8. Repeat the steps for region 1 until $2r_y^2 x \geq 2r_x^2 y$.

Example 3-3 Midpoint Ellipse Drawing

Given input ellipse parameters $r_x = 8$ and $r_y = 6$, we illustrate the steps in the midpoint ellipse algorithm by determining raster positions along the ellipse path in the first quadrant. Initial values and increments for the decision parameter calculations are

$$\begin{aligned} 2r_y^2x &= 0 && \text{(with increment } 2r_y^2 = 72) \\ 2r_x^2y &= 2r_x^2r_y && \text{(with increment } -2r_x^2 = -128) \end{aligned}$$

For region 1: The initial point for the ellipse centered on the origin is $(x_0, y_0) = (0, 6)$, and the initial decision parameter value is

$$p1_0 = r_y^2 - r_x^2r_y + \frac{1}{4}r_x^2 = -332$$

Successive decision parameter values and positions along the ellipse path are calculated using the midpoint method as

k	$p1_k$	(x_{k+1}, y_{k+1})	$2r_y^2x_{k+1}$	$2r_x^2y_{k+1}$
0	-332	(1, 6)	72	768
1	-224	(2, 6)	144	768
2	-44	(3, 6)	216	768
3	208	(4, 5)	288	640
4	-108	(5, 5)	360	640
5	288	(6, 4)	432	512
6	244	(7, 3)	504	384

We now move out of region 1, since $2r_y^2x > 2r_x^2y$.

For region 2, the initial point is $(x_0, y_0) = (7, 3)$ and the initial decision parameter is

$$p2_0 = f\left(7 + \frac{1}{2}, 2\right) = -151$$

The remaining positions along the ellipse path in the first quadrant are then calculated as

k	$p2_k$	(x_{k+1}, y_{k+1})	$2r_y^2x_{k+1}$	$2r_x^2y_{k+1}$
0	-151	(8, 2)	576	256
1	233	(8, 1)	576	128
2	745	(8, 0)	—	—

A plot of the selected positions around the ellipse boundary within the first quadrant is shown in Fig. 3-23.

In the following procedure, the midpoint algorithm is used to display an ellipse with input parameters R_x , R_y , $xCenter$, and $yCenter$. Positions along the

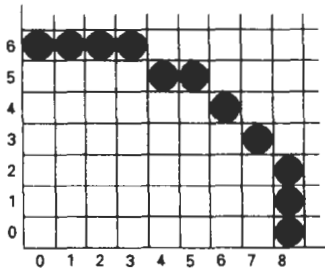


Figure 3-23

Positions along an elliptical path centered on the origin with $r_x = 8$ and $r_y = 6$ using the midpoint algorithm to calculate pixel addresses in the first quadrant.

curve in the first quadrant are generated and then shifted to their proper screen positions. Intensities for these positions and the symmetry positions in the other three quadrants are loaded into the frame buffer using the `setPixel` routine.

```
#include "device.h"

#define ROUND(a) ((int)(a+0.5))

void ellipseMidpoint (int xCenter, int yCenter, int Rx, int Ry)
{
    int Rx2 = Rx*Rx;
    int Ry2 = Ry*Ry;
    int twoRx2 = 2*Rx2;
    int twoRy2 = 2*Ry2;
    int p;
    int x = 0;
    int y = Ry;
    int px = 0;
    int py = twoRx2 * y;
    void ellipsePlotPoints (int, int, int, int);

    /* Plot the first set of points */
    ellipsePlotPoints (xCenter, yCenter, x, y);

    /* Region 1 */
    p = ROUND (Ry2 - (Rx2 * Ry) + (0.25 * Rx2));
    while (px < py) {
        x++;
        px += twoRy2;
        if (p < 0)
            p += Ry2 + px;
        else {
            y--;
            py -= twoRx2;
            p += Ry2 + px - py;
        }
        ellipsePlotPoints (xCenter, yCenter, x, y);
    }

    /* Region 2 */
    p = ROUND (Ry2*(x+0.5)*(x+0.5) + Rx2*(y-1)*(y-1) - Rx2*Ry2);
    while (y > 0) {
        y--;
        py -= twoRx2;
        if (p > 0)
            p += Rx2 - py;
        else {
            x++;
            px += twoRy2;
            p += Rx2 - py + px;
        }
    }
}
```

```
    }  
    ellipsePlotPoints (xCenter, yCenter, x, y);  
  }  
}  
  
void ellipsePlotPoints (int xCenter, int yCenter, int x, int y)  
{  
  setPixel (xCenter + x, yCenter + y);  
  setPixel (xCenter - x, yCenter + y);  
  setPixel (xCenter + x, yCenter - y);  
  setPixel (xCenter - x, yCenter - y);  
}
```

3-7

OTHER CURVES

Various curve functions are useful in object modeling, animation path specifications, data and function graphing, and other graphics applications. Commonly encountered curves include conics, trigonometric and exponential functions, probability distributions, general polynomials, and spline functions. Displays of these curves can be generated with methods similar to those discussed for the circle and ellipse functions. We can obtain positions along curve paths directly from explicit representations $y = f(x)$ or from parametric forms. Alternatively, we could apply the incremental midpoint method to plot curves described with implicit functions $f(x, y) = 0$.

A straightforward method for displaying a specified curve function is to approximate it with straight line segments. Parametric representations are useful in this case for obtaining equally spaced line endpoint positions along the curve path. We can also generate equally spaced positions from an explicit representation by choosing the independent variable according to the slope of the curve. Where the slope of $y = f(x)$ has a magnitude less than 1, we choose x as the independent variable and calculate y values at equal x increments. To obtain equal spacing where the slope has a magnitude greater than 1, we use the inverse function, $x = f^{-1}(y)$, and calculate values of x at equal y steps.

Straight-line or curve approximations are used to graph a data set of discrete coordinate points. We could join the discrete points with straight line segments, or we could use linear regression (least squares) to approximate the data set with a single straight line. A nonlinear least-squares approach is used to display the data set with some approximating function, usually a polynomial.

As with circles and ellipses, many functions possess symmetries that can be exploited to reduce the computation of coordinate positions along curve paths. For example, the normal probability distribution function is symmetric about a center position (the mean), and all points along one cycle of a sine curve can be generated from the points in a 90° interval.

Conic Sections

In general, we can describe a **conic section** (or **conic**) with the second-degree equation:

$$Ax^2 + By^2 + Cxy + Dx + Ey + F = 0 \quad (3-50)$$

where values for parameters A, B, C, D, E , and F determine the kind of curve we are to display. Given this set of coefficients, we can determine the particular conic that will be generated by evaluating the discriminant $B^2 - 4AC$:

$$B^2 - 4AC \begin{cases} < 0, & \text{generates an ellipse (or circle)} \\ = 0, & \text{generates a parabola} \\ > 0, & \text{generates a hyperbola} \end{cases} \quad (3-51)$$

For example, we get the circle equation 3-24 when $A = B = 1, C = 0, D = -2x_c, E = -2y_c$, and $F = x_c^2 + y_c^2 - r^2$. Equation 3-50 also describes the "degenerate" conics: points and straight lines.

Ellipses, hyperbolas, and parabolas are particularly useful in certain animation applications. These curves describe orbital and other motions for objects subjected to gravitational, electromagnetic, or nuclear forces. Planetary orbits in the solar system, for example, are ellipses; and an object projected into a uniform gravitational field travels along a parabolic trajectory. Figure 3-24 shows a parabolic path in standard position for a gravitational field acting in the negative y direction. The explicit equation for the parabolic trajectory of the object shown can be written as

$$y = y_0 + a(x - x_0)^2 + b(x - x_0) \quad (3-52)$$

with constants a and b determined by the initial velocity v_0 of the object and the acceleration g due to the uniform gravitational force. We can also describe such parabolic motions with parametric equations using a time parameter t , measured in seconds from the initial projection point:

$$\begin{aligned} x &= x_0 + v_{x0}t \\ y &= y_0 + v_{y0}t - \frac{1}{2}gt^2 \end{aligned} \quad (3-53)$$

Here, v_{x0} and v_{y0} are the initial velocity components, and the value of g near the surface of the earth is approximately 980cm/sec^2 . Object positions along the parabolic path are then calculated at selected time steps.

Hyperbolic motions (Fig. 3-25) occur in connection with the collision of charged particles and in certain gravitational problems. For example, comets or meteorites moving around the sun may travel along hyperbolic paths and escape to outer space, never to return. The particular branch (left or right, in Fig. 3-25) describing the motion of an object depends on the forces involved in the problem. We can write the standard equation for the hyperbola centered on the origin in Fig. 3-25 as

$$\left(\frac{x}{r_x}\right)^2 - \left(\frac{y}{r_y}\right)^2 = 1 \quad (3-54)$$

with $x \leq -r_x$ for the left branch and $x \geq r_x$ for the right branch. Since this equation differs from the standard ellipse equation 3-35 only in the sign between the x^2 and y^2 terms, we can generate points along a hyperbolic path with a slightly modified ellipse algorithm. We will return to the discussion of animation applications and methods in more detail in Chapter 16. And in Chapter 10, we discuss applications of computer graphics in scientific visualization.

Section 3-7

Other Curves

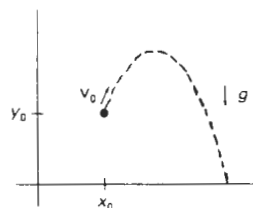


Figure 3-24
Parabolic path of an object tossed into a downward gravitational field at the initial position (x_0, y_0) .

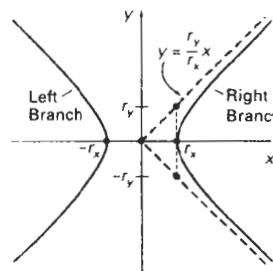


Figure 3-25
Left and right branches of a hyperbola in standard position with symmetry axis along the x axis.

Parabolas and hyperbolas possess a symmetry axis. For example, the parabola described by Eq. 3-53 is symmetric about the axis:

$$x = x_0 + v_{x0}v_{y0}/g$$

The methods used in the midpoint ellipse algorithm can be directly applied to obtain points along one side of the symmetry axis of hyperbolic and parabolic paths in the two regions: (1) where the magnitude of the curve slope is less than 1, and (2) where the magnitude of the slope is greater than 1. To do this, we first select the appropriate form of Eq. 3-50 and then use the selected function to set up expressions for the decision parameters in the two regions.

Polynomials and Spline Curves

A polynomial function of n th degree in x is defined as

$$y = \sum_{k=0}^n a_k x^k \quad (3-55)$$

$$= a_0 + a_1 x + \cdots + a_{n-1} x^{n-1} + a_n x^n$$

where n is a nonnegative integer and the a_k are constants, with $a_n \neq 0$. We get a quadratic when $n = 2$; a cubic polynomial when $n = 3$; a quartic when $n = 4$; and so forth. And we have a straight line when $n = 1$. Polynomials are useful in a number of graphics applications, including the design of object shapes, the specification of animation paths, and the graphing of data trends in a discrete set of data points.

Designing object shapes or motion paths is typically done by specifying a few points to define the general curve contour, then fitting the selected points with a polynomial. One way to accomplish the curve fitting is to construct a cubic polynomial curve section between each pair of specified points. Each curve section is then described in parametric form as

$$x = a_{x0} + a_{x1}u + a_{x2}u^2 + a_{x3}u^3 \quad (3-56)$$

$$y = a_{y0} + a_{y1}u + a_{y2}u^2 + a_{y3}u^3 \quad (3-57)$$



Figure 3-26
A spline curve formed with individual cubic polynomial sections between specified coordinate points.

where parameter u varies over the interval 0 to 1. Values for the coefficients of u in the parametric equations are determined from boundary conditions on the curve sections. One boundary condition is that two adjacent curve sections have the same coordinate position at the boundary, and a second condition is to match the two curve slopes at the boundary so that we obtain one continuous, smooth curve (Fig. 3-26). Continuous curves that are formed with polynomial pieces are called **spline curves**, or simply **splines**. There are other ways to set up spline curves, and the various spline-generating methods are explored in Chapter 10.

3-8

PARALLEL CURVE ALGORITHMS

Methods for exploiting parallelism in curve generation are similar to those used in displaying straight line segments. We can either adapt a sequential algorithm by allocating processors according to curve partitions, or we could devise other

methods and assign processors to screen partitions.

A parallel midpoint method for displaying circles is to divide the circular arc from 90° to 45° into equal subarcs and assign a separate processor to each subarc. As in the parallel Bresenham line algorithm, we then need to set up computations to determine the beginning y value and decision parameter p_k value for each processor. Pixel positions are then calculated throughout each subarc, and positions in the other circle octants are then obtained by symmetry. Similarly, a parallel ellipse midpoint method divides the elliptical arc over the first quadrant into equal subarcs and parcels these out to separate processors. Pixel positions in the other quadrants are determined by symmetry. A screen-partitioning scheme for circles and ellipses is to assign each scan line crossing the curve to a separate processor. In this case, each processor uses the circle or ellipse equation to calculate curve-intersection coordinates.

For the display of elliptical arcs or other curves, we can simply use the scan-line partitioning method. Each processor uses the curve equation to locate the intersection positions along its assigned scan line. With processors assigned to individual pixels, each processor would calculate the distance (or distance squared) from the curve to its assigned pixel. If the calculated distance is less than a predefined value, the pixel is plotted.

3-9

CURVE FUNCTIONS

Routines for circles, splines, and other commonly used curves are included in many graphics packages. The PHIGS standard does not provide explicit functions for these curves, but it does include the following general curve function:

```
generalizedDrawingPrimitive (n, wcPoints, id, datalist)
```

where *wcPoints* is a list of n coordinate positions, *datalist* contains noncoordinate data values, and parameter *id* selects the desired function. At a particular installation, a circle might be referenced with *id* = 1, an ellipse with *id* = 2, and so on.

As an example of the definition of curves through this PHIGS function, a circle (*id* = 1, say) could be specified by assigning the two center coordinate values to *wcpoints* and assigning the radius value to *datalist*. The generalized drawing primitive would then reference the appropriate algorithm, such as the midpoint method, to generate the circle. With interactive input, a circle could be defined with two coordinate points: the center position and a point on the circumference. Similarly, interactive specification of an ellipse can be done with three points: the two foci and a point on the ellipse boundary, all stored in *wc-points*. For an ellipse in standard position, *wcpoints* could be assigned only the center coordinates, with *datalist* assigned the values for r_x and r_y . Splines defined with control points would be generated by assigning the control point coordinates to *wcpoints*.

Functions to generate circles and ellipses often include the capability of drawing curve sections by specifying parameters for the line endpoints. Expanding the parameter list allows specification of the beginning and ending angular values for an arc, as illustrated in Fig. 3-27. Another method for designating a cir-

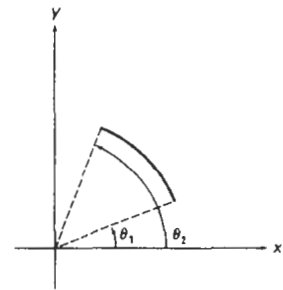


Figure 3-27
Circular arc specified by beginning and ending angles. Circle center is at the coordinate origin.

cular or elliptical arc is to input the beginning and ending coordinate positions of the arc.

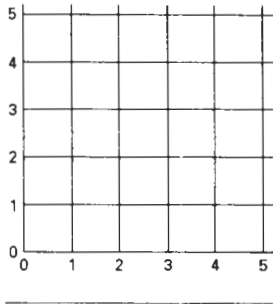


Figure 3-28
Lower-left section of the screen grid referencing integer coordinate positions.

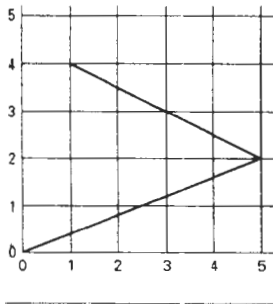


Figure 3-29
Line path for a series of connected line segments between screen grid coordinate positions.

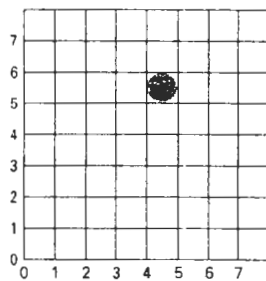


Figure 3-30
Illuminated pixel at raster position (4, 5).

3-10 PIXEL ADDRESSING AND OBJECT GEOMETRY

So far we have assumed that all input positions were given in terms of scan-line number and pixel-position number across the scan line. As we saw in Chapter 2, there are, in general, several coordinate references associated with the specification and generation of a picture. Object descriptions are given in a world-reference frame, chosen to suit a particular application, and input world coordinates are ultimately converted to screen display positions. World descriptions of objects are given in terms of precise coordinate positions, which are infinitesimally small mathematical points. Pixel coordinates, however, reference finite screen areas. If we want to preserve the specified geometry of world objects, we need to compensate for the mapping of mathematical input points to finite pixel areas. One way to do this is simply to adjust the dimensions of displayed objects to account for the amount of overlap of pixel areas with the object boundaries. Another approach is to map world coordinates onto screen positions between pixels, so that we align object boundaries with pixel boundaries instead of pixel centers.

Screen Grid Coordinates

An alternative to addressing display positions in terms of pixel centers is to reference screen coordinates with respect to the grid of horizontal and vertical pixel boundary lines spaced one unit apart (Fig. 3-28). A screen coordinate position is then the pair of integer values identifying a grid intersection position between two pixels. For example, the mathematical line path for a polyline with screen endpoints (0, 0), (5, 2), and (1, 4) is shown in Fig. 3-29.

With the coordinate origin at the lower left of the screen, each pixel area can be referenced by the integer grid coordinates of its lower left corner. Figure 3-30 illustrates this convention for an 8 by 8 section of a raster, with a single illuminated pixel at screen coordinate position (4, 5). In general, we identify the area occupied by a pixel with screen coordinates (x, y) as the unit square with diagonally opposite corners at (x, y) and $(x + 1, y + 1)$. This pixel-addressing scheme has several advantages: It avoids half-integer pixel boundaries, it facilitates precise object representations, and it simplifies the processing involved in many scan-conversion algorithms and in other raster procedures.

The algorithms for line drawing and curve generation discussed in the preceding sections are still valid when applied to input positions expressed as screen grid coordinates. Decision parameters in these algorithms are now simply a measure of screen grid separation differences, rather than separation differences from pixel centers.

Maintaining Geometric Properties of Displayed Objects

When we convert geometric descriptions of objects into pixel representations, we transform mathematical points and lines into finite screen areas. If we are to maintain the original geometric measurements specified by the input coordinates

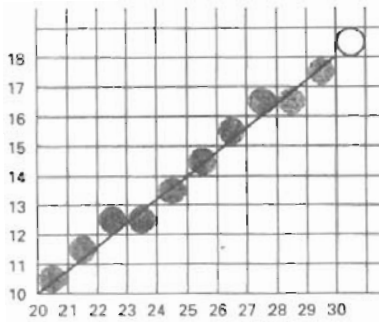


Figure 3-31
Line path and corresponding pixel
display for input screen grid
endpoint coordinates (20, 10) and
(30, 18).

for an object, we need to account for the finite size of pixels when we transform the object definition to a screen display.

Figure 3-31 shows the line plotted in the Bresenham line-algorithm example of Section 3-2. Interpreting the line endpoints (20, 10) and (30, 18) as precise grid crossing positions, we see that the line should not extend past screen grid position (30, 18). If we were to plot the pixel with screen coordinates (30, 18), as in the example given in Section 3-2, we would display a line that spans 11 horizontal units and 9 vertical units. For the mathematical line, however, $\Delta x = 10$ and $\Delta y = 8$. If we are addressing pixels by their center positions, we can adjust the length of the displayed line by omitting one of the endpoint pixels. If we think of screen coordinates as addressing pixel boundaries, as shown in Fig. 3-31, we plot a line using only those pixels that are "interior" to the line path; that is, only those pixels that are between the line endpoints. For our example, we would plot the leftmost pixel at (20, 10) and the rightmost pixel at (29, 17). This displays a line that

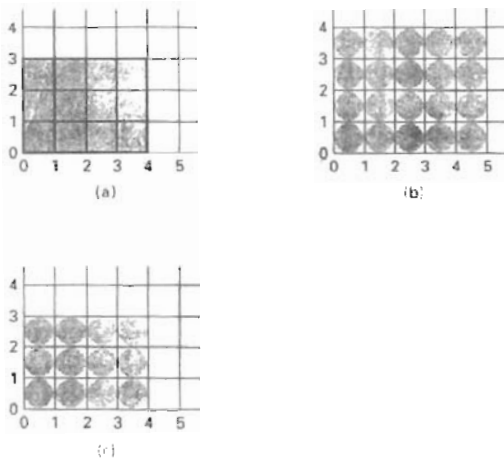


Figure 3-32
Conversion of rectangle (a) with vertices at screen
coordinates (0, 0), (4, 0), (4, 3), and (0, 3) into display
(b) that includes the right and top boundaries and into
display (c) that maintains geometric magnitudes.

has the same geometric magnitudes as the mathematical line from (20, 10) to (30, 18).

For an enclosed area, input geometric properties are maintained by displaying the area only with those pixels that are interior to the object boundaries. The rectangle defined with the screen coordinate vertices shown in Fig. 3-32(a), for example, is larger when we display it filled with pixels up to and including the border pixel lines joining the specified vertices. As defined, the area of the rectangle is 12 units, but as displayed in Fig. 3-32(b), it has an area of 20 units. In Fig. 3-32(c), the original rectangle measurements are maintained by displaying

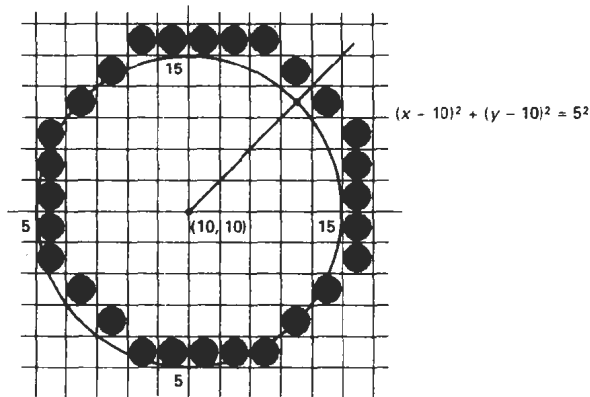


Figure 3-33
Circle path and midpoint circle algorithm plot of a circle with radius 5 in screen coordinates.

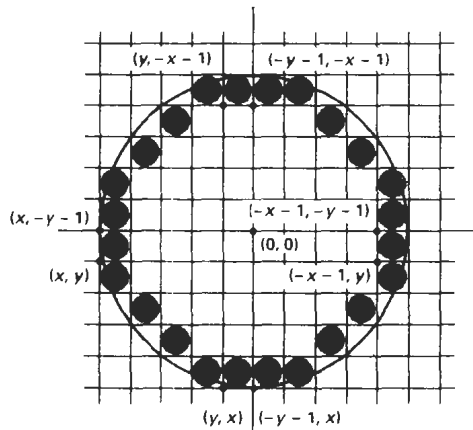


Figure 3-34
Modification of the circle plot in Fig. 3-33 to maintain the specified circle diameter of 10.

only the internal pixels. The right boundary of the input rectangle is at $x = 4$. To maintain this boundary in the display, we set the rightmost pixel grid coordinate at $x = 3$. The pixels in this vertical column then span the interval from $x = 3$ to $x = 4$. Similarly, the mathematical top boundary of the rectangle is at $y = 3$, so we set the top pixel row for the displayed rectangle at $y = 2$.

These compensations for finite pixel width along object boundaries can be applied to other polygons and to curved figures so that the raster display maintains the input object specifications. A circle of radius 5 and center position (10, 10), for instance, would be displayed as in Fig. 3-33 by the midpoint circle algorithm using screen grid coordinate positions. But the plotted circle has a diameter of 11. To plot the circle with the defined diameter of 10, we can modify the circle algorithm to shorten each pixel scan line and each pixel column, as in Fig. 3-34. One way to do this is to generate points clockwise along the circular arc in the third quadrant, starting at screen coordinates (10, 5). For each generated point, the other seven circle symmetry points are generated by decreasing the x coordinate values by 1 along scan lines and decreasing the y coordinate values by 1 along pixel columns. Similar methods are applied in ellipse algorithms to maintain the specified proportions in the display of an ellipse.

3-11

FILLED-AREA PRIMITIVES

A standard output primitive in general graphics packages is a solid-color or patterned polygon area. Other kinds of area primitives are sometimes available, but polygons are easier to process since they have linear boundaries.

There are two basic approaches to area filling on raster systems. One way to fill an area is to determine the overlap intervals for scan lines that cross the area. Another method for area filling is to start from a given interior position and paint outward from this point until we encounter the specified boundary conditions. The scan-line approach is typically used in general graphics packages to fill polygons, circles, ellipses, and other simple curves. Fill methods starting from an interior point are useful with more complex boundaries and in interactive painting systems. In the following sections, we consider methods for solid fill of specified areas. Other fill options are discussed in Chapter 4.

Scan-Line Polygon Fill Algorithm

Figure 3-35 illustrates the scan-line procedure for solid filling of polygon areas. For each scan line crossing a polygon, the area-fill algorithm locates the intersection points of the scan line with the polygon edges. These intersection points are then sorted from left to right, and the corresponding frame-buffer positions between each intersection pair are set to the specified fill color. In the example of Fig. 3-35, the four pixel intersection positions with the polygon boundaries define two stretches of interior pixels from $x = 10$ to $x = 14$ and from $x = 18$ to $x = 24$.

Some scan-line intersections at polygon vertices require special handling. A scan line passing through a vertex intersects two polygon edges at that position, adding two points to the list of intersections for the scan line. Figure 3-36 shows two scan lines at positions y and y' that intersect edge endpoints. Scan line y intersects five polygon edges. Scan line y' , however, intersects an even number of edges although it also passes through a vertex. Intersection points along scan line

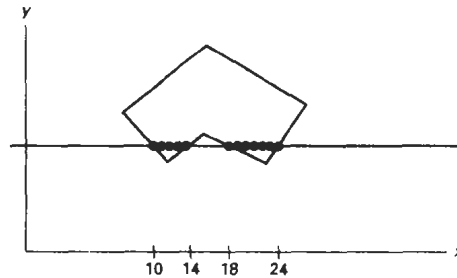


Figure 3-35
Interior pixels along a scan line
passing through a polygon area.

y' correctly identify the interior pixel spans. But with scan line y , we need to do some additional processing to determine the correct interior points.

The topological difference between scan line y and scan line y' in Fig. 3-36 is identified by noting the position of the intersecting edges relative to the scan line. For scan line y , the two intersecting edges sharing a vertex are on opposite sides of the scan line. But for scan line y' , the two intersecting edges are both above the scan line. Thus, the vertices that require additional processing are those that have connecting edges on opposite sides of the scan line. We can identify these vertices by tracing around the polygon boundary either in clockwise or counterclockwise order and observing the relative changes in vertex y coordinates as we move from one edge to the next. If the endpoint y values of two consecutive edges monotonically increase or decrease, we need to count the middle vertex as a single intersection point for any scan line passing through that vertex. Otherwise, the shared vertex represents a local extremum (minimum or maximum) on the polygon boundary, and the two edge intersections with the scan line passing through that vertex can be added to the intersection list.

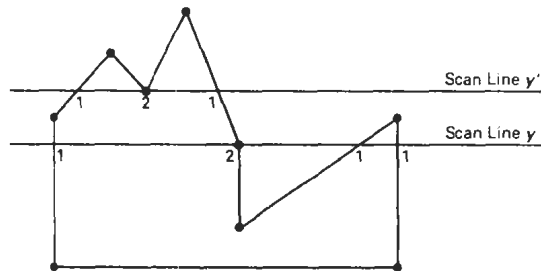


Figure 3-36
Intersection points along scan lines that intersect polygon vertices. Scan line y generates an odd number of intersections, but scan line y' generates an even number of intersections that can be paired to identify correctly the interior pixel spans.

One way to resolve the question as to whether we should count a vertex as one intersection or two is to shorten some polygon edges to split those vertices that should be counted as one intersection. We can process nonhorizontal edges around the polygon boundary in the order specified, either clockwise or counter-clockwise. As we process each edge, we can check to determine whether that edge and the next nonhorizontal edge have either monotonically increasing or decreasing endpoint y values. If so, the lower edge can be shortened to ensure that only one intersection point is generated for the scan line going through the common vertex joining the two edges. Figure 3-37 illustrates shortening of an edge. When the endpoint y coordinates of the two edges are increasing, the y value of the upper endpoint of the current edge is decreased by 1, as in Fig. 3-37(a). When the endpoint y values are monotonically decreasing, as in Fig. 3-37(b), we decrease the y coordinate of the upper endpoint of the edge following the current edge.

Calculations performed in scan-conversion and other graphics algorithms typically take advantage of various **coherence** properties of a scene that is to be displayed. What we mean by coherence is simply that the properties of one part of a scene are related in some way to other parts of the scene so that the relationship can be used to reduce processing. Coherence methods often involve incremental calculations applied along a single scan line or between successive scan lines. In determining edge intersections, we can set up incremental coordinate calculations along any edge by exploiting the fact that the slope of the edge is constant from one scan line to the next. Figure 3-38 shows two successive scan lines crossing a left edge of a polygon. The slope of this polygon boundary line can be expressed in terms of the scan-line intersection coordinates:

$$m = \frac{y_{k+1} - y_k}{x_{k+1} - x_k} \quad (3-58)$$

Since the change in y coordinates between the two scan lines is simply

$$y_{k+1} - y_k = 1 \quad (3-59)$$

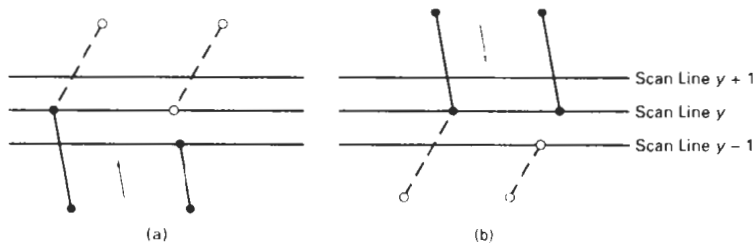


Figure 3-37

Adjusting endpoint y values for a polygon, as we process edges in order around the polygon perimeter. The edge currently being processed is indicated as a solid line. In (a), the y coordinate of the upper endpoint of the current edge is decreased by 1. In (b), the y coordinate of the upper endpoint of the next edge is decreased by 1.

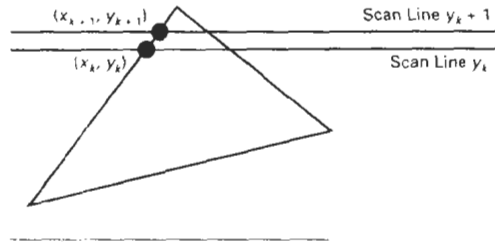


Figure 3-38
Two successive scan lines
intersecting a polygon boundary.

the x -intersection value x_{k+1} on the upper scan line can be determined from the x -intersection value x_k on the preceding scan line as

$$x_{k+1} = x_k + \frac{1}{m} \quad (3-60)$$

Each successive x intercept can thus be calculated by adding the inverse of the slope and rounding to the nearest integer.

An obvious parallel implementation of the fill algorithm is to assign each scan line crossing the polygon area to a separate processor. Edge-intersection calculations are then performed independently. Along an edge with slope m , the intersection x_k value for scan line k above the initial scan line can be calculated as

$$x_k = x_0 + \frac{k}{m} \quad (3-61)$$

In a sequential fill algorithm, the increment of x values by the amount $1/m$ along an edge can be accomplished with integer operations by recalling that the slope m is the ratio of two integers:

$$m = \frac{\Delta y}{\Delta x}$$

where Δx and Δy are the differences between the edge endpoint x and y coordinate values. Thus, incremental calculations of x intercepts along an edge for successive scan lines can be expressed as

$$x_{k+1} = x_k + \frac{\Delta x}{\Delta y} \quad (3-62)$$

Using this equation, we can perform integer evaluation of the x intercepts by initializing a counter to 0, then incrementing the counter by the value of Δx each time we move up to a new scan line. Whenever the counter value becomes equal to or greater than Δy , we increment the current x intersection value by 1 and decrease the counter by the value Δy . This procedure is equivalent to maintaining integer and fractional parts for x intercepts and incrementing the fractional part until we reach the next integer value.

As an example of integer incrementing, suppose we have an edge with slope $m = 7/3$. At the initial scan line, we set the counter to 0 and the counter in-

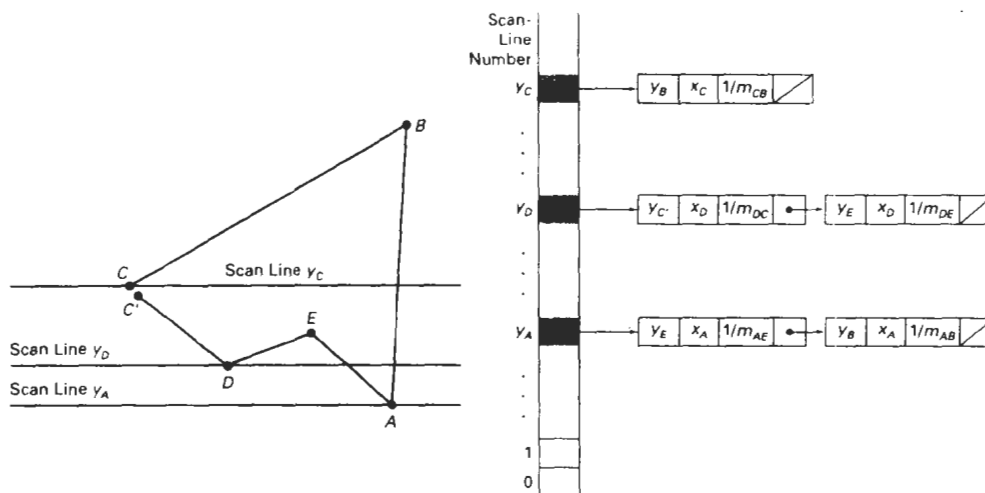


Figure 3-39

A polygon and its sorted edge table, with edge \overline{DC} shortened by one unit in the y direction.

crement to 3. As we move up to the next three scan lines along this edge, the counter is successively assigned the values 3, 6, and 9. On the third scan line above the initial scan line, the counter now has a value greater than 7. So we increment the x -intersection coordinate by 1, and reset the counter to the value $9 - 7 = 2$. We continue determining the scan-line intersections in this way until we reach the upper endpoint of the edge. Similar calculations are carried out to obtain intersections for edges with negative slopes.

We can round to the nearest pixel x -intersection value, instead of truncating to obtain integer positions, by modifying the edge-intersection algorithm so that the increment is compared to $\Delta y/2$. This can be done with integer arithmetic by incrementing the counter with the value $2\Delta x$ at each step and comparing the increment to Δy . When the increment is greater than or equal to Δy , we increase the x value by 1 and decrement the counter by the value of $2\Delta y$. In our previous example with $m = 7/3$, the counter values for the first few scan lines above the initial scan line on this edge would now be 6, 12 (reduced to -2), 4, 10 (reduced to -4), 2, 8 (reduced to -6), 0, 6, and 12 (reduced to -2). Now x would be incremented on scan lines 2, 4, 6, 9, etc., above the initial scan line for this edge. The extra calculations required for each edge are $2\Delta x = \Delta x + \Delta x$ and $2\Delta y = \Delta y + \Delta y$.

To efficiently perform a polygon fill, we can first store the polygon boundary in a *sorted edge table* that contains all the information necessary to process the scan lines efficiently. Proceeding around the edges in either a clockwise or a counterclockwise order, we can use a bucket sort to store the edges, sorted on the smallest y value of each edge, in the correct scan-line positions. Only nonhorizontal edges are entered into the sorted edge table. As the edges are processed, we can also shorten certain edges to resolve the vertex-intersection question. Each entry in the table for a particular scan line contains the maximum y value for that edge, the x -intercept value (at the lower vertex) for the edge, and the inverse slope of the edge. For each scan line, the edges are in sorted order from left to right. Figure 3-39 shows a polygon and the associated sorted edge table.

Next, we process the scan lines from the bottom of the polygon to its top, producing an *active edge list* for each scan line crossing the polygon boundaries. The active edge list for a scan line contains all edges crossed by that scan line, with iterative coherence calculations used to obtain the edge intersections.

Implementation of edge-intersection calculations can also be facilitated by storing Δx and Δy values in the sorted edge table. Also, to ensure that we correctly fill the interior of specified polygons, we can apply the considerations discussed in Section 3-10. For each scan line, we fill in the pixel spans for each pair of x -intercepts starting from the leftmost x -intercept value and ending at one position before the rightmost x intercept. And each polygon edge can be shortened by one unit in the y direction at the top endpoint. These measures also guarantee that pixels in adjacent polygons will not overlap each other.

The following procedure performs a solid-fill scan conversion for an input set of polygon vertices. For each scan line within the vertical extents of the polygon, an active edge list is set up and edge intersections are calculated. Across each scan line, the interior fill is then applied between successive pairs of edge intersections, processed from left to right.

```
#include "device.h"

typedef struct tEdge {
    int yUpper;
    float xIntersect, dxPerScan;
    struct tEdge * next;
} Edge;

/* Inserts edge into list in order of increasing xIntersect field. */
void insertEdge (Edge * list, Edge * edge)
{
    Edge * p, * q = list;

    p = q->next;
    while (p != NULL) {
        if (edge->xIntersect < p->xIntersect)
            p = NULL;
        else {
            q = p;
            p = p->next;
        }
    }
    edge->next = q->next;
    q->next = edge;
}

/* For an index, return y-coordinate of next nonhorizontal line */
int yNext (int k, int cnt, dcPt * pts)
{
    int j;

    if ((k+1) > (cnt-1))
        j = 0;
    else
        j = k + 1;
    while (pts[k].y == pts[j].y)
        if ((j+1) > (cnt-1))
            j = 0;
        else
```



```

        j++;
        return (pts[j].y);
    }

/* Store lower-y coordinate and inverse slope for each edge. Adjust
   and store upper-y coordinate for edges that are the lower member
   of a monotonically increasing or decreasing pair of edges */
void makeEdgeRec
(dcPt lower, dcPt upper, int yComp, Edge * edge, Edge * edges[])
{
    edge->dxPerScan =
        (float) (upper.x - lower.x) / (upper.y - lower.y);
    edge->xIntersect = lower.x;
    if (upper.y < yComp)
        edge->yUpper = upper.y - 1;
    else
        edge->yUpper = upper.y;
    insertEdge (edges[lower.y], edge);
}

void buildEdgeList (int cnt, dcPt * pts, Edge * edges[])
{
    Edge * edge;
    dcPt v1, v2;
    int i, yPrev = pts[cnt - 2].y;

    v1.x = pts[cnt-1].x; v1.y = pts[cnt-1].y;
    for (i=0; i<cnt; i++) {
        v2 = pts[i];
        if (v1.y != v2.y) {
            /* nonhorizontal line */
            edge = (Edge *) malloc (sizeof (Edge));
            if (v1.y < v2.y) /* up-going edge */
                makeEdgeRec (v1, v2, yNext (i, cnt, pts), edge, edges);
            else /* down-going edge */
                makeEdgeRec (v2, v1, yPrev, edge, edges);
        }
        yPrev = v1.y;
        v1 = v2;
    }
}

void buildActiveList (int scan, Edge * active, Edge * edges[])
{
    Edge * p, * q;

    p = edges[scan]->next;
    while (p) {
        q = p->next;
        insertEdge (active, p);
        p = q;
    }
}

void fillScan (int scan, Edge * active)
{
    Edge * p1, * p2;
    int i;

    p1 = active->next;
    while (p1) {
        p2 = p1->next;

```

```

        for (i=p1->xIntersect; i<p2->xIntersect; i++)
            setPixel ((int) i, scan);
        p1 = p2->next;
    }
}

void deleteAfter (Edge * q)
{
    Edge * p = q->next;

    q->next = p->next;
    free (p);
}

/* Delete completed edges. Update 'xIntersect' field for others */
void updateActiveList (int scan, Edge * active)
{
    Edge * q = active, * p = active->next;

    while (p)
        if (scan >= p->yUpper) {
            p = p->next;
            deleteAfter (q);
        }
        else {
            p->xIntersect = p->xIntersect + p->dxPerScan;
            q = p;
            p = p->next;
        }
    }

void resortActiveList (Edge * active)
{
    Edge * q, * p = active->next;

    active->next = NULL;
    while (p) {
        q = p->next;
        insertEdge (active, p);
        p = q;
    }
}

void scanFill (int cnt, dcPt * pts)
{
    Edge * edges[WINDOW_HEIGHT], * active;
    int i, scan;

    for (i=0; i<WINDOW_HEIGHT; i++) {
        edges[i] = (Edge *) malloc (sizeof (Edge));
        edges[i]->next = NULL;
    }
    buildEdgeList (cnt, pts, edges);
    active = (Edge *) malloc (sizeof (Edge));
    active->next = NULL;

    for (scan=0; scan<WINDOW_HEIGHT; scan++) {
        buildActiveList (scan, active, edges);
        if (active->next) {
            fillScan (scan, active);
            updateActiveList (scan, active);
            resortActiveList (active);
        }
    }
}

```

```

    }
    /* Free edge records that have been malloc'ed ... */
}

```

Inside-Outside Tests

Area-filling algorithms and other graphics processes often need to identify interior regions of objects. So far, we have discussed area filling only in terms of standard polygon shapes. In elementary geometry, a polygon is usually defined as having no self-intersections. Examples of standard polygons include triangles, rectangles, octagons, and decagons. The component edges of these objects are joined only at the vertices, and otherwise the edges have no common points in the plane. Identifying the interior regions of standard polygons is generally a straightforward process. But in most graphics applications, we can specify any sequence for the vertices of a fill area, including sequences that produce intersecting edges, as in Fig. 3-40. For such shapes, it is not always clear which regions of the xy plane we should call "interior" and which regions we should designate as "exterior" to the object. Graphics packages normally use either the odd-even rule or the nonzero winding number rule to identify interior regions of an object.

We apply the **odd-even rule**, also called the **odd parity rule** or the **even-odd rule**, by conceptually drawing a line from any position P to a distant point outside the coordinate extents of the object and counting the number of edge crossings along the line. If the number of polygon edges crossed by this line is odd, then P is an *interior* point. Otherwise, P is an *exterior* point. To obtain an accurate edge count, we must be sure that the line path we choose does not intersect any polygon vertices. Figure 3-40(a) shows the interior and exterior regions obtained from the odd-even rule for a self-intersecting set of edges. The scan-line polygon fill algorithm discussed in the previous section is an example of area filling using the odd-even rule.

Another method for defining interior regions is the **nonzero winding number rule**, which counts the number of times the polygon edges wind around a particular point in the counterclockwise direction. This count is called the **winding number**, and the interior points of a two-dimensional object are defined to be

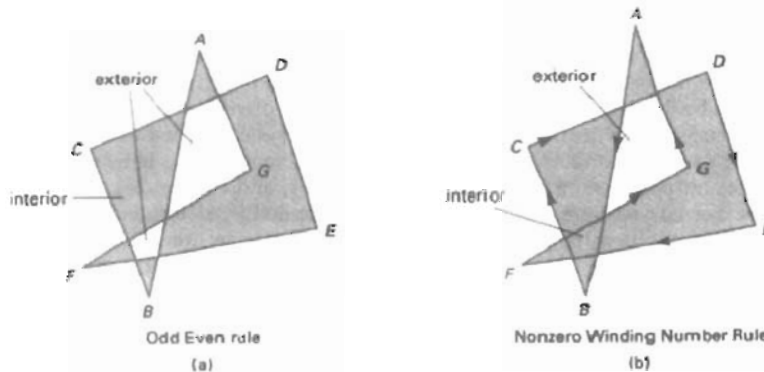


Figure 3-40
Identifying interior and exterior regions for a self-intersecting polygon.

those that have a nonzero value for the winding number. We apply the nonzero winding number rule to polygons by initializing the winding number to 0 and again imagining a line drawn from any position P to a distant point beyond the coordinate extents of the object. The line we choose must not pass through any vertices. As we move along the line from position P to the distant point, we count the number of edges that cross the line in each direction. We add 1 to the winding number every time we intersect a polygon edge that crosses the line from right to left, and we subtract 1 every time we intersect an edge that crosses from left to right. The final value of the winding number, after all edge crossings have been counted, determines the relative position of P . If the winding number is nonzero, P is defined to be an interior point. Otherwise, P is taken to be an exterior point. Figure 3-40(b) shows the interior and exterior regions defined by the nonzero winding number rule for a self-intersecting set of edges. For standard polygons and other simple shapes, the nonzero winding number rule and the odd-even rule give the same results. But for more complicated shapes, the two methods may yield different interior and exterior regions, as in the example of Fig. 3-40.

One way to determine directional edge crossings is to take the vector cross product of a vector u along the line from P to a distant point with the edge vector E for each edge that crosses the line. If the z component of the cross product $u \times E$ for a particular edge is positive, that edge crosses from right to left and we add 1 to the winding number. Otherwise, the edge crosses from left to right and we subtract 1 from the winding number. An edge vector is calculated by subtracting the starting vertex position for that edge from the ending vertex position. For example, the edge vector for the first edge in the example of Fig. 3-40 is

$$E_{AB} = V_B - V_A$$

where V_A and V_B represent the point vectors for vertices A and B. A somewhat simpler way to compute directional edge crossings is to use vector dot products instead of cross products. To do this, we set up a vector that is perpendicular to u and that points from right to left as we look along the line from P in the direction of u . If the components of u are (u_x, u_y) , then this perpendicular to u has components $(-u_y, u_x)$ (Appendix A). Now, if the dot product of the perpendicular and an edge vector is positive, that edge crosses the line from right to left and we add 1 to the winding number. Otherwise, the edge crosses the line from left to right, and we subtract 1 from the winding number.

Some graphics packages use the nonzero winding number rule to implement area filling, since it is more versatile than the odd-even rule. In general, objects can be defined with multiple, unconnected sets of vertices or disjoint sets of closed curves, and the direction specified for each set can be used to define the interior regions of objects. Examples include characters, such as letters of the alphabet and punctuation symbols, nested polygons, and concentric circles or ellipses. For curved lines, the odd-even rule is applied by determining intersections with the curve path, instead of finding edge intersections. Similarly, with the nonzero winding number rule, we need to calculate tangent vectors to the curves at the crossover intersection points with the line from position P .

Scan-Line Fill of Curved Boundary Areas

In general, scan-line fill of regions with curved boundaries requires more work than polygon filling, since intersection calculations now involve nonlinear boundaries. For simple curves such as circles or ellipses, performing a scan-line fill is a straightforward process. We only need to calculate the two scan-line inter-

sections on opposite sides of the curve. This is the same as generating pixel positions along the curve boundary, and we can do that with the midpoint method. Then we simply fill in the horizontal pixel spans between the boundary points on opposite sides of the curve. Symmetries between quadrants (and between octants for circles) are used to reduce the boundary calculations.

Similar methods can be used to generate a fill area for a curve section. An elliptical arc, for example, can be filled as in Fig. 3-41. The interior region is bounded by the ellipse section and a straight-line segment that closes the curve by joining the beginning and ending positions of the arc. Symmetries and incremental calculations are exploited whenever possible to reduce computations.



Figure 3-41
Interior fill of an elliptical arc.

Boundary-Fill Algorithm

Another approach to area filling is to start at a point inside a region and paint the interior outward toward the boundary. If the boundary is specified in a single color, the fill algorithm proceeds outward pixel by pixel until the boundary color is encountered. This method, called the **boundary-fill algorithm**, is particularly useful in interactive painting packages, where interior points are easily selected. Using a graphics tablet or other interactive device, an artist or designer can sketch a figure outline, select a fill color or pattern from a color menu, and pick an interior point. The system then paints the figure interior. To display a solid color region (with no border), the designer can choose the fill color to be the same as the boundary color.

A boundary-fill procedure accepts as input the coordinates of an interior point (x, y) , a fill color, and a boundary color. Starting from (x, y) , the procedure tests neighboring positions to determine whether they are of the boundary color. If not, they are painted with the fill color, and their neighbors are tested. This process continues until all pixels up to the boundary color for the area have been tested. Both inner and outer boundaries can be set up to specify an area, and some examples of defining regions for boundary fill are shown in Fig. 3-42.

Figure 3-43 shows two methods for proceeding to neighboring pixels from the current test position. In Fig. 3-43(a), four neighboring points are tested. These are the pixel positions that are right, left, above, and below the current pixel. Areas filled by this method are called **4-connected**. The second method, shown in Fig. 3-43(b), is used to fill more complex figures. Here the set of neighboring positions to be tested includes the four diagonal pixels. Fill methods using this approach are called **8-connected**. An 8-connected boundary-fill algorithm would correctly fill the interior of the area defined in Fig. 3-44, but a 4-connected boundary-fill algorithm produces the partial fill shown.

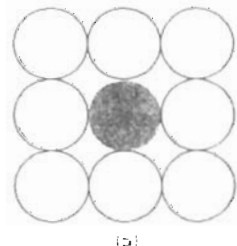
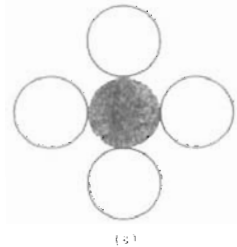


Figure 3-43
Fill methods applied to a 4-connected area (a) and to an 8-connected area (b). Open circles represent pixels to be tested from the current test position, shown as a solid color



Figure 3-42
Example color boundaries for a boundary-fill procedure.

The following procedure illustrates a recursive method for filling a 4-connected area with an intensity specified in parameter `fill` up to a boundary color specified with parameter `boundary`. We can extend this procedure to fill an 8-connected region by including four additional statements to test diagonal positions, such as $(x + 1, y + 1)$.

```
void boundaryFill4 (int x, int y, int fill, int boundary)
{
    int current;

    current = getPixel (x, y);
    if ((current != boundary) && (current != fill)) {
        setColor (fill);
        setPixel (x, y);
        boundaryFill4 (x+1, y, fill, boundary);
        boundaryFill4 (x-1, y, fill, boundary);
        boundaryFill4 (x, y+1, fill, boundary);
        boundaryFill4 (x, y-1, fill, boundary);
    }
}
```

Recursive boundary-fill algorithms may not fill regions correctly if some interior pixels are already displayed in the fill color. This occurs because the algorithm checks next pixels both for boundary color and for fill color. Encountering a pixel with the fill color can cause a recursive branch to terminate, leaving other interior pixels unfilled. To avoid this, we can first change the color of any interior pixels that are initially set to the fill color before applying the boundary-fill procedure.

Also, since this procedure requires considerable stacking of neighboring points, more efficient methods are generally employed. These methods fill horizontal pixel spans across scan lines, instead of proceeding to 4-connected or 8-connected neighboring points. Then we need only stack a beginning position for each horizontal pixel span, instead of stacking all unprocessed neighboring positions around the current position. Starting from the initial interior point with this method, we first fill in the contiguous span of pixels on this starting scan line. Then we locate and stack starting positions for spans on the adjacent scan lines, where spans are defined as the contiguous horizontal string of positions

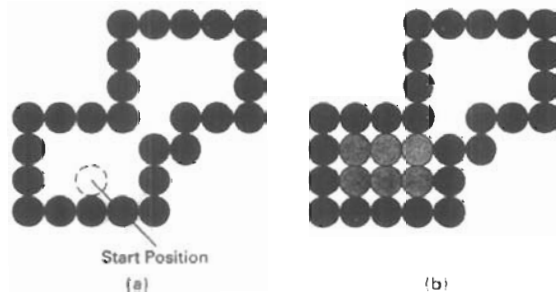


Figure 3-44

The area defined within the color boundary (a) is only partially filled in (b) using a 4-connected boundary-fill algorithm.

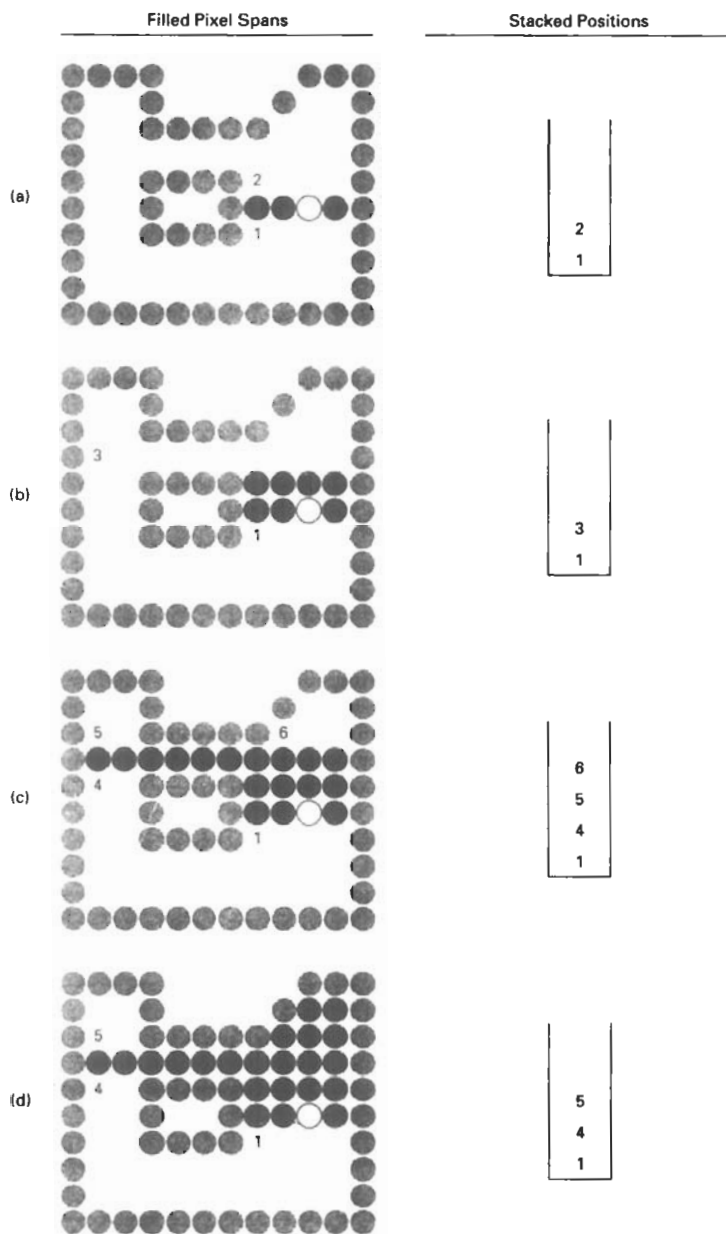


Figure 3-45
Boundary fill across pixel spans for a 4-connected area. (a) The filled initial pixel span, showing the position of the initial point (open circle) and the stacked positions for pixel spans on adjacent scan lines. (b) Filled pixel span on the first scan line above the initial scan line and the current contents of the stack. (c) Filled pixel spans on the first two scan lines above the initial scan line and the current contents of the stack. (d) Completed pixel spans for the upper-right portion of the defined region and the remaining stacked positions to be processed.



Figure 3-46
An area defined within
multiple color boundaries.

bounded by pixels displayed in the area border color. At each subsequent step, we unstack the next start position and repeat the process.

An example of how pixel spans could be filled using this approach is illustrated for the 4-connected fill region in Fig. 3-45. In this example, we first process scan lines successively from the start line to the top boundary. After all upper scan lines are processed, we fill in the pixel spans on the remaining scan lines in order down to the bottom boundary. The leftmost pixel position for each horizontal span is located and stacked, in left to right order across successive scan lines, as shown in Fig. 3-45. In (a) of this figure, the initial span has been filled, and starting positions 1 and 2 for spans on the next scan lines (below and above) are stacked. In Fig. 3-45(b), position 2 has been unstacked and processed to produce the filled span shown, and the starting pixel (position 3) for the single span on the next scan line has been stacked. After position 3 is processed, the filled spans and stacked positions are as shown in Fig. 3-45(c). And Fig. 3-45(d) shows the filled pixels after processing all spans in the upper right of the specified area. Position 5 is next processed, and spans are filled in the upper left of the region; then position 4 is picked up to continue the processing for the lower scan lines.

Flood-Fill Algorithm

Sometimes we want to fill in (or recolor) an area that is not defined within a single color boundary. Figure 3-46 shows an area bordered by several different color regions. We can paint such areas by replacing a specified interior color instead of searching for a boundary color value. This approach is called a **flood-fill algorithm**. We start from a specified interior point (x, y) and reassign all pixel values that are currently set to a given interior color with the desired fill color. If the area we want to paint has more than one interior color, we can first reassign pixel values so that all interior points have the same color. Using either a 4-connected or 8-connected approach, we then step through pixel positions until all interior points have been repainted. The following procedure flood fills a 4-connected region recursively, starting from the input position.

```
void floodFill4 (int x, int y, int fillColor, int oldColor)
{
    if (getPixel (x, y) == oldColor) {
        setColor (fillColor);
        setPixel (x, y);
        floodFill4 (x+1, y, fillColor, oldColor);
        floodFill4 (x-1, y, fillColor, oldColor);
        floodFill4 (x, y+1, fillColor, oldColor);
        floodFill4 (x, y-1, fillColor, oldColor);
    }
}
```

We can modify procedure `floodFill4` to reduce the storage requirements of the stack by filling horizontal pixel spans, as discussed for the boundary-fill algorithm. In this approach, we stack only the beginning positions for those pixel spans having the value `oldColor`. The steps in this modified flood-fill algorithm are similar to those illustrated in Fig. 3-45 for a boundary fill. Starting at the first position of each span, the pixel values are replaced until a value other than `oldColor` is encountered.

We display a filled polygon in PHIGS and GKS with the function

```
fillArea (n, wcVertices)
```

The displayed polygon area is bounded by a series of n straight line segments connecting the set of vertex positions specified in `wcVertices`. These packages do not provide fill functions for objects with curved boundaries.

Implementation of the `fillArea` function depends on the selected type of interior fill. We can display the polygon boundary surrounding a hollow interior, or we can choose a solid color or pattern fill with no border for the display of the polygon. For solid fill, the `fillArea` function is implemented with the scan-line fill algorithm to display a single color area. The various attribute options for displaying polygon fill areas in PHIGS are discussed in the next chapter.

Another polygon primitive available in PHIGS is `fillAreaSet`. This function allows a series of polygons to be displayed by specifying the list of vertices for each polygon. Also, in other graphics packages, functions are often provided for displaying a variety of commonly used fill areas besides general polygons. Some examples are `fillRectangle`, `fillCircle`, `fillCircleArc`, `fillEllipse`, and `fillEllipseArc`.

3-13

CELL ARRAY

The **cell array** is a primitive that allows users to display an arbitrary shape defined as a two-dimensional grid pattern. A predefined matrix of color values is mapped by this function onto a specified rectangular coordinate region. The PHIGS version of this function is

```
cellArray (wcPoints, n, m, colorArray)
```

where `colorArray` is the n by m matrix of integer color values and `wcPoints` lists the limits of the rectangular coordinate region: (x_{\min}, y_{\min}) and (x_{\max}, y_{\max}) . Figure 3-47 shows the distribution of the elements of the color matrix over the coordinate rectangle.

Each coordinate cell in Fig. 3-47 has width $(x_{\max} - x_{\min})/n$ and height $(y_{\max} - y_{\min})/m$. Pixel color values are assigned according to the relative positions of the pixel center coordinates. If the center of a pixel lies within one of the n by m coordinate cells, that pixel is assigned the color of the corresponding element in the matrix `colorArray`.

3-14

CHARACTER GENERATION

Letters, numbers, and other characters can be displayed in a variety of sizes and styles. The overall design style for a set (or family) of characters is called a **type-**

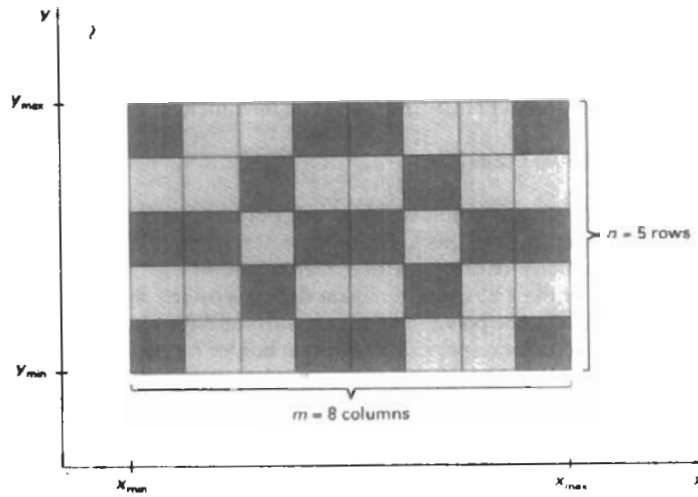


Figure 3-47

Mapping an n by m cell array into a rectangular coordinate region.

face. Today, there are hundreds of typefaces available for computer applications. Examples of a few common typefaces are Courier, Helvetica, New York, Palatino, and Zapf Chancery. Originally, the term *font* referred to a set of cast metal character forms in a particular size and format, such as 10-point Courier Italic or 12-point Palatino Bold. Now, the terms *font* and *typeface* are often used interchangeably, since printing is no longer done with cast metal forms.

Typefaces (or fonts) can be divided into two broad groups: *serif* and *sans serif*. Serif type has small lines or accents at the ends of the main character strokes, while sans-serif type does not have accents. For example, the text in this book is set in a serif font (Palatino). But this sentence is printed in a sans-serif font (Optima). Serif type is generally more *readable*; that is, it is easier to read in longer blocks of text. On the other hand, the individual characters in sans-serif type are easier to *recognize*. For this reason, sans-serif type is said to be more *legible*. Since sans-serif characters can be quickly recognized, this typeface is good for labeling and short headings.

Two different representations are used for storing computer fonts. A simple method for representing the character shapes in a particular typeface is to use rectangular grid patterns. The set of characters are then referred to as a **bitmap font** (or **bitmapped font**). Another, more flexible, scheme is to describe character shapes using straight-line and curve sections, as in PostScript, for example. In this case, the set of characters is called an **outline font**. Figure 3-48 illustrates the two methods for character representation. When the pattern in Fig. 3-48(a) is copied to an area of the frame buffer, the 1 bits designate which pixel positions are to be displayed on the monitor. To display the character shape in Fig. 3-48(b), the interior of the character outline must be filled using the scan-line fill procedure (Section 3-11).

Bitmap fonts are the simplest to define and display: The character grid only needs to be mapped to a frame-buffer position. In general, however, bitmap fonts

require more space, because each variation (size and format) must be stored in a *font cache*. It is possible to generate different sizes and other variations, such as bold and italic, from one set, but this usually does not produce good results.

In contrast to bitmap fonts, outline fonts require less storage since each variation does not require a distinct font cache. We can produce boldface, italic, or different sizes by manipulating the curve definitions for the character outlines. But it does take more time to process the outline fonts, because they must be scan converted into the frame buffer.

A character string is displayed in PHIGS with the following function:

```
text (wcPoint, string)
```

Parameter *string* is assigned a character sequence, which is then displayed at coordinate position *wcPoint* = (*x*, *y*). For example, the statement

```
text (wcPoint, 'Population Distribution')
```

along with the coordinate specification for *wcPoint*, could be used as a label on a distribution graph.

Just how the string is positioned relative to coordinates (*x*, *y*) is a user option. The default is that (*x*, *y*) sets the coordinate location for the lower left corner of the first character of the horizontal string to be displayed. Other string orientations, such as vertical, horizontal, or slanting, are set as attribute options and will be discussed in the next chapter.

Another convenient character function in PHIGS is one that places a designated character, called a **marker symbol**, at one or more selected positions. This function is defined with the same parameter list as in the line function:

```
polymarker (n, wcPoints)
```

A predefined character is then centered at each of the *n* coordinate positions in the list *wcPoints*. The default symbol displayed by *polymarker* depends on the

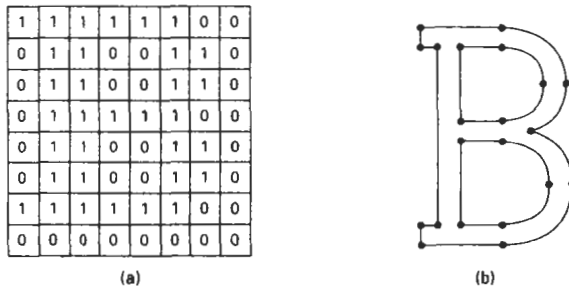


Figure 3-48

The letter B represented in (a) with an 8 by 8 bilevel bitmap pattern and in (b) with an outline shape defined with straight-line and curve segments.

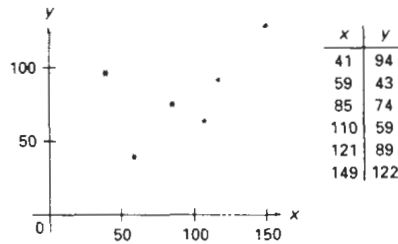


Figure 3-49
Sequence of data values plotted
with the polymarker function.

particular implementation, but we assume for now that an asterisk is to be used. Figure 3-49 illustrates plotting of a data set with the statement

```
polymarker (6, wcPoints)
```

SUMMARY

The output primitives discussed in this chapter provide the basic tools for constructing pictures with straight lines, curves, filled areas, cell-array patterns, and text. Examples of pictures generated with these primitives are given in Figs. 3-50 and 3-51.

Three methods that can be used to plot pixel positions along a straight-line path are the DDA algorithm, Bresenham's algorithm, and the midpoint method. For straight lines, Bresenham's algorithm and the midpoint method are identical and are the most efficient. Frame-buffer access in these methods can also be performed efficiently by incrementally calculating memory addresses. Any of the line-generating algorithms can be adapted to a parallel implementation by partitioning line segments.

Circles and ellipses can be efficiently and accurately scan converted using midpoint methods and taking curve symmetry into account. Other conic sections, parabolas and hyperbolas, can be plotted with similar methods. Spline curves, which are piecewise continuous polynomials, are widely used in design applications. Parallel implementation of curve generation can be accomplished by partitioning the curve paths.

To account for the fact that displayed lines and curves have finite widths, we must adjust the pixel dimensions of objects to coincide to the specified geometric dimensions. This can be done with an addressing scheme that references pixel positions at their lower left corner, or by adjusting line lengths.

Filled area primitives in many graphics packages refer to filled polygons. A common method for providing polygon fill on raster systems is the scan-line fill algorithm, which determines interior pixel spans across scan lines that intersect the polygon. The scan-line algorithm can also be used to fill the interior of objects with curved boundaries. Two other methods for filling the interior regions of objects are the boundary-fill algorithm and the flood-fill algorithm. These two fill procedures paint the interior, one pixel at a time, outward from a specified interior point.

The scan-line fill algorithm is an example of filling object interiors using the odd-even rule to locate the interior regions. Other methods for defining object interiors are also useful, particularly with unusual, self-intersecting objects. A common example is the nonzero winding number rule. This rule is more flexible than the odd-even rule for handling objects defined with multiple boundaries.

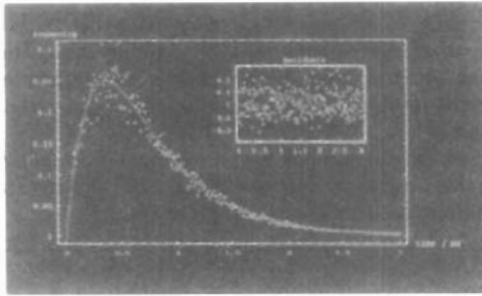


Figure 3-50
A data plot generated with straight line segments, a curve, circles (or markers), and text. (Courtesy of Wolfram Research, Inc., The Maker of Mathematica.)

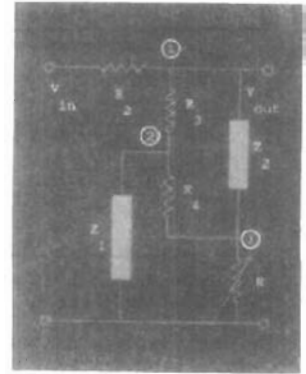


Figure 3-51
An electrical diagram drawn with straight line sections, circles, filled rectangles, and text. (Courtesy of Wolfram Research, Inc., The Maker of Mathematica.)

Additional primitives available in graphics packages include cell arrays, character strings, and marker symbols. Cell arrays are used to define and store color patterns. Character strings are used to provide picture and graph labeling. And marker symbols are useful for plotting the position of data points.

Table 3-1 lists implementations for some of the output primitives discussed in this chapter.

TABLE 3-1
OUTPUT PRIMITIVE IMPLEMENTATIONS

<code>typedef struct { float x, y; } wcPt2;</code>	Defines a location in 2-dimensional world-coordinates.
<code>pPolyline (int n, wcPt2 * pts)</code>	Draw a connected sequence of n-1 line segments, specified in pts.
<code>pCircle (wcPt2 center, float r)</code>	Draw a circle of radius r at center.
<code>pFillarea (int n, wcPt2 * pts)</code>	Draw a filled polygon with n vertices, specified in pts.
<code>pCellArray (wcPt2 * pts, int n, int m, int colors)</code>	Map an n by m array of colors onto a rectangular area defined by pts.
<code>pText (wcPt2 position, char * txt)</code>	Draw the character string txt at position.
<code>pPolymarker (int n, wcPt2 * pts)</code>	Draw a collection of n marker symbols at pts.

Here, we present a few example programs illustrating applications of output primitives. Functions listed in Table 3-1 are defined in the header file `graphics.h`, along with the routines `openGraphics`, `closeGraphics`, `setColor`, and `setBackground`.

The first program produces a line graph for monthly data over a period of one year. Output of this procedure is drawn in Fig. 3-52. This data set is also used by the second program to produce the bar graph in Fig. 3-53.

```
#include <stdio.h>
#include "graphics.h"

#define WINDOW_WIDTH 600
#define WINDOW_HEIGHT 500
/* Amount of space to leave on each side of the chart */
#define MARGIN_WIDTH 0.05 * WINDOW_WIDTH
#define N_DATA 12

typedef enum
{ Jan, Feb, Mar, Apr, May, Jun, Jul, Aug, Sep, Oct, Nov, Dec } Months;

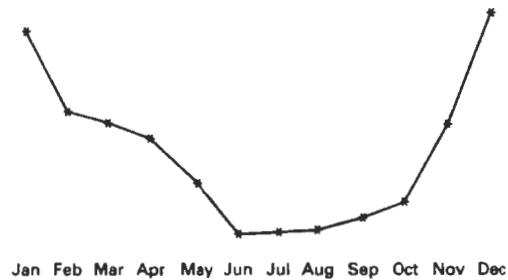
char * monthNames[N_DATA] = { "Jan", "Feb", "Mar", "Apr", "May", "Jun",
                               "Jul", "Aug", "Sep", "Oct", "Nov", "Dec" };

int readData (char * inFile, float * data)
{
    int fileError = FALSE;
    FILE * fp;
    Months month;

    if ((fp = fopen (inFile, "r")) == NULL)
        fileError = TRUE;
    else {
        for (month = Jan; month <= Dec; month++)
            fscanf (fp, "%f", &data[month]);
        fclose (fp);
    }
    return (fileError);
}

void lineChart (float * data)
{
    wcPt2 dataPos[N_DATA], labelPos;
    Months m;
    float mWidth = (WINDOW_WIDTH - 2 * MARGIN_WIDTH) / N_DATA;
    int chartBottom = 0.1 * WINDOW_HEIGHT;
    int offset = 0.05 * WINDOW_HEIGHT; /* Space between data and labels */
    int labelLength = 24; /* Assuming fixed-width 8-pixel characters */

    labelPos.y = chartBottom;
    for (m = Jan; m <= Dec; m++) {
        /* Calculate x and y positions for data markers */
        dataPos[m].x = MARGIN_WIDTH + m * mWidth + 0.5 * mWidth;
        dataPos[m].y = chartBottom + offset + data[m];
        /* Shift the label to the left by one-half its length */
        labelPos.x = dataPos[m].x - 0.5 * labelLength;
        pText (labelPos, monthNames[m]);
    }
    pPolyline (N_DATA, dataPos);
    pPolymarker (N_DATA, dataPos);
}
```



Summary

Figure 3-52
A line plot of data points output by
the lineChart procedure.

```

}

void main (int argc, char ** argv)
{
    float data[N_DATA];
    int dataError = FALSE;
    long windowID;

    if (argc < 2) {
        fprintf (stderr, "Usage: %s dataFileName\n", argv[0]);
        exit ();
    }
    dataError = readData (argv[1], data);
    if (dataError) {
        fprintf (stderr, "%s error. Can't read file %s\n", argv[1]);
        exit ();
    }
    windowID = openGraphics (*argv, WINDOW_WIDTH, WINDOW_HEIGHT);
    setBackground (WHITE);
    setColor (BLACK);
    lineChart (data);
    sleep (10);
    closeGraphics (windowID);
}

```

```

void barChart (float * data)
{
    wcPt2 dataPos[4], labelPos;
    Months m;
    float x, mWidth = (WINDOW_WIDTH - 2 * MARGIN_WIDTH) / N_DATA;
    int chartBottom = 0.1 * WINDOW_HEIGHT;
    int offset = 0.05 * WINDOW_HEIGHT; /* Space between data and labels */
    int labelLength = 24; /* Assuming fixed-width 8-pixel characters */

    labelPos.y = chartBottom;
    for (m = Jan; m <= Dec; m++) {
        /* Find the center of this month's bar */
        x = MARGIN_WIDTH + m * mWidth + 0.5 * mWidth;

        /* Shift the label to the left by one-half its assumed length */
        labelPos.x = x - 0.5 * labelLength;
    }
}

```

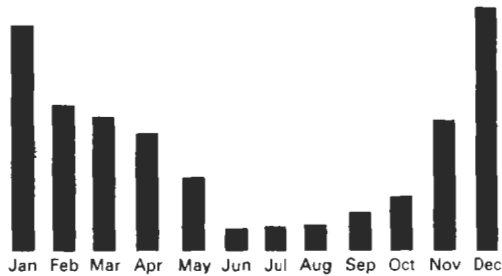


Figure 3-53
A bar-chart plot output by the
barChart procedure.

```
pText (labelPos, monthNames[m]);

/* Get the coordinates for this month's bar */
dataPos[0].x = dataPos[3].x = x - 0.5 * labelLength;
dataPos[1].x = dataPos[2].x = x + 0.5 * labelLength;
dataPos[0].y = dataPos[1].y = chartBottom + offset;
dataPos[2].y = dataPos[3].y = chartBottom + offset + data[m];
pFillArea (4, dataPos);
}
```

Pie charts are used to show the percentage contribution of individual parts to the whole. The next procedure constructs a pie chart, with the number and relative size of the slices determined by input. A sample output from this procedure appears in Fig. 3-54.

```
#define TWO_PI 6.28

void pieChart (float * data)
{
    wcPt2 pts[2], center;
    float radius = WINDOW_HEIGHT / 4.0;
    float newSlice, total = 0.0, lastSlice = 0.0;
    Months month;

    center.x = WINDOW_WIDTH / 2;
    center.y = WINDOW_HEIGHT / 2;
    pCircle (center, radius);
    for (month = Jan; month <= Dec; month++)
        total += data[month];
    pts[0].x = center.x; pts[0].y = center.y;
    for (month = Jan; month <= Dec; month++) {
        newSlice = TWO_PI * data[month] / total + lastSlice;
        pts[1].x = center.x + radius * cosf (newSlice);
        pts[1].y = center.y + radius * sinf (newSlice);
        pPolyline (2, pts);
        lastSlice = newSlice;
    }
}
```


Some variations on the circle equations are output by this next procedure. The shapes shown in Fig. 3-55 are generated by varying the radius r of a circle. Depending on how we vary r , we can produce a spiral, cardioid, limaçon, or other similar figure.

```
#include <stdio.h>
#include <math.h>
#include "graphics.h"

#define TWO_PI 6.28

/* Limaçon equation is  $r = a * \cos(\theta) + b$ . Cardioid is the same,
   with  $a = b$ , so  $r = a * (1 + \cos(\theta))$ .
*/
typedef enum { spiral, cardioid, threeLeaf, fourLeaf, limaçon } Fig;

void drawCurlyFig (Fig figure, wcPt2 pos, int * p)
{
    float r, theta = 0.0, dtheta = 1.0 / (float) p[0];
    int nPoints = (int) ceilf (TWO_PI * p[0]) + 1;
    wcPt2 * pt;

    if ((pt = (wcPt2 *) malloc (nPoints * sizeof (wcPt2))) == NULL) {
        fprintf (stderr, "Couldn't allocate points\n");
        return;
    }

    /* Set first point for figure */
    pt[0].y = pos.y;
    switch (figure) {
        case spiral:    pt[0].x = pos.x;                break;
        case limaçon:   pt[0].x = pos.x + p[0] + p[1]; break;
        case cardioid:  pt[0].x = pos.x + p[0] * 2;     break;
        case threeLeaf: pt[0].x = pos.x + p[0];         break;
        case fourLeaf:  pt[0].x = pos.x + p[0];         break;
    }
    nPoints = 1;
    while (theta < TWO_PI) {
        switch (figure) {
            case spiral:    r = p[0] * theta;                break;
            case limaçon:   r = p[0] * cosf (theta) + p[1]; break;
            case cardioid:  r = p[0] * (1 + cosf (theta));    break;
            case threeLeaf: r = p[0] * cosf (3 * theta);      break;
            case fourLeaf:  r = p[0] * cosf (2 * theta);      break;
        }
        pt[nPoints].x = pos.x + r * cosf (theta);
        pt[nPoints].y = pos.y + r * sinf (theta);
        nPoints++;
        theta += dtheta;
    }

    pPolyline (nPoints, pt);
    free (pt);
}

void main (int argc, char ** argv)
{

```

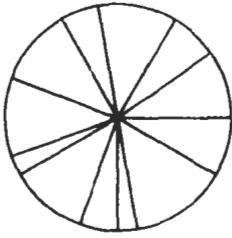


Figure 3-54
Output generated from the
pieChart procedure.



Figure 3-55
Curved figures produced with the drawShape procedure.

```
long windowID = openGraphics (*argv, 400, 100);
Fig f;
/* Center positions for each figure */
wcPt2 center[] = { 50, 50, 100, 50, 175, 50, 250, 50, 300, 50 };

/* Parameters to define each figure. First four need one parameter.
   Fifth figure (limaçon) needs two. */
int p[5][2] = { 5, -1, 20, -1, 30, -1, 30, -1, 40, 10 };

setBackground (WHITE);
setColor (BLACK);
for (f=spiral; f<=limaçon; f++)
    drawCurlyFig (f, center[f], p[f]);
sleep (10);
closeGraphics (windowID);
}
```

REFERENCES

Information on Bresenham's algorithms can be found in Bresenham (1965, 1977). For mid-point methods, see Kappel (1985). Parallel methods for generating lines and circles are discussed in Pang (1990) and in Wright (1990).

Additional programming examples and information on PHIGS primitives can be found in Howard, et al. (1991), Hopgood and Duce (1991), Gaskins (1992), and Blake (1993). For information on GKS output primitive functions, see Hopgood et al. (1983) and Enderle, Kansy, and Pfaff (1984).

EXERCISES

- 3-1. Implement the polyline function using the DDA algorithm, given any number (n) of input points. A single point is to be plotted when $n = 1$.
- 3-2. Extend Bresenham's line algorithm to generate lines with any slope, taking symmetry between quadrants into account. Implement the polyline function using this algorithm as a routine that displays the set of straight lines connecting the n input points. For $n = 1$, the routine displays a single point.

- 3-3. Devise a consistent scheme for implementing the `polyline` function, for any set of input line endpoints, using a modified Bresenham line algorithm so that geometric magnitudes are maintained (Section 3-10).
- 3-4. Use the midpoint method to derive decision parameters for generating points along a straight-line path with slope in the range $0 < m < 1$. Show that the midpoint decision parameters are the same as those in the Bresenham line algorithm.
- 3-5. Use the midpoint method to derive decision parameters that can be used to generate straight line segments with any slope.
- 3-6. Set up a parallel version of Bresenham's line algorithm for slopes in the range $0 < m < 1$.
- 3-7. Set up a parallel version of Bresenham's algorithm for straight lines of any slope.
- 3-8. Suppose you have a system with an 8-inch by 10-inch video monitor that can display 100 pixels per inch. If memory is organized in one-byte words, the starting frame-buffer address is 0, and each pixel is assigned one byte of storage, what is the frame-buffer address of the pixel with screen coordinates (x, y) ?
- 3-9. Suppose you have a system with an 8-inch by 10-inch video monitor that can display 100 pixels per inch. If memory is organized in one-byte words, the starting frame-buffer address is 0, and each pixel is assigned 6 bits of storage, what is the frame-buffer address (or addresses) of the pixel with screen coordinates (x, y) ?
- 3-10. Implement the `setPixel` routine in Bresenham's line algorithm using iterative techniques for calculating frame-buffer addresses (Section 3-3).
- 3-11. Revise the midpoint circle algorithm to display so that geometric magnitudes are maintained (Section 3-10).
- 3-12. Set up a procedure for a parallel implementation of the midpoint circle algorithm.
- 3-13. Derive decision parameters for the midpoint ellipse algorithm assuming the start position is $(r_x, 0)$ and points are to be generated along the curve path in counterclockwise order.
- 3-14. Set up a procedure for a parallel implementation of the midpoint ellipse algorithm.
- 3-15. Devise an efficient algorithm that takes advantage of symmetry properties to display a sine function.
- 3-16. Devise an efficient algorithm, taking function symmetry into account, to display a plot of damped harmonic motion:

$$y = Ae^{-kx} \sin(\omega x + \theta)$$

where ω is the angular frequency and θ is the phase of the sine function. Plot y as a function of x for several cycles of the sine function or until the maximum amplitude is reduced to $A/10$.

- 3-17. Using the midpoint method, and taking symmetry into account, develop an efficient algorithm for scan conversion of the following curve over the interval $-10 \leq x \leq 10$:

$$y = \frac{1}{12} x^3$$

- 3-18. Use the midpoint method and symmetry considerations to scan convert the parabola

$$y = 100 - x^2$$

over the interval $-10 \leq x \leq 10$.

- 3-19. Use the midpoint method and symmetry considerations to scan convert the parabola

$$x = y^2$$

for the interval $-10 \leq y \leq 10$.

- 3-20. Set up a midpoint algorithm, taking symmetry considerations into account to scan convert any parabola of the form

$$y = ax^2 - b$$

with input values for parameters a , b , and the range of x .

- 3-21. Write a program to scan convert the interior of a specified ellipse into a solid color.
- 3-22. Devise an algorithm for determining interior regions for any input set of vertices using the nonzero winding number rule and cross-product calculations to identify the direction of edge crossings.
- 3-23. Devise an algorithm for determining interior regions for any input set of vertices using the nonzero winding number rule and dot-product calculations to identify the direction of edge crossings.
- 3-24. Write a procedure for filling the interior of any specified set of "polygon" vertices using the nonzero winding number rule to identify interior regions.
- 3-25. Modify the boundary-fill algorithm for a 4-connected region to avoid excessive stacking by incorporating scan-line methods.
- 3-26. Write a boundary-fill procedure to fill an 8-connected region.
- 3-27. Explain how an ellipse displayed with the midpoint method could be properly filled with a boundary-fill algorithm.
- 3-28. Develop and implement a flood-fill algorithm to fill the interior of any specified area.
- 3-29. Write a routine to implement the `text` function.
- 3-30. Write a routine to implement the `polymarker` function.
- 3-31. Write a program to display a bar graph using the `polyline` function. Input to the program is to include the data points and the labeling required for the x and y axes. The data points are to be scaled by the program so that the graph is displayed across the full screen area.
- 3-32. Write a program to display a bar graph in any selected screen area. Use the `polyline` function to draw the bars.
- 3-33. Write a procedure to display a line graph for any input set of data points in any selected area of the screen, with the input data set scaled to fit the selected screen area. Data points are to be displayed as asterisks joined with straight line segments, and the x and y axes are to be labeled according to input specifications. (Instead of asterisks, small circles or some other symbols could be used to plot the data points.)
- 3-34. Using a `circle` function, write a routine to display a pie chart with appropriate labeling. Input to the routine is to include a data set giving the distribution of the data over some set of intervals, the name of the pie chart, and the names of the intervals. Each section label is to be displayed outside the boundary of the pie chart near the corresponding pie section.

4

Attributes of Output Primitives



In general, any parameter that affects the way a primitive is to be displayed is referred to as an **attribute parameter**. Some attribute parameters, such as color and size, determine the fundamental characteristics of a primitive. Others specify how the primitive is to be displayed under special conditions. Examples of attributes in this class include depth information for three-dimensional viewing and visibility or detectability options for interactive object-selection programs. These special-condition attributes will be considered in later chapters. Here, we consider only those attributes that control the basic display properties of primitives, without regard for special situations. For example, lines can be dotted or dashed, fat or thin, and blue or orange. Areas might be filled with one color or with a multicolor pattern. Text can appear reading from left to right, slanted diagonally across the screen, or in vertical columns. Individual characters can be displayed in different fonts, colors, and sizes. And we can apply intensity variations at the edges of objects to smooth out the raster staircase effect.

One way to incorporate attribute options into a graphics package is to extend the parameter list associated with each output primitive function to include the appropriate attributes. A line-drawing function, for example, could contain parameters to set color, width, and other properties, in addition to endpoint coordinates. Another approach is to maintain a system list of current attribute values. Separate functions are then included in the graphics package for setting the current values in the attribute list. To generate an output primitive, the system checks the relevant attributes and invokes the display routine for that primitive using the current attribute settings. Some packages provide users with a combination of attribute functions and attribute parameters in the output primitive commands. With the GKS and PHIGS standards, attribute settings are accomplished with separate functions that update a system attribute list.

4-1

LINE ATTRIBUTES

Basic attributes of a straight line segment are its type, its width, and its color. In some graphics packages, lines can also be displayed using selected pen or brush options. In the following sections, we consider how line-drawing routines can be modified to accommodate various attribute specifications.

Line Type

Possible selections for the line-type attribute include solid lines, dashed lines, and dotted lines. We modify a line-drawing algorithm to generate such lines by setting the length and spacing of displayed solid sections along the line path. A dashed line could be displayed by generating an interdash spacing that is equal to the length of the solid sections. Both the length of the dashes and the interdash spacing are often specified as user options. A dotted line can be displayed by

generating very short dashes with the spacing equal to or greater than the dash size. Similar methods are used to produce other line-type variations.

To set line type attributes in a PHIGS application program, a user invokes the function

```
setLinetype (lt)
```

where parameter *lt* is assigned a positive integer value of 1, 2, 3, or 4 to generate lines that are, respectively, solid, dashed, dotted, or dash-dotted. Other values for the line-type parameter *lt* could be used to display variations in the dot-dash patterns. Once the line-type parameter has been set in a PHIGS application program, all subsequent line-drawing commands produce lines with this line type. The following program segment illustrates use of the *linetype* command to display the data plots in Fig. 4-1.

```
#include <stdio.h>
#include "graphics.h"

#define MARGIN_WIDTH 0.05 * WINDOW_WIDTH

int readData (char * inFile, float * data)
{
    int fileError = FALSE;
    FILE * fp;
    int month;

    if ((fp = fopen (inFile, "r")) == NULL)
        fileError = TRUE;
    else {
        for (month=0; month<12; month++)
            fscanf (fp, "%f", &data[month]);
        fclose (fp);
    }
    return (fileError);
}

void chartData (float * data, pLineType lineType)
{
    wcPt2 pts[12];
    float monthWidth = (WINDOW_WIDTH - 2 * MARGIN_WIDTH) / 12;
    int i;

    for (i=0; i<12; i++) {
        pts[i].x = MARGIN_WIDTH + i * monthWidth + 0.5 * monthWidth;
        pts[i].y = data[i];
    }
    pSetLineType (lineType);
    pPolyline (12, pts);
}

int main (int argc, char ** argv)
{
    long windowID = openGraphics (*argv, WINDOW_WIDTH, WINDOW_HEIGHT);
    float data[12];

    setBackground (WHITE);
    setColor (BLUE);
    readData ("../data/data1960", data);
    chartData (data, SOLID);
    readData ("../data/data1970", data);
    chartData (data, DASHED);
    readData ("../data/data1980", data);
    chartData (data, DOTTED);
    sleep (10);
    closeGraphics (windowID);
}
```

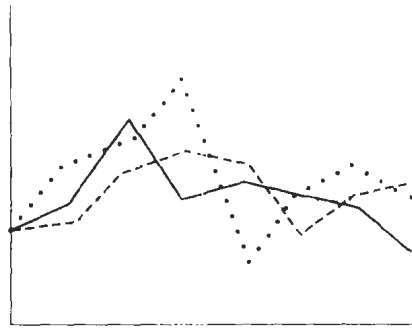


Figure 4-1
Plotting three data sets with three different line types, as output by the `chartData` procedure.

Raster line algorithms display line-type attributes by plotting pixel spans. For the various dashed, dotted, and dot-dashed patterns, the line-drawing procedure outputs sections of contiguous pixels along the line path, skipping over a number of intervening pixels between the solid spans. Pixel counts for the span length and interspan spacing can be specified in a pixel `mask`, which is a string containing the digits 1 and 0 to indicate which positions to plot along the line path. The mask 1111000, for instance, could be used to display a dashed line with a dash length of four pixels and an interdash spacing of three pixels. On a bilevel system, the mask gives the bit values that should be loaded into the frame buffer along the line path to display the selected line type.

Plotting dashes with a fixed number of pixels results in unequal-length dashes for different line orientations, as illustrated in Fig. 4-2. Both dashes shown are plotted with four pixels, but the diagonal dash is longer by a factor of $\sqrt{2}$. For precision drawings, dash lengths should remain approximately constant for any line orientation. To accomplish this, we can adjust the pixel counts for the solid spans and interspan spacing according to the line slope. In Fig. 4-2, we can display approximately equal-length dashes by reducing the diagonal dash to three pixels. Another method for maintaining dash length is to treat dashes as individual line segments. Endpoint coordinates for each dash are located and passed to the line routine, which then calculates pixel positions along the dash path.

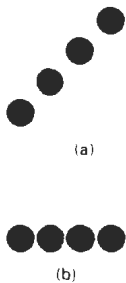


Figure 4-2
Unequal-length dashes displayed with the same number of pixels.

Line Width

Implementation of line-width options depends on the capabilities of the output device. A heavy line on a video monitor could be displayed as adjacent parallel lines, while a pen plotter might require pen changes. As with other PHIGS attributes, a line-width command is used to set the current line-width value in the attribute list. This value is then used by line-drawing algorithms to control the thickness of lines generated with subsequent output primitive commands.

We set the line-width attribute with the command:

```
setLinewidthScaleFactor (lw)
```

Line-width parameter `lw` is assigned a positive number to indicate the relative width of the line to be displayed. A value of 1 specifies a standard-width line. On a pen plotter, for instance, a user could set `lw` to a value of 0.5 to plot a line whose width is half that of the standard line. Values greater than 1 produce lines thicker than the standard.

For raster implementation, a standard-width line is generated with single pixels at each sample position, as in the Bresenham algorithm. Other-width lines are displayed as positive integer multiples of the standard line by plotting additional pixels along adjacent parallel line paths. For lines with slope magnitude less than 1, we can modify a line-drawing routine to display thick lines by plotting a vertical span of pixels at each x position along the line. The number of pixels in each span is set equal to the integer magnitude of parameter $1w$. In Fig. 4-3, we plot a double-width line by generating a parallel line above the original line path. At each x sampling position, we calculate the corresponding y coordinate and plot pixels with screen coordinates (x, y) and $(x, y+1)$. We display lines with $1w \geq 3$ by alternately plotting pixels above and below the single-width line path.

For lines with slope magnitude greater than 1, we can plot thick lines with horizontal spans, alternately picking up pixels to the right and left of the line path. This scheme is demonstrated in Fig. 4-4, where a line width of 4 is plotted with horizontal pixel spans.

Although thick lines are generated quickly by plotting horizontal or vertical pixel spans, the displayed width of a line (measured perpendicular to the line path) is dependent on its slope. A 45° line will be displayed thinner by a factor of $1/\sqrt{2}$ compared to a horizontal or vertical line plotted with the same-length pixel spans.

Another problem with implementing width options using horizontal or vertical pixel spans is that the method produces lines whose ends are horizontal or vertical regardless of the slope of the line. This effect is more noticeable with very thick lines. We can adjust the shape of the line ends to give them a better appearance by adding **line caps** (Fig. 4-5). One kind of line cap is the *butt cap* obtained by adjusting the end positions of the component parallel lines so that the thick line is displayed with square ends that are perpendicular to the line path. If the specified line has slope m , the square end of the thick line has slope $-1/m$. Another line cap is the *round cap* obtained by adding a filled semicircle to each butt cap. The circular arcs are centered on the line endpoints and have a diameter equal to the line thickness. A third type of line cap is the *projecting square cap*. Here, we simply extend the line and add butt caps that are positioned one-half of the line width beyond the specified endpoints.

Other methods for producing thick lines include displaying the line as a filled rectangle or generating the line with a selected pen or brush pattern, as discussed in the next section. To obtain a rectangle representation for the line

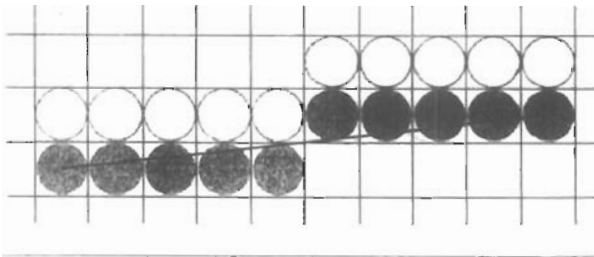


Figure 4-3

Double-wide raster line with slope $|m| < 1$ generated with vertical pixel spans.

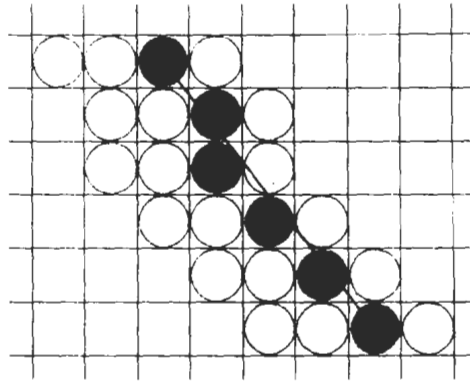


Figure 4-4

Raster line with slope $|m| > 1$
and line-width parameter $lw = 4$
plotted with horizontal pixel spans.

boundary, we calculate the position of the rectangle vertices along perpendiculars to the line path so that vertex coordinates are displaced from the line endpoints by one-half the line width. The rectangular line then appears as in Fig. 4-5(a). We could then add round caps to the filled rectangle or extend its length to display projecting square caps.

Generating thick polylines requires some additional considerations. In general, the methods we have considered for displaying a single line segment will not produce a smoothly connected series of line segments. Displaying thick lines using horizontal and vertical pixel spans, for example, leaves pixel gaps at the boundaries between lines of different slopes where there is a shift from horizontal spans to vertical spans. We can generate thick polylines that are smoothly joined at the cost of additional processing at the segment endpoints. Figure 4-6 shows three possible methods for smoothly joining two line segments. A *miter join* is accomplished by extending the outer boundaries of each of the two lines until they meet. A *round join* is produced by capping the connection between the two segments with a circular boundary whose diameter is equal to the line

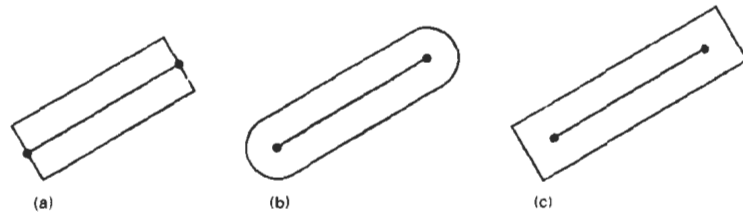


Figure 4-5

Thick lines drawn with (a) butt caps, (b) round caps, and (c) projecting square caps.

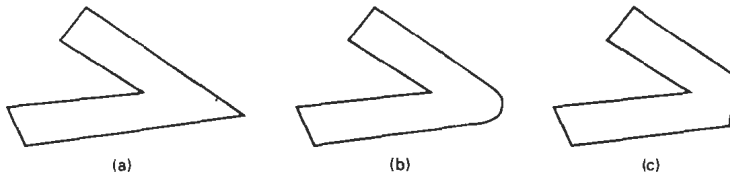


Figure 4-6
Thick line segments connected with (a) miter join, (b) round join, and (c) bevel join.

width. And a *bevel join* is generated by displaying the line segments with butt caps and filling in the triangular gap where the segments meet. If the angle between two connected line segments is very small, a miter join can generate a long spike that distorts the appearance of the polyline. A graphics package can avoid this effect by switching from a miter join to a bevel join, say, when any two consecutive segments meet at a small enough angle.

Pen and Brush Options

With some packages, lines can be displayed with pen or brush selections. Options in this category include shape, size, and pattern. Some possible pen or brush shapes are given in Fig. 4-7. These shapes can be stored in a pixel mask that identifies the array of pixel positions that are to be set along the line path. For example, a rectangular pen can be implemented with the mask shown in Fig. 4-8 by moving the center (or one corner) of the mask along the line path, as in Fig. 4-9. To avoid setting pixels more than once in the frame buffer, we can simply accumulate the horizontal spans generated at each position of the mask and keep track of the beginning and ending x positions for the spans across each scan line.

Lines generated with pen (or brush) shapes can be displayed in various widths by changing the size of the mask. For example, the rectangular pen line in Fig. 4-9 could be narrowed with a 2×2 rectangular mask or widened with a 4×4 mask. Also, lines can be displayed with selected patterns by superimposing the pattern values onto the pen or brush mask. Some examples of line patterns are shown in Fig. 4-10. An additional pattern option that can be provided in a paint package is the display of simulated brush strokes. Figure 4-11 illustrates some patterns that can be displayed by modeling different types of brush strokes.

Line Color

When a system provides color (or intensity) options, a parameter giving the current color index is included in the list of system-attribute values. A polyline routine displays a line in the current color by setting this color value in the frame buffer at pixel locations along the line path using the `setPixel` procedure. The number of color choices depends on the number of bits available per pixel in the frame buffer.

We set the line color value in PHIGS with the function

```
setPolylineColourIndex (lc)
```

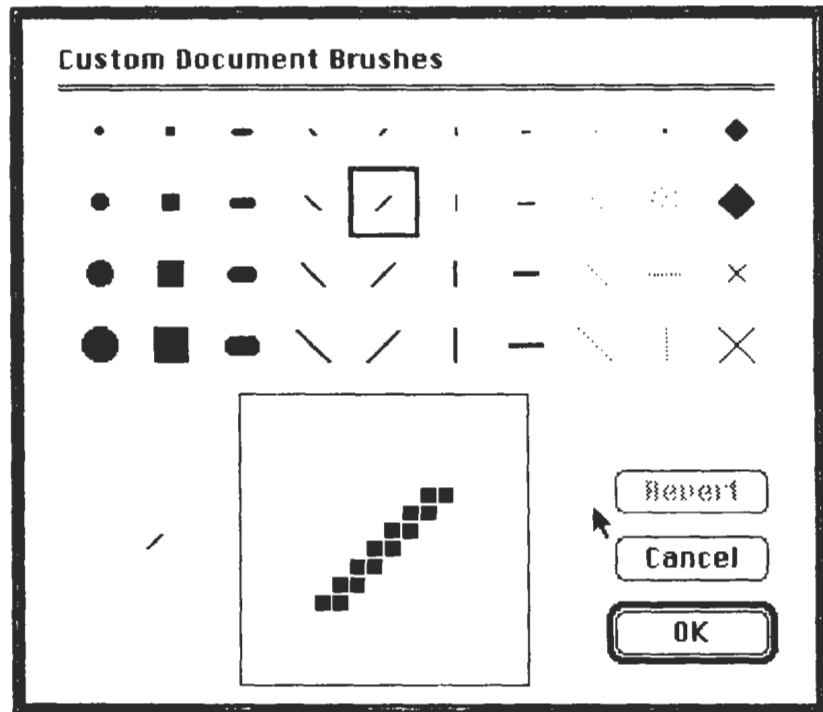


Figure 4-7
Pen and brush shapes for line display.

Nonnegative integer values, corresponding to allowed color choices, are assigned to the line color parameter *lc*. A line drawn in the background color is invisible, and a user can erase a previously displayed line by respecifying it in the background color (assuming the line does not overlap more than one background color area).

An example of the use of the various line attribute commands in an applications program is given by the following sequence of statements:

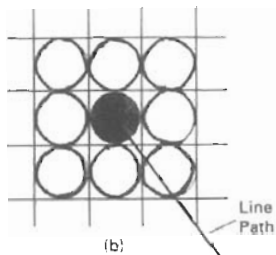
```
setLinetype (2);
setLinewidthScaleFactor (2);
setPolylineColourIndex (5);
polyline (n1, wcpoints1);

setPolylineColourIndex (6);
polyline (n2, wcpoints2);
```

This program segment would display two figures, drawn with double-wide dashed lines. The first is displayed in a color corresponding to code 5, and the second in color 6.

$$\begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$$

(a)



(b)

Figure 4-8

(a) A pixel mask for a rectangular pen, and (b) the associated array of pixels displayed by centering the mask over a specified pixel position.

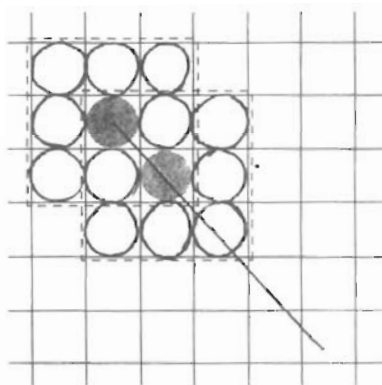


Figure 4-9

Generating a line with the pen shape of Fig. 4-8.



Figure 4-10

Curved lines drawn with a paint program using various shapes and patterns. From left to right, the brush shapes are square, round, diagonal line, dot pattern, and faded airbrush.



Figure 4-11
A daruma doll, a symbol of good fortune in Japan, drawn by computer artist Koichi Kozaki using a paintbrush system. Daruma dolls actually come without eyes. One eye is painted in when a wish is made, and the other is painted in when the wish comes true.
(Courtesy of Wacom Technology, Inc.)

4-2 CURVE ATTRIBUTES

Parameters for curve attributes are the same as those for line segments. We can display curves with varying colors, widths, dot-dash patterns, and available pen or brush options. Methods for adapting curve-drawing algorithms to accommodate attribute selections are similar to those for line drawing.

The pixel masks discussed for implementing line-type options are also used in raster curve algorithms to generate dashed and dotted patterns. For example, the mask 11100 produces the dashed circle shown in Fig. 4-12. We can generate the dashes in the various octants using circle symmetry, but we must shift the pixel positions to maintain the correct sequence of dashes and spaces as we move from one octant to the next. Also, as in line algorithms, pixel masks display dashes and interdash spaces that vary in length according to the slope of the curve. If we want to display constant-length dashes, we need to adjust the number of pixels plotted in each dash as we move around the circle circumference. Instead of applying a pixel mask with constant spans, we plot pixels along equal angular arcs to produce equal length dashes.

Raster curves of various widths can be displayed using the method of horizontal or vertical pixel spans. Where the magnitude of the curve slope is less than 1, we plot vertical spans; where the slope magnitude is greater than 1, we plot horizontal spans. Figure 4-13 demonstrates this method for displaying a circular arc of width 4 in the first quadrant. Using circle symmetry, we generate the circle path with vertical spans in the octant from $x = 0$ to $x = y$, and then reflect pixel positions about the line $y = x$ to obtain the remainder of the curve shown. Circle sections in the other quadrants are obtained by reflecting pixel positions in the

first quadrant about the coordinate axes. The thickness of curves displayed with this method is again a function of curve slope. Circles, ellipses, and other curves will appear thinnest where the slope has a magnitude of 1.

Another method for displaying thick curves is to fill in the area between two parallel curve paths, whose separation distance is equal to the desired width. We could do this using the specified curve path as one boundary and setting up the second boundary either inside or outside the original curve path. This approach, however, shifts the original curve path either inward or outward, depending on which direction we choose for the second boundary. We can maintain the original curve position by setting the two boundary curves at a distance of one-half the width on either side of the specified curve path. An example of this approach is shown in Fig. 4-14 for a circle segment with radius 16 and a specified width of 4. The boundary arcs are then set at a separation distance of 2 on either side of the radius of 16. To maintain the proper dimensions of the circular arc, as discussed in Section 3-10, we can set the radii for the concentric boundary arcs at $r = 14$ and $r = 17$. Although this method is accurate for generating thick circles, in general, it provides only an approximation to the true area of other thick

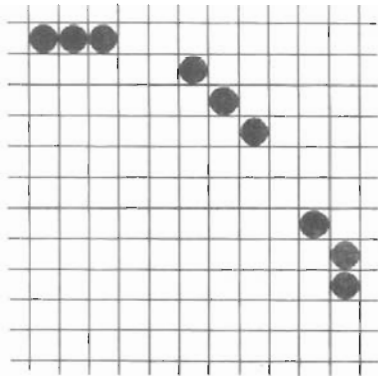


Figure 4-12
A dashed circular arc displayed with a dash span of 3 pixels and an interdash spacing of 2 pixels.

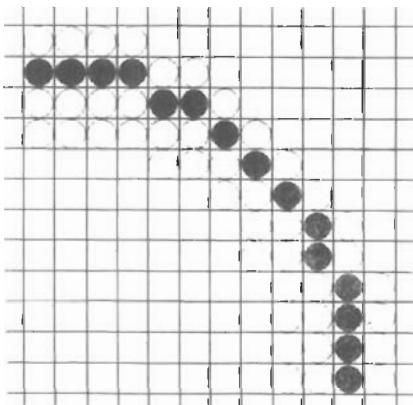


Figure 4-13
Circular arc of width 4 plotted with pixel spans.

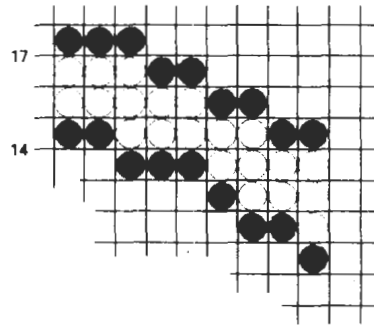


Figure 4-14
A circular arc of width 4 and radius 16 displayed by filling the region between two concentric arcs.

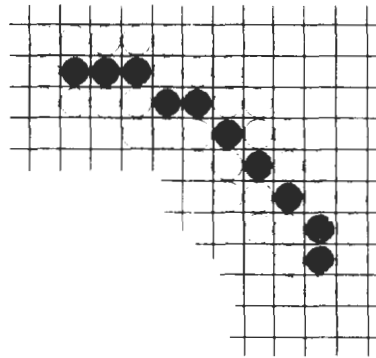


Figure 4-15
Circular arc displayed with a rectangular pen.

curves. For example, the inner and outer boundaries of a fat ellipse generated with this method do not have the same foci.

Pen (or brush) displays of curves are generated using the same techniques discussed for straight line segments. We replicate a pen shape along the line path, as illustrated in Fig. 4-15 for a circular arc in the first quadrant. Here, the center of the rectangular pen is moved to successive curve positions to produce the curve shape shown. Curves displayed with a rectangular pen in this manner will be thicker where the magnitude of the curve slope is 1. A uniform curve thickness can be displayed by rotating the rectangular pen to align it with the slope direction as we move around the curve or by using a circular pen shape. Curves drawn with pen and brush shapes can be displayed in different sizes and with superimposed patterns or simulated brush strokes.

4-3 COLOR AND GRAYSCALE LEVELS

Various color and intensity-level options can be made available to a user, depending on the capabilities and design objectives of a particular system. General-purpose raster-scan systems, for example, usually provide a wide range of colors, while random-scan monitors typically offer only a few color choices, if any. Color

options are numerically coded with values ranging from 0 through the positive integers. For CRT monitors, these color codes are then converted to intensity-level settings for the electron beams. With color plotters, the codes could control ink-jet deposits or pen selections.

In a color raster system, the number of color choices available depends on the amount of storage provided per pixel in the frame buffer. Also, color information can be stored in the frame buffer in two ways: We can store color codes directly in the frame buffer, or we can put the color codes in a separate table and use pixel values as an index into this table. With the direct storage scheme, whenever a particular color code is specified in an application program, the corresponding binary value is placed in the frame buffer for each component pixel in the output primitives to be displayed in that color. A minimum number of colors can be provided in this scheme with 3 bits of storage per pixel, as shown in Table 4-1. Each of the three bit positions is used to control the intensity level (either on or off) of the corresponding electron gun in an RGB monitor. The leftmost bit controls the red gun, the middle bit controls the green gun, and the rightmost bit controls the blue gun. Adding more bits per pixel to the frame buffer increases the number of color choices. With 6 bits per pixel, 2 bits can be used for each gun. This allows four different intensity settings for each of the three color guns, and a total of 64 color values are available for each screen pixel. With a resolution of 1024 by 1024, a full-color (24-bit per pixel) RGB system needs 3 megabytes of storage for the frame buffer. Color tables are an alternate means for providing extended color capabilities to a user without requiring large frame buffers. Lower-cost personal computer systems, in particular, often use color tables to reduce frame-buffer storage requirements.

Color Tables

Figure 4-16 illustrates a possible scheme for storing color values in a **color lookup table** (or **video lookup table**), where frame-buffer values are now used as indices into the color table. In this example, each pixel can reference any one of the 256 table positions, and each entry in the table uses 24 bits to specify an RGB color. For the color code 2081, a combination green-blue color is displayed for pixel location (x, y). Systems employing this particular lookup table would allow

TABLE 4-1
THE EIGHT COLOR CODES FOR A THREE-BIT
PER PIXEL FRAME BUFFER

Color Code	Stored Color Values in Frame Buffer			Displayed Color
	RED	GREEN	BLUE	
0	0	0	0	Black
1	0	0	1	Blue
2	0	1	0	Green
3	0	1	1	Cyan
4	1	0	0	Red
5	1	0	1	Magenta
6	1	1	0	Yellow
7	1	1	1	White

a user to select any 256 colors for simultaneous display from a palette of nearly 17 million colors. Compared to a full-color system, this scheme reduces the number of simultaneous colors that can be displayed, but it also reduces the frame-buffer storage requirements to 1 megabyte. Some graphics systems provide 9 bits per pixel in the frame buffer, permitting a user to select 512 colors that could be used in each display.

A user can set color-table entries in a PHIGS applications program with the function

```
setColourRepresentation (ws, ci, colorptr)
```

Parameter *ws* identifies the workstation output device; parameter *ci* specifies the color index, which is the color-table position number (0 to 255 for the example in Fig. 4-16); and parameter *colorptr* points to a trio of RGB color values (*r*, *g*, *b*) each specified in the range from 0 to 1. An example of possible table entries for color monitors is given in Fig. 4-17.

There are several advantages in storing color codes in a lookup table. Use of a color table can provide a "reasonable" number of simultaneous colors without requiring large frame buffers. For most applications, 256 or 512 different colors are sufficient for a single picture. Also, table entries can be changed at any time, allowing a user to be able to experiment easily with different color combinations in a design, scene, or graph without changing the attribute settings for the graphics data structure. Similarly, visualization applications can store values for some physical quantity, such as energy, in the frame buffer and use a lookup table to try out various color encodings without changing the pixel values. And in visualization and image-processing applications, color tables are a convenient means for setting color thresholds so that all pixel values above or below a specified threshold can be set to the same color. For these reasons, some systems provide both capabilities for color-code storage, so that a user can elect either to use color tables or to store color codes directly in the frame buffer.

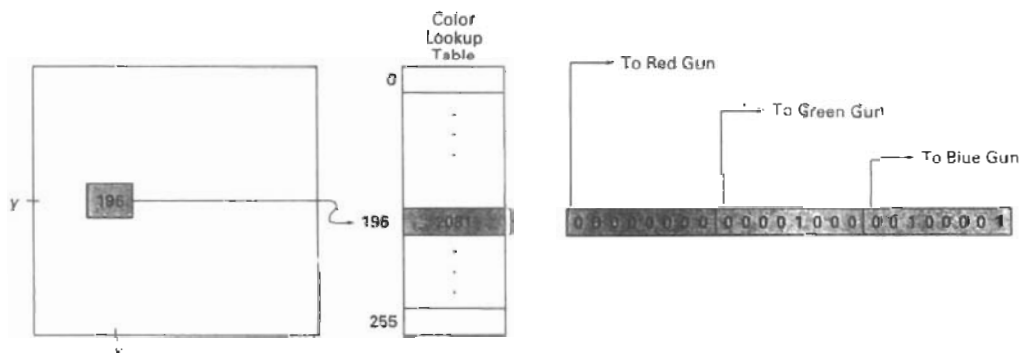


Figure 4-16

A color lookup table with 24 bits per entry accessed from a frame buffer with 8 bits per pixel. A value of 196 stored at pixel position (*x*, *y*) references the location in this table containing the value 2081. Each 8-bit segment of this entry controls the intensity level of one of the three electron guns in an RGB monitor.

WS = 1		WS = 2	
Ci	Color	Ci	Color
0	(0, 0, 0)	0	(1, 1, 1)
1	(0, 0, 0.2)	1	(0.9, 1, 1)
.	.	2	(0.8, 1, 1)
.	.	.	.
.	.	.	.
192	(0, 0.03, 0.13)	.	.
.	.	.	.
.	.	.	.
.	.	.	.

Figure 4-17
Workstation color tables.

Grayscale

With monitors that have no color capability, color functions can be used in an application program to set the shades of gray, or **grayscale**, for displayed primitives. Numeric values over the range from 0 to 1 can be used to specify grayscale levels, which are then converted to appropriate binary codes for storage in the raster. This allows the intensity settings to be easily adapted to systems with differing grayscale capabilities.

Table 4-2 lists the specifications for intensity codes for a four-level grayscale system. In this example, any intensity input value near 0.33 would be stored as the binary value 01 in the frame buffer, and pixels with this value would be displayed as dark gray. If additional bits per pixel are available in the frame buffer, the value of 0.33 would be mapped to the nearest level. With 3 bits per pixel, we can accommodate 8 gray levels; while 8 bits per pixel would give us 256 shades of gray. An alternative scheme for storing the intensity information is to convert each intensity code directly to the voltage value that produces this grayscale level on the output device in use.

When multiple output devices are available at an installation, the same color-table interface may be used for all monitors. In this case, a color table for a monochrome monitor can be set up using a range of RGB values as in Fig. 4-17, with the display intensity corresponding to a given color index *ci* calculated as

$$\text{intensity} = 0.5[\min(r, g, b) + \max(r, g, b)]$$

TABLE 4-2
INTENSITY CODES FOR A FOUR-LEVEL
GRAYSCALE SYSTEM

Intensity Codes	Stored Intensity Values In The Frame Buffer (Binary Code)		Displayed Grayscale
0.0	0	(00)	Black
0.33	1	(01)	Dark gray
0.67	2	(10)	Light gray
1.0	3	(11)	White

4-4

AREA-FILL ATTRIBUTES

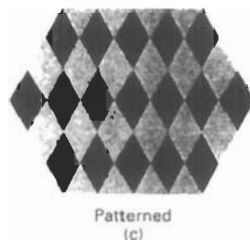
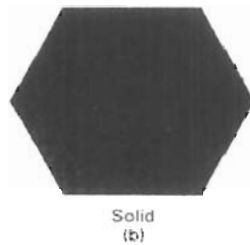
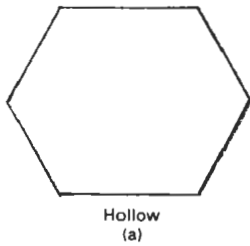


Figure 4-18
Polygon fill styles.

Options for filling a defined region include a choice between a solid color or a patterned fill and choices for the particular colors and patterns. These fill options can be applied to polygon regions or to areas defined with curved boundaries, depending on the capabilities of the available package. In addition, areas can be painted using various brush styles, colors, and transparency parameters.

Fill Styles

Areas are displayed with three basic fill styles: hollow with a color border, filled with a solid color, or filled with a specified pattern or design. A basic fill style is selected in a PHIGS program with the function

```
setInteriorStyle (fs)
```

Values for the fill-style parameter *fs* include *hollow*, *solid*, and *pattern* (Fig. 4-18). Another value for fill style is *hatch*, which is used to fill an area with selected hatching patterns—parallel lines or crossed lines—as in Fig. 4-19. As with line attributes, a selected fill-style value is recorded in the list of system attributes and applied to fill the interiors of subsequently specified areas. Fill selections for parameter *fs* are normally applied to polygon areas, but they can also be implemented to fill regions with curved boundaries.

Hollow areas are displayed using only the boundary outline, with the interior color the same as the background color. A solid fill is displayed in a single color up to and including the borders of the region. The color for a solid interior or for a hollow area outline is chosen with

```
setInteriorColourIndex (fc)
```

where fill-color parameter *fc* is set to the desired color code. A polygon hollow fill is generated with a line-drawing routine as a closed polyline. Solid fill of a region can be accomplished with the scan-line procedures discussed in Section 3-11.

Other fill options include specifications for the edge type, edge width, and edge color of a region. These attributes are set independently of the fill style or fill color, and they provide for the same options as the line-attribute parameters (line type, line width, and line color). That is, we can display area edges dotted or dashed, fat or thin, and in any available color regardless of how we have filled the interior.



Figure 4-19
Polygon fill using hatch patterns.

Pattern Fill

We select fill patterns with

```
setInteriorStyleIndex (pi)
```

where pattern index parameter *pi* specifies a table position. For example, the following set of statements would fill the area defined in the *fillArea* command with the second pattern type stored in the pattern table:

```
setInteriorStyle (pattern);
setInteriorStyleIndex (2);
fillArea (n, points);
```

Separate tables are set up for hatch patterns. If we had selected *hatch* fill for the interior style in this program segment, then the value assigned to parameter *pi* is an index to the stored patterns in the hatch table.

For fill style *pattern*, table entries can be created on individual output devices with

```
setPatternRepresentation (ws, pi, nx, ny, cp)
```

Parameter *pi* sets the pattern index number for workstation code *ws*, and *cp* is a two-dimensional array of color codes with *nx* columns and *ny* rows. The following program segment illustrates how this function could be used to set the first entry in the pattern table for workstation 1.

```
cp[1,1] := 4;          cp[2,2] := 4;
cp[1,2] := 0;          cp[2,1] := 0;
setPatternRepresentation (1, 1, 2, 2, cp);
```

Table 4-3 shows the first two entries for this color table. Color array *cp* in this example specifies a pattern that produces alternate red and black diagonal pixel lines on an eight-color system.

When a color array *cp* is to be applied to fill a region, we need to specify the size of the area that is to be covered by each element of the array. We do this by setting the rectangular coordinate extents of the pattern:

```
setPatternSize (dx, dy)
```

where parameters *dx* and *dy* give the coordinate width and height of the array mapping. An example of the coordinate size associated with a pattern array is given in Fig. 4-20. If the values for *dx* and *dy* in this figure are given in screen coordinates, then each element of the color array would be applied to a 2 by 2 screen grid containing four pixels.

A reference position for starting a *pattern* fill is assigned with the statement

```
setPatternReferencePoint (position)
```

Parameter *position* is a pointer to coordinates (*x_p*, *y_p*) that fix the lower left corner of the rectangular pattern. From this starting position, the pattern is then replicated in the *x* and *y* directions until the defined area is covered by nonover-

TABLE 4-3
A WORKSTATION
PATTERN TABLE WITH
TWO ENTRIES, USING
THE COLOR CODES OF
TABLE 4-1

Index (<i>pi</i>)	Pattern (<i>cp</i>)
1	$\begin{bmatrix} 4 & 0 \\ 0 & 4 \end{bmatrix}$
2	$\begin{bmatrix} 2 & 1 & 2 \\ 1 & 2 & 1 \\ 2 & 1 & 2 \end{bmatrix}$

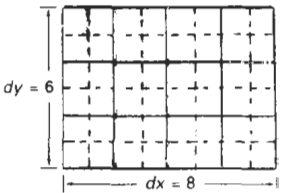


Figure 4-20
A pattern array with 4
columns and 3 rows mapped
to an 8 by 6 coordinate
rectangle.

lapping copies of the pattern array. The process of filling an area with a rectangular pattern is called **tiling** and rectangular fill patterns are sometimes referred to as **tiling patterns**. Figure 4-21 demonstrates tiling of a triangular fill area starting from a pattern reference point.

To illustrate the use of the pattern commands, the following program example displays a black-and-white pattern in the interior of a parallelogram fill area (Fig. 4-22). The pattern size in this program is set to map each array element to a single pixel.

```
#define WS 1

void patternFill ()
{
    wcPt2 pts[4];
    int bwPattern[3][3] = { 1, 0, 0, 0, 1, 1, 1, 0, 0 };

    pSetPatternRepresentation (WS, 8, 3, 3, bwPattern);

    pts[0].x = 10; pts[0].y = 10;
    pts[1].x = 20; pts[1].y = 10;
    pts[2].x = 28; pts[2].y = 18;
    pts[3].x = 18; pts[3].y = 18;

    pSetFillAreaInteriorStyle (PATTERN);
    pSetFillAreaPatternIndex (8);
    pSetPatternReferencePoint (14, 11);

    pFillArea (4, pts);
}
```

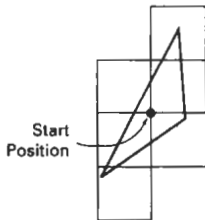


Figure 4-21

Tiling an area from a designated start position. Nonoverlapping adjacent patterns are laid out to cover all scan lines passing through the defined area.

Pattern fill can be implemented by modifying the scan-line procedures discussed in Chapter 3 so that a selected pattern is superimposed onto the scan lines. Beginning from a specified start position for a pattern fill, the rectangular patterns would be mapped vertically to scan lines between the top and bottom of the fill area and horizontally to interior pixel positions across these scan lines. Horizontally, the pattern array is repeated at intervals specified by the value of size parameter *dx*. Similarly, vertical repeats of the pattern are separated by intervals set with parameter *dy*. This scan-line pattern procedure applies both to polygons and to areas bounded by curves.

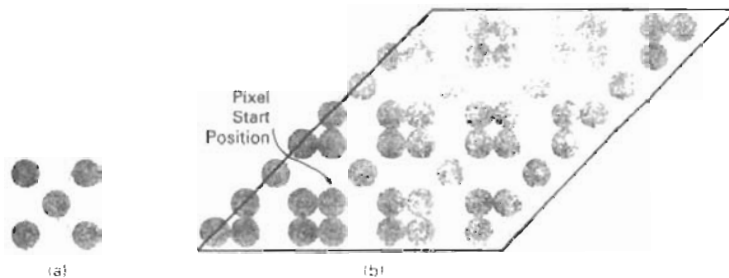


Figure 4-22

A pattern array (a) superimposed on a parallelogram fill area to produce the display (b).

Hatch fill is applied to regions by displaying sets of parallel lines. The fill procedures are implemented to draw either single hatching or cross hatching. Spacing and slope for the hatch lines can be set as parameters in the hatch table. On raster systems, a hatch fill can be specified as a pattern array that sets color values for groups of diagonal pixels.

In many systems, the pattern reference point (xp, yp) is assigned by the system. For instance, the reference point could be set automatically at a polygon vertex. In general, for any fill region, the reference point can be chosen as the lower left corner of the *bounding rectangle* (or *bounding box*) determined by the coordinate extents of the region (Fig. 4-23). To simplify selection of the reference coordinates, some packages always use the screen coordinate origin as the pattern start position, and window systems often set the reference point at the coordinate origin of the window. Always setting (xp, yp) at the coordinate origin also simplifies the tiling operations when each color-array element of a pattern is to be mapped to a single pixel. For example, if the row positions in the pattern array are referenced in reverse (that is, from bottom to top starting at 1), a pattern value is then assigned to pixel position (x, y) in screen or window coordinates as

```
setPixel ( x, y, cp(y mod ny + 1, x mod nx + 1) )
```

where ny and nx specify the number of rows and number of columns in the pattern array. Setting the pattern start position at the coordinate origin, however, effectively attaches the pattern fill to the screen or window background, rather than to the fill regions. Adjacent or overlapping areas filled with the same pattern would show no apparent boundary between the areas. Also, repositioning and refilling an object with the same pattern can result in a shift in the assigned pixel values over the object interior. A moving object would appear to be transparent against a stationary pattern background, instead of moving with a fixed interior pattern.

It is also possible to combine a fill pattern with background colors (including grayscale) in various ways. With a bitmap pattern containing only the digits 1 and 0, the 0 values could be used as transparency indicators to let the background show through. Alternatively, the 1 and 0 digits can be used to fill an interior with two-color patterns. In general, color-fill patterns can be combined in several other ways with background colors. The pattern and background colors can be combined using Boolean operations, or the pattern colors can simply replace the background colors. Figure 4-24 demonstrates how the Boolean and replace operations for a 2 by 2 fill pattern would set pixel values on a binary (black-and-white) system against a particular background pattern.

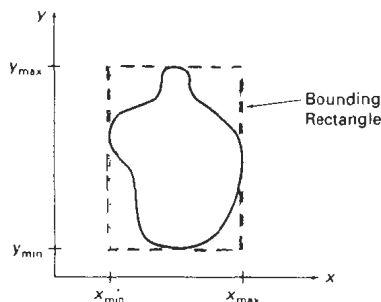


Figure 4-23
Bounding rectangle for a region
with coordinate extents x_{min} , x_{max} ,
 y_{min} , and y_{max} in the x and y
directions.

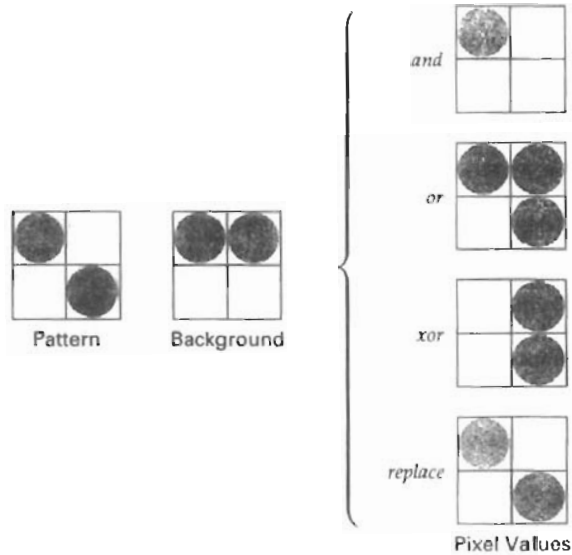


Figure 4-24

Combining a fill pattern with a background pattern using Boolean operations, *and*, *or*, and *xor* (exclusive or), and using simple replacement.

Soft Fill

Modified boundary-fill and flood-fill procedures that are applied to repaint areas so that the fill color is combined with the background colors are referred to as **soft-fill** or **tint-fill** algorithms. One use for these fill methods is to soften the fill colors at object borders that have been blurred to antialias the edges. Another is to allow repainting of a color area that was originally filled with a semitransparent brush, where the current color is then a mixture of the brush color and the background colors "behind" the area. In either case, we want the new fill color to have the same variations over the area as the current fill color.

As an example of this type of fill, the **linear soft-fill** algorithm repaints an area that was originally painted by merging a foreground color F with a single background color B , where $F \neq B$. Assuming we know the values for F and B , we can determine how these colors were originally combined by checking the current color contents of the frame buffer. The current RGB color P of each pixel within the area to be refilled is some linear combination of F and B :

$$P = tF + (1 - t)B \quad (4-1)$$

where the "transparency" factor t has a value between 0 and 1 for each pixel. For values of t less than 0.5, the background color contributes more to the interior color of the region than does the fill color. Vector Equation 4-1 holds for each

RGB component of the colors, with

$$\mathbf{P} = (P_R, P_G, P_B), \quad \mathbf{F} = (F_R, F_G, F_B), \quad \mathbf{B} = (B_R, B_G, B_B) \quad (4-2)$$

We can thus calculate the value of parameter t using one of the RGB color components as

$$t = \frac{P_k - B_k}{F_k - B_k} \quad (4-3)$$

where $k = R, G, \text{ or } B$; and $F_k \neq B_k$. Theoretically, parameter t has the same value for each RGB component, but roundoff to integer codes can result in different values of t for different components. We can minimize this roundoff error by selecting the component with the largest difference between \mathbf{F} and \mathbf{B} . This value of t is then used to mix the new fill color \mathbf{NF} with the background color, using either a modified flood-fill or boundary-fill procedure.

Similar soft-fill procedures can be applied to an area whose foreground color is to be merged with multiple background color areas, such as a checkerboard pattern. When two background colors B_1 and B_2 are mixed with foreground color \mathbf{F} , the resulting pixel color \mathbf{P} is

$$\mathbf{P} = t_0\mathbf{F} + t_1\mathbf{B}_1 + (1 - t_0 - t_1)\mathbf{B}_2 \quad (4-4)$$

where the sum of the coefficients t_0 , t_1 , and $(1 - t_0 - t_1)$ on the color terms must equal 1. We can set up two simultaneous equations using two of the three RGB color components to solve for the two proportionality parameters, t_0 and t_1 . These parameters are then used to mix the new fill color with the two background colors to obtain the new pixel color. With three background colors and one foreground color, or with two background and two foreground colors, we need all three RGB equations to obtain the relative amounts of the four colors. For some foreground and background color combinations, however, the system of two or three RGB equations cannot be solved. This occurs when the color values are all very similar or when they are all proportional to each other.

4-5

CHARACTER ATTRIBUTES

The appearance of displayed characters is controlled by attributes such as font, size, color, and orientation. Attributes can be set both for entire character strings (text) and for individual characters defined as marker symbols.

Text Attributes

There are a great many text options that can be made available to graphics programmers. First of all, there is the choice of font (or typeface), which is a set of characters with a particular design style such as New York, Courier, Helvetica, London, Times Roman, and various special symbol groups. The characters in a selected font can also be displayed with assorted underlining styles (solid, dotted, double), in **boldface**, in *italics*, and in outline or shadow styles. A particular

font and associated style is selected in a PHIGS program by setting an integer code for the text font parameter `tf` in the function

```
setTextFont (tf)
```

Font options can be made available as predefined sets of grid patterns or as character sets designed with polylines and spline curves.

Color settings for displayed text are stored in the system attribute list and used by the procedures that load character definitions into the frame buffer. When a character string is to be displayed, the current color is used to set pixel values in the frame buffer corresponding to the character shapes and positions. Control of text color (or intensity) is managed from an application program with

```
setTextColourIndex (tc)
```

where text color parameter `tc` specifies an allowable color code.

We can adjust text size by scaling the overall dimensions (height and width) of characters or by scaling only the character width. Character size is specified by printers and compositors in *points*, where 1 point is 0.013837 inch (or approximately 1/72 inch). For example, the text you are now reading is a 10-point font. Point measurements specify the size of the *body* of a character (Fig. 4-25), but different fonts with the same point specifications can have different character sizes, depending on the design of the typeface. The distance between the *baseline* and the *topline* of the character body is the same for all characters in a particular size and typeface, but the body width may vary. *Proportionally spaced fonts* assign a smaller body width to narrow characters such as *i*, *j*, *l*, and *f* compared to broad characters such as *W* or *M*. *Character height* is defined as the distance between the *baseline* and the *capline* of characters. *Kerned* characters, such as *f* and *j* in Fig. 4-25, typically extend beyond the character-body limits, and letters with descenders (*g*, *j*, *p*, *q*, *y*) extend below the baseline. Each character is positioned within the character body by a font designer to allow suitable spacing along and between print lines when text is displayed with character bodies touching.

Text size can be adjusted without changing the width-to-height ratio of characters with

```
setCharacterHeight (ch)
```

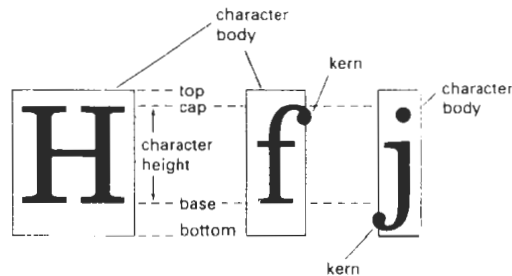


Figure 4-25
Character body.

Height 1
Height 2
Height 3

Figure 4-26
The effect of different character-height settings on displayed text.

Parameter `ch` is assigned a real value greater than 0 to set the coordinate height of capital letters: the distance between baseline and capline in user coordinates. This setting also affects character-body size, so that the width and spacing of characters is adjusted to maintain the same text proportions. For instance, doubling the height also doubles the character width and the spacing between characters. Figure 4-26 shows a character string displayed with three different character heights.

The width only of text can be set with the function

```
setCharacterExpansionFactor (cw)
```

where the character-width parameter `cw` is set to a positive real value that scales the body width of characters. Text height is unaffected by this attribute setting. Examples of text displayed with different character expansions is given in Fig. 4-27.

Spacing between characters is controlled separately with

```
setCharacterSpacing (cs)
```

where the character-spacing parameter `cs` can be assigned any real value. The value assigned to `cs` determines the spacing between character bodies along print lines. Negative values for `cs` overlap character bodies; positive values insert space to spread out the displayed characters. Assigning the value 0 to `cs` causes text to be displayed with no space between character bodies. The amount of spacing to be applied is determined by multiplying the value of `cs` by the character height (distance between baseline and capline). In Fig. 4-28, a character string is displayed with three different settings for the character-spacing parameter.

The orientation for a displayed character string is set according to the direction of the **character up vector**:

```
setCharacterUpVector (upvect)
```

Parameter `upvect` in this function is assigned two values that specify the *x* and *y* vector components. Text is then displayed so that the orientation of characters from baseline to capline is in the direction of the up vector. For example, with `upvect = (1, 1)`, the direction of the up vector is 45° and text would be displayed as shown in Fig. 4-29. A procedure for orienting text rotates characters so that the sides of character bodies, from baseline to capline, are aligned with the up vector. The rotated character shapes are then scan converted into the frame buffer.

width 0.5
width 1.0
width 2.0

Figure 4-27
The effect of different character-width settings on displayed text.

Spacing 0.0
Spacing 0.5
Spacing 1.0

Figure 4-28
The effect of different character spacings on displayed text.

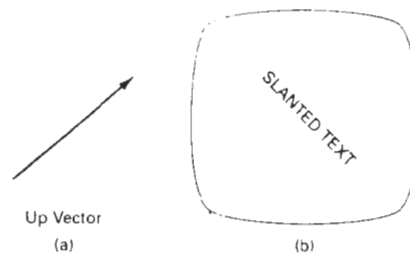


Figure 4-29
Direction of the up vector (a) controls the orientation of displayed text (b).

It is useful in many applications to be able to arrange character strings vertically or horizontally (Fig. 4-30). An attribute parameter for this option is set with the statement

```
setTextPath (tp)
```

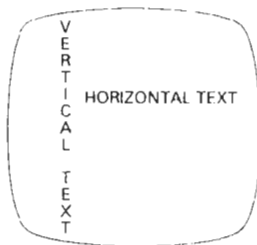


Figure 4-30
Text path attributes can be set to produce horizontal or vertical arrangements of character strings.

where the text-path parameter *tp* can be assigned the value: *right*, *left*, *up*, or *down*. Examples of text displayed with these four options are shown in Fig. 4-31. A procedure for implementing this option must transform the character patterns into the specified orientation before transferring them to the frame buffer.

Character strings can also be oriented using a combination of up-vector and text-path specifications to produce slanted text. Figure 4-32 shows the directions of character strings generated by the various text-path settings for a 45° up vector. Examples of text generated for text-path values *down* and *right* with this up vector are illustrated in Fig. 4-33.

Another handy attribute for character strings is alignment. This attribute specifies how text is to be positioned with respect to the start coordinates. Alignment attributes are set with

```
setTextAlignment (h, v)
```

where parameters *h* and *v* control horizontal and vertical alignment, respectively. Horizontal alignment is set by assigning *h* a value of *left*, *centre*, or *right*. Vertical alignment is set by assigning *v* a value of *top*, *cap*, *half*, *base*, or *bottom*. The interpretation of these alignment values depends on the current setting for the text path. Figure 4-34 shows the position of the alignment settings when text is to be displayed horizontally to the right or vertically down. Similar interpretations apply to text path values of *left* and *up*. The "most natural" alignment for a particular text path is chosen by assigning the value *normal* to the *h* and *v* parameters. Figure 4-35 illustrates common alignment positions for horizontal and vertical text labels.

A precision specification for text display is given with

```
setTextPrecision (tpr)
```

where text precision parameter *tpr* is assigned one of the values: *string*, *char*, or *stroke*. The highest-quality text is displayed when the precision parameter is set to the value *stroke*. For this precision setting, greater detail would be used in defining the character shapes, and the processing of attribute selections and other



Figure 4-31
Text displayed with the four text-path options.

string-manipulation procedures would be carried out to the highest possible accuracy. The lowest-quality precision setting, *string*, is used for faster display of character strings. At this precision, many attribute selections such as text path are ignored, and string-manipulation procedures are simplified to reduce processing time.

Marker Attributes

A marker symbol is a single character that can be displayed in different colors and in different sizes. Marker attributes are implemented by procedures that load the chosen character into the raster at the defined positions with the specified color and size.

We select a particular character to be the marker symbol with

```
setMarkerType (mt)
```

where marker type parameter *mt* is set to an integer code. Typical codes for marker type are the integers 1 through 5, specifying, respectively, a dot (·), a vertical cross (+), an asterisk (*), a circle (o), and a diagonal cross (×). Displayed marker types are centered on the marker coordinates.

We set the marker size with

```
setMarkerSizeScaleFactor (ms)
```

with parameter marker size *ms* assigned a positive number. This scaling parameter is applied to the nominal size for the particular marker symbol chosen. Values greater than 1 produce character enlargement; values less than 1 reduce the marker size.

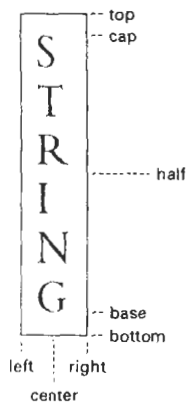


Figure 4-34
Alignment attribute values for horizontal and vertical strings.

Section 4-5

Character Attributes

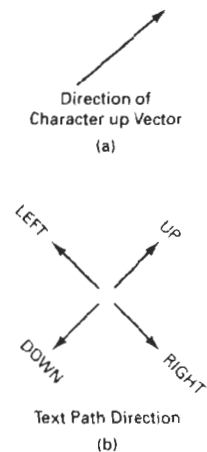


Figure 4-32
An up-vector specification (a) controls the direction of the text path (b).

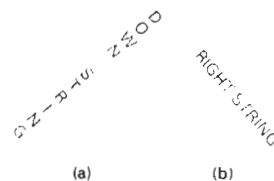


Figure 4-33
The 45° up vector in Fig. 4-32 produces the display (a) for a down path and the display (b) for a right path.

RIGHT	
ALIGNMENT	
CENTER	
ALIGNMENT	
LEFT	
ALIGNMENT	
BOTTOM	
ALIGNMENT	
TOP	
ALIGNMENT	

Figure 4-35
Character-string alignments.

Marker color is specified with

```
setPolymarkerColourIndex (mc)
```

A selected color code for parameter *mc* is stored in the current attribute list and used to display subsequently specified marker primitives.

4-6 BUNDLED ATTRIBUTES

With the procedures we have considered so far, each function references a single attribute that specifies exactly how a primitive is to be displayed with that attribute setting. These specifications are called **individual** (or **unbundled**) attributes, and they are meant to be used with an output device that is capable of displaying primitives in the way specified. If an application program, employing individual attributes, is interfaced to several output devices, some of the devices may not have the capability to display the intended attributes. A program using individual color attributes, for example, may have to be modified to produce acceptable output on a monochromatic monitor.

Individual attribute commands provide a simple and direct method for specifying attributes when a single output device is used. When several kinds of output devices are available at a graphics installation, it is convenient for a user to be able to say how attributes are to be interpreted on each of the different devices. This is accomplished by setting up tables for each output device that lists sets of attribute values that are to be used on that device to display each primitive type. A particular set of attribute values for a primitive on each output device is then chosen by specifying the appropriate table index. Attributes specified in this manner are called **bundled** attributes. The table for each primitive that defines groups of attribute values to be used when displaying that primitive on a particular output device is called a **bundle table**.

Attributes that may be bundled into the workstation table entries are those that do not involve coordinate specifications, such as color and line type. The choice between a bundled or an unbundled specification is made by setting a switch called the **aspect source flag** for each of these attributes:

```
setIndividualASF (attributeptr, flagptr)
```

where parameter *attributeptr* points to a list of attributes, and parameter *flagptr* points to the corresponding list of aspect source flags. Each aspect source flag can be assigned a value of *individual* or *bundled*. Attributes that may be bundled are listed in the following sections.

Bundled Line Attributes

Entries in the bundle table for line attributes on a specified workstation are set with the function

```
setPolylineRepresentation (ws, li, lt, lw, lc)
```

Parameter *ws* is the workstation identifier, and line index parameter *li* defines the bundle table position. Parameters *lt*, *lw*, and *lc* are then bundled and assigned values to set the line type, line width, and line color specifications, respectively, for the designated table index. For example, the following statements define groups of line attributes that are to be referenced as index number 3 on two different workstations:

```
setPolylineRepresentation (1, 3, 2, 0.5, 1);  
setPolylineRepresentation (4, 3, 1, 1, 7);
```

A polyline that is assigned a table index value of 3 would then be displayed using dashed lines at half thickness in a blue color on workstation 1; while on workstation 4, this same index generates solid, standard-sized white lines.

Once the bundle tables have been set up, a group of bundled line attributes is chosen for each workstation by specifying the table index value:

```
setPolylineIndex (li)
```

Subsequent *polyline* commands would then generate lines on each workstation according to the set of bundled attribute values defined at the table position specified by the value of the line index parameter *li*.

Bundled Area-Fill Attributes

Table entries for bundled area-fill attributes are set with

```
setInteriorRepresentation (ws, fi, fs, pi, fc)
```

which defines the attribute list corresponding to fill index *fi* on workstation *ws*. Parameters *fs*, *pi*, and *fc* are assigned values for the fill style, pattern index, and fill color, respectively, on the designated workstation. Similar bundle tables can also be set up for edge attributes of polygon fill areas.

A particular attribute bundle is then selected from the table with the function

```
setInteriorIndex (fi)
```

Subsequently defined fill areas are then displayed on each active workstation according to the table entry specified by the fill index parameter *fi*. Other fill-area attributes, such as pattern reference point and pattern size, are independent of the workstation designation and are set with the functions previously described.

Bundled Text Attributes

The function

```
setTextRepresentation (ws, ti, tf, tp, te, ts, tc)
```

bundles values for text font, precision, expansion factor, size, and color in a table position for workstation *ws* that is specified by the value assigned to text index

parameter `ti`. Other text attributes, including character up vector, text path, character height, and text alignment are set individually.

A particular text index value is then chosen with the function

```
setTextIndex (ti)
```

Each text function that is then invoked is displayed on each workstation with the set of attributes referenced by this table position.

Bundled Marker Attributes

Table entries for bundled marker attributes are set up with

```
setPolymarkerRepresentation (ws, mi, mt, ms, mc)
```

This defines the marker type, marker scale factor, and marker color for index `mi` on workstation `ws`. Bundle table selections are then made with the function

```
setPolymarkerIndex (mi)
```

4-7

INQUIRY FUNCTIONS

Current settings for attributes and other parameters, such as workstation types and status, in the system lists can be retrieved with inquiry functions. These functions allow current values to be copied into specified parameters, which can then be saved for later reuse or used to check the current state of the system if an error occurs.

We check current attribute values by stating the name of the attribute in the inquiry function. For example, the functions

```
inquirePolylineIndex (lastli)
```

and

```
inquireInteriorColourIndex (lastfc)
```

copy the current values for line index and fill color into parameters `lastli` and `lastfc`. The following program segment illustrates reusing the current line type value after a set of lines are drawn with a new line type.

```
inquireLinetype (oldlt);  
setLinetype (newlt);  
.  
.  
.  
setLinetype (oldlt);
```


Displayed primitives generated by the raster algorithms discussed in Chapter 3 have a jagged, or stairstep, appearance because the sampling process digitizes coordinate points on an object to discrete integer pixel positions. This distortion of information due to low-frequency sampling (undersampling) is called **aliasing**. We can improve the appearance of displayed raster lines by applying **antialiasing** methods that compensate for the undersampling process.

An example of the effects of undersampling is shown in Fig. 4-36. To avoid losing information from such periodic objects, we need to set the sampling frequency to at least twice that of the highest frequency occurring in the object, referred to as the **Nyquist sampling frequency** (or Nyquist sampling rate) f_s :

$$f_s = 2f_{\max} \quad (4-5)$$

Another way to state this is that the sampling interval should be no larger than one-half the cycle interval (called the **Nyquist sampling interval**). For x -interval sampling, the Nyquist sampling interval Δx_s is

$$\Delta x_s = \frac{\Delta x_{\text{cycle}}}{2} \quad (4-6)$$

where $\Delta x_{\text{cycle}} = 1/f_{\max}$. In Fig. 4-36, our sampling interval is one and one-half times the cycle interval, so the sampling interval is at least three times too big. If we want to recover all the object information for this example, we need to cut the sampling interval down to one-third the size shown in the figure.

One way to increase sampling rate with raster systems is simply to display objects at higher resolution. But even at the highest resolution possible with current technology, the jaggies will be apparent to some extent. There is a limit to how big we can make the frame buffer and still maintain the refresh rate at 30 to 60 frames per second. And to represent objects accurately with continuous parameters, we need arbitrarily small sampling intervals. Therefore, unless hardware technology is developed to handle arbitrarily large frame buffers, increased screen resolution is not a complete solution to the aliasing problem.

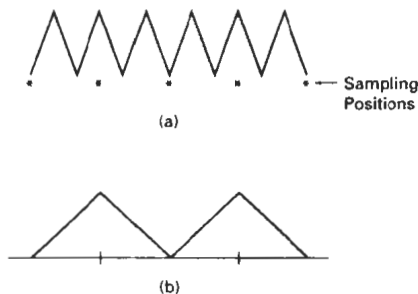


Figure 4-36
Sampling the periodic shape in (a) at the marked positions produces the aliased lower-frequency representation in (b).

With raster systems that are capable of displaying more than two intensity levels (color or gray scale), we can apply antialiasing methods to modify pixel intensities. By appropriately varying the intensities of pixels along the boundaries of primitives, we can smooth the edges to lessen the jagged appearance.

A straightforward antialiasing method is to increase sampling rate by treating the screen as if it were covered with a finer grid than is actually available. We can then use multiple sample points across this finer grid to determine an appropriate intensity level for each screen pixel. This technique of sampling object characteristics at a high resolution and displaying the results at a lower resolution is called **supersampling** (or **postfiltering**, since the general method involves computing intensities at subpixel grid positions, then combining the results to obtain the pixel intensities). Displayed pixel positions are spots of light covering a finite area of the screen, and not infinitesimal mathematical points. Yet in the line and fill-area algorithms we have discussed, the intensity of each pixel is determined by the location of a single point on the object boundary. By supersampling, we obtain intensity information from multiple points that contribute to the overall intensity of a pixel.

An alternative to supersampling is to determine pixel intensity by calculating the areas of overlap of each pixel with the objects to be displayed. Antialiasing by computing overlap areas is referred to as **area sampling** (or **prefiltering**, since the intensity of the pixel as a whole is determined without calculating subpixel intensities). Pixel overlap areas are obtained by determining where object boundaries intersect individual pixel boundaries.

Raster objects can also be antialiased by shifting the display location of pixel areas. This technique, called **pixel phasing**, is applied by "micropositioning" the electron beam in relation to object geometry.

Supersampling Straight Line Segments

Supersampling straight lines can be performed in several ways. For the gray-scale display of a straight-line segment, we can divide each pixel into a number of subpixels and count the number of subpixels that are along the line path. The intensity level for each pixel is then set to a value that is proportional to this subpixel count. An example of this method is given in Fig. 4-37. Each square pixel area is divided into nine equal-sized square subpixels, and the shaded regions show the subpixels that would be selected by Bresenham's algorithm. This scheme provides for three intensity settings above zero, since the maximum number of subpixels that can be selected within any pixel is three. For this example, the pixel at position (10, 20) is set to the maximum intensity (level 3); pixels at (11, 21) and (12, 21) are each set to the next highest intensity (level 2); and pixels at (11, 20) and (12, 22) are each set to the lowest intensity above zero (level 1). Thus the line intensity is spread out over a greater number of pixels, and the stairstep effect is smoothed by displaying a somewhat blurred line path in the vicinity of the stair steps (between horizontal runs). If we want to use more intensity levels to antialias the line with this method, we increase the number of sampling positions across each pixel. Sixteen subpixels gives us four intensity levels above zero; twenty-five subpixels gives us five levels; and so on.

In the supersampling example of Fig. 4-37, we considered pixel areas of finite size, but we treated the line as a mathematical entity with zero width. Actually, displayed lines have a width approximately equal to that of a pixel. If we take the finite width of the line into account, we can perform supersampling by setting each pixel intensity proportional to the number of subpixels inside the

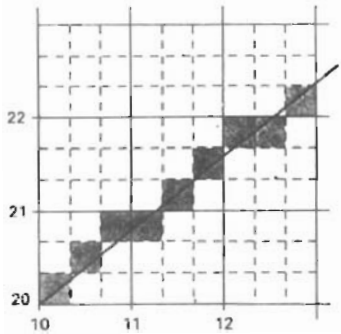


Figure 4-37
Supersampling subpixel positions
along a straight line segment whose
left endpoint is at screen
coordinates (10, 20).

polygon representing the line area. A subpixel can be considered to be inside the line if its lower left corner is inside the polygon boundaries. An advantage of this supersampling procedure is that the number of possible intensity levels for each pixel is equal to the total number of subpixels within the pixel area. For the example in Fig. 4-37, we can represent this line with finite width by positioning the polygon boundaries parallel to the line path as in Fig. 4-38. And each pixel can now be set to one of nine possible brightness levels above zero.

Another advantage of supersampling with a finite-width line is that the total line intensity is distributed over more pixels. In Fig. 4-38, we now have the pixel at grid position (10, 21) turned on (at intensity level 2), and we also pick up contributions from pixels immediately below and immediately to the left of position (10, 21). Also, if we have a color display, we can extend the method to take background colors into account. A particular line might cross several different color areas, and we can average subpixel intensities to obtain pixel color settings. For instance, if five subpixels within a particular pixel area are determined to be inside the boundaries for a red line and the remaining four pixels fall within a blue background area, we can calculate the color for this pixel as

$$\text{pixel}_{\text{color}} = (5 \cdot \text{red} + 4 \cdot \text{blue})/9$$

The trade-off for these gains from supersampling a finite-width line is that identifying interior subpixels requires more calculations than simply determining which subpixels are along the line path. These calculations are also complicated by the positioning of the line boundaries in relation to the line path. This posi-

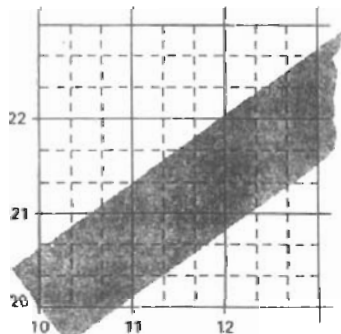


Figure 4-38
Supersampling subpixel positions
in relation to the interior of a line of
finite width.

tioning depends on the slope of the line. For a 45° line, the line path is centered on the polygon area; but for either a horizontal or a vertical line, we want the line path to be one of the polygon boundaries. For instance, a horizontal line passing through grid coordinates (10, 20) would be represented as the polygon bounded by horizontal grid lines $y = 20$ and $y = 21$. Similarly, the polygon representing a vertical line through (10, 20) would have vertical boundaries along grid lines $x = 10$ and $x = 11$. For lines with slope $|m| < 1$, the mathematical line path is positioned proportionately closer to the lower polygon boundary; and for lines with slope $|m| > 1$, this line path is placed closer to the upper polygon boundary.

Pixel-Weighting Masks

Supersampling algorithms are often implemented by giving more weight to subpixels near the center of a pixel area, since we would expect these subpixels to be more important in determining the overall intensity of a pixel. For the 3 by 3 pixel subdivisions we have considered so far, a weighting scheme as in Fig. 4-39 could be used. The center subpixel here is weighted four times that of the corner subpixels and twice that of the remaining subpixels. Intensities calculated for each grid of nine subpixels would then be averaged so that the center subpixel is weighted by a factor of $1/4$; the top, bottom, and side subpixels are each weighted by a factor of $1/8$; and the corner subpixels are each weighted by a factor of $1/16$. An array of values specifying the relative importance of subpixels is sometimes referred to as a “mask” of subpixel weights. Similar masks can be set up for larger subpixel grids. Also, these masks are often extended to include contributions from subpixels belonging to neighboring pixels, so that intensities can be averaged over adjacent pixels.

Area Sampling Straight Line Segments

We perform area sampling for a straight line by setting each pixel intensity proportional to the area of overlap of the pixel with the finite-width line. The line can be treated as a rectangle, and the section of the line area between two adjacent vertical (or two adjacent horizontal) screen grid lines is then a trapezoid. Overlap areas for pixels are calculated by determining how much of the trapezoid overlaps each pixel in that vertical column (or horizontal row). In Fig. 4-38, the pixel with screen grid coordinates (10, 20) is about 90 percent covered by the line area, so its intensity would be set to 90 percent of the maximum intensity. Similarly, the pixel at (10, 21) would be set to an intensity of about 15 percent of maximum. A method for estimating pixel overlap areas is illustrated by the supersampling example in Fig. 4-38. The total number of subpixels within the line boundaries is approximately equal to the overlap area, and this estimation is improved by using finer subpixel grids. With color displays, the areas of pixel overlap with different color regions is calculated and the final pixel color is taken as the average color of the various overlap areas.

1	2	1
2	4	2
1	2	1

Figure 4-39
Relative weights for a grid of
3 by 3 subpixels.

Filtering Techniques

A more accurate method for antialiasing lines is to use **filtering** techniques. The method is similar to applying a weighted pixel mask, but now we imagine a continuous *weighting surface* (or *filter function*) covering the pixel. Figure 4-40 shows examples of rectangular, conical, and Gaussian filter functions. Methods for applying the filter function are similar to applying a weighting mask, but now we

integrate over the pixel surface to obtain the weighted average intensity. To reduce computation, table lookups are commonly used to evaluate the integrals.

Pixel Phasing

On raster systems that can address subpixel positions within the screen grid, pixel phasing can be used to antialias objects. Stairsteps along a line path or object boundary are smoothed out by moving (micropositioning) the electron beam to more nearly approximate positions specified by the object geometry. Systems incorporating this technique are designed so that individual pixel positions can be shifted by a fraction of a pixel diameter. The electron beam is typically shifted by $1/4$, $1/2$, or $3/4$ of a pixel diameter to plot points closer to the true path of a line or object edge. Some systems also allow the size of individual pixels to be adjusted as an additional means for distributing intensities. Figure 4-41 illustrates the antialiasing effects of pixel phasing on a variety of line paths.

Compensating for Line Intensity Differences

Antialiasing a line to soften the stairstep effect also compensates for another raster effect, illustrated in Fig. 4-42. Both lines are plotted with the same number of pixels, yet the diagonal line is longer than the horizontal line by a factor of $\sqrt{2}$. The visual effect of this is that the diagonal line appears less bright than the horizontal line, because the diagonal line is displayed with a lower intensity per unit length. A line-drawing algorithm could be adapted to compensate for this effect by adjusting the intensity of each line according to its slope. Horizontal and vertical lines would be displayed with the lowest intensity, while 45° lines would be given the highest intensity. But if antialiasing techniques are applied to a display,

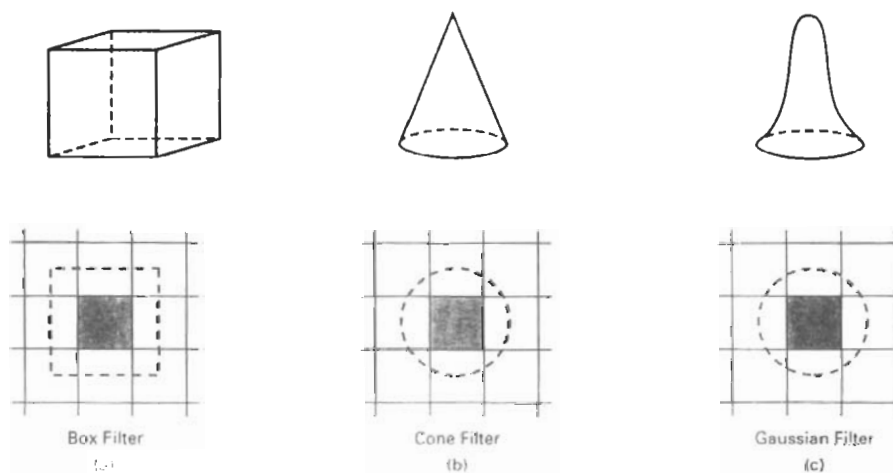


Figure 4-40

Common filter functions used to antialias line paths. The volume of each filter is normalized to 1, and the height gives the relative weight at any subpixel position.

intensities are automatically compensated. When the finite width of lines is taken into account, pixel intensities are adjusted so that lines display a total intensity proportional to their length.

Antialiasing Area Boundaries

The antialiasing concepts we have discussed for lines can also be applied to the boundaries of areas to remove their jagged appearance. We can incorporate these procedures into a scan-line algorithm to smooth the area outline as the area is generated.

If system capabilities permit the repositioning of pixels, area boundaries can be smoothed by adjusting boundary pixel positions so that they are along the line defining an area boundary. Other methods adjust each pixel intensity at a boundary position according to the percent of pixel area that is inside the boundary. In Fig. 4-43, the pixel at position (x, y) has about half its area inside the polygon boundary. Therefore, the intensity at that position would be adjusted to one-half its assigned value. At the next position $(x + 1, y + 1)$ along the boundary, the intensity is adjusted to about one-third the assigned value for that point. Similar adjustments, based on the percent of pixel area coverage, are applied to the other intensity values around the boundary.

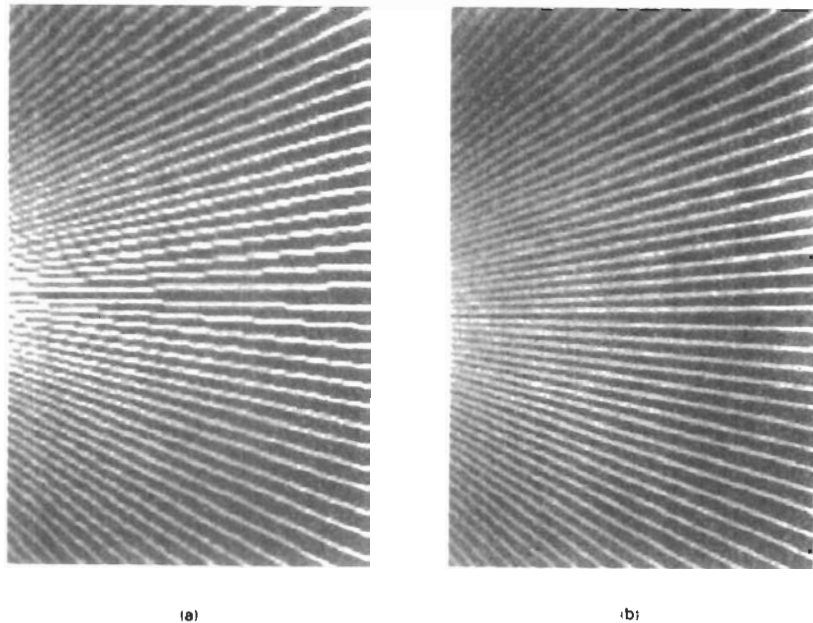


Figure 4-41

Jagged lines (a), plotted on the Merlin 9200 system, are smoothed (b) with an antialiasing technique called pixel phasing. This technique increases the number of addressable points on the system from 768×576 to 3072×2304 . (Courtesy of Megatek Corp.)

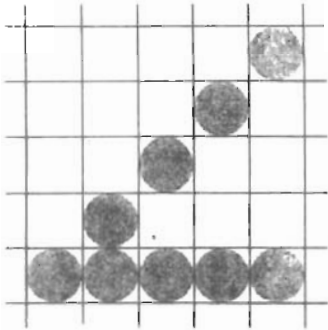


Figure 4-42
Unequal-length lines displayed
with the same number of pixels in
each line.

Supersampling methods can be applied by subdividing the total area and determining the number of subpixels inside the area boundary. A pixel partitioning into four subareas is shown in Fig. 4-44. The original 4 by 4 grid of pixels is turned into an 8 by 8 grid, and we now process eight scan lines across this grid instead of four. Figure 4-45 shows one of the pixel areas in this grid that overlaps an object boundary. Along the two scan lines we determine that three of the subpixel areas are inside the boundary. So we set the pixel intensity at 75 percent of its maximum value.

Another method for determining the percent of pixel area within a boundary, developed by Pitteway and Watkinson, is based on the midpoint line algorithm. This algorithm selects the next pixel along a line by determining which of two pixels is closer to the line by testing the location of the midpoint between the two pixels. As in the Bresenham algorithm, we set up a decision parameter p whose sign tells us which of the next two candidate pixels is closer to the line. By slightly modifying the form of p , we obtain a quantity that also gives the percent of the current pixel area that is covered by an object.

We first consider the method for a line with slope m in the range from 0 to 1. In Fig. 4-46, a straight line path is shown on a pixel grid. Assuming that the pixel at position (x_k, y_k) has been plotted, the next pixel nearest the line at $x = x_k + 1$ is either the pixel at y_k or the one at $y_k + 1$. We can determine which pixel is nearer with the calculation

$$y - y_{\text{mid}} = [m(x_k + 1) + b] - (y_k + 0.5) \quad (4-7)$$

This gives the vertical distance from the actual y coordinate on the line to the halfway point between pixels at position y_k and $y_k + 1$. If this difference calculation is negative, the pixel at y_k is closer to the line. If the difference is positive, the

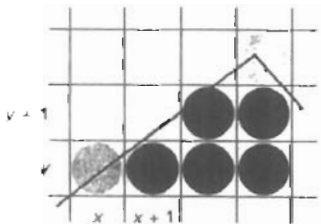


Figure 4-43
Adjusting pixel intensities along an
area boundary.

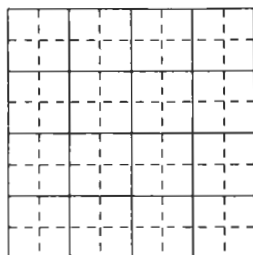


Figure 4-44
A 4 by 4 pixel section of a raster display subdivided into an 8 by 8 grid.

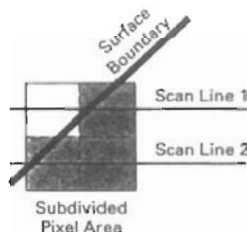


Figure 4-45
A subdivided pixel area with three subdivisions inside an object boundary line.

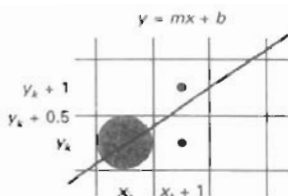


Figure 4-46
Boundary edge of an area passing through a pixel grid section.

pixel at $y_k + 1$ is closer. We can adjust this calculation so that it produces a positive number in the range from 0 to 1 by adding the quantity $1 - m$:

$$p = [m(x_k + 1) + b] - (y_k + 0.5) + (1 - m) \quad (4-8)$$

Now the pixel at y_k is nearer if $p < 1 - m$, and the pixel at $y_k + 1$ is nearer if $p > 1 - m$.

Parameter p also measures the amount of the current pixel that is overlapped by the area. For the pixel at (x_k, y_k) in Fig. 4-47, the interior part of the pixel has an area that can be calculated as

$$\text{area} = mx_k + b - y_k + 0.5 \quad (4-9)$$

This expression for the overlap area of the pixel at (x_k, y_k) is the same as that for parameter p in Eq. 4-8. Therefore, by evaluating p to determine the next pixel position along the polygon boundary, we also determine the percent of area coverage for the current pixel.

We can generalize this algorithm to accommodate lines with negative slopes and lines with slopes greater than 1. This calculation for parameter p could then be incorporated into a midpoint line algorithm to locate pixel positions and an object edge and to concurrently adjust pixel intensities along the boundary lines. Also, we can adjust the calculations to reference pixel coordinates at their lower left coordinates and maintain area proportions as discussed in Section 3-10.

At polygon vertices and for very skinny polygons, as shown in Fig. 4-48, we have more than one boundary edge passing through a pixel area. For these cases, we need to modify the Pitteway-Watkinson algorithm by processing all edges passing through a pixel and determining the correct interior area.

Filtering techniques discussed for line antialiasing can also be applied to area edges. Also, the various antialiasing methods can be applied to polygon areas or to regions with curved boundaries. Boundary equations are used to estimate area overlap of pixel regions with the area to be displayed. And coherence techniques are used along and between scan lines to simplify the calculations.

SUMMARY

In this chapter, we have explored the various attributes that control the appearance of displayed primitives. Procedures for displaying primitives use attribute settings to adjust the output of algorithms for line-generation, area-filling, and text-string displays.

The basic line attributes are line type, line color, and line width. Specifications for line type include solid, dashed, and dotted lines. Line-color specifications can be given in terms of RGB components, which control the intensity of the three electron guns in an RGB monitor. Specifications for line width are given in terms of multiples of a standard, one-pixel-wide line. These attributes can be applied to both straight lines and curves.

To reduce the size of the frame buffer, some raster systems use a separate color lookup table. This limits the number of colors that can be displayed to the size of the lookup table. Full-color systems are those that provide 24 bits per pixel and no separate color lookup table.

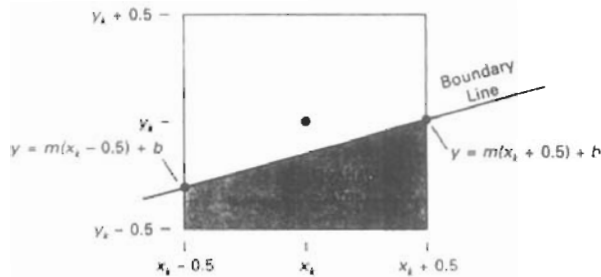


Figure 4-47

Overlap area of a pixel rectangle, centered at position (x_k, y_k) , with the interior of a polygon area.

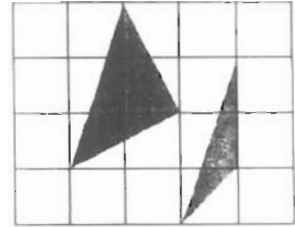


Figure 4-48

Polygons with more than one boundary line passing through individual pixel regions.

Fill-area attributes include the fill style and the fill color or the fill pattern. When the fill style is to be solid, the fill color specifies the color for the solid fill of the polygon interior. A hollow-fill style produces an interior in the background color and a border in the fill color. The third type of fill is patterned. In this case, a selected array pattern is used to fill the polygon interior.

An additional fill option provided in some packages is soft fill. This fill has applications in antialiasing and in painting packages. Soft-fill procedures provide a new fill color for a region that has the same variations as the previous fill color. One example of this approach is the linear soft-fill algorithm that assumes that the previous fill was a linear combination of foreground and background colors. This same linear relationship is then determined from the frame-buffer settings and used to repaint the area in a new color.

Characters, defined as pixel grid patterns or as outline fonts, can be displayed in different colors, sizes, and orientations. To set the orientation of a character string, we select a direction for the character up vector and a direction for the text path. In addition, we can set the alignment of a text string in relation to the start coordinate position. Marker symbols can be displayed using selected characters of various sizes and colors.

Graphics packages can be devised to handle both unbundled and bundled attribute specifications. Unbundled attributes are those that are defined for only one type of output device. Bundled attribute specifications allow different sets of attributes to be used on different devices, but accessed with the same index number in a bundle table. Bundle tables may be installation-defined, user-defined, or both. Functions to set the bundle table values specify workstation type and the attribute list for a given attribute index.

To determine current settings for attributes and other parameters, we can invoke inquiry functions. In addition to retrieving color and other attribute information, we can obtain workstation codes and status values with inquiry functions.

Because scan conversion is a digitizing process on raster systems, displayed primitives have a jagged appearance. This is due to the undersampling of information which rounds coordinate values to pixel positions. We can improve the appearance of raster primitives by applying antialiasing procedures that adjust pixel intensities. One method for doing this is to supersample. That is, we consider each pixel to be composed of subpixels and we calculate the intensity of the

subpixels and average the values of all subpixels. Alternatively, we can perform area sampling and determine the percentage of area coverage for a screen pixel, then set the pixel intensity proportional to this percentage. We can also weight the subpixel contributions according to position, giving higher weights to the central subpixels. Another method for antialiasing is to build special hardware configurations that can shift pixel positions.

Table 4-4 lists the attributes discussed in this chapter for the output primitive classifications: line, fill area, text, and marker. The attribute functions that can be used in graphics packages are listed for each category.

TABLE 4-4
SUMMARY OF ATTRIBUTES

<i>Output Primitive Type</i>	<i>Associated Attributes</i>	<i>Attribute-Setting Functions</i>	<i>Bundled- Attribute Functions</i>
Line	Type	setLinetype	setPolylineIndex
	Width	setLineWidthScaleFactor	setPolylineRepresentation
	Color	setPolylineColourIndex	
Fill Area	Fill Style	setInteriorStyle	setInteriorIndex
	Fill Color	setInteriorColorIndex	setInteriorRepresentation
	Pattern	setInteriorStyleIndex	
		setPatternRepresentation	
		setPatternSize	
Text	Font	setPatternReferencePoint	
		setTextFont	setTextIndex
		setTextColourIndex	setTextRepresentation
		setCharacterHeight	
		setCharacterExpansionFactor	
	Orientation	setCharacterUpVector	
		setTextPath	
Marker	Type	setTextAlignment	
		setMarkerType	setPolymarkerIndex
		setMarkerSizeScaleFactor	setPolymarkerRepresentation
		setPolymarkerColourIndex	

REFERENCES

Color and grayscale considerations are discussed in Crow (1978) and in Heckbert (1982). Soft-fill techniques are given in Fishkin and Barsky (1984). Antialiasing techniques are discussed in Pitteway and Watkinson (1980), Crow (1981), Turkowski (1982), Korein and Badler (1983), and Kirk and Avro, Schilling, and Wu (1991). Attribute functions in PHIGS are discussed in Howard et al. (1991), Hopgood and Duce (1991), Gaskins (1992), and Blake (1993). For information on GKS workstations and attributes, see Hopgood et al. (1983) and Enderle, Kansy, and Pfaff (1984).

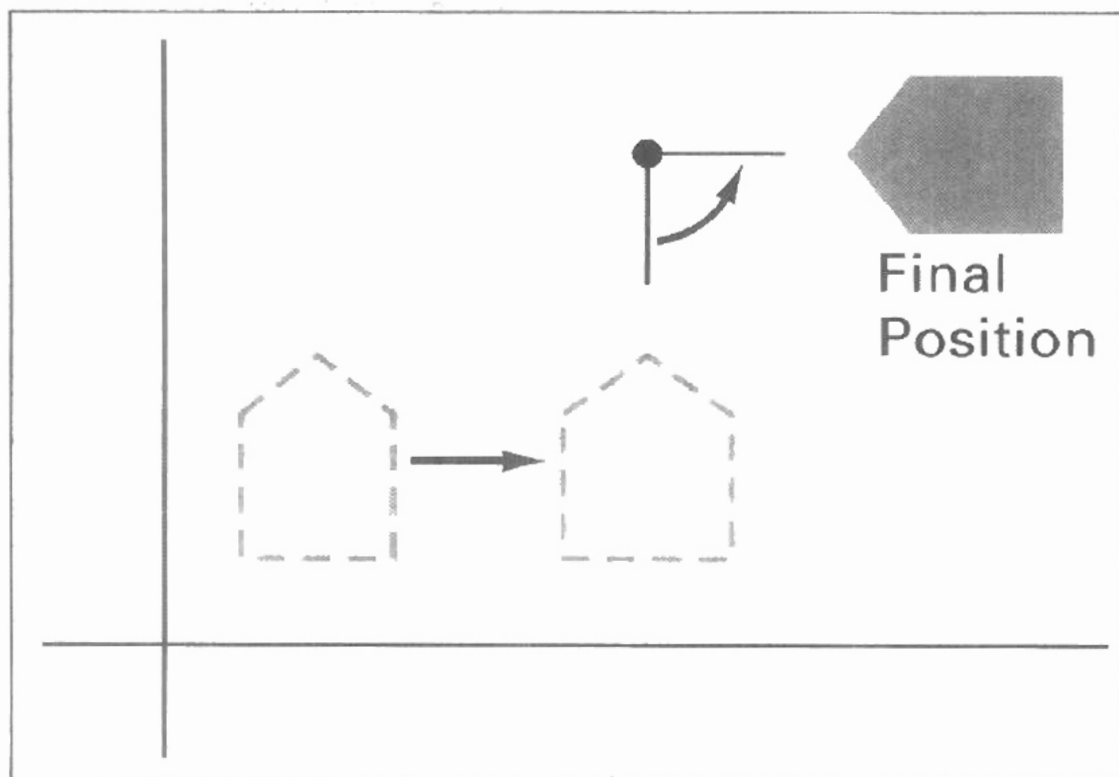
EXERCISES

- 4-1. Implement the line-type function by modifying Bresenham's line-drawing algorithm to display either solid, dashed, or dotted lines.

- 4-2. Implement the line-type function with a midpoint line algorithm to display either solid, dashed, or dotted lines.
- 4-3. Devise a parallel method for implementing the line-type function.
- 4-4. Devise a parallel method for implementing the line-width function.
- 4-5. A line specified by two endpoints and a width can be converted to a rectangular polygon with four vertices and then displayed using a scan-line method. Develop an efficient algorithm for computing the four vertices needed to define such a rectangle using the line endpoints and line width.
- 4-6. Implement the line-width function in a line-drawing program so that any one of three line widths can be displayed.
- 4-7. Write a program to output a line graph of three data sets defined over the same x coordinate range. Input to the program is to include the three sets of data values, labeling for the axes, and the coordinates for the display area on the screen. The data sets are to be scaled to fit the specified area, each plotted line is to be displayed in a different line type (solid, dashed, dotted), and the axes are to be labeled. (Instead of changing the line type, the three data sets can be plotted in different colors.)
- 4-8. Set up an algorithm for displaying thick lines with either butt caps, round caps, or projecting square caps. These options can be provided in an option menu.
- 4-9. Devise an algorithm for displaying thick polylines with either a miter join, a round join, or a bevel join. These options can be provided in an option menu.
- 4-10. Implement pen and brush menu options for a line-drawing procedure, including at least two options: round and square shapes.
- 4-11. Modify a line-drawing algorithm so that the intensity of the output line is set according to its slope. That is, by adjusting pixel intensities according to the value of the slope, all lines are displayed with the same intensity per unit length.
- 4-12. Define and implement a function for controlling the line type (solid, dashed, dotted) of displayed ellipses.
- 4-13. Define and implement a function for setting the width of displayed ellipses.
- 4-14. Write a routine to display a bar graph in any specified screen area. Input is to include the data set, labeling for the coordinate axes, and the coordinates for the screen area. The data set is to be scaled to fit the designated screen area, and the bars are to be displayed in designated colors or patterns.
- 4-15. Write a procedure to display two data sets defined over the same x-coordinate range, with the data values scaled to fit a specified region of the display screen. The bars for one of the data sets are to be displaced horizontally to produce an overlapping bar pattern for easy comparison of the two sets of data. Use a different color or a different fill pattern for the two sets of bars.
- 4-16. Devise an algorithm for implementing a color lookup table and the `setColourRepresentation` operation.
- 4-17. Suppose you have a system with an 8-inch by 10-inch video screen that can display 100 pixels per inch. If a color lookup table with 64 positions is used with this system, what is the smallest possible size (in bytes) for the frame buffer?
- 4-18. Consider an RGB raster system that has a 512-by-512 frame buffer with a 20 bits per pixel and a color lookup table with 24 bits per pixel. (a) How many distinct gray levels can be displayed with this system? (b) How many distinct colors (including gray levels) can be displayed? (c) How many colors can be displayed at any one time? (d) What is the total memory size? (e) Explain two methods for reducing memory size while maintaining the same color capabilities.
- 4-19. Modify the scan-line algorithm to apply any specified rectangular fill pattern to a polygon interior, starting from a designated pattern position.
- 4-20. Write a procedure to fill the interior of a given ellipse with a specified pattern.
- 4-21. Write a procedure to implement the `setPatternRepresentation` function.

- 4-22. Define and implement a procedure for changing the size of an existing rectangular fill pattern.
- 4-23. Write a procedure to implement a soft-fill algorithm. Carefully define what the soft-fill algorithm is to accomplish and how colors are to be combined.
- 4-24. Devise an algorithm for adjusting the height and width of characters defined as rectangular grid patterns
- 4-25. Implement routines for setting the character up vector and the text path for controlling the display of character strings.
- 4-26. Write a program to align text as specified by input values for the alignment parameters.
- 4-27. Develop procedures for implementing the marker attribute functions.
- 4-28. Compare attribute-implementation procedures needed by systems that employ bundled attributes to those needed by systems using unbundled attributes.
- 4-29. Develop procedures for storing and accessing attributes in unbundled system attribute tables. The procedures are to be designed to store designated attribute values in the system tables, to pass attributes to the appropriate output routines, and to pass attributes to memory locations specified in inquiry commands.
- 4-30. Set up the same procedures described in the previous exercise for bundled system attribute tables.
- 4-31. Implement an antialiasing procedure by extending Bresenham's line algorithm to adjust pixel intensities in the vicinity of a line path.
- 4-32. Implement an antialiasing procedure for the midpoint line algorithm.
- 4-33. Develop an algorithm for antialiasing elliptical boundaries.
- 4-34. Modify the scan-line algorithm for area fill to incorporate antialiasing. Use coherence techniques to reduce calculations on successive scan lines
- 4-35. Write a program to implement the Pitteway-Watkinson antialiasing algorithm as a scan-line procedure to fill a polygon interior. Use the routine `setPixel (x, y, intensity)` to load the intensity value into the frame buffer at location (x, y).

5

Two-Dimensional
Geometric Transformations

With the procedures for displaying output primitives and their attributes, we can create a variety of pictures and graphs. In many applications, there is also a need for altering or manipulating displays. Design applications and facility layouts are created by arranging the orientations and sizes of the component parts of the scene. And animations are produced by moving the “camera” or the objects in a scene along animation paths. Changes in orientation, size, and shape are accomplished with **geometric transformations** that alter the coordinate descriptions of objects. The basic geometric transformations are translation, rotation, and scaling. Other transformations that are often applied to objects include reflection and shear. We first discuss methods for performing geometric transformations and then consider how transformation functions can be incorporated into graphics packages.

5-1 BASIC TRANSFORMATIONS

Here, we first discuss general procedures for applying translation, rotation, and scaling parameters to reposition and resize two-dimensional objects. Then, in Section 5-2, we consider how transformation equations can be expressed in a more convenient matrix formulation that allows efficient combination of object transformations.

Translation

A **translation** is applied to an object by repositioning it along a straight-line path from one coordinate location to another. We translate a two-dimensional point by adding **translation distances**, t_x and t_y , to the original coordinate position (x, y) to move the point to a new position (x', y') (Fig. 5-1).

$$x' = x + t_x, \quad y' = y + t_y \quad (5-1)$$

The translation distance pair (t_x, t_y) is called a **translation vector** or **shift vector**.

We can express the translation equations 5-1 as a single matrix equation by using column vectors to represent coordinate positions and the translation vector:

$$\mathbf{P} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \quad \mathbf{P}' = \begin{bmatrix} x'_1 \\ x'_2 \end{bmatrix}, \quad \mathbf{T} = \begin{bmatrix} t_x \\ t_y \end{bmatrix} \quad (5-2)$$

This allows us to write the two-dimensional translation equations in the matrix form:

$$\mathbf{P}' = \mathbf{P} + \mathbf{T} \quad (5-3)$$

Sometimes matrix-transformation equations are expressed in terms of coordinate row vectors instead of column vectors. In this case, we would write the matrix representations as $\mathbf{P} = [x \ y]$ and $\mathbf{T} = [t_x \ t_y]$. Since the column-vector representation for a point is standard mathematical notation, and since many graphics packages, for example, GKS and PHIGS, also use the column-vector representation, we will follow this convention.

Translation is a *rigid-body transformation* that moves objects without deformation. That is, every point on the object is translated by the same amount. A straight line segment is translated by applying the transformation equation 5-3 to each of the line endpoints and redrawing the line between the new endpoint positions. Polygons are translated by adding the translation vector to the coordinate position of each vertex and regenerating the polygon using the new set of vertex coordinates and the current attribute settings. Figure 5-2 illustrates the application of a specified translation vector to move an object from one position to another.

Similar methods are used to translate curved objects. To change the position of a circle or ellipse, we translate the center coordinates and redraw the figure in the new location. We translate other curves (for example, splines) by displacing the coordinate positions defining the objects, then we reconstruct the curve paths using the translated coordinate points.

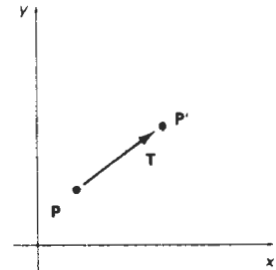


Figure 5-1
Translating a point from position \mathbf{P} to position \mathbf{P}' with translation vector \mathbf{T} .

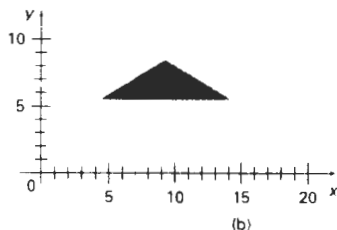
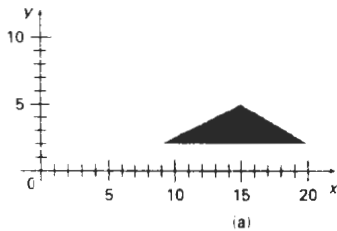


Figure 5-2
Moving a polygon from position (a) to position (b) with the translation vector $(-5.50, 3.75)$.

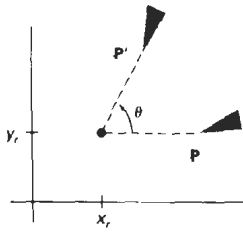


Figure 5-3
Rotation of an object through angle θ about the pivot point (x_r, y_r) .

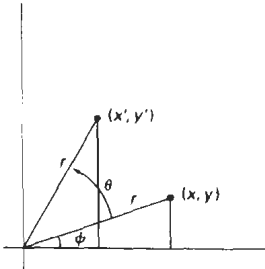


Figure 5-4
Rotation of a point from position (x, y) to position (x', y') through an angle θ relative to the coordinate origin. The original angular displacement of the point from the x axis is ϕ .

Rotation

A two-dimensional **rotation** is applied to an object by repositioning it along a circular path in the xy plane. To generate a rotation, we specify a **rotation angle** θ and the position (x_r, y_r) of the **rotation point** (or **pivot point**) about which the object is to be rotated (Fig. 5-3). Positive values for the rotation angle define counterclockwise rotations about the pivot point, as in Fig. 5-3, and negative values rotate objects in the clockwise direction. This transformation can also be described as a rotation about a **rotation axis** that is perpendicular to the xy plane and passes through the pivot point.

We first determine the transformation equations for rotation of a point position P when the pivot point is at the coordinate origin. The angular and coordinate relationships of the original and transformed point positions are shown in Fig. 5-4. In this figure, r is the constant distance of the point from the origin, angle ϕ is the original angular position of the point from the horizontal, and θ is the rotation angle. Using standard trigonometric identities, we can express the transformed coordinates in terms of angles θ and ϕ as

$$\begin{aligned} x' &= r \cos(\phi + \theta) = r \cos \phi \cos \theta - r \sin \phi \sin \theta \\ y' &= r \sin(\phi + \theta) = r \cos \phi \sin \theta + r \sin \phi \cos \theta \end{aligned} \quad (5-4)$$

The original coordinates of the point in polar coordinates are

$$x = r \cos \phi, \quad y = r \sin \phi \quad (5-5)$$

Substituting expressions 5-5 into 5-4, we obtain the transformation equations for rotating a point at position (x, y) through an angle θ about the origin:

$$\begin{aligned} x' &= x \cos \theta - y \sin \theta \\ y' &= x \sin \theta + y \cos \theta \end{aligned} \quad (5-6)$$

With the column-vector representations 5-2 for coordinate positions, we can write the rotation equations in the matrix form:

$$\mathbf{P}' = \mathbf{R} \cdot \mathbf{P} \quad (5-7)$$

where the rotation matrix is

$$\mathbf{R} = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix} \quad (5-8)$$

When coordinate positions are represented as row vectors instead of column vectors, the matrix product in rotation equation 5-7 is transposed so that the transformed row coordinate vector $[x' y']$ is calculated as

$$\begin{aligned} \mathbf{P}'^T &= (\mathbf{R} \cdot \mathbf{P})^T \\ &= \mathbf{P}^T \cdot \mathbf{R}^T \end{aligned}$$

where $\mathbf{P}^T = [x \ y]$, and the transpose \mathbf{R}^T of matrix \mathbf{R} is obtained by interchanging rows and columns. For a rotation matrix, the transpose is obtained by simply changing the sign of the sine terms.

Rotation of a point about an arbitrary pivot position is illustrated in Fig. 5-5. Using the trigonometric relationships in this figure, we can generalize Eqs. 5-6 to obtain the transformation equations for rotation of a point about any specified rotation position (x_r, y_r) :

$$\begin{aligned}x' &= x_r + (x - x_r) \cos \theta - (y - y_r) \sin \theta \\y' &= y_r + (x - x_r) \sin \theta + (y - y_r) \cos \theta\end{aligned}\quad (5-9)$$

These general rotation equations differ from Eqs. 5-6 by the inclusion of additive terms, as well as the multiplicative factors on the coordinate values. Thus, the matrix expression 5-7 could be modified to include pivot coordinates by matrix addition of a column vector whose elements contain the additive (translational) terms in Eqs. 5-9. There are better ways, however, to formulate such matrix equations, and we discuss in Section 5-2 a more consistent scheme for representing the transformation equations.

As with translations, rotations are rigid-body transformations that move objects without deformation. Every point on an object is rotated through the same angle. A straight line segment is rotated by applying the rotation equations 5-9 to each of the line endpoints and redrawing the line between the new endpoint positions. Polygons are rotated by displacing each vertex through the specified rotation angle and regenerating the polygon using the new vertices. Curved lines are rotated by repositioning the defining points and redrawing the curves. A circle or an ellipse, for instance, can be rotated about a noncentral axis by moving the center position through the arc that subtends the specified rotation angle. An ellipse can be rotated about its center coordinates by rotating the major and minor axes.

Scaling

A **scaling** transformation alters the size of an object. This operation can be carried out for polygons by multiplying the coordinate values (x, y) of each vertex by **scaling factors** s_x and s_y to produce the transformed coordinates (x', y') :

$$x' = x \cdot s_x, \quad y' = y \cdot s_y \quad (5-10)$$

Scaling factor s_x scales objects in the x direction, while s_y scales in the y direction. The transformation equations 5-10 can also be written in the matrix form:

$$\begin{bmatrix} x' \\ y' \end{bmatrix} = \begin{bmatrix} s_x & 0 \\ 0 & s_y \end{bmatrix} \cdot \begin{bmatrix} x \\ y \end{bmatrix} \quad (5-11)$$

or

$$P' = S \cdot P \quad (5-12)$$

where S is the 2 by 2 scaling matrix in Eq. 5-11.

Any positive numeric values can be assigned to the scaling factors s_x and s_y . Values less than 1 reduce the size of objects; values greater than 1 produce an enlargement. Specifying a value of 1 for both s_x and s_y leaves the size of objects unchanged. When s_x and s_y are assigned the same value, a **uniform scaling** is pro-

Section 5-1

Basic Transformations

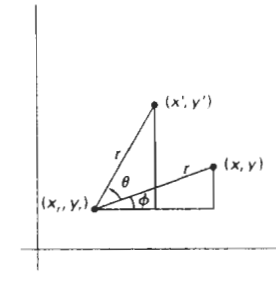


Figure 5-5
Rotating a point from position (x, y) to position (x', y') through an angle θ about rotation point (x_r, y_r) .

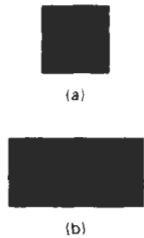


Figure 5-6
Turning a square (a) into a rectangle (b) with scaling factors $s_x = 2$ and $s_y = 1$.

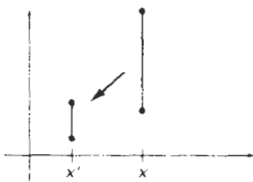


Figure 5-7
A line scaled with Eq. 5-12 using $s_x = s_y = 0.5$ is reduced in size and moved closer to the coordinate origin.

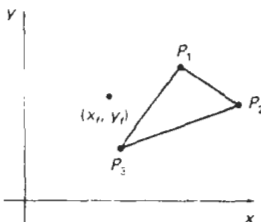


Figure 5-8
Scaling relative to a chosen fixed point (x_f, y_f) . Distances from each polygon vertex to the fixed point are scaled by transformation equations 5-13.

duced that maintains relative object proportions. Unequal values for s_x and s_y result in a **differential scaling** that is often used in design applications, where pictures are constructed from a few basic shapes that can be adjusted by scaling and positioning transformations (Fig. 5-6).

Objects transformed with Eq. 5-11 are both scaled and repositioned. Scaling factors with values less than 1 move objects closer to the coordinate origin, while values greater than 1 move coordinate positions farther from the origin. Figure 5-7 illustrates scaling a line by assigning the value 0.5 to both s_x and s_y in Eq. 5-11. Both the line length and the distance from the origin are reduced by a factor of 1/2.

We can control the location of a scaled object by choosing a position, called the **fixed point**, that is to remain unchanged after the scaling transformation. Coordinates for the fixed point (x_f, y_f) can be chosen as one of the vertices, the object centroid, or any other position (Fig. 5-8). A polygon is then scaled relative to the fixed point by scaling the distance from each vertex to the fixed point. For a vertex with coordinates (x, y) , the scaled coordinates (x', y') are calculated as

$$x' = x_f + (x - x_f)s_x, \quad y' = y_f + (y - y_f)s_y \quad (5-13)$$

We can rewrite these scaling transformations to separate the multiplicative and additive terms:

$$\begin{aligned} x' &= x \cdot s_x + x_f(1 - s_x) \\ y' &= y \cdot s_y + y_f(1 - s_y) \end{aligned} \quad (5-14)$$

where the additive terms $x_f(1 - s_x)$ and $y_f(1 - s_y)$ are constant for all points in the object.

Including coordinates for a fixed point in the scaling equations is similar to including coordinates for a pivot point in the rotation equations. We can set up a column vector whose elements are the constant terms in Eqs. 5-14, then we add this column vector to the product $S \cdot P$ in Eq. 5-12. In the next section, we discuss a matrix formulation for the transformation equations that involves only matrix multiplication.

Polygons are scaled by applying transformations 5-14 to each vertex and then regenerating the polygon using the transformed vertices. Other objects are scaled by applying the scaling transformation equations to the parameters defining the objects. An ellipse in standard position is resized by scaling the semimajor and semiminor axes and redrawing the ellipse about the designated center coordinates. Uniform scaling of a circle is done by simply adjusting the radius. Then we redisplay the circle about the center coordinates using the transformed radius.

5-2

MATRIX REPRESENTATIONS AND HOMOGENEOUS COORDINATES

Many graphics applications involve sequences of geometric transformations. An animation, for example, might require an object to be translated and rotated at each increment of the motion. In design and picture construction applications,

we perform translations, rotations, and scalings to fit the picture components into their proper positions. Here we consider how the matrix representations discussed in the previous sections can be reformulated so that such transformation sequences can be efficiently processed.

We have seen in Section 5-1 that each of the basic transformations can be expressed in the general matrix form

$$\mathbf{P}' = \mathbf{M}_1 \cdot \mathbf{P} + \mathbf{M}_2 \quad (5-15)$$

with coordinate positions \mathbf{P} and \mathbf{P}' represented as column vectors. Matrix \mathbf{M}_1 is a 2 by 2 array containing multiplicative factors, and \mathbf{M}_2 is a two-element column matrix containing translational terms. For translation, \mathbf{M}_1 is the identity matrix. For rotation or scaling, \mathbf{M}_2 contains the translational terms associated with the pivot point or scaling fixed point. To produce a sequence of transformations with these equations, such as scaling followed by rotation then translation, we must calculate the transformed coordinates one step at a time. First, coordinate positions are scaled, then these scaled coordinates are rotated, and finally the rotated coordinates are translated. A more efficient approach would be to combine the transformations so that the final coordinate positions are obtained directly from the initial coordinates, thereby eliminating the calculation of intermediate coordinate values. To be able to do this, we need to reformulate Eq. 5-15 to eliminate the matrix addition associated with the translation terms in \mathbf{M}_2 .

We can combine the multiplicative and translational terms for two-dimensional geometric transformations into a single matrix representation by expanding the 2 by 2 matrix representations to 3 by 3 matrices. This allows us to express all transformation equations as matrix multiplications, providing that we also expand the matrix representations for coordinate positions. To express any two-dimensional transformation as a matrix multiplication, we represent each Cartesian coordinate position (x, y) with the **homogeneous coordinate triple** (x_h, y_h, h) , where

$$x = \frac{x_h}{h}, \quad y = \frac{y_h}{h} \quad (5-16)$$

Thus, a general homogeneous coordinate representation can also be written as $(h \cdot x, h \cdot y, h)$. For two-dimensional geometric transformations, we can choose the homogeneous parameter h to be any nonzero value. Thus, there is an infinite number of equivalent homogeneous representations for each coordinate point (x, y) . A convenient choice is simply to set $h = 1$. Each two-dimensional position is then represented with homogeneous coordinates $(x, y, 1)$. Other values for parameter h are needed, for example, in matrix formulations of three-dimensional viewing transformations.

The term *homogeneous coordinates* is used in mathematics to refer to the effect of this representation on Cartesian equations. When a Cartesian point (x, y) is converted to a homogeneous representation (x_h, y_h, h) , equations containing x and y , such as $f(x, y) = 0$, become homogeneous equations in the three parameters x_h, y_h , and h . This just means that if each of the three parameters is replaced by any value v times that parameter, the value v can be factored out of the equations.

Expressing positions in homogeneous coordinates allows us to represent all geometric transformation equations as matrix multiplications. Coordinates are

represented with three-element column vectors, and transformation operations are written as 3 by 3 matrices. For translation, we have

$$\begin{bmatrix} x' \\ y' \\ 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 & t_x \\ 0 & 1 & t_y \\ 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} \quad (5-17)$$

which we can write in the abbreviated form

$$\mathbf{P}' = \mathbf{T}(t_x, t_y) \cdot \mathbf{P} \quad (5-18)$$

with $\mathbf{T}(t_x, t_y)$ as the 3 by 3 translation matrix in Eq. 5-17. The inverse of the translation matrix is obtained by replacing the translation parameters t_x and t_y with their negatives: $-t_x$ and $-t_y$.

Similarly, rotation transformation equations about the coordinate origin are now written as

$$\begin{bmatrix} x' \\ y' \\ 1 \end{bmatrix} = \begin{bmatrix} \cos\theta & -\sin\theta & 0 \\ \sin\theta & \cos\theta & 0 \\ 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} \quad (5-19)$$

or as

$$\mathbf{P}' = \mathbf{R}(\theta) \cdot \mathbf{P} \quad (5-20)$$

The rotation transformation operator $\mathbf{R}(\theta)$ is the 3 by 3 matrix in Eq. 5-19 with rotation parameter θ . We get the inverse rotation matrix when θ is replaced with $-\theta$.

Finally, a scaling transformation relative to the coordinate origin is now expressed as the matrix multiplication

$$\begin{bmatrix} x' \\ y' \\ 1 \end{bmatrix} = \begin{bmatrix} s_x & 0 & 0 \\ 0 & s_y & 0 \\ 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} \quad (5-21)$$

or

$$\mathbf{P}' = \mathbf{S}(s_x, s_y) \cdot \mathbf{P} \quad (5-22)$$

where $\mathbf{S}(s_x, s_y)$ is the 3 by 3 matrix in Eq. 5-21 with parameters s_x and s_y . Replacing these parameters with their multiplicative inverses ($1/s_x$ and $1/s_y$) yields the inverse scaling matrix.

Matrix representations are standard methods for implementing transformations in graphics systems. In many systems, rotation and scaling functions produce transformations with respect to the coordinate origin, as in Eqs. 5-19 and 5-21. Rotations and scalings relative to other reference positions are then handled as a succession of transformation operations. An alternate approach in a graphics package is to provide parameters in the transformation functions for the scaling fixed-point coordinates and the pivot-point coordinates. General rotation and scaling matrices that include the pivot or fixed point are then set up directly without the need to invoke a succession of transformation functions.

With the matrix representations of the previous section, we can set up a matrix for any sequence of transformations as a **composite transformation matrix** by calculating the matrix product of the individual transformations. Forming products of transformation matrices is often referred to as a **concatenation**, or **composition**, of matrices. For column-matrix representation of coordinate positions, we form composite transformations by multiplying matrices in order from right to left. That is, each successive transformation matrix premultiplies the product of the preceding transformation matrices.

Translations

If two successive translation vectors (t_{x1}, t_{y1}) and (t_{x2}, t_{y2}) are applied to a coordinate position \mathbf{P} , the final transformed location \mathbf{P}' is calculated as

$$\begin{aligned}\mathbf{P}' &= T(t_{x2}, t_{y2}) \cdot \{T(t_{x1}, t_{y1}) \cdot \mathbf{P}\} \\ &= \{T(t_{x2}, t_{y2}) \cdot T(t_{x1}, t_{y1})\} \cdot \mathbf{P}\end{aligned}\quad (5-23)$$

where \mathbf{P} and \mathbf{P}' are represented as homogeneous-coordinate column vectors. We can verify this result by calculating the matrix product for the two associative groupings. Also, the composite transformation matrix for this sequence of translations is

$$\begin{bmatrix} 1 & 0 & t_{x2} \\ 0 & 1 & t_{y2} \\ 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} 1 & 0 & t_{x1} \\ 0 & 1 & t_{y1} \\ 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 & t_{x1} + t_{x2} \\ 0 & 1 & t_{y1} + t_{y2} \\ 0 & 0 & 1 \end{bmatrix}\quad (5-24)$$

or

$$T(t_{x2}, t_{y2}) \cdot T(t_{x1}, t_{y1}) = T(t_{x1} + t_{x2}, t_{y1} + t_{y2})\quad (5-25)$$

which demonstrates that two successive translations are additive.

Rotations

Two successive rotations applied to point \mathbf{P} produce the transformed position

$$\begin{aligned}\mathbf{P}' &= \mathbf{R}(\theta_2) \cdot \{\mathbf{R}(\theta_1) \cdot \mathbf{P}\} \\ &= \{\mathbf{R}(\theta_2) \cdot \mathbf{R}(\theta_1)\} \cdot \mathbf{P}\end{aligned}\quad (5-26)$$

By multiplying the two rotation matrices, we can verify that two successive rotations are additive:

$$\mathbf{R}(\theta_2) \cdot \mathbf{R}(\theta_1) = \mathbf{R}(\theta_1 + \theta_2)\quad (5-27)$$

so that the final rotated coordinates can be calculated with the composite rotation matrix as

$$\mathbf{P}' = \mathbf{R}(\theta_1 + \theta_2) \cdot \mathbf{P}\quad (5-28)$$

Scalings

Concatenating transformation matrices for two successive scaling operations produces the following composite scaling matrix:

$$\begin{bmatrix} s_{x2} & 0 & 0 \\ 0 & s_{y2} & 0 \\ 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} s_{x1} & 0 & 0 \\ 0 & s_{y1} & 0 \\ 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} s_{x1} \cdot s_{x2} & 0 & 0 \\ 0 & s_{y1} \cdot s_{y2} & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad (5-29)$$

or

$$S(s_{x2}, s_{y2}) \cdot S(s_{x1}, s_{y1}) = S(s_{x1} \cdot s_{x2}, s_{y1} \cdot s_{y2}) \quad (5-30)$$

The resulting matrix in this case indicates that successive scaling operations are multiplicative. That is, if we were to triple the size of an object twice in succession, the final size would be nine times that of the original.

General Pivot-Point Rotation

With a graphics package that only provides a rotate function for revolving objects about the coordinate origin, we can generate rotations about any selected pivot point (x_p, y_p) by performing the following sequence of translate-rotate-translate operations:

1. Translate the object so that the pivot-point position is moved to the coordinate origin.
2. Rotate the object about the coordinate origin.
3. Translate the object so that the pivot point is returned to its original position.

This transformation sequence is illustrated in Fig. 5-9. The composite transforma-

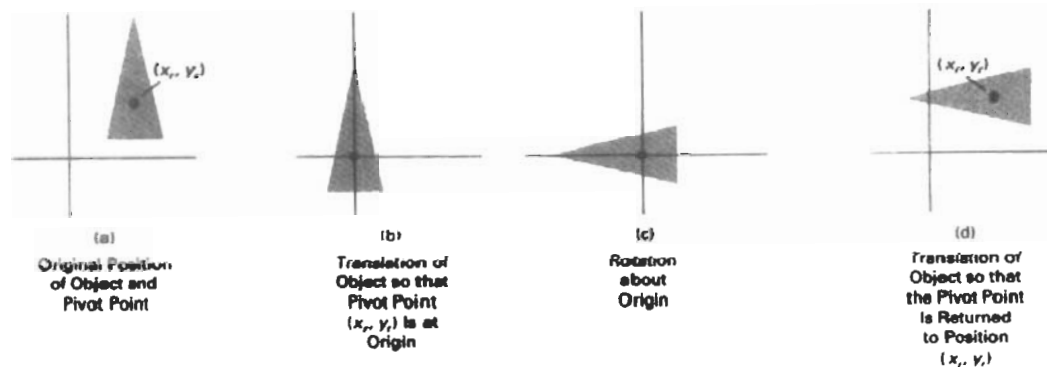


Figure 5-9

A transformation sequence for rotating an object about a specified pivot point using the rotation matrix $R(\theta)$ of transformation 5-19.

tion matrix for this sequence is obtained with the concatenation

$$\begin{bmatrix} 1 & 0 & x_r \\ 0 & 1 & y_r \\ 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} \cos \theta & -\sin \theta & 0 \\ \sin \theta & \cos \theta & 0 \\ 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} 1 & 0 & -x_r \\ 0 & 1 & -y_r \\ 0 & 0 & 1 \end{bmatrix} \\ = \begin{bmatrix} \cos \theta & -\sin \theta & x_r(1 - \cos \theta) + y_r \sin \theta \\ \sin \theta & \cos \theta & y_r(1 - \cos \theta) - x_r \sin \theta \\ 0 & 0 & 1 \end{bmatrix} \quad (5.31)$$

which can be expressed in the form

$$\mathbf{T}(x_r, y_r) \cdot \mathbf{R}(\theta) \cdot \mathbf{T}(-x_r, -y_r) = \mathbf{R}(x_r, y_r, \theta) \quad (5.32)$$

where $\mathbf{T}(-x_r, -y_r) = \mathbf{T}^{-1}(x_r, y_r)$. In general, a rotate function can be set up to accept parameters for pivot-point coordinates, as well as the rotation angle, and to generate automatically the rotation matrix of Eq. 5-31.

General Fixed-Point Scaling

Figure 5-10 illustrates a transformation sequence to produce scaling with respect to a selected fixed position (x_f, y_f) using a scaling function that can only scale relative to the coordinate origin.

1. Translate object so that the fixed point coincides with the coordinate origin.
2. Scale the object with respect to the coordinate origin.
3. Use the inverse translation of step 1 to return the object to its original position.

Concatenating the matrices for these three operations produces the required scaling matrix

$$\begin{bmatrix} 1 & 0 & x_f \\ 0 & 1 & y_f \\ 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} s_x & 0 & 0 \\ 0 & s_y & 0 \\ 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} 1 & 0 & -x_f \\ 0 & 1 & -y_f \\ 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} s_x & 0 & x_f(1 - s_x) \\ 0 & s_y & y_f(1 - s_y) \\ 0 & 0 & 1 \end{bmatrix} \quad (5.33)$$

or

$$\mathbf{T}(x_f, y_f) \cdot \mathbf{S}(s_x, s_y) \cdot \mathbf{T}(-x_f, -y_f) = \mathbf{S}(x_f, y_f, s_x, s_y) \quad (5.34)$$

This transformation is automatically generated on systems that provide a scale function that accepts coordinates for the fixed point.

General Scaling Directions

Parameters s_x and s_y scale objects along the x and y directions. We can scale an object in other directions by rotating the object to align the desired scaling directions with the coordinate axes before applying the scaling transformation.

Suppose we want to apply scaling factors with values specified by parameters s_1 and s_2 in the directions shown in Fig. 5-11. To accomplish the scaling with-

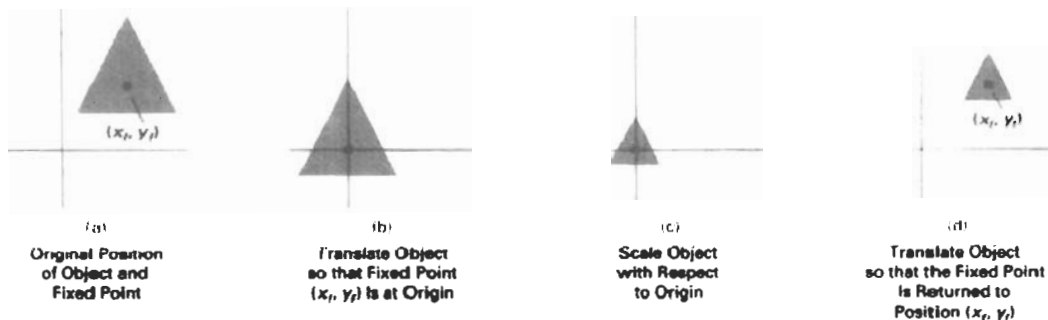


Figure 5-10

A transformation sequence for scaling an object with respect to a specified fixed position using the scaling matrix $S(s_x, s_y)$ of transformation 5-21.

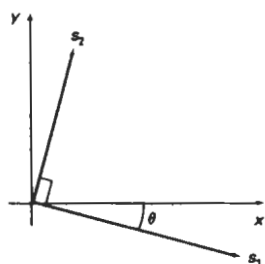


Figure 5-11

Scaling parameters s_1 and s_2 are to be applied in orthogonal directions defined by the angular displacement θ .

out changing the orientation of the object, we first perform a rotation so that the directions for s_1 and s_2 coincide with the x and y axes, respectively. Then the scaling transformation is applied, followed by an opposite rotation to return points to their original orientations. The composite matrix resulting from the product of these three transformations is

$$\begin{aligned} & \mathbf{R}^{-1}(\theta) \cdot \mathbf{S}(s_1, s_2) \cdot \mathbf{R}(\theta) \\ &= \begin{bmatrix} s_1 \cos^2 \theta + s_2 \sin^2 \theta & (s_2 - s_1) \cos \theta \sin \theta & 0 \\ (s_2 - s_1) \cos \theta \sin \theta & s_1 \sin^2 \theta + s_2 \cos^2 \theta & 0 \\ 0 & 0 & 1 \end{bmatrix} \end{aligned} \quad (5-35)$$

As an example of this scaling transformation, we turn a unit square into a parallelogram (Fig. 5-12) by stretching it along the diagonal from $(0, 0)$ to $(1, 1)$. We rotate the diagonal onto the y axis and double its length with the transformation parameters $\theta = 45^\circ$, $s_1 = 1$, and $s_2 = 2$.

In Eq. 5-35, we assumed that scaling was to be performed relative to the origin. We could take this scaling operation one step further and concatenate the matrix with translation operators, so that the composite matrix would include parameters for the specification of a scaling fixed position.

Concatenation Properties

Matrix multiplication is associative. For any three matrices, A , B , and C , the matrix product $A \cdot B \cdot C$ can be performed by first multiplying A and B or by first multiplying B and C :

$$A \cdot B \cdot C = (A \cdot B) \cdot C = A \cdot (B \cdot C) \quad (5-36)$$

Therefore, we can evaluate matrix products using either a left-to-right or a right-to-left associative grouping.

On the other hand, transformation products may not be commutative: The matrix product $A \cdot B$ is not equal to $B \cdot A$, in general. This means that if we want

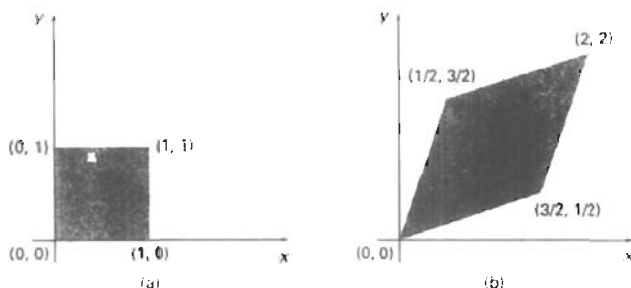


Figure 5-12
A square (a) is converted to a parallelogram (b) using the composite transformation matrix 5-35, with $s_1 = 1$, $s_2 = 2$, and $\theta = 45^\circ$.

to translate and rotate an object, we must be careful about the order in which the composite matrix is evaluated (Fig. 5-13). For some special cases, such as a sequence of transformations all of the same kind, the multiplication of transformation matrices is commutative. As an example, two successive rotations could be performed in either order and the final position would be the same. This commutative property holds also for two successive translations or two successive scalings. Another commutative pair of operations is rotation and uniform scaling ($s_x = s_y$).

General Composite Transformations and Computational Efficiency

A general two-dimensional transformation, representing a combination of translations, rotations, and scalings, can be expressed as

$$\begin{bmatrix} x' \\ y' \\ 1 \end{bmatrix} = \begin{bmatrix} rs_{xx} & rs_{xy} & trs_x \\ rs_{yx} & rs_{yy} & trs_y \\ 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} \quad (5-37)$$

The four elements rs_{ij} are the multiplicative rotation-scaling terms in the transformation that involve only rotation angles and scaling factors. Elements trs_x and trs_y are the translational terms containing combinations of translation distances, pivot-point and fixed-point coordinates, and rotation angles and scaling parameters. For example, if an object is to be scaled and rotated about its centroid coordinates (x_c, y_c) and then translated, the values for the elements of the composite transformation matrix are

$$\begin{aligned} & T(t_x, t_y) \cdot R(x_c, y_c, \theta) \cdot S(x_c, y_c, s_x, s_y) \\ &= \begin{bmatrix} s_x \cos \theta & -s_y \sin \theta & x_c(1 - s_x \cos \theta) + y_c s_y \sin \theta + t_x \\ s_x \sin \theta & s_y \cos \theta & y_c(1 - s_y \cos \theta) - x_c s_x \sin \theta + t_y \\ 0 & 0 & 1 \end{bmatrix} \end{aligned} \quad (5-38)$$

Although matrix equation 5-37 requires nine multiplications and six additions, the explicit calculations for the transformed coordinates are

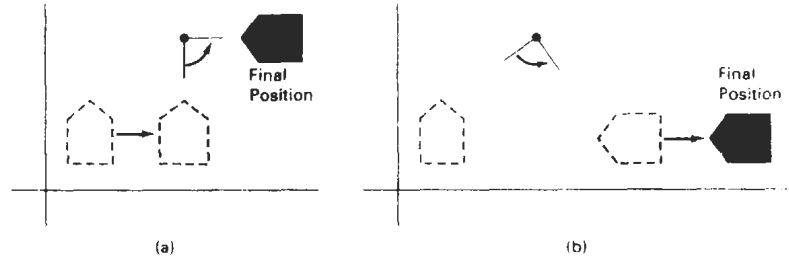


Figure 5-13

Reversing the order in which a sequence of transformations is performed may affect the transformed position of an object. In (a), an object is first translated, then rotated. In (b), the object is rotated first, then translated.

$$x' = x \cdot rs_{xx} + y \cdot rs_{xy} + trs_x, \quad y' = x \cdot rs_{yx} + y \cdot rs_{yy} + trs_y \quad (5.39)$$

Thus, we actually only need to perform four multiplications and four additions to transform coordinate positions. This is the maximum number of computations required for any transformation sequence, once the individual matrices have been concatenated and the elements of the composite matrix evaluated. Without concatenation, the individual transformations would be applied one at a time and the number of calculations could be significantly increased. An efficient implementation for the transformation operations, therefore, is to formulate transformation matrices, concatenate any transformation sequence, and calculate transformed coordinates using Eq. 5-39. On parallel systems, direct matrix multiplications with the composite transformation matrix of Eq. 5-37 can be equally efficient.

A general **rigid-body transformation matrix**, involving only translations and rotations, can be expressed in the form

$$\begin{bmatrix} r_{xx} & r_{xy} & tr_x \\ r_{yx} & r_{yy} & tr_y \\ 0 & 0 & 1 \end{bmatrix} \quad (5.40)$$

where the four elements r_{ij} are the multiplicative rotation terms, and elements tr_x and tr_y are the translational terms. A rigid-body change in coordinate position is also sometimes referred to as a **rigid-motion** transformation. All angles and distances between coordinate positions are unchanged by the transformation. In addition, matrix 5-40 has the property that its upper-left 2-by-2 submatrix is an orthogonal matrix. This means that if we consider each row of the submatrix as a vector, then the two vectors (r_{xx}, r_{xy}) and (r_{yx}, r_{yy}) form an orthogonal set of unit vectors: Each vector has unit length

$$r_{xx}^2 + r_{xy}^2 = r_{yx}^2 + r_{yy}^2 = 1 \quad (5.41)$$

and the vectors are perpendicular (their dot product is 0):

$$r_{xx}r_{yx} + r_{xy}r_{yy} = 0 \quad (5.42)$$

Therefore, if these unit vectors are transformed by the rotation submatrix, (r_{xx}, r_{xy}) is converted to a unit vector along the x axis and (r_{yx}, r_{yy}) is transformed into a unit vector along the y axis of the coordinate system:

$$\begin{bmatrix} r_{xx} & r_{xy} & 0 \\ r_{yx} & r_{yy} & 0 \\ 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} r_{xx} \\ r_{xy} \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} \quad (5-43)$$

$$\begin{bmatrix} r_{xx} & r_{xy} & 0 \\ r_{yx} & r_{yy} & 0 \\ 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} r_{yx} \\ r_{yy} \\ 1 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix} \quad (5-44)$$

As an example, the following rigid-body transformation first rotates an object through an angle θ about a pivot point (x_r, y_r) and then translates:

$$\begin{aligned} & \mathbf{T}(t_x, t_y) \cdot \mathbf{R}(x_r, y_r, \theta) \\ &= \begin{bmatrix} \cos \theta & -\sin \theta & x_r(1 - \cos \theta) + y_r \sin \theta + t_x \\ \sin \theta & \cos \theta & y_r(1 - \cos \theta) - x_r \sin \theta + t_y \\ 0 & 0 & 1 \end{bmatrix} \end{aligned} \quad (5-45)$$

Here, orthogonal unit vectors in the upper-left 2-by-2 submatrix are $(\cos \theta, -\sin \theta)$ and $(\sin \theta, \cos \theta)$, and

$$\begin{bmatrix} \cos \theta & -\sin \theta & 0 \\ \sin \theta & \cos \theta & 0 \\ 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} \cos \theta \\ -\sin \theta \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} \quad (5-46)$$

Similarly, unit vector $(\sin \theta, \cos \theta)$ is converted by the transformation matrix in Eq. 5-46 to the unit vector $(0, 1)$ in the y direction.

The orthogonal property of rotation matrices is useful for constructing a rotation matrix when we know the final orientation of an object rather than the amount of angular rotation necessary to put the object into that position. Directions for the desired orientation of an object could be determined by the alignment of certain objects in a scene or by selected positions in the scene. Figure 5-14 shows an object that is to be aligned with the unit direction vectors \mathbf{u}' and \mathbf{v}' . Assuming that the original object orientation, as shown in Fig. 5-14(a), is aligned with the coordinate axes, we construct the desired transformation by assigning the elements of \mathbf{u}' to the first row of the rotation matrix and the elements of \mathbf{v}' to the second row. This can be a convenient method for obtaining the transformation matrix for rotation within a local (or "object") coordinate system when we know the final orientation vectors. A similar transformation is the conversion of object descriptions from one coordinate system to another, and in Section 5-5, we consider how to set up transformations to accomplish this coordinate conversion.

Since rotation calculations require trigonometric evaluations and several multiplications for each transformed point, computational efficiency can become an important consideration in rotation transformations. In animations and other applications that involve many repeated transformations and small rotation angles, we can use approximations and iterative calculations to reduce computa-

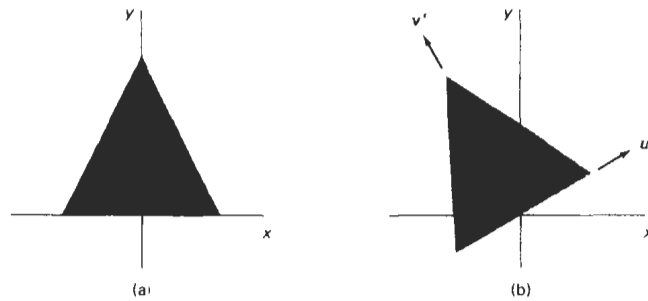


Figure 5-14

The rotation matrix for revolving an object from position (a) to position (b) can be constructed with the values of the unit orientation vectors u' and v' relative to the original orientation.

tions in the composite transformation equations. When the rotation angle is small, the trigonometric functions can be replaced with approximation values based on the first few terms of their power-series expansions. For small enough angles (less than 10°), $\cos \theta$ is approximately 1 and $\sin \theta$ has a value very close to the value of θ in radians. If we are rotating in small angular steps about the origin, for instance, we can set $\cos \theta$ to 1 and reduce transformation calculations at each step to two multiplications and two additions for each set of coordinates to be rotated:

$$x' = x - y \sin \theta, \quad y' = x \sin \theta + y \quad (5.47)$$

where $\sin \theta$ is evaluated once for all steps, assuming the rotation angle does not change. The error introduced by this approximation at each step decreases as the rotation angle decreases. But even with small rotation angles, the accumulated error over many steps can become quite large. We can control the accumulated error by estimating the error in x' and y' at each step and resetting object positions when the error accumulation becomes too great.

Composite transformations often involve inverse matrix calculations. Transformation sequences for general scaling directions and for reflections and shears (Section 5-4), for example, can be described with inverse rotation components. As we have noted, the inverse matrix representations for the basic geometric transformations can be generated with simple procedures. An inverse translation matrix is obtained by changing the signs of the translation distances, and an inverse rotation matrix is obtained by performing a matrix transpose (or changing the sign of the sine terms). These operations are much simpler than direct inverse matrix calculations.

An implementation of composite transformations is given in the following procedure. Matrix M is initialized to the identity matrix. As each individual transformation is specified, it is concatenated with the total transformation matrix M . When all transformations have been specified, this composite transformation is applied to a given object. For this example, a polygon is scaled and rotated about a given reference point. Then the object is translated. Figure 5-15 shows the original and final positions of the polygon transformed by this sequence.

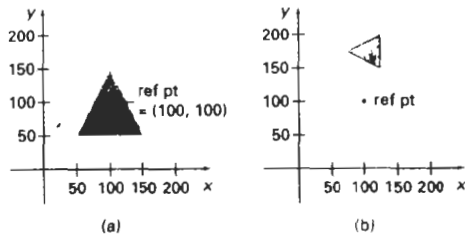


Figure 5-15

A polygon (a) is transformed into (b) by the composite operations in the following procedure.

```
#include <math.h>
#include "graphics.h"

typedef float Matrix3x3[3][3];
Matrix3x3 theMatrix;

void matrix3x3SetIdentity (Matrix3x3 m)
{
    int i,j;

    for (i=0; i<3; i++) for (j=0; j<3; j++) m[i][j] = (i == j);
}

/* Multiplies matrix a times b, putting result in b */
void matrix3x3PreMultiply (Matrix3x3 a, Matrix3x3 b)
{
    int r,c;
    Matrix3x3 tmp;

    for (r = 0; r < 3; r++)
        for (c = 0; c < 3; c++)
            tmp[r][c] =
                a[r][0]*b[0][c] + a[r][1]*b[1][c] + a[r][2]*b[2][c];

    for (r = 0; r < 3; r++)
        for (c = 0; c < 3; c++)
            b[r][c] = tmp[r][c];
}

void translate2 (int tx, int ty)
{
    Matrix3x3 m;

    matrix3x3SetIdentity (m);
    m[0][2] = tx;
    m[1][2] = ty;
    matrix3x3PreMultiply (m, theMatrix);
}
```

```

}

void scale2 (float sx, float sy, wcPt2 refPt)
{
    Matrix3x3 m;

    matrix3x3SetIdentity (m);
    m[0][0] = sx;
    m[0][2] = (1 - sx) * refPt.x;
    m[1][1] = sy;
    m[1][2] = (1 - sy) * refPt.y;
    matrix3x3PreMultiply (m, theMatrix);
}

void rotate2 (float a, wcPt2 refPt)
{
    Matrix3x3 m;

    matrix3x3SetIdentity (m);
    a = pToRadians (a);
    m[0][0] = cosf (a);
    m[0][1] = -sinf (a);
    m[0][2] = refPt.x * (1 - cosf (a)) + refPt.y * sinf (a);
    m[1][0] = sinf (a);
    m[1][1] = cosf (a);
    m[1][2] = refPt.y * (1 - cosf (a)) - refPt.x * sinf (a);
    matrix3x3PreMultiply (m, theMatrix);
}

void transformPoints2 (int npts, wcPt2 *pts)
{
    int k;
    float tmp;

    for (k = 0; k < npts; k++) {
        tmp = theMatrix[0][0] * pts[k].x + theMatrix[0][1] *
            pts[k].y + theMatrix[0][2];
        pts[k].y = theMatrix[1][0] * pts[k].x + theMatrix[1][1] *
            pts[k].y + theMatrix[1][2];
        pts[k].x = tmp;
    }
}

void main (int argc, char ** argv)
{
    wcPt2 pts[3] = { 50.0, 50.0, 150.0, 50.0, 100.0, 150.0};
    wcPt2 refPt = {100.0, 100.0};
    long windowID = openGraphics (*argv, 200, 350);

    setBackground (WHITE);
    setColor (BLUE);
    pFillArea (3, pts);
    matrix3x3SetIdentity (theMatrix);
    scale2 (0.5, 0.5, refPt);
    rotate2 (90.0, refPt);
    translate2 (0, 150);
    transformPoints2 (3, pts);
    pFillArea (3, pts);
    sleep (10);
    closeGraphics (windowID);
}

```

Basic transformations such as translation, rotation, and scaling are included in most graphics packages. Some packages provide a few additional transformations that are useful in certain applications. Two such transformations are reflection and shear.

Reflection

A **reflection** is a transformation that produces a mirror image of an object. The mirror image for a two-dimensional reflection is generated relative to an **axis of reflection** by rotating the object 180° about the reflection axis. We can choose an axis of reflection in the xy plane or perpendicular to the xy plane. When the reflection axis is a line in the xy plane, the rotation path about this axis is in a plane perpendicular to the xy plane. For reflection axes that are perpendicular to the xy plane, the rotation path is in the xy plane. Following are examples of some common reflections.

Reflection about the line $y = 0$, the x axis, is accomplished with the transformation matrix

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & -1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad (5-48)$$

This transformation keeps x values the same, but “flips” the y values of coordinate positions. The resulting orientation of an object after it has been reflected about the x axis is shown in Fig. 5-16. To envision the rotation transformation path for this reflection, we can think of the flat object moving out of the xy plane and rotating 180° through three-dimensional space about the x axis and back into the xy plane on the other side of the x axis.

A reflection about the y axis flips x coordinates while keeping y coordinates the same. The matrix for this transformation is

$$\begin{bmatrix} -1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad (5-49)$$

Figure 5-17 illustrates the change in position of an object that has been reflected about the line $x = 0$. The equivalent rotation in this case is 180° through three-dimensional space about the y axis.

We flip both the x and y coordinates of a point by reflecting relative to an axis that is perpendicular to the xy plane and that passes through the coordinate origin. This transformation, referred to as a reflection relative to the coordinate origin, has the matrix representation:

$$\begin{bmatrix} -1 & 0 & 0 \\ 0 & -1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad (5-50)$$

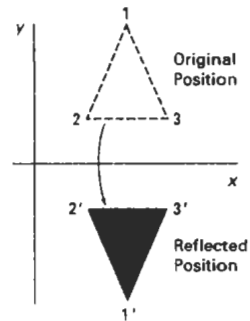


Figure 5-16
Reflection of an object about the x axis.

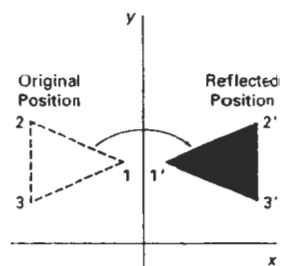


Figure 5-17
Reflection of an object about the y axis.

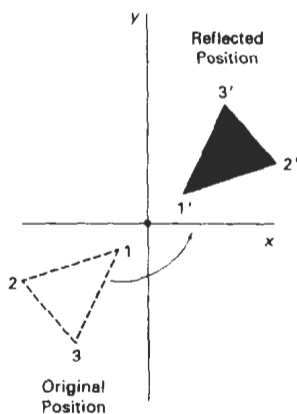


Figure 5-18
Reflection of an object relative to an axis perpendicular to the xy plane and passing through the coordinate origin.

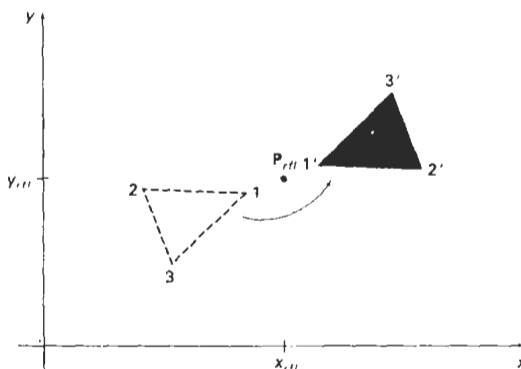


Figure 5-19
Reflection of an object relative to an axis perpendicular to the xy plane and passing through point P_{refl} .

An example of reflection about the origin is shown in Fig. 5-18. The reflection matrix 5-50 is the rotation matrix $R(\theta)$ with $\theta = 180^\circ$. We are simply rotating the object in the xy plane half a revolution about the origin.

Reflection 5-50 can be generalized to any reflection point in the xy plane (Fig. 5-19). This reflection is the same as a 180° rotation in the xy plane using the reflection point as the pivot point.

If we chose the reflection axis as the diagonal line $y = x$ (Fig. 5-20), the reflection matrix is

$$\begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad (5-51)$$

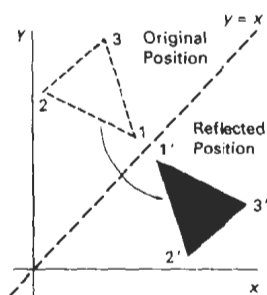


Figure 5-20
Reflection of an object with respect to the line $y = x$.

We can derive this matrix by concatenating a sequence of rotation and coordinate-axis reflection matrices. One possible sequence is shown in Fig. 5-21. Here, we first perform a clockwise rotation through a 45° angle, which rotates the line $y = x$ onto the x axis. Next, we perform a reflection with respect to the x axis. The final step is to rotate the line $y = x$ back to its original position with a counterclockwise rotation through 45° . An equivalent sequence of transformations is first to reflect the object about the x axis, and then to rotate counterclockwise 90° .

To obtain a transformation matrix for reflection about the diagonal $y = -x$, we could concatenate matrices for the transformation sequence: (1) clockwise rotation by 45° , (2) reflection about the y axis, and (3) counterclockwise rotation by 45° . The resulting transformation matrix is

$$\begin{bmatrix} 0 & -1 & 0 \\ -1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad (5-52)$$

Figure 5-22 shows the original and final positions for an object transformed with this reflection matrix.

Reflections about any line $y = mx + b$ in the xy plane can be accomplished with a combination of translate-rotate-reflect transformations. In general, we first translate the line so that it passes through the origin. Then we can rotate the line onto one of the coordinate axes and reflect about that axis. Finally, we restore the line to its original position with the inverse rotation and translation transformations.

We can implement reflections with respect to the coordinate axes or coordinate origin as scaling transformations with negative scaling factors. Also, elements of the reflection matrix can be set to values other than ± 1 . Values whose magnitudes are greater than 1 shift the mirror image farther from the reflection axis, and values with magnitudes less than 1 bring the mirror image closer to the reflection axis.

Shear

A transformation that distorts the shape of an object such that the transformed shape appears as if the object were composed of internal layers that had been caused to slide over each other is called a **shear**. Two common shearing transformations are those that shift coordinate x values and those that shift y values.

An x -direction shear relative to the x axis is produced with the transformation matrix

$$\begin{bmatrix} 1 & sh_x & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad (5-53)$$

which transforms coordinate positions as

$$x' = x + sh_x \cdot y, \quad y' = y \quad (5-54)$$

Any real number can be assigned to the shear parameter sh_x . A coordinate position (x, y) is then shifted horizontally by an amount proportional to its distance (y value) from the x axis ($y = 0$). Setting sh_x to 2, for example, changes the square in Fig. 5-23 into a parallelogram. Negative values for sh_x shift coordinate positions to the left.

We can generate x -direction shears relative to other reference lines with

$$\begin{bmatrix} 1 & sh_x & -sh_x \cdot y_{ref} \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad (5-55)$$

with coordinate positions transformed as

$$x' = x + sh_x(y - y_{ref}), \quad y' = y \quad (5-56)$$

An example of this shearing transformation is given in Fig. 5-24 for a shear parameter value of $1/2$ relative to the line $y_{ref} = -1$.

Section 5-4

Other Transformations

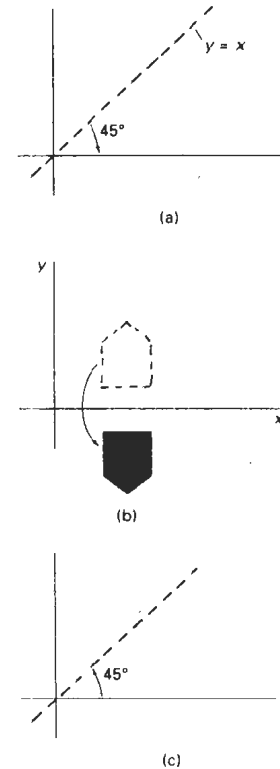


Figure 5-21
Sequence of transformations to produce reflection about the line $y = x$: (a) clockwise rotation of 45° , (b) reflection about the x axis, and (c) counterclockwise rotation by 45° .

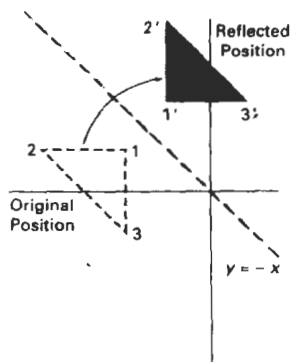


Figure 5-22
Reflection with respect to the line $y = -x$.

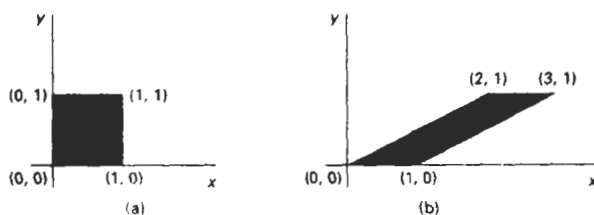


Figure 5-23
A unit square (a) is converted to a parallelogram (b) using the x -direction shear matrix 5-53 with $sh_x = 2$.

A y -direction shear relative to the line $x = x_{ref}$ is generated with the transformation matrix

$$\begin{bmatrix} 1 & 0 & 0 \\ sh_y & 1 & -sh_y \cdot x_{ref} \\ 0 & 0 & 1 \end{bmatrix} \quad (5-57)$$

which generates transformed coordinate positions

$$x' = x, \quad y' = sh_y(x - x_{ref}) + y \quad (5-58)$$

This transformation shifts a coordinate position vertically by an amount proportional to its distance from the reference line $x = x_{ref}$. Figure 5-25 illustrates the conversion of a square into a parallelogram with $sh_y = 1/2$ and $x_{ref} = -1$.

Shearing operations can be expressed as sequences of basic transformations. The x -direction shear matrix 5-53, for example, can be written as a composite transformation involving a series of rotation and scaling matrices that would scale the unit square of Fig. 5-23 along its diagonal, while maintaining the original lengths and orientations of edges parallel to the x axis. Shifts in the positions of objects relative to shearing reference lines are equivalent to translations.

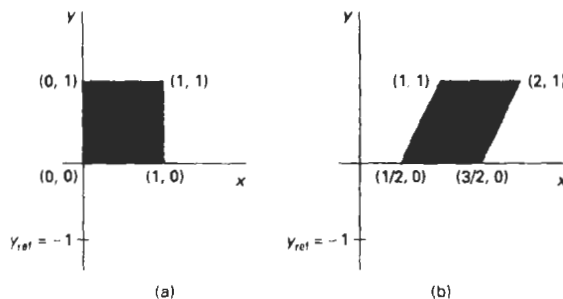
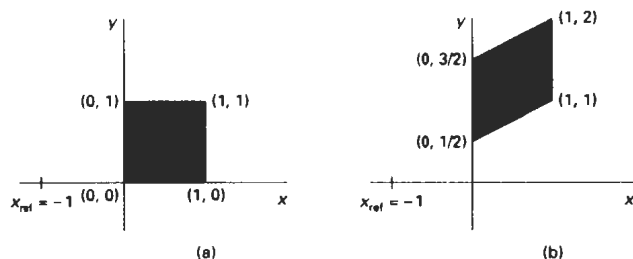


Figure 5-24
A unit square (a) is transformed to a shifted parallelogram (b) with $sh_y = 1/2$ and $y_{ref} = -1$ in the shear matrix 5-55.

**Figure 5-25**

A unit square (a) is turned into a shifted parallelogram (b) with parameter values $sh_y = 1/2$ and $x_{ref} = -1$ in the y -direction using shearing transformation 5-57.

5-5

TRANSFORMATIONS BETWEEN COORDINATE SYSTEMS

Graphics applications often require the transformation of object descriptions from one coordinate system to another. Sometimes objects are described in non-Cartesian reference frames that take advantage of object symmetries. Coordinate descriptions in these systems must then be converted to Cartesian device coordinates for display. Some examples of two-dimensional non-Cartesian systems are polar coordinates, elliptical coordinates, and parabolic coordinates. In other cases, we need to transform between two Cartesian systems. For modeling and design applications, individual objects may be defined in their own local Cartesian references, and the local coordinates must then be transformed to position the objects within the overall scene coordinate system. A facility management program for office layouts, for instance, has individual coordinate reference descriptions for chairs and tables and other furniture that can be placed into a floor plan, with multiple copies of the chairs and other items in different positions. In other applications, we may simply want to reorient the coordinate reference for displaying a scene. Relationships between Cartesian reference systems and some common non-Cartesian systems are given in Appendix A. Here, we consider transformations between two Cartesian frames of reference.

Figure 5-26 shows two Cartesian systems, with the coordinate origins at $(0, 0)$ and (x_0, y_0) and with an orientation angle θ between the x and x' axes. To transform object descriptions from xy coordinates to $x'y'$ coordinates, we need to set up a transformation that superimposes the $x'y'$ axes onto the xy axes. This is done in two steps:

1. Translate so that the origin (x_0, y_0) of the $x'y'$ system is moved to the origin of the xy system.
2. Rotate the x' axis onto the x axis.

Translation of the coordinate origin is expressed with the matrix operation

$$T(-x_0, -y_0) = \begin{bmatrix} 1 & 0 & -x_0 \\ 0 & 1 & -y_0 \\ 0 & 0 & 1 \end{bmatrix} \quad (5-59)$$

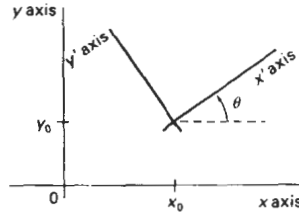


Figure 5-26
A Cartesian $x'y'$ system positioned at (x_0, y_0) with orientation θ in an xy Cartesian system.

and the orientation of the two systems after the translation operation would appear as in Fig. 5-27. To get the axes of the two systems into coincidence, we then perform the clockwise rotation

$$R(-\theta) = \begin{bmatrix} \cos \theta & \sin \theta & 0 \\ -\sin \theta & \cos \theta & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad (5-60)$$

Concatinating these two transformations matrices gives us the complete composite matrix for transforming object descriptions from the xy system to the $x'y'$ system:

$$\mathbf{M}_{xy,x'y'} = \mathbf{R}(-\theta) \cdot \mathbf{T}(-x_0, -y_0) \quad (5-61)$$

An alternate method for giving the orientation of the second coordinate system is to specify a vector \mathbf{V} that indicates the direction for the positive y' axis, as shown in Fig. 5-28. Vector \mathbf{V} is specified as a point in the xy reference frame relative to the origin of the xy system. A unit vector in the y' direction can then be obtained as

$$\mathbf{v} = \frac{\mathbf{V}}{|\mathbf{V}|} = (v_x, v_y) \quad (5-62)$$

And we obtain the unit vector \mathbf{u} along the x' axis by rotating \mathbf{v} 90° clockwise:

$$\begin{aligned} \mathbf{u} &= (v_y, -v_x) \\ &= (u_x, u_y) \end{aligned} \quad (5-63)$$

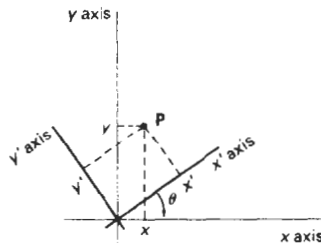


Figure 5-27
Position of the reference frames shown in Fig. 5-26 after translating the origin of the $x'y'$ system to the coordinate origin of the xy system.

Section 5-5

Transformations between Coordinate Systems

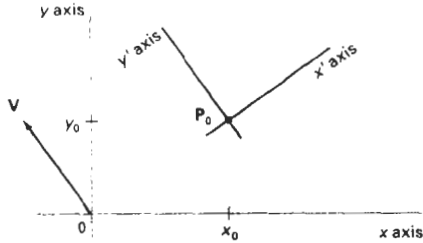


Figure 5-28
Cartesian system $x'y'$ with origin at $P_0 = (x_0, y_0)$ and y' axis parallel to vector V .

In Section 5-3, we noted that the elements of any rotation matrix could be expressed as elements of a set of orthogonal unit vectors. Therefore, the matrix to rotate the $x'y'$ system into coincidence with the xy system can be written as

$$R = \begin{bmatrix} u_x & u_y & 0 \\ v_x & v_y & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad (5-64)$$

As an example, suppose we choose the orientation for the y' axis as $V = (-1, 0)$, then the x' axis is in the positive y direction and the rotation transformation matrix is

$$\begin{bmatrix} 0 & 1 & 0 \\ -1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

Equivalently, we can obtain this rotation matrix from 5-60 by setting the orientation angle as $\theta = 90^\circ$.

In an interactive application, it may be more convenient to choose the direction for V relative to position P_0 than it is to specify it relative to the xy -coordinate origin. Unit vectors u and v would then be oriented as shown in Fig. 5-29. The components of v are now calculated as

$$v = \frac{P_1 - P_0}{|P_1 - P_0|} \quad (5-65)$$

and u is obtained as the perpendicular to v that forms a right-handed Cartesian system.

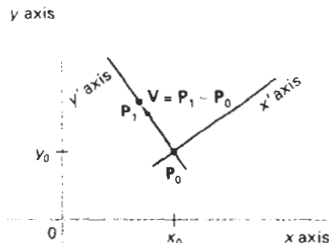


Figure 5-29
A Cartesian $x'y'$ system defined with two coordinate positions, P_0 and P_1 , within an xy reference frame.

5-6 AFFINE TRANSFORMATIONS

A coordinate transformation of the form

$$x' = a_{xx}x + a_{xy}y + b_x, \quad y' = a_{yx}x + a_{yy}y + b_y \quad (5-66)$$

is called a two-dimensional **affine transformation**. Each of the transformed coordinates x' and y' is a linear function of the original coordinates x and y , and parameters a_{ij} and b_k are constants determined by the transformation type. Affine transformations have the general properties that parallel lines are transformed into parallel lines and finite points map to finite points.

Translation, rotation, scaling, reflection, and shear are examples of two-dimensional affine transformations. Any general two-dimensional affine transformation can always be expressed as a composition of these five transformations. Another affine transformation is the conversion of coordinate descriptions from one reference system to another, which can be described as a combination of translation and rotation. An affine transformation involving only rotation, translation, and reflection preserves angles and lengths, as well as parallel lines. For these three transformations, the lengths and angle between two lines remains the same after the transformation.

5-7 TRANSFORMATION FUNCTIONS

Graphics packages can be structured so that separate commands are provided to a user for each of the basic transformation operations, as in procedure `transformObject`. A composite transformation is then set up by referencing individual functions in the order required for the transformation sequence. An alternate formulation is to provide users with a single transformation function that includes parameters for each of the basic transformations. The output of this function is the composite transformation matrix for the specified parameter values. Both options are useful. Separate functions are convenient for simple transformation operations, and a composite function can provide an expedient method for specifying complex transformation sequences.

The PHIGS library provides users with both options. Individual commands for generating the basic transformation matrices are

```
translate (translateVector, matrixTranslate)
rotate (theta, matrixRotate)
scale (scaleVector, matrixScale)
```

Each of these functions produces a 3 by 3 transformation matrix that can then be used to transform coordinate positions expressed as homogeneous column vectors. Parameter `translateVector` is a pointer to the pair of translation distances t_x and t_y . Similarly, parameter `scaleVector` specifies the pair of scaling values s_x and s_y . Rotate and scale matrices (`matrixTranslate` and `matrixScale`) transform with respect to the coordinate origin.

We concatenate transformation matrices that have been previously set up with the function

```
composeMatrix (matrix2, matrix1, matrixOut)
```

where elements of the composite output matrix are calculated by postmultiplying `matrix2` by `matrix1`. A composite transformation matrix to perform a combination scaling, rotation, and translation is produced with the function

```
buildTransformationMatrix (referencePoint, translateVector,
                           theta, scaleVector, matrix)
```

Rotation and scaling are carried out with respect to the coordinate position specified by parameter `referencePoint`. The order for the transformation sequence is assumed to be (1) scale, (2) rotate, and (3) translate, with the elements for the composite transformation stored in parameter `matrix`. We can use this function to generate a single transformation matrix or a composite matrix for two or three transformations (in the order stated). We could generate a translation matrix by setting `scaleVector` = (1, 1), `theta` = 0, and assigning `x` and `y` shift values to parameter `translateVector`. Any coordinate values could be assigned to parameter `referencePoint`, since the transformation calculations are unaffected by this parameter when no scaling or rotation takes place. But if we only want to set up a translation matrix, we can use function `translate` and simply specify the translation vector. A rotation or scaling transformation matrix is specified by setting `translateVector` = (0, 0) and assigning appropriate values to parameters `referencePoint`, `theta`, and `scaleVector`. To obtain a rotation matrix, we set `scaleVector` = (1, 1); and for scaling only, we set `theta` = 0. If we want to rotate or scale with respect to the coordinate origin, it is simpler to set up the matrix using either the `rotate` or `scale` function.

Since the function `buildTransformationMatrix` always generates the transformation sequence in the order (1) scale, (2) rotate, and (3) translate, the following function is provided to allow specification of other sequences:

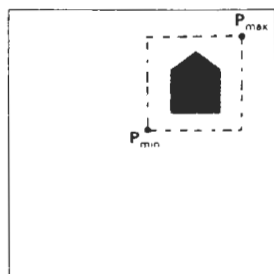
```
composeTransformationMatrix (matrixIn, referencePoint,
                             translateVector, theta, scaleVector, matrixOut)
```

We can use this function in combination with the `buildTransformationMatrix` function or with any of the other matrix-construction functions to compose any transformation sequence. For example, we could set up a scale matrix about a fixed point with the `buildTransformationMatrix` function, then we could use the `composeTransformationMatrix` function to concatenate this scale matrix with a rotation about a specified pivot point. The composite rotate-scale sequence is then stored in `matrixOut`.

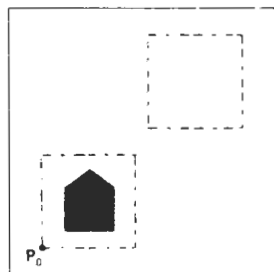
After we have set up a transformation matrix, we can apply the matrix to individual coordinate positions of an object with the function

```
transformPoint (inPoint, matrix, outPoint)
```

where parameter `inPoint` gives the initial `xy`-coordinate position of an object point, and parameter `outPoint` contains the corresponding transformed coordinates. Additional functions, discussed in Chapter 7, are available for performing two-dimensional modeling transformations.



(a)



(b)

Figure 5-30

Translating an object from screen position (a) to position (b) by moving a rectangular block of pixel values.

Coordinate positions P_{min} and P_{max} specify the limits of the rectangular block to be moved, and P_0 is the destination reference position.

The particular capabilities of raster systems suggest an alternate method for transforming objects. Raster systems store picture information as pixel patterns in the frame buffer. Therefore, some simple transformations can be carried out rapidly by simply moving rectangular arrays of stored pixel values from one location to another within the frame buffer. Few arithmetic operations are needed, so the pixel transformations are particularly efficient.

Raster functions that manipulate rectangular pixel arrays are generally referred to as **raster ops**. Moving a block of pixels from one location to another is also called a **block transfer** of pixel values. On a bilevel system, this operation is called a **bitBlt (bit-block transfer)**, particularly when the function is hardware implemented. The term **pixBlt** is sometimes used for block transfers on multi-level systems (multiple bits per pixel).

Figure 5-30 illustrates translation performed as a block transfer of a raster area. All bit settings in the rectangular area shown are copied as a block into another part of the raster. We accomplish this translation by first reading pixel intensities from a specified rectangular area of a raster into an array, then we copy the array back into the raster at the new location. The original object could be erased by filling its rectangular area with the background intensity (assuming the object does not overlap other objects in the scene).

Typical raster functions often provided in graphics packages are:

- *copy* - move a pixel block from one raster area to another.
- *read* - save a pixel block in a designated array.
- *write* - transfer a pixel array to a position in the frame buffer.

Some implementations provide options for combining pixel values. In *replace* mode, pixel values are simply transferred to the destination positions. Other options for combining pixel values include Boolean operations (*and*, *or*, and *exclusive or*) and binary arithmetic operations. With the *exclusive or* mode, two successive copies of a block to the same raster area restores the values that were originally present in that area. This technique can be used to move an object across a scene without destroying the background. Another option for adjusting pixel values is to combine the source pixels with a specified mask. This allows only selected positions within a block to be transferred or shaded by the patterns defined in the mask.

1	2	3
4	5	6
7	8	9
10	11	12

(a)

3	6	9	12
2	5	8	11
1	4	7	10

(b)

12	11	10
9	8	7
6	5	4
3	2	1

(c)

Figure 5-31

Rotating an array of pixel values. The original array orientation is shown in (a), the array orientation after a 90° counterclockwise rotation is shown in (b), and the array orientation after a 180° rotation is shown in (c).

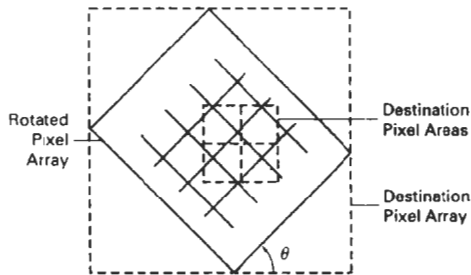


Figure 5-32

A raster rotation for a rectangular block of pixels is accomplished by mapping the destination pixel areas onto the rotated block.

Rotations in 90-degree increments are easily accomplished with block transfers. We can rotate an object 90° counterclockwise by first reversing the pixel values in each row of the array, then we interchange rows and columns. A 180° rotation is obtained by reversing the order of the elements in each row of the array, then reversing the order of the rows. Figure 5-31 demonstrates the array manipulations necessary to rotate a pixel block by 90° and by 180°.

For array rotations that are not multiples of 90°, we must perform more computations. The general procedure is illustrated in Fig. 5-32. Each destination pixel area is mapped onto the rotated array and the amount of overlap with the rotated pixel areas is calculated. An intensity for the destination pixel is then computed by averaging the intensities of the overlapped source pixels, weighted by their percentage of area overlap.

Raster scaling of a block of pixels is analogous to the cell-array mapping discussed in Section 3-13. We scale the pixel areas in the original block using specified values for s_x and s_y and map the scaled rectangle onto a set of destination pixels. The intensity of each destination pixel is then assigned according to its area of overlap with the scaled pixel areas (Fig. 5-33).

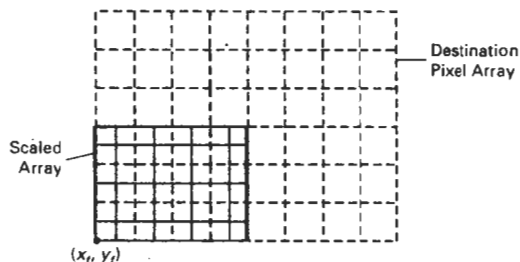


Figure 5-33

Mapping destination pixel areas onto a scaled array of pixel values. Scaling factors $s_x = s_y = 0.5$ are applied relative to fixed point (x_f, y_f) .

SUMMARY

The basic geometric transformations are translation, rotation, and scaling. Translation moves an object in a straight-line path from one position to another. Rotation moves an object from one position to another in a circular path around a specified pivot point (rotation point). Scaling changes the dimensions of an object relative to a specified fixed point.

We can express two-dimensional geometric transformations as 3 by 3 matrix operators, so that sequences of transformations can be concatenated into a single composite matrix. This is an efficient formulation, since it allows us to reduce computations by applying the composite matrix to the initial coordinate positions of an object to obtain the final transformed positions. To do this, we also need to express two-dimensional coordinate positions as three-element column or row matrices. We choose a column-matrix representation for coordinate points because this is the standard mathematical convention and because many graphics packages also follow this convention. For two-dimensional transformations, coordinate positions are then represented with three-element homogeneous coordinates with the third (homogeneous) coordinate assigned the value 1.

Composite transformations are formed as multiplications of any combination of translation, rotation, and scaling matrices. We can use combinations of translation and rotation for animation applications, and we can use combinations of rotation and scaling to scale objects in any specified direction. In general, matrix multiplications are not commutative. We obtain different results, for example, if we change the order of a translate-rotate sequence. A transformation sequence involving only translations and rotations is a rigid-body transformation, since angles and distances are unchanged. Also, the upper-left submatrix of a rigid-body transformation is an orthogonal matrix. Thus, rotation matrices can be formed by setting the upper-left 2-by-2 submatrix equal to the elements of two orthogonal unit vectors. Computations in rotational transformations can be reduced by using approximations for the sine and cosine functions when the rotation angle is small. Over many rotational steps, however, the approximation error can accumulate to a significant value.

Other transformations include reflections and shears. Reflections are transformations that rotate an object 180° about a reflection axis. This produces a mirror image of the object with respect to that axis. When the reflection axis is perpendicular to the xy plane, the reflection is obtained as a rotation in the xy plane. When the reflection axis is in the xy plane, the reflection is obtained as a rotation in a plane that is perpendicular to the xy plane. Shear transformations distort the shape of an object by shifting x or y coordinate values by an amount proportional to the coordinate distance from a shear reference line.

Transformations between Cartesian coordinate systems are accomplished with a sequence of translate-rotate transformations. One way to specify a new coordinate reference frame is to give the position of the new coordinate origin and the direction of the new y axis. The direction of the new x axis is then obtained by rotating the y direction vector 90° clockwise. Coordinate descriptions of objects in the old reference frame are transferred to the new reference with the transformation matrix that superimposes the new coordinate axes onto the old coordinate axes. This transformation matrix can be calculated as the concatenation of a translation that moves the new origin to the old coordinate origin and a rotation to align the two sets of axes. The rotation matrix is obtained from unit vectors in the x and y directions for the new system.

Two-dimensional geometric transformations are affine transformations. That is, they can be expressed as a linear function of coordinates x and y . Affine transformations transform parallel lines to parallel lines and transform finite points to finite points. Geometric transformations that do not involve scaling or shear also preserve angles and lengths.

Transformation functions in graphics packages are usually provided only for translation, rotation, and scaling. These functions include individual procedures for creating a translate, rotate, or scale matrix, and functions for generating a composite matrix given the parameters for a transformation sequence.

Fast raster transformations can be performed by moving blocks of pixels. This avoids calculating transformed coordinates for an object and applying scan-conversion routines to display the object at the new position. Three common raster operations (bitBlts or pixBlts) are copy, read, and write. When a block of pixels is moved to a new position in the frame buffer, we can simply replace the old pixel values or we can combine the pixel values using Boolean or arithmetic operations. Raster translations are carried out by copying a pixel block to a new location in the frame buffer. Raster rotations in multiples of 90° are obtained by manipulating row and column positions of the pixel values in a block. Other rotations are performed by first mapping rotated pixel areas onto destination positions in the frame buffer, then calculating overlap areas. Scaling in raster transformations is also accomplished by mapping transformed pixel areas to the frame-buffer destination positions.

REFERENCES

For additional information on homogeneous coordinates in computer graphics, see Blinn (1977 and 1978).

Transformation functions in PHIGS are discussed in Hopgood and Duce (1991), Howard et al. (1991), Gaskins (1992), and Blake (1993). For information on GKS transformation functions, see Hopgood et al. (1983) and Enderle, Kansy, and Pfaff (1984).

EXERCISES

- 5-1 Write a program to continuously rotate an object about a pivot point. Small angles are to be used for each successive rotation, and approximations to the sine and cosine functions are to be used to speed up the calculations. The rotation angle for each step is to be chosen so that the object makes one complete revolution in less than 30 seconds. To avoid accumulation of coordinate errors, reset the original coordinate values for the object at the start of each new revolution.
- 5-2 Show that the composition of two rotations is additive by concatenating the matrix representations for $R(\theta_1)$ and $R(\theta_2)$ to obtain

$$R(\theta_1) \cdot R(\theta_2) = R(\theta_1 + \theta_2)$$

- 5-3 Write a set of procedures to implement the `buildTransformationMatrix` and the `composeTransformationMatrix` functions to produce a composite transformation matrix for any set of input transformation parameters.
- 5-4 Write a program that applies any specified sequence of transformations to a displayed object. The program is to be designed so that a user selects the transformation sequence and associated parameters from displayed menus, and the composite transfor-

mation is then calculated and used to transform the object. Display the original object and the transformed object in different colors or different fill patterns.

- 5-5 Modify the transformation matrix (5-35), for scaling in an arbitrary direction, to include coordinates for any specified scaling fixed point (x_0, y_0) .
- 5-6 Prove that the multiplication of transformation matrices for each of the following sequence of operations is commutative:
 - (a) Two successive rotations.
 - (b) Two successive translations.
 - (c) Two successive scalings.
- 5-7 Prove that a uniform scaling ($s_x = s_y$) and a rotation form a commutative pair of operations but that, in general, scaling and rotation are not commutative operations.
- 5-8 Multiply the individual scale, rotate, and translate matrices in Eq. 5-38 to verify the elements in the composite transformation matrix.
- 5-9 Show that transformation matrix (5-51), for a reflection about the line $y = x$, is equivalent to a reflection relative to the x axis followed by a counterclockwise rotation of 90° .
- 5-10 Show that transformation matrix (5-52), for a reflection about the line $y = -x$, is equivalent to a reflection relative to the y axis followed by a counterclockwise rotation of 90° .
- 5-11 Show that two successive reflections about either of the coordinate axes is equivalent to a single rotation about the coordinate origin.
- 5-12 Determine the form of the transformation matrix for a reflection about an arbitrary line with equation $y = mx + b$.
- 5-13 Show that two successive reflections about any line passing through the coordinate origin is equivalent to a single rotation about the origin.
- 5-14 Determine a sequence of basic transformations that are equivalent to the x -direction shearing matrix (5-53).
- 5-15 Determine a sequence of basic transformations that are equivalent to the y -direction shearing matrix (5-57).
- 5-16 Set up a shearing procedure to display italic characters, given a vector font definition. That is, all character shapes in this font are defined with straight-line segments, and italic characters are formed with shearing transformations. Determine an appropriate value for the shear parameter by comparing italics and plain text in some available font. Define a simple vector font for input to your routine.
- 5-17 Derive the following equations for transforming a coordinate point $P = (x, y)$ in one Cartesian system to the coordinate values (x', y') in another Cartesian system that is rotated by an angle θ , as in Fig. 5-27. Project point P onto each of the four axes and analyse the resulting right triangles.

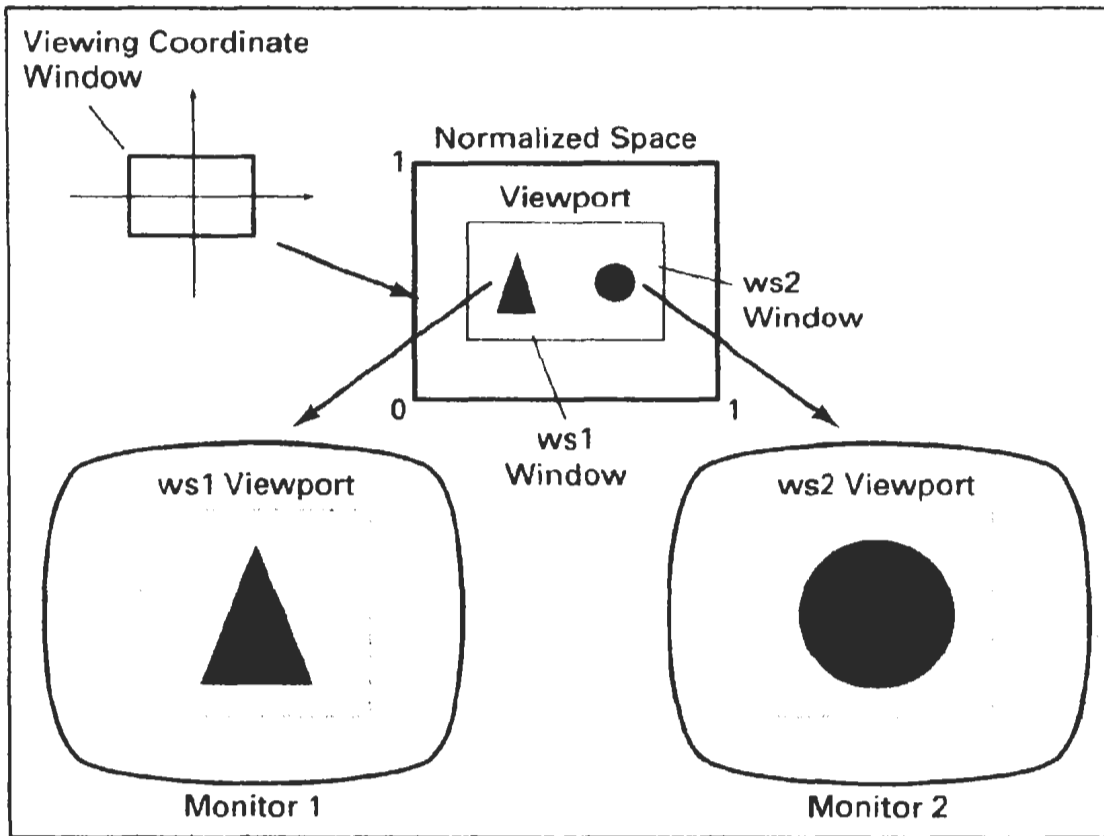
$$x' = x \cos \theta + y \sin \theta, \quad y' = -x \sin \theta + y \cos \theta$$

- 5-18 Write a procedure to compute the elements of the matrix for transforming object descriptions from one Cartesian coordinate system to another. The second coordinate system is to be defined with an origin point P_0 and a vector V that gives the direction for the positive y' axis of this system.
- 5-19 Set up procedures for implementing a block transfer of a rectangular area of a frame buffer, using one function to read the area into an array and another function to copy the array into the designated transfer area.
- 5-20 Determine the results of performing two successive block transfers into the same area of a frame buffer using the various Boolean operations.
- 5-21 What are the results of performing two successive block transfers into the same area of a frame buffer using the binary arithmetic operations?

- 5-22 Implement a routine to perform block transfers in a frame buffer using any specified Boolean operation or a replacement (copy) operation
- 5-23 Write a routine to implement rotations in increments of 90° in frame-buffer block transfers.
- 5-24 Write a routine to implement rotations by any specified angle in a frame-buffer block transfer.
- 5-25 Write a routine to implement scaling as a raster transformation of a pixel block.

Exercises

Two-Dimensional Viewing



We now consider the formal mechanism for displaying views of a picture on an output device. Typically, a graphics package allows a user to specify which part of a defined picture is to be displayed and where that part is to be placed on the display device. Any convenient Cartesian coordinate system, referred to as the world-coordinate reference frame, can be used to define the picture. For a two-dimensional picture, a view is selected by specifying a subarea of the total picture area. A user can select a single area for display, or several areas could be selected for simultaneous display or for an animated panning sequence across a scene. The picture parts within the selected areas are then mapped onto specified areas of the device coordinates. When multiple view areas are selected, these areas can be placed in separate display locations, or some areas could be inserted into other, larger display areas. Transformations from world to device coordinates involve translation, rotation, and scaling operations, as well as procedures for deleting those parts of the picture that are outside the limits of a selected display area.

6-1

THE VIEWING PIPELINE

A world-coordinate area selected for display is called a **window**. An area on a display device to which a window is mapped is called a **viewport**. The window defines *what* is to be viewed; the viewport defines *where* it is to be displayed. Often, windows and viewports are rectangles in standard position, with the rectangle edges parallel to the coordinate axes. Other window or viewport geometries, such as general polygon shapes and circles, are used in some applications, but these shapes take longer to process. In general, the mapping of a part of a world-coordinate scene to device coordinates is referred to as a **viewing transformation**. Sometimes the two-dimensional viewing transformation is simply referred to as the *window-to-viewport transformation* or the *windowing transformation*. But, in general, viewing involves more than just the transformation from the window to the viewport. Figure 6-1 illustrates the mapping of a picture section that falls within a rectangular window onto a designated rectangular viewport.

In computer graphics terminology, the term *window* originally referred to an area of a picture that is selected for viewing, as defined at the beginning of this section. Unfortunately, the same term is now used in window-manager systems to refer to any rectangular screen area that can be moved about, resized, and made active or inactive. In this chapter, we will only use the term window to

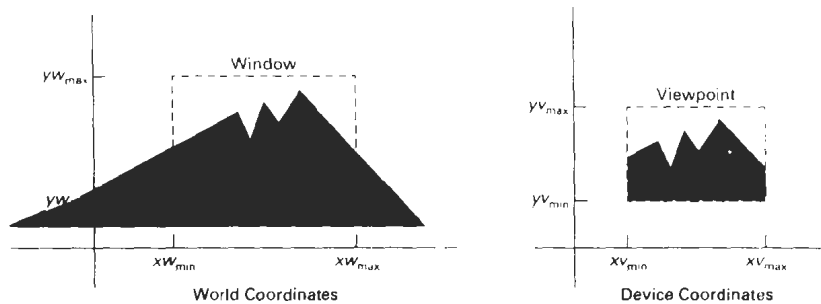


Figure 6-1

A viewing transformation using standard rectangles for the window and viewport.

refer to an area of a world-coordinate scene that has been selected for display. When we consider graphical user interfaces in Chapter 8, we will discuss screen windows and window-manager systems.

Some graphics packages that provide window and viewport operations allow only standard rectangles, but a more general approach is to allow the rectangular window to have any orientation. In this case, we carry out the viewing transformation in several steps, as indicated in Fig. 6-2. First, we construct the scene in world coordinates using the output primitives and attributes discussed in Chapters 3 and 4. Next, to obtain a particular orientation for the window, we can set up a two-dimensional **viewing-coordinate system** in the world-coordinate plane, and define a window in the viewing-coordinate system. The viewing-coordinate reference frame is used to provide a method for setting up arbitrary orientations for rectangular windows. Once the viewing reference frame is established, we can transform descriptions in world coordinates to viewing coordinates. We then define a viewport in normalized coordinates (in the range from 0 to 1) and map the viewing-coordinate description of the scene to normalized coordinates. At the final step, all parts of the picture that lie outside the viewport are clipped, and the contents of the viewport are transferred to device coordinates. Figure 6-3 illustrates a rotated viewing-coordinate reference frame and the mapping to normalized coordinates.

By changing the position of the viewport, we can view objects at different positions on the display area of an output device. Also, by varying the size of viewports, we can change the size and proportions of displayed objects. We achieve zooming effects by successively mapping different-sized windows on a

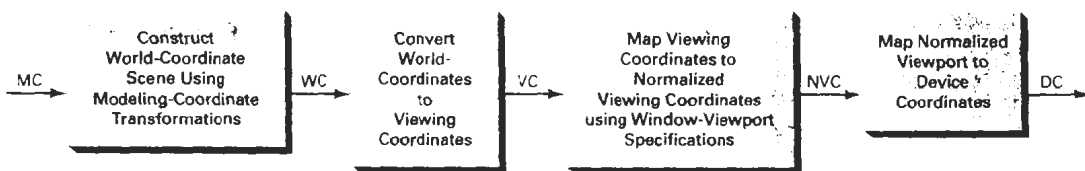


Figure 6-2

The two-dimensional viewing-transformation pipeline.

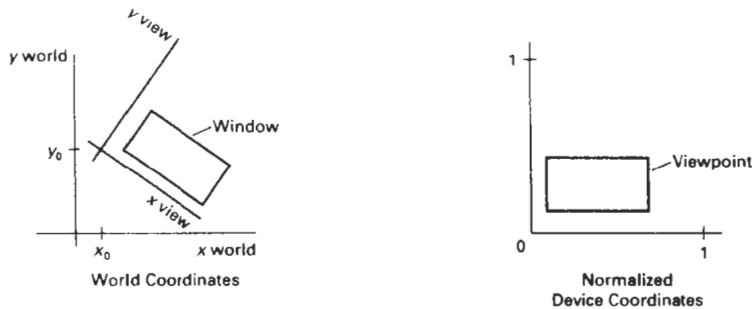


Figure 6-3

Setting up a rotated world window in viewing coordinates and the corresponding normalized-coordinate viewport.

fixed-size viewport. As the windows are made smaller, we zoom in on some part of a scene to view details that are not shown with larger windows. Similarly, more overview is obtained by zooming out from a section of a scene with successively larger windows. Panning effects are produced by moving a fixed-size window across the various objects in a scene.

Viewports are typically defined within the unit square (normalized coordinates). This provides a means for separating the viewing and other transformations from specific output-device requirements, so that the graphics package is largely device-independent. Once the scene has been transferred to normalized coordinates, the unit square is simply mapped to the display area for the particular output device in use at that time. Different output devices can be used by providing the appropriate device drivers.

When all coordinate transformations are completed, viewport clipping can be performed in normalized coordinates or in device coordinates. This allows us to reduce computations by concatenating the various transformation matrices. Clipping procedures are of fundamental importance in computer graphics. They are used not only in viewing transformations, but also in window-manager systems, in painting and drawing packages to eliminate parts of a picture inside or outside of a designated screen area, and in many other applications.

6-2

VIEWING COORDINATE REFERENCE FRAME

This coordinate system provides the reference frame for specifying the world-coordinate window. We set up the viewing coordinate system using the procedures discussed in Section 5-5. First, a viewing-coordinate origin is selected at some world position: $P_0 = (x_0, y_0)$. Then we need to establish the orientation, or rotation, of this reference frame. One way to do this is to specify a world vector \mathbf{V} that defines the viewing y_v direction. Vector \mathbf{V} is called the **view up vector**.

Given \mathbf{V} , we can calculate the components of unit vectors $\mathbf{v} = (v_x, v_y)$ and $\mathbf{u} = (u_x, u_y)$ for the viewing y_v and x_v axes, respectively. These unit vectors are used to form the first and second rows of the rotation matrix \mathbf{R} that aligns the viewing x_v, y_v axes with the world x_w, y_w axes.



Figure 6-4
A viewing-coordinate frame is moved into coincidence with the world frame in two steps: (a) translate the viewing origin to the world origin, then (b) rotate to align the axes of the two systems.

We obtain the matrix for converting world-coordinate positions to viewing coordinates as a two-step composite transformation: First, we translate the viewing origin to the world origin, then we rotate to align the two coordinate reference frames. The composite two-dimensional transformation to convert world coordinates to viewing coordinates is

$$M_{wc,vc} = R \cdot T \quad (6-1)$$

where T is the translation matrix that takes the viewing origin point P_0 to the world origin, and R is the rotation matrix that aligns the axes of the two reference frames. Figure 6-4 illustrates the steps in this coordinate transformation.

6-3

WINDOW-TO-VIEWPORT COORDINATE TRANSFORMATION

Once object descriptions have been transferred to the viewing reference frame, we choose the window extents in viewing coordinates and select the viewport limits in normalized coordinates (Fig. 6-3). Object descriptions are then transferred to normalized device coordinates. We do this using a transformation that maintains the same relative placement of objects in normalized space as they had in viewing coordinates. If a coordinate position is at the center of the viewing window, for instance, it will be displayed at the center of the viewport.

Figure 6-5 illustrates the window-to-viewport mapping. A point at position (xw, yw) in the window is mapped into position (xv, yv) in the associated viewport. To maintain the same relative placement in the viewport as in the window, we require that

$$\begin{aligned} \frac{xv - xv_{min}}{xv_{max} - xv_{min}} &= \frac{xw - xw_{min}}{xw_{max} - xw_{min}} \\ \frac{yv - yv_{min}}{yv_{max} - yv_{min}} &= \frac{yw - yw_{min}}{yw_{max} - yw_{min}} \end{aligned} \quad (6-2)$$

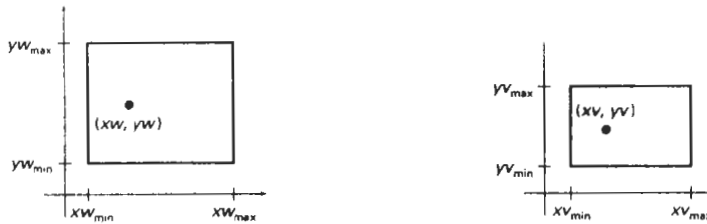


Figure 6-5

A point at position (xw, yw) in a designated window is mapped to viewport coordinates (xv, yv) so that relative positions in the two areas are the same.

Solving these expressions for the viewport position (xv, yv) , we have

$$\begin{aligned} xv &= xv_{\min} + (xw - xw_{\min})sx \\ yv &= yv_{\min} + (yw - yw_{\min})sy \end{aligned} \quad (6-3)$$

where the scaling factors are

$$\begin{aligned} sx &= \frac{xv_{\max} - xv_{\min}}{xw_{\max} - xw_{\min}} \\ sy &= \frac{yv_{\max} - yv_{\min}}{yw_{\max} - yw_{\min}} \end{aligned} \quad (6-4)$$

Equations 6-3 can also be derived with a set of transformations that converts the window area into the viewport area. This conversion is performed with the following sequence of transformations:

1. Perform a scaling transformation using a fixed-point position of (xw_{\min}, yw_{\min}) that scales the window area to the size of the viewport.
2. Translate the scaled window area to the position of the viewport.

Relative proportions of objects are maintained if the scaling factors are the same ($sx = sy$). Otherwise, world objects will be stretched or contracted in either the x or y direction when displayed on the output device.

Character strings can be handled in two ways when they are mapped to a viewport. The simplest mapping maintains a constant character size, even though the viewport area may be enlarged or reduced relative to the window. This method would be employed when text is formed with standard character fonts that cannot be changed. In systems that allow for changes in character size, string definitions can be windowed the same as other primitives. For characters formed with line segments, the mapping to the viewport can be carried out as a sequence of line transformations.

From normalized coordinates, object descriptions are mapped to the various display devices. Any number of output devices can be open in a particular application, and another window-to-viewport transformation can be performed for each open output device. This mapping, called the **workstation transforma-**

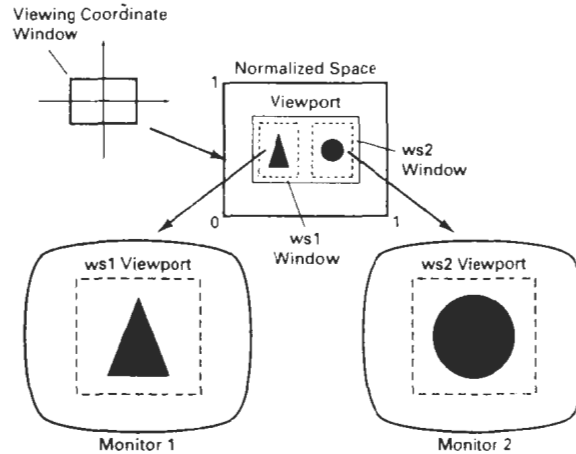


Figure 6-6
Mapping selected parts of a scene in normalized coordinates to different video monitors with workstation transformations.

tion, is accomplished by selecting a window area in normalized space and a viewport area in the coordinates of the display device. With the workstation transformation, we gain some additional control over the positioning of parts of a scene on individual output devices. As illustrated in Fig. 6-6, we can use workstation transformations to partition a view so that different parts of normalized space can be displayed on different output devices.

6-4

TWO-DIMENSIONAL VIEWING FUNCTIONS

We define a viewing reference system in a PHIGS application program with the following function:

```
evaluateViewOrientationMatrix (x0, y0, xV, yV,  
                             error, viewMatrix)
```

where parameters x_0 and y_0 are the coordinates of the viewing origin, and parameters x_V and y_V are the world-coordinate positions for the view up vector. An integer error code is generated if the input parameters are in error; otherwise, the `viewMatrix` for the world-to-viewing transformation is calculated. Any number of viewing transformation matrices can be defined in an application.

To set up the elements of a window-to-viewport mapping matrix, we invoke the function

```
evaluateViewMappingMatrix (xwmin, xwmax, ywmin, ywmax,  
                          xvmin, xvmax, yvmin, yvmax, error, viewMappingMatrix)
```

Here, the window limits in viewing coordinates are chosen with parameters $xwmin$, $xwmax$, $ywmin$, and $ywmax$; and the viewport limits are set with the nor-

malized coordinate positions *xvmin*, *xvmax*, *yvmin*, *yvmax*. As with the viewing-transformation matrix, we can construct several window-viewport pairs and use them for projecting various parts of the scene to different areas of the unit square.

Next, we can store combinations of viewing and window-viewport mappings for various workstations in a *viewing table* with

```
setViewRepresentation (ws, viewIndex, viewMatrix,
    viewMappingMatrix, xclipmin, xclipmax, yclipmin,
    yclipmax, clipxy)
```

where parameter *ws* designates the output device (workstation), and parameter *viewIndex* sets an integer identifier for this particular window-viewport pair. The matrices *viewMatrix* and *viewMappingMatrix* can be concatenated and referenced by the *viewIndex*. Additional clipping limits can also be specified here, but they are usually set to coincide with the viewport boundaries. And parameter *clipxy* is assigned either the value *noclip* or the value *clip*. This allows us to turn off clipping if we want to view the parts of the scene outside the viewport. We can also select *noclip* to speed up processing when we know that all of the scene is included within the viewport limits.

The function

```
setViewIndex (viewIndex)
```

selects a particular set of options from the viewing table. This view-index selection is then applied to subsequently specified output primitives and associated attributes and generates a display on each of the active workstations.

At the final stage, we apply a workstation transformation by selecting a workstation window-viewport pair:

```
setWorkstationWindow (ws, xwsWindmir, xwsWindmax,
    ywsWindmin, ywsWindmax)
setWorkstationViewport (ws, xwsVPortmin, xwsVPortmax,
    ywsVPortmin, ywsVPortmax)
```

where parameter *ws* gives the workstation number. Window-coordinate extents are specified in the range from 0 to 1 (normalized space), and viewport limits are in integer device coordinates.

If a workstation viewport is not specified, the unit square of the normalized reference frame is mapped onto the largest square area possible on an output device. The coordinate origin of normalized space is mapped to the origin of device coordinates, and the aspect ratio is retained by transforming the unit square onto a square area on the output device.

Example 6-1 Two-Dimensional Viewing Example

As an example of the use of viewing functions, the following sequence of statements sets up a rotated window in world coordinates and maps its contents to the upper right corner of workstation 2. We keep the viewing coordinate origin at the world origin, and we choose the view up direction for the window as (1, 1). This gives us a viewing-coordinate system that is rotated 45° clockwise in the world-coordinate reference frame. The view index is set to the value 5.

```
evaluateViewOrientationMatrix (0, 0, 1, 1,  
                               viewError, viewMat);  
evaluateViewMappingMatrix (-60.5, 41.24, -20.75, 82.5, 0.5,  
                           0.8, 0.7, 1.0, viewMapError, viewMapMat);  
setViewRepresentation (2, 5, viewMat, viewMapMat, 0.5, 0.8,  
                       0.7, 1.0, clip);  
setViewIndex (5);
```

Similarly, we could set up an additional transformation with view index 6 that would map a specified window into a viewport at the lower left of the screen. Two graphs, for example, could then be displayed at opposite screen corners with the following statements.

```
setViewIndex (5);  
polyline (3, axes);  
polyline (15, data1);  
setViewIndex (6);  
polyline (3, axes);  
polyline (25, data2);
```

View index 5 selects a viewport in the upper right of the screen display, and view index 6 selects a viewport in the lower left corner. The function `polyline (3, axes)` produces the horizontal and vertical coordinate reference for the data plot in each graph.

6-5

CLIPPING OPERATIONS

Generally, any procedure that identifies those portions of a picture that are either inside or outside of a specified region of space is referred to as a **clipping algorithm**, or simply **clipping**. The region against which an object is to be clipped is called a **clip window**.

Applications of clipping include extracting part of a defined scene for viewing; identifying visible surfaces in three-dimensional views; antialiasing line segments or object boundaries; creating objects using solid-modeling procedures; displaying a multiwindow environment; and drawing and painting operations that allow parts of a picture to be selected for copying, moving, erasing, or duplicating. Depending on the application, the clip window can be a general polygon or it can even have curved boundaries. We first consider clipping methods using rectangular clip regions, then we discuss methods for other clip-region shapes.

For the viewing transformation, we want to display only those picture parts that are within the window area (assuming that the clipping flags have not been set to `noclip`). Everything outside the window is discarded. Clipping algorithms can be applied in world coordinates, so that only the contents of the window interior are mapped to device coordinates. Alternatively, the complete world-coordinate picture can be mapped first to device coordinates, or normalized device coordinates, then clipped against the viewport boundaries. World-coordinate clipping removes those primitives outside the window from further consideration, thus eliminating the processing necessary to transform those primitives to device space. Viewport clipping, on the other hand, can reduce calculations by allowing concatenation of viewing and geometric transformation matrices. But

viewport clipping does require that the transformation to device coordinates be performed for all objects, including those outside the window area. On raster systems, clipping algorithms are often combined with scan conversion.

In the following sections, we consider algorithms for clipping the following primitive types

- Point Clipping
- Line Clipping (straight-line segments)
- Area Clipping (polygons)
- Curve Clipping
- Text Clipping

Line and polygon clipping routines are standard components of graphics packages, but many packages accommodate curved objects, particularly spline curves and conics, such as circles and ellipses. Another way to handle curved objects is to approximate them with straight-line segments and apply the line- or polygon-clipping procedure.

6-6

POINT CLIPPING

Assuming that the clip window is a rectangle in standard position, we save a point $P = (x, y)$ for display if the following inequalities are satisfied:

$$\begin{aligned} xw_{\min} &\leq x \leq xw_{\max} \\ yw_{\min} &\leq y \leq yw_{\max} \end{aligned} \quad (6-5)$$

where the edges of the clip window $(xw_{\min}, xw_{\max}, yw_{\min}, yw_{\max})$ can be either the world-coordinate window boundaries or viewport boundaries. If any one of these four inequalities is not satisfied, the point is clipped (not saved for display).

Although point clipping is applied less often than line or polygon clipping, some applications may require a point-clipping procedure. For example, point clipping can be applied to scenes involving explosions or sea foam that are modeled with particles (points) distributed in some region of the scene.

6-7

LINE CLIPPING

Figure 6-7 illustrates possible relationships between line positions and a standard rectangular clipping region. A line-clipping procedure involves several parts. First, we can test a given line segment to determine whether it lies completely inside the clipping window. If it does not, we try to determine whether it lies completely outside the window. Finally, if we cannot identify a line as completely inside or completely outside, we must perform intersection calculations with one or more clipping boundaries. We process lines through the “inside-outside” tests by checking the line endpoints. A line with both endpoints inside all clipping boundaries, such as the line from P_1 to P_2 , is saved. A line with both endpoints outside any one of the clip boundaries (line P_3P_4 in Fig. 6-7) is outside the win-

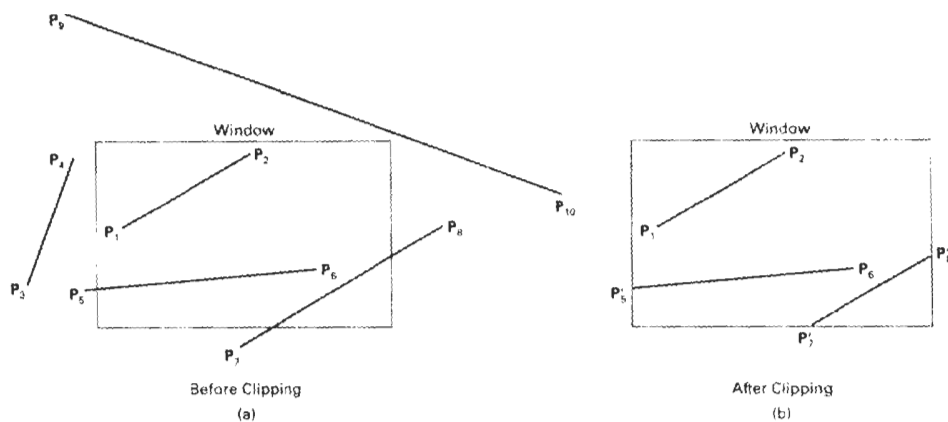


Figure 6-7
Line clipping against a rectangular clip window.

dow. All other lines cross one or more clipping boundaries, and may require calculation of multiple intersection points. To minimize calculations, we try to devise clipping algorithms that can efficiently identify outside lines and reduce intersection calculations.

For a line segment with endpoints (x_1, y_1) and (x_2, y_2) and one or both endpoints outside the clipping rectangle, the parametric representation

$$\begin{aligned} x &= x_1 + u(x_2 - x_1) \\ y &= y_1 + u(y_2 - y_1), \quad 0 \leq u \leq 1 \end{aligned} \quad (6-6)$$

could be used to determine values of parameter u for intersections with the clipping boundary coordinates. If the value of u for an intersection with a rectangle boundary edge is outside the range 0 to 1, the line does not enter the interior of the window at that boundary. If the value of u is within the range from 0 to 1, the line segment does indeed cross into the clipping area. This method can be applied to each clipping boundary edge in turn to determine whether any part of the line segment is to be displayed. Line segments that are parallel to window edges can be handled as special cases.

Clipping line segments with these parametric tests requires a good deal of computation, and faster approaches to clipping are possible. A number of efficient line clippers have been developed, and we survey the major algorithms in the next sections. Some algorithms are designed explicitly for two-dimensional pictures and some are easily adapted to three-dimensional applications.

Cohen-Sutherland Line Clipping

This is one of the oldest and most popular line-clipping procedures. Generally, the method speeds up the processing of line segments by performing initial tests that reduce the number of intersections that must be calculated. Every line end-

point in a picture is assigned a four-digit binary code, called a **region code**, that identifies the location of the point relative to the boundaries of the clipping rectangle. Regions are set up in reference to the boundaries as shown in Fig. 6-8. Each bit position in the region code is used to indicate one of the four relative coordinate positions of the point with respect to the clip window: to the left, right, top, or bottom. By numbering the bit positions in the region code as 1 through 4 from right to left, the coordinate regions can be correlated with the bit positions as

- bit 1: left
- bit 2: right
- bit 3: below
- bit 4: above

A value of 1 in any bit position indicates that the point is in that relative position; otherwise, the bit position is set to 0. If a point is within the clipping rectangle, the region code is 0000. A point that is below and to the left of the rectangle has a region code of 0101.

Bit values in the region code are determined by comparing endpoint coordinate values (x , y) to the clip boundaries. Bit 1 is set to 1 if $x < xw_{\min}$. The other three bit values can be determined using similar comparisons. For languages in which bit manipulation is possible, region-code bit values can be determined with the following two steps: (1) Calculate differences between endpoint coordinates and clipping boundaries. (2) Use the resultant sign bit of each difference calculation to set the corresponding value in the region code. Bit 1 is the sign bit of $x - xw_{\min}$; bit 2 is the sign bit of $xw_{\max} - x$; bit 3 is the sign bit of $y - yw_{\min}$; and bit 4 is the sign bit of $yw_{\max} - y$.

Once we have established region codes for all line endpoints, we can quickly determine which lines are completely inside the clip window and which are clearly outside. Any lines that are completely contained within the window boundaries have a region code of 0000 for both endpoints, and we trivially accept these lines. Any lines that have a 1 in the same bit position in the region codes for each endpoint are completely outside the clipping rectangle, and we trivially reject these lines. We would discard the line that has a region code of 1001 for one endpoint and a code of 0101 for the other endpoint. Both endpoints of this line are left of the clipping rectangle, as indicated by the 1 in the first bit position of each region code. A method that can be used to test lines for total clipping is to perform the logical *and* operation with both region codes. If the result is not 0000, the line is completely outside the clipping region.

Lines that cannot be identified as completely inside or completely outside a clip window by these tests are checked for intersection with the window boundaries. As shown in Fig. 6-9, such lines may or may not cross into the window interior. We begin the clipping process for a line by comparing an outside endpoint to a clipping boundary to determine how much of the line can be discarded. Then the remaining part of the line is checked against the other boundaries, and we continue until either the line is totally discarded or a section is found inside the window. We set up our algorithm to check line endpoints against clipping boundaries in the order left, right, bottom, top.

To illustrate the specific steps in clipping lines against rectangular boundaries using the Cohen-Sutherland algorithm, we show how the lines in Fig. 6-9 could be processed. Starting with the bottom endpoint of the line from P_1 to P_2 ,

Section 6-7

Line Clipping

1001	1000	1010
0001	0000	0010
0101	0100	0110

Figure 6-8

Binary region codes assigned to line endpoints according to relative position with respect to the clipping rectangle.

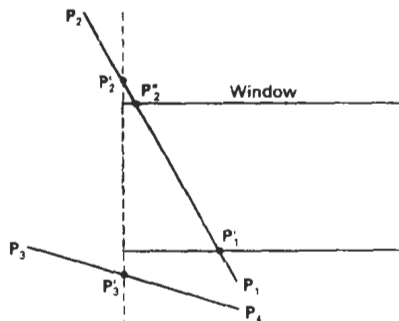


Figure 6-9
Lines extending from one coordinate region to another may pass through the clip window, or they may intersect clipping boundaries without entering the window.

we check P_1 against the left, right, and bottom boundaries in turn and find that this point is below the clipping rectangle. We then find the intersection point P'_1 with the bottom boundary and discard the line section from P_1 to P'_1 . The line now has been reduced to the section from P'_1 to P_2 . Since P_2 is outside the clip window, we check this endpoint against the boundaries and find that it is to the left of the window. Intersection point P'_2 is calculated, but this point is above the window. So the final intersection calculation yields P''_2 , and the line from P'_1 to P''_2 is saved. This completes processing for this line, so we save this part and go on to the next line. Point P_3 in the next line is to the left of the clipping rectangle, so we determine the intersection P'_3 and eliminate the line section from P_3 to P'_3 . By checking region codes for the line section from P'_3 to P_4 , we find that the remainder of the line is below the clip window and can be discarded also.

Intersection points with a clipping boundary can be calculated using the slope-intercept form of the line equation. For a line with endpoint coordinates (x_1, y_1) and (x_2, y_2) , the y coordinate of the intersection point with a vertical boundary can be obtained with the calculation

$$y = y_1 + m(x - x_1) \quad (6-7)$$

where the x value is set either to xw_{\min} or to xw_{\max} , and the slope of the line is calculated as $m = (y_2 - y_1)/(x_2 - x_1)$. Similarly, if we are looking for the intersection with a horizontal boundary, the x coordinate can be calculated as

$$x = x_1 + \frac{y - y_1}{m} \quad (6-8)$$

with y set either to yw_{\min} or to yw_{\max} .

The following procedure demonstrates the Cohen-Sutherland line-clipping algorithm. Codes for each endpoint are stored as bytes and processed using bit manipulations.

```
#define ROUND(a)    ((int)(a+0.5))

/* Bit masks encode a point's position relative to the clip edges. A
   point's status is encoded by OR'ing together appropriate bit masks.
*/
#define LEFT_EDGE    0x1
```

```

#define RIGHT_EDGE 0x2
#define BOTTOM_EDGE 0x4
#define TOP_EDGE 0x8

/* Points encoded as 0000 are completely Inside the clip rectangle;
   all others are outside at least one edge. If OR'ing two codes is
   FALSE (no bits are set in either code), the line can be Accepted. If
   the AND operation between two codes is TRUE, the line defined by those
   endpoints is completely outside the clip region and can be Rejected.
*/
#define INSIDE(a) (!a)
#define REJECT(a,b) (a&b)
#define ACCEPT(a,b) (!(a|b))

unsigned char encode (wcPt2 pt, dcPt winMin, dcPt winMax)
{
    unsigned char code=0x00;

    if (pt.x < winMin.x)
        code = code | LEFT_EDGE;
    if (pt.x > winMax.x)
        code = code | RIGHT_EDGE;
    if (pt.y < winMin.y)
        code = code | BOTTOM_EDGE;
    if (pt.y > winMax.y)
        code = code | TOP_EDGE;
    return (code);
}

void swapPts (wcPt2 * p1, wcPt2 * p2)
{
    wcPt2 tmp;

    tmp = *p1; *p1 = *p2; *p2 = tmp;
}

void swapCodes (unsigned char * c1, unsigned char * c2)
{
    unsigned char tmp;

    tmp = *c1; *c1 = *c2; *c2 = tmp;
}

void clipLine (dcPt winMin, dcPt winMax, wcPt2 p1, wcPt2 p2)
{
    unsigned char code1, code2;
    int done = FALSE, draw = FALSE;
    float m;

    while (!done) {
        code1 = encode (p1, winMin, winMax);
        code2 = encode (p2, winMin, winMax);
        if (ACCEPT (code1, code2)) {
            done = TRUE;
            draw = TRUE;
        }
        else
            if (REJECT (code1, code2))
                done = TRUE;
            else {
                /* Ensure that p1 is outside window */
                if (INSIDE (code1)) {

```

```

    swapPts (&p1, &p2);
    swapCodes (&code1, &code2);
}
/* Use slope (m) to find line-clipEdge intersections */
if (p2.x != p1.x)
    m = (p2.y - p1.y) / (p2.x - p1.x);
if (code1 & LEFT_EDGE) {
    p1.y += (winMin.x - p1.x) * m;
    p1.x = winMin.x;
}
else
    if (code1 & RIGHT_EDGE) {
        p1.y += (winMax.x - p1.x) * m;
        p1.x = winMax.x;
    }
    else
        if (code1 & BOTTOM_EDGE) {
            /* Need to update p1.x for non-vertical lines only */
            if (p2.x != p1.x)
                p1.x += (winMin.y - p1.y) / m;
            p1.y = winMin.y;
        }
        else
            if (code1 & TOP_EDGE) {
                if (p2.x != p1.x)
                    p1.x += (winMax.y - p1.y) / m;
                p1.y = winMax.y;
            }
    }
}
if (draw)
    lineEDA (ROUND(p1.x), ROUND(p1.y), ROUND(p2.x), ROUND(p2.y));
}

```

Liang-Barsky Line Clipping

Faster line clippers have been developed that are based on analysis of the parametric equation of a line segment, which we can write in the form

$$\begin{aligned} x &= x_1 + u\Delta x \\ y &= y_1 + u\Delta y, \quad 0 \leq u \leq 1 \end{aligned} \quad (6-9)$$

where $\Delta x = x_2 - x_1$ and $\Delta y = y_2 - y_1$. Using these parametric equations, Cyrus and Beck developed an algorithm that is generally more efficient than the Cohen-Sutherland algorithm. Later, Liang and Barsky independently devised an even faster parametric line-clipping algorithm. Following the Liang-Barsky approach, we first write the point-clipping conditions 6-5 in the parametric form:

$$\begin{aligned} xw_{\min} &\leq x_1 + u\Delta x \leq xw_{\max} \\ yw_{\min} &\leq y_1 + u\Delta y \leq yw_{\max} \end{aligned} \quad (6-10)$$

Each of these four inequalities can be expressed as

$$up_k \leq q_k, \quad k = 1, 2, 3, 4 \quad (6-11)$$

where parameters p and q are defined as

Section 6-7

Line Clipping

$$\begin{aligned} p_1 &= -\Delta x, & q_1 &= x_1 - xu_{\min} \\ p_2 &= \Delta x, & q_2 &= xu_{\max} - x_1 \\ p_3 &= -\Delta y, & q_3 &= y_1 - yu_{\min} \\ p_4 &= \Delta y, & q_4 &= yu_{\max} - y_1 \end{aligned} \quad (6-12)$$

Any line that is parallel to one of the clipping boundaries has $p_k = 0$ for the value of k corresponding to that boundary ($k = 1, 2, 3$, and 4 correspond to the left, right, bottom, and top boundaries, respectively). If, for that value of k , we also find $q_k < 0$, then the line is completely outside the boundary and can be eliminated from further consideration. If $q_k \geq 0$, the line is inside the parallel clipping boundary.

When $p_k < 0$, the infinite extension of the line proceeds from the outside to the inside of the infinite extension of this particular clipping boundary. If $p_k > 0$, the line proceeds from the inside to the outside. For a nonzero value of p_k , we can calculate the value of u that corresponds to the point where the infinitely extended line intersects the extension of boundary k as

$$u = \frac{q_k}{p_k} \quad (6-13)$$

For each line, we can calculate values for parameters u_1 and u_2 that define that part of the line that lies within the clip rectangle. The value of u_1 is determined by looking at the rectangle edges for which the line proceeds from the outside to the inside ($p < 0$). For these edges, we calculate $r_k = q_k/p_k$. The value of u_1 is taken as the largest of the set consisting of 0 and the various values of r . Conversely, the value of u_2 is determined by examining the boundaries for which the line proceeds from inside to outside ($p > 0$). A value of r_k is calculated for each of these boundaries, and the value of u_2 is the minimum of the set consisting of 1 and the calculated r values. If $u_1 > u_2$, the line is completely outside the clip window and it can be rejected. Otherwise, the endpoints of the clipped line are calculated from the two values of parameter u .

This algorithm is presented in the following procedure. Line intersection parameters are initialized to the values $u_1 = 0$ and $u_2 = 1$. For each clipping boundary, the appropriate values for p and q are calculated and used by the function *clipTest* to determine whether the line can be rejected or whether the intersection parameters are to be adjusted. When $p < 0$, the parameter r is used to update u_1 ; when $p > 0$, parameter r is used to update u_2 . If updating u_1 or u_2 results in $u_1 > u_2$, we reject the line. Otherwise, we update the appropriate u parameter only if the new value results in a shortening of the line. When $p = 0$ and $q < 0$, we can discard the line since it is parallel to and outside of this boundary. If the line has not been rejected after all four values of p and q have been tested, the endpoints of the clipped line are determined from values of u_1 and u_2 .

```
#include "graphics.h"

#define ROUND(a) ((int)(a+0.5))

int clipTest (float p, float q, float * u1, float * u2)
```

```

{
    float r;
    int retVal = TRUE;

    if (p < 0.0) {
        r = q / p;
        if (r > *u2)
            retVal = FALSE;
        else
            if (r > *u1)
                *u1 = r;
    }
    else
        if (p > 0.0) {
            r = q / p;
            if (r < *u1)
                retVal = FALSE;
            else if (r < *u2)
                *u2 = r;
        }
    else
        /* p = 0, so line is parallel to this clipping edge */
        if (q < 0.0)
            /* Line is outside clipping edge */
            retVal = FALSE;

    return (retVal);
}

void clipLine (dcPt winMin, dcPt winMax, wcPt2 p1, wcPt2 p2)
{
    float u1 = 0.0, u2 = 1.0, dx = p2.x - p1.x, dy;

    if (clipTest (-dx, p1.x - winMin.x, &u1, &u2))
        if (clipTest (dx, winMax.x - p1.x, &u1, &u2)) {
            dy = p2.y - p1.y;
            if (clipTest (-dy, p1.y - winMin.y, &u1, &u2))
                if (clipTest (dy, winMax.y - p1.y, &u1, &u2)) {
                    if (u2 < 1.0) {
                        p2.x = p1.x + u2 * dx;
                        p2.y = p1.y + u2 * dy;
                    }
                    if (u1 > 0.0) {
                        p1.x += u1 * dx;
                        p1.y += u1 * dy;
                    }
                    lineDDA (ROUND(p1.x), ROUND(p1.y), ROUND(p2.x), ROUND(p2.y));
                }
            }
        }
}

```

In general, the Liang-Barsky algorithm is more efficient than the Cohen-Sutherland algorithm, since intersection calculations are reduced. Each update of parameters u_1 and u_2 requires only one division; and window intersections of the line are computed only once, when the final values of u_1 and u_2 have been computed. In contrast, the Cohen-Sutherland algorithm can repeatedly calculate intersections along a line path, even though the line may be completely outside the clip window. And, each intersection calculation requires both a division and a multiplication. Both the Cohen-Sutherland and the Liang-Barsky algorithms can be extended to three-dimensional clipping (Chapter 12).

By creating more regions around the clip window, the Nicholl-Lee-Nicholl (or NLN) algorithm avoids multiple clipping of an individual line segment. In the Cohen-Sutherland method, for example, multiple intersections may be calculated along the path of a single line before an intersection on the clipping rectangle is located or the line is completely rejected. These extra intersection calculations are eliminated in the NLN algorithm by carrying out more region testing before intersection positions are calculated. Compared to both the Cohen-Sutherland and the Liang-Barsky algorithms, the Nicholl-Lee-Nicholl algorithm performs fewer comparisons and divisions. The trade-off is that the NLN algorithm can only be applied to two-dimensional clipping, whereas both the Liang-Barsky and the Cohen-Sutherland methods are easily extended to three-dimensional scenes.

For a line with endpoints P_1 and P_2 , we first determine the position of point P_1 for the nine possible regions relative to the clipping rectangle. Only the three regions shown in Fig. 6-10 need be considered. If P_1 lies in any one of the other six regions, we can move it to one of the three regions in Fig. 6-10 using a symmetry transformation. For example, the region directly above the clip window can be transformed to the region left of the clip window using a reflection about the line $y = -x$, or we could use a 90° counterclockwise rotation.

Next, we determine the position of P_2 relative to P_1 . To do this, we create some new regions in the plane, depending on the location of P_1 . Boundaries of the new regions are half-infinite line segments that start at the position of P_1 and pass through the window corners. If P_1 is inside the clip window and P_2 is outside, we set up the four regions shown in Fig. 6-11. The intersection with the appropriate window boundary is then carried out, depending on which one of the four regions (L , T , R , or B) contains P_2 . Of course, if both P_1 and P_2 are inside the clipping rectangle, we simply save the entire line.

If P_1 is in the region to the left of the window, we set up the four regions, L , LT , LR , and LB , shown in Fig. 6-12. These four regions determine a unique boundary for the line segment. For instance, if P_2 is in region L , we clip the line at the left boundary and save the line segment from this intersection point to P_2 . But if P_2 is in region LT , we save the line segment from the left window boundary to the top boundary. If P_2 is not in any of the four regions, L , LT , LR , or LB , the entire line is clipped.

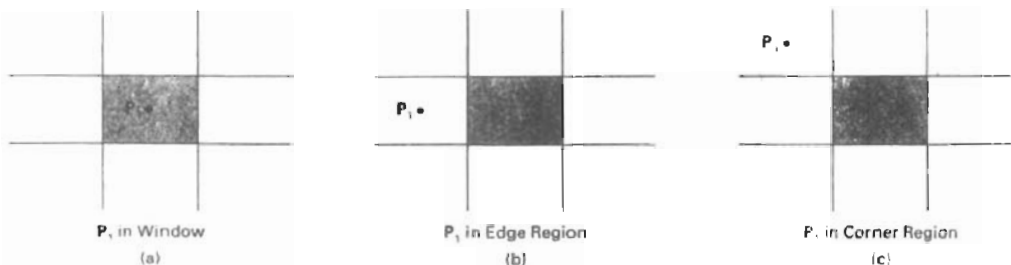


Figure 6-10

Three possible positions for a line endpoint P_1 in the NLN line-clipping algorithm.

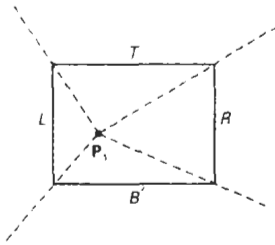


Figure 6-11
The four clipping regions used in the NLN algorithm when P_1 is inside the clip window and P_2 is outside.

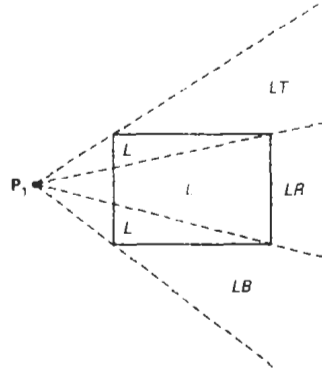


Figure 6-12
The four clipping regions used in the NLN algorithm when P_1 is directly left of the clip window.

For the third case, when P_1 is to the left and above the clip window, we use the clipping regions in Fig. 6-13. In this case, we have the two possibilities shown, depending on the position of P_1 relative to the top left corner of the window. If P_2 is in one of the regions T , L , TR , TB , LR , or LB , this determines a unique clip-window edge for the intersection calculations. Otherwise, the entire line is rejected.

To determine the region in which P_2 is located, we compare the slope of the line to the slopes of the boundaries of the clip regions. For example, if P_1 is left of the clipping rectangle (Fig. 6-12), then P_2 is in region LT if

$$\text{slope } \overline{P_1 P_{TR}} < \text{slope } \overline{P_1 P_2} < \text{slope } \overline{P_1 P_{TL}} \quad (6-14)$$

or

$$\frac{y_T - y_1}{x_R - x_1} < \frac{y_2 - y_1}{x_2 - x_1} < \frac{y_T - y_1}{x_L - x_1} \quad (6-15)$$

And we clip the entire line if

$$(y_1 - y_1)(x_2 - x_1) < (x_L - x_1)(y_2 - y_1) \quad (6-16)$$

The coordinate difference and product calculations used in the slope tests are saved and also used in the intersection calculations. From the parametric equations

$$x = x_1 + (x_2 - x_1)u$$

$$y = y_1 + (y_2 - y_1)u$$

an x -intersection position on the left window boundary is $x = x_L$, with $u = (x_L - x_1)/(x_2 - x_1)$, so that the y -intersection position is

$$y = y_1 + \frac{y_2 - y_1}{x_2 - x_1}(x_L - x_1) \quad (6-17)$$

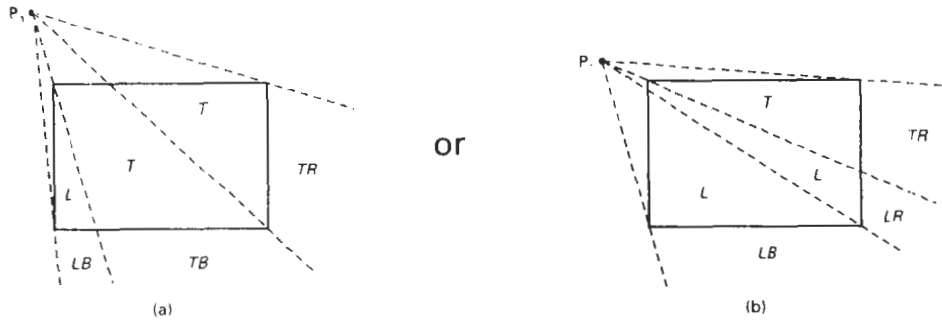


Figure 6-13

The two possible sets of clipping regions used in the NLN algorithm when P_1 is above and to the left of the clip window.

And an intersection position on the top boundary has $y = y_T$ and $u = (y_T - y_1)/(y_2 - y_1)$, with

$$x = x_1 + \frac{x_2 - x_1}{y_2 - y_1}(y_T - y_1) \quad (6-18)$$

Line Clipping Using Nonrectangular Clip Windows

In some applications, it is often necessary to clip lines against arbitrarily shaped polygons. Algorithms based on parametric line equations, such as the Liang-Barsky method and the earlier Cyrus-Beck approach, can be extended easily to convex polygon windows. We do this by modifying the algorithm to include the parametric equations for the boundaries of the clip region. Preliminary screening of line segments can be accomplished by processing lines against the coordinate extents of the clipping polygon. For concave polygon-clipping regions, we can still apply these parametric clipping procedures if we first split the concave polygon into a set of convex polygons.

Circles or other curved-boundary clipping regions are also possible, but less commonly used. Clipping algorithms for these areas are slower because intersection calculations involve nonlinear curve equations. At the first step, lines can be clipped against the bounding rectangle (coordinate extents) of the curved clipping region. Lines that can be identified as completely outside the bounding rectangle are discarded. To identify inside lines, we can calculate the distance of line endpoints from the circle center. If the square of this distance for both endpoints of a line is less than or equal to the radius squared, we can save the entire line. The remaining lines are then processed through the intersection calculations, which must solve simultaneous circle-line equations.

Splitting Concave Polygons

We can identify a concave polygon by calculating the cross products of successive edge vectors in order around the polygon perimeter. If the z component of

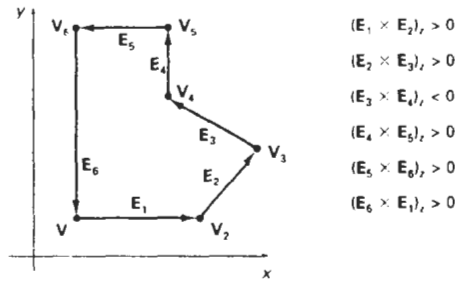


Figure 6-14
Identifying a concave polygon by calculating cross products of successive pairs of edge vectors.

some cross products is positive while others have a negative z component, we have a concave polygon. Otherwise, the polygon is convex. This is assuming that no series of three successive vertices are collinear, in which case the cross product of the two edge vectors for these vertices is zero. If all vertices are collinear, we have a degenerate polygon (a straight line). Figure 6-14 illustrates the edge-vector cross-product method for identifying concave polygons.

A *vector method* for splitting a concave polygon in the xy plane is to calculate the edge-vector cross products in a counterclockwise order and to note the sign of the z component of the cross products. If any z component turns out to be negative (as in Fig. 6-14), the polygon is concave and we can split it along the line of the first edge vector in the cross-product pair. The following example illustrates this method for splitting a concave polygon.

Example 6-2: Vector Method for Splitting Concave Polygons

Figure 6-15 shows a concave polygon with six edges. Edge vectors for this polygon can be expressed as

$$\begin{aligned}
 E_1 &= (1, 0, 0), & E_2 &= (1, 1, 0) \\
 E_3 &= (1, -1, 0), & E_4 &= (0, 2, 0) \\
 E_5 &= (-3, 0, 0), & E_6 &= (0, -2, 0)
 \end{aligned}$$

where the z component is 0, since all edges are in the xy plane. The cross product $E_i \times E_j$ for two successive edge vectors is a vector perpendicular to the xy plane with z component equal to $E_{ix}E_{jy} - E_{iy}E_{jx}$.

$$\begin{aligned}
 E_1 \times E_2 &= (0, 0, 1), & E_2 \times E_3 &= (0, 0, -2) \\
 E_3 \times E_4 &= (0, 0, 2), & E_4 \times E_5 &= (0, 0, 6) \\
 E_5 \times E_6 &= (0, 0, 6), & E_6 \times E_1 &= (0, 0, 2)
 \end{aligned}$$

Since the cross product $E_2 \times E_3$ has a negative z component, we split the polygon along the line of vector E_2 . The line equation for this edge has a slope of 1 and a y intercept of -1 . We then determine the intersection of this line and the other

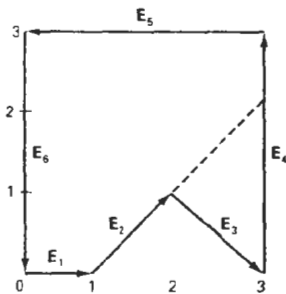


Figure 6-15
Splitting a concave polygon using the vector method.

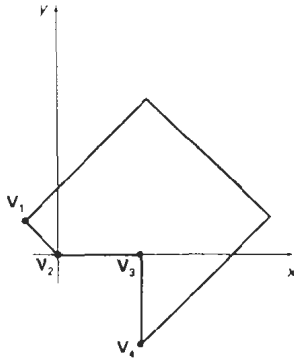


Figure 6-16
Splitting a concave polygon using the rotational method. After rotating V_3 onto the x axis, we find that V_4 is below the x axis. So we split the polygon along the line of V_2V_3 .

polygon edges to split the polygon into two pieces. No other edge cross products are negative, so the two new polygons are both convex.

We can also split a concave polygon using a *rotational method*. Proceeding counterclockwise around the polygon edges, we translate each polygon vertex V_k in turn to the coordinate origin. We then rotate in a clockwise direction so that the next vertex V_{k+1} is on the x axis. If the next vertex, V_{k+2} , is below the x axis, the polygon is concave. We then split the polygon into two new polygons along the x axis and repeat the concave test for each of the two new polygons. Otherwise, we continue to rotate vertices on the x axis and to test for negative y vertex values. Figure 6-16 illustrates the rotational method for splitting a concave polygon.

6-8

POLYGON CLIPPING

To clip polygons, we need to modify the line-clipping procedures discussed in the previous section. A polygon boundary processed with a line clipper may be displayed as a series of unconnected line segments (Fig. 6-17), depending on the orientation of the polygon to the clipping window. What we really want to display is a bounded area after clipping, as in Fig. 6-18. For polygon clipping, we require an algorithm that will generate one or more closed areas that are then scan converted for the appropriate area fill. The output of a polygon clipper should be a sequence of vertices that defines the clipped polygon boundaries.

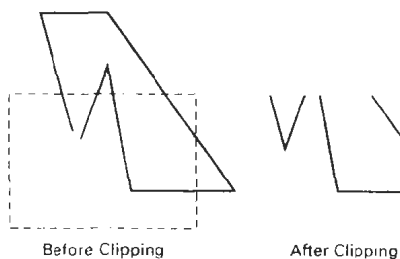


Figure 6-17
Display of a polygon processed by a line-clipping algorithm

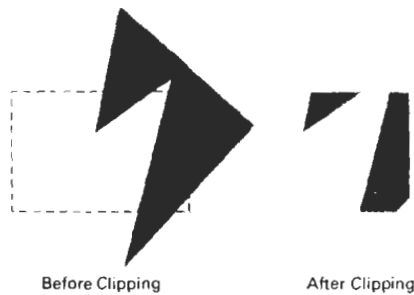


Figure 6-18
Display of a correctly clipped
polygon.

Sutherland-Hodgeman Polygon Clipping

We can correctly clip a polygon by processing the polygon boundary as a whole against each window edge. This could be accomplished by processing all polygon vertices against each clip rectangle boundary in turn. Beginning with the initial set of polygon vertices, we could first clip the polygon against the left rectangular boundary to produce a new sequence of vertices. The new set of vertices could then be successively passed to a right boundary clipper, a bottom boundary clipper, and a top boundary clipper, as in Fig. 6-19. At each step, a new sequence of output vertices is generated and passed to the next window boundary clipper.

There are four possible cases when processing vertices in sequence around the perimeter of a polygon. As each pair of adjacent polygon vertices is passed to a window boundary clipper, we make the following tests: (1) If the first vertex is outside the window boundary and the second vertex is inside, both the intersection point of the polygon edge with the window boundary and the second vertex are added to the output vertex list. (2) If both input vertices are inside the window boundary, only the second vertex is added to the output vertex list. (3) If the first vertex is inside the window boundary and the second vertex is outside, only the edge intersection with the window boundary is added to the output vertex list. (4) If both input vertices are outside the window boundary, nothing is added to the output list. These four cases are illustrated in Fig. 6-20 for successive pairs of polygon vertices. Once all vertices have been processed for one clip window boundary, the output list of vertices is clipped against the next window boundary.

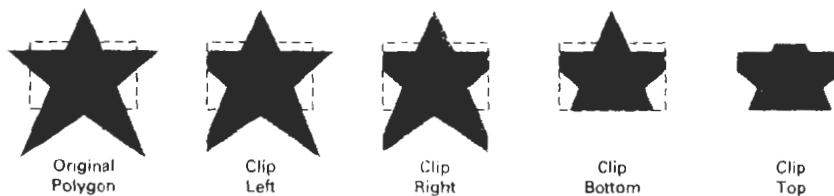


Figure 6-19
Clipping a polygon against successive window boundaries.

```

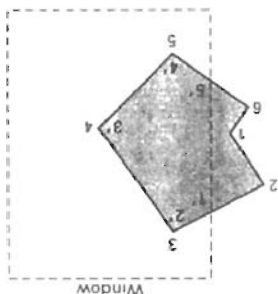
typedef enum { Left, Right, Bottom, Top } Edge;
#define N_EDGES 4

int inside (wcp2 p, Edge b, dcpt wMin, dcpt wMax)
{
    switch (b) {

```

Clipping a polygon against the left boundary of a window, starting with vertex 1. Primed numbers are used to label the points in the output vertex list for this window boundary.

Figure 6-21

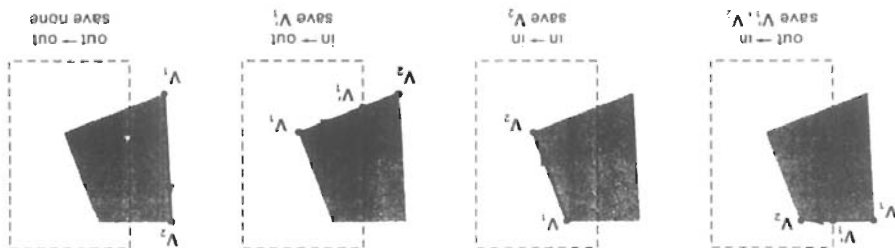


The following procedure demonstrates the pipeline clipping approach. An array, *s*, records the most recent point that was clipped for each clip-window boundary. The main routine passes each vertex *p* to the `clipPoint` routine for clipping against the first window boundary. If the line defined by endpoints *p* and *s* [boundary] crosses this window boundary, the intersection is calculated and passed to the next clipping stage. If *p* is inside the window, it is passed to the next clipping stage. Any point that survives clipping against all window boundaries is then entered into the output array of points. The array `firstPoint` stores for each window boundary the first point clipped against that boundary. After all polygon vertices have been processed, a closing routine clips lines defined by the first and last points clipped against each boundary.

illustrate the progression of the polygon vertices in Fig. 6-22 through a pipeline shows a polygon and its intersection points with a clip window. In Fig. 6-23, we any clippers. Otherwise, the point does not continue in the pipeline. Figure 6-22 been determined to be inside or on a window boundary by all four bound- calculated intersection point) is added to the output vertex list only after it has a processor and a pipeline of clipping routines. A point (either an input vertex or a next boundary clipper. This can be done with parallel processors or a single ping individual vertices at each step and passing the clipped vertices on to the boundary. We can eliminate the intermediate output vertex lists by simply clip- storage for an output list of vertices as a polygon is clipped against each window implementing the algorithm as we have just described requires setting up

repeat the process for the next window boundary. find and save the intersection point. Using the five saved points, we would re- save both the intersection point and vertex 3. Vertices 4 and 5 are determined to save along to vertex 3, which is inside, we calculate the intersection and window boundary. Vertices 1 and 2 are found to be on the outside of the bound- We illustrate this method by processing the area in Fig. 6-21 against the left

Figure 6-20
Successive processing of pairs of polygon vertices against the left window boundary.



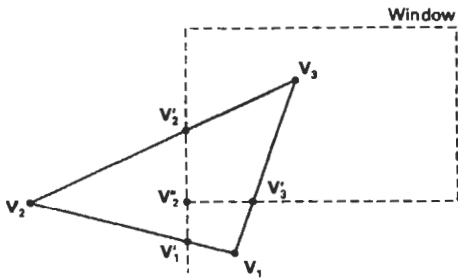


Figure 6-22
A polygon overlapping a rectangular clip window.

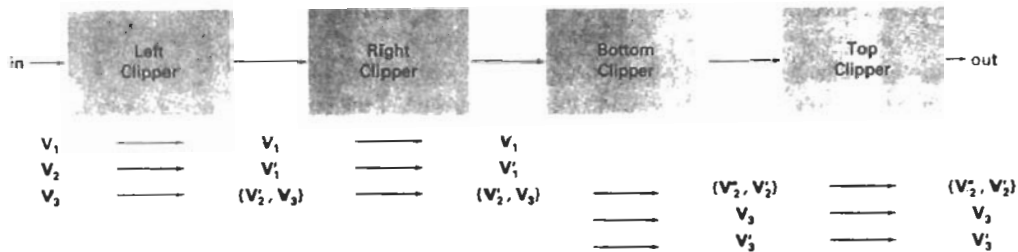


Figure 6-23
Processing the vertices of the polygon in Fig. 6-22 through a boundary-clipping pipeline. After all vertices are processed through the pipeline, the vertex list for the clipped polygon is $\{V'_2, V'_3, V_3, V'_1\}$.

```

case Left:   if (p.x < wMin.x) return (FALSE); break;
case Right:  if (p.x > wMax.x) return (FALSE); break;
case Bottom: if (p.y < wMin.y) return (FALSE); break;
case Top:    if (p.y > wMax.y) return (FALSE); break;
}
return (TRUE);
}

int cross (wcPt2 p1, wcPt2 p2, Edge b, dcPt wMin, dcPt wMax)
{
    if (inside (p1, b, wMin, wMax) == inside (p2, b, wMin, wMax))
        return (FALSE);
    else return (TRUE);
}

wcPt2 intersect (wcPt2 p1, wcPt2 p2, Edge b, dcPt wMin, dcPt wMax)
{
    wcPt2 iPt;
    float m;

    if (p1.x != p2.x) m = (p1.y - p2.y) / (p1.x - p2.x);
    switch (b) {
    case Left:
        iPt.x = wMin.x;
        iPt.y = p2.y + (wMin.x - p2.x) * m;
        break;
    case Right:
        iPt.x = wMax.x;

```

```

    iPt.y = p2.y + (wMax.x - p2.x) * m;
    break;
case Bottom:
    iPt.y = wMin.y;
    if (p1.x != p2.x) iPt.x = p2.x + (wMin.y - p2.y) / m;
    else iPt.x = p2.x;
    break;
case Top:
    iPt.y = wMax.y;
    if (p1.x != p2.x) iPt.x = p2.x + (wMax.y - p2.y) / m;
    else iPt.x = p2.x;
    break;
}
return (iPt);
}

void clipPoint (wcPt2 p, Edge b, dcPt wMin, dcPt wMax,
               wcPt2 * pOut, int * cnt, wcPt2 * first[], wcPt2 * s)
{
    wcPt2 iPt;

    /* If no previous point exists for this edge, save this point. */
    if (!first[b])
        first[b] = &p;
    else
        /* Previous point exists. If 'p' and previous point cross edge,
           find intersection. Clip against next boundary, if any. If
           no more edges, add intersection to output list. */
        if (cross (p, s[b], b, wMin, wMax)) {
            iPt = intersect (p, s[b], b, wMin, wMax);
            if (b < Top)
                clipPoint (iPt, b+1, wMin, wMax, pOut, cnt, first, s);
            else {
                pOut[*cnt] = iPt; (*cnt)++;
            }
        }

    s[b] = p; /* Save 'p' as most recent point for this edge */

    /* For all, if point is 'inside' proceed to next clip edge, if any */
    if (inside (p, b, wMin, wMax))
        if (b < Top)
            clipPoint (p, b+1, wMin, wMax, pOut, cnt, first, s);
        else {
            pOut[*cnt] = p; (*cnt)++;
        }
}

void closeClip (dcPt wMin, dcPt wMax, wcPt2 * pOut,
               int * cnt, wcPt2 * first[], wcPt2 * s)
{
    wcPt2 i;
    Edge b;

    for (b = Left; b <= Top; b++) {
        if (cross (s[b], *first[b], b, wMin, wMax)) {
            i = intersect (s[b], *first[b], b, wMin, wMax);
            if (b < Top)
                clipPoint (i, b+1, wMin, wMax, pOut, cnt, first, s);
            else {
                pOut[*cnt] = i; (*cnt)++;
            }
        }
    }
}

```

```

    }
}

int clipPolygon (dcPt wMin, dcPt wMax, int n, wcPt2 * pIn, wcPt2 * pOut)
{
    /* 'first' holds pointer to first point processed against a clip
       edge. 's' holds most recent point processed against an edge */
    wcPt2 * first[N_EDGE] = { 0, 0, 0, 0 } s[N_EDGE];
    int i, cnt = 0;

    for (i=0; i<n; i++)
        clipPoint (pIn[i], Left, wMin, wMax, pOut, &cnt, first, s);
    closeClip (wMin, wMax, pOut, &cnt, first, s);
    return (cnt);
}

```

Convex polygons are correctly clipped by the Sutherland-Hodgeman algorithm, but concave polygons may be displayed with extraneous lines, as demonstrated in Fig. 6-24. This occurs when the clipped polygon should have two or more separate sections. But since there is only one output vertex list, the last vertex in the list is always joined to the first vertex. There are several things we could do to correctly display concave polygons. For one, we could split the concave polygon into two or more convex polygons and process each convex polygon separately. Another possibility is to modify the Sutherland-Hodgeman approach to check the final vertex list for multiple vertex points along any clip window boundary and correctly join pairs of vertices. Finally, we could use a more general polygon clipper, such as either the Weiler-Atherton algorithm or the Weiler algorithm described in the next section.

Weiler-Atherton Polygon Clipping

Here, the vertex-processing procedures for window boundaries are modified so that concave polygons are displayed correctly. This clipping procedure was developed as a method for identifying visible surfaces, and so it can be applied with arbitrary polygon-clipping regions.

The basic idea in this algorithm is that instead of always proceeding around the polygon edges as vertices are processed, we sometimes want to follow the window boundaries. Which path we follow depends on the polygon-processing direction (clockwise or counterclockwise) and whether the pair of polygon vertices currently being processed represents an outside-to-inside pair or an inside-

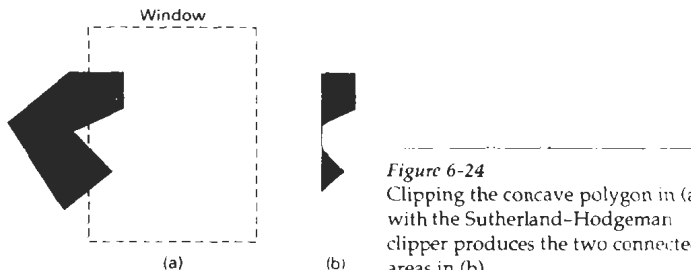


Figure 6-24
Clipping the concave polygon in (a) with the Sutherland-Hodgeman clipper produces the two connected areas in (b).

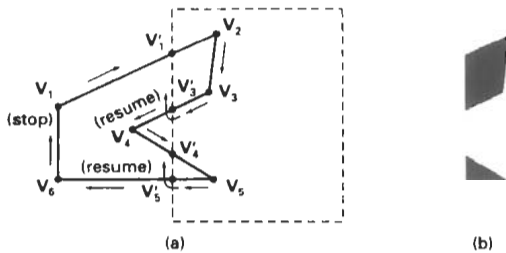


Figure 6-25
Clipping a concave polygon (a) with the Weiler-Atherton algorithm generates the two separate polygon areas in (b).

to-outside pair. For clockwise processing of polygon vertices, we use the following rules:

- For an outside-to-inside pair of vertices, follow the polygon boundary.
- For an inside-to-outside pair of vertices, follow the window boundary in a clockwise direction.

In Fig. 6-25, the processing direction in the Weiler-Atherton algorithm and the resulting clipped polygon is shown for a rectangular clipping window.

An improvement on the Weiler-Atherton algorithm is the Weiler algorithm, which applies constructive solid geometry ideas to clip an arbitrary polygon against any polygon-clipping region. Figure 6-26 illustrates the general idea in this approach. For the two polygons in this figure, the correctly clipped polygon is calculated as the intersection of the clipping polygon and the polygon object.

Other Polygon-Clipping Algorithms

Various parametric line-clipping methods have also been adapted to polygon clipping. And they are particularly well suited for clipping against convex polygon-clipping windows. The Liang-Barsky Line Clipper, for example, can be extended to polygon clipping with a general approach similar to that of the Sutherland-Hodgeman method. Parametric line representations are used to process polygon edges in order around the polygon perimeter using region-testing procedures similar to those used in line clipping.



Figure 6-26
Clipping a polygon by determining the intersection of two polygon areas.

6-9 CURVE CLIPPING

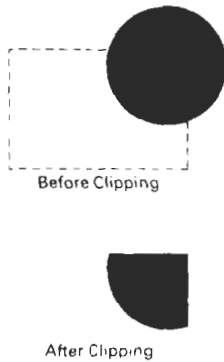


Figure 6-27
Clipping a filled circle.

Areas with curved boundaries can be clipped with methods similar to those discussed in the previous sections. Curve-clipping procedures will involve nonlinear equations, however, and this requires more processing than for objects with linear boundaries.

The bounding rectangle for a circle or other curved object can be used first to test for overlap with a rectangular clip window. If the bounding rectangle for the object is completely inside the window, we save the object. If the rectangle is determined to be completely outside the window, we discard the object. In either case, there is no further computation necessary. But if the bounding rectangle test fails, we can look for other computation-saving approaches. For a circle, we can use the coordinate extents of individual quadrants and then octants for preliminary testing before calculating curve-window intersections. For an ellipse, we can test the coordinate extents of individual quadrants. Figure 6-27 illustrates circle clipping against a rectangular window.

Similar procedures can be applied when clipping a curved object against a general polygon clip region. On the first pass, we can clip the bounding rectangle of the object against the bounding rectangle of the clip region. If the two regions overlap, we will need to solve the simultaneous line-curve equations to obtain the clipping intersection points.

6-10 TEXT CLIPPING

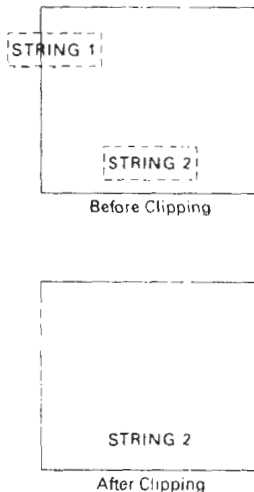


Figure 6-28
Text clipping using a
bounding rectangle about the
entire string.

There are several techniques that can be used to provide text clipping in a graphics package. The clipping technique used will depend on the methods used to generate characters and the requirements of a particular application.

The simplest method for processing character strings relative to a window boundary is to use the *all-or-none string-clipping* strategy shown in Fig. 6-28. If all of the string is inside a clip window, we keep it. Otherwise, the string is discarded. This procedure is implemented by considering a bounding rectangle around the text pattern. The boundary positions of the rectangle are then compared to the window boundaries, and the string is rejected if there is any overlap. This method produces the fastest text clipping.

An alternative to rejecting an entire character string that overlaps a window boundary is to use the *all-or-none character-clipping* strategy. Here we discard only those characters that are not completely inside the window (Fig. 6-29). In this case, the boundary limits of individual characters are compared to the window. Any character that either overlaps or is outside a window boundary is clipped.

A final method for handling text clipping is to clip the components of individual characters. We now treat characters in much the same way that we treated lines. If an individual character overlaps a clip window boundary, we clip off the parts of the character that are outside the window (Fig. 6-30). Outline character fonts formed with line segments can be processed in this way using a line-clipping algorithm. Characters defined with bit maps would be clipped by comparing the relative position of the individual pixels in the character grid patterns to the clipping boundaries.

6-11

EXTERIOR CLIPPING

So far, we have considered only procedures for clipping a picture to the interior of a region by eliminating everything outside the clipping region. What is saved by these procedures is *inside* the region. In some cases, we want to do the reverse, that is, we want to clip a picture to the exterior of a specified region. The picture parts to be saved are those that are *outside* the region. This is referred to as **exterior clipping**.

A typical example of the application of exterior clipping is in multiple-window systems. To correctly display the screen windows, we often need to apply both internal and external clipping. Figure 6-31 illustrates a multiple-window display. Objects within a window are clipped to the interior of that window. When other higher-priority windows overlap these objects, the objects are also clipped to the exterior of the overlapping windows.

Exterior clipping is used also in other applications that require overlapping pictures. Examples here include the design of page layouts in advertising or publishing applications or for adding labels or design patterns to a picture. The technique can also be used for combining graphs, maps, or schematics. For these applications, we can use exterior clipping to provide a space for an insert into a larger picture.

Procedures for clipping objects to the interior of concave polygon windows can also make use of external clipping. Figure 6-32 shows a line $\overline{P_1P_2}$ that is to be clipped to the interior of a concave window with vertices $V_1V_2V_3V_4V_5$. Line $\overline{P_1P_2}$ can be clipped in two passes: (1) First, $\overline{P_1P_2}$ is clipped to the interior of the convex polygon $V_1V_2V_3V_4$ to yield the clipped segment $\overline{P'_1P'_2}$ (Fig. 6-32(b)). (2) Then an external clip of $\overline{P'_1P'_2}$ is performed against the convex polygon $V_1V_5V_4$ to yield the final clipped line segment $\overline{P''_1P''_2}$.

SUMMARY

In this chapter, we have seen how we can map a two-dimensional world-coordinate scene to a display device. The viewing-transformation pipeline in-

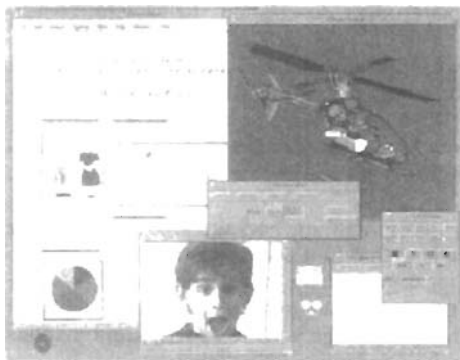
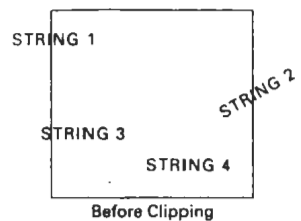


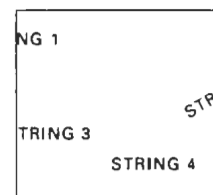
Figure 6-31

A multiple-window screen display showing examples of both interior and exterior clipping. (Courtesy of Sun Microsystems).

Summary



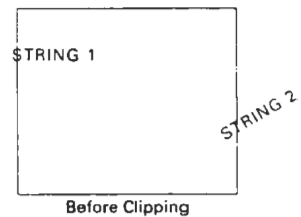
Before Clipping



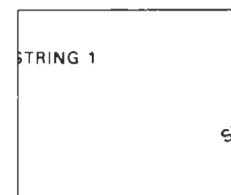
After Clipping

Figure 6-29

Text clipping using a bounding rectangle about individual characters.



Before Clipping



After Clipping

Figure 6-30

Text clipping performed on the components of individual characters.

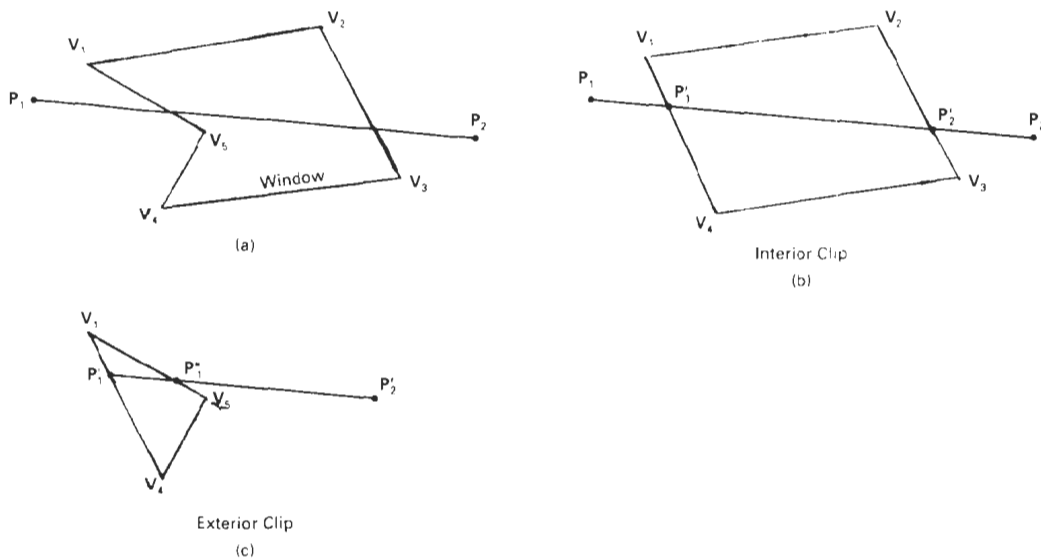


Figure 6-32
Clipping line $\overline{P_1P_2}$ to the interior of a concave polygon with vertices V_1, V_2, V_3, V_4, V_5 : (a), using convex polygons $V_1V_2V_3V_4$ (b) and $V_1V_5V_4$ (c), to produce the clipped line $\overline{P'_1P'_2}$.

cludes constructing the world-coordinate scene using modeling transformations transferring world-coordinates to viewing coordinates, mapping the viewing-coordinate descriptions of objects to normalized device coordinates, and finally mapping to device coordinates. Normalized coordinates are specified in the range from 0 to 1, and they are used to make viewing packages independent of particular output devices.

Viewing coordinates are specified by giving the world-coordinate position of the viewing origin and the view up vector that defines the direction of the viewing y axis. These parameters are used to construct the viewing transformation matrix that maps world-coordinate object descriptions to viewing coordinates.

A window is then set up in viewing coordinates, and a viewport is specified in normalized device coordinates. Typically, the window and viewport are rectangles in standard position (rectangle boundaries are parallel to the coordinate axes). The mapping from viewing coordinates to normalized device coordinates is then carried out so that relative positions in the window are maintained in the viewport.

Viewing functions in a graphics programming package are used to create one or more sets of viewing parameters. One function is typically provided to calculate the elements of the matrix for transforming world coordinates to viewing coordinates. Another function is used to set up the window-to-viewport transformation matrix, and a third function can be used to specify combinations of viewing transformations and window mapping in a viewing table. We can

then select different viewing combinations by specifying particular view indices listed in the viewing table.

When objects are displayed on the output device, all parts of a scene outside the window (and the viewport) are clipped off unless we set clip parameters to turn off clipping. In many packages, clipping is done in normalized device coordinates so that all transformations can be concatenated into a single transformation operation before applying the clipping algorithms. The clipping region is commonly referred to as the clipping window, or as the clipping rectangle when the window and viewport are standard rectangles. Several algorithms have been developed for clipping objects against the clip-window boundaries.

Line-clipping algorithms include the Cohen-Sutherland method, the Liang-Barsky method, and the Nicholl-Lee-Nicholl method. The Cohen-Sutherland method is widely used, since it was one of the first line-clipping algorithms to be developed. Region codes are used to identify the position of line endpoints relative to the rectangular, clipping window boundaries. Lines that cannot be immediately identified as completely inside the window or completely outside are then clipped against window boundaries. Liang and Barsky use a parametric line representation, similar to that of the earlier Cyrus-Beck algorithm, to set up a more efficient line-clipping procedure that reduces intersection calculations. The Nicholl-Lee-Nicholl algorithm uses more region testing in the xy plane to reduce intersection calculations even further. Parametric line clipping is easily extended to convex clipping windows and to three-dimensional clipping windows.

Line clipping can also be carried out for concave, polygon clipping windows and for clipping windows with curved boundaries. With concave polygons, we can use either the vector method or the rotational method to split a concave polygon into a number of convex polygons. With curved clipping windows, we calculate line intersections using the curve equations.

Polygon-clipping algorithms include the Sutherland-Hodgeman method, the Liang-Barsky method, and the Weiler-Atherton method. In the Sutherland-Hodgeman clipper, vertices of a convex polygon are processed in order against the four rectangular window boundaries to produce an output vertex list for the clipped polygon. Liang and Barsky use parametric line equations to represent the convex polygon edges, and they use similar testing to that performed in line clipping to produce an output vertex list for the clipped polygon. Both the Weiler-Atherton method and the Weiler method correctly clip both convex and concave polygons, and these polygon clippers also allow the clipping window to be a general polygon. The Weiler-Atherton algorithm processes polygon vertices in order to produce one or more lists of output polygon vertices. The Weiler method performs clipping by finding the intersection region of the two polygons.

Objects with curved boundaries are processed against rectangular clipping windows by calculating intersections using the curve equations. These clipping procedures are slower than line clippers or polygon clippers, because the curve equations are nonlinear.

The fastest text-clipping method is to completely clip a string if any part of the string is outside any window boundary. Another method for text clipping is to use the all-or-none approach with the individual characters in a string. A third method is to apply either point, line, polygon, or curve clipping to the individual characters in a string, depending on whether characters are defined as point grids or as outline fonts.

In some applications, such as creating picture insets and managing multiple-screen windows, exterior clipping is performed. In this case, all parts of a scene that are inside a window are clipped and the exterior parts are saved.

REFERENCES

- Line-clipping algorithms are discussed in Sproull and Sutherland (1968), Cyrus and Beck (1978), and Liang and Barsky (1984). Methods for improving the speed of the Cohen–Sutherland line-clipping algorithm are given in Duvanenko (1990).
- Polygon-clipping methods are presented in Sutherland and Hodgeman (1974) and in Liang and Barsky (1983). General techniques for clipping arbitrarily shaped polygons against each other are given in Weiler and Atherton (1977) and in Weiler (1980).
- Two-dimensional viewing operations in PHIGS are discussed in Howard et al. (1991), Gaskins (1992), Hopgood and Duce (1991), and Blake (1993). For information on GKS viewing operations, see Hopgood et al. (1983) and Enderle et al. (1984).

EXERCISES

- 6-1. Write a procedure to implement the `evaluateViewOrientationMatrix` function that calculates the elements of the matrix for transforming world coordinates to viewing coordinates, given the viewing coordinate origin P_0 and the view up vector V .
- 6-2. Derive the window-to-viewport transformation equations 6-3 by first scaling the window to the size of the viewport and then translating the scaled window to the viewport position.
- 6-3. Write a procedure to implement the `evaluateViewMappingMatrix` function that calculates the elements of a matrix for performing the window-to-viewport transformation.
- 6-4. Write a procedure to implement the `setViewRepresentation` function to concatenate `viewMatrix` and `viewMappingMatrix` and to store the result, referenced by a specified view index, in a viewing table.
- 6-5. Write a set of procedures to implement the viewing pipeline without clipping and without the workstation transformation. Your program should allow a scene to be constructed with modeling-coordinate transformations, a specified viewing system, and a specified window–viewport pair. As an option, a viewing table can be implemented to store different sets of viewing transformation parameters.
- 6-6. Derive the matrix representation for a workstation transformation.
- 6-7. Write a set of procedures to implement the viewing pipeline without clipping, but including the workstation transformation. Your program should allow a scene to be constructed with modeling-coordinate transformations, a specified viewing system, a specified window–viewport pair, and workstation transformation parameters. For a given world-coordinate scene, the composite viewing transformation matrix should transform the scene to an output device for display.
- 6-8. Implement the Cohen–Sutherland line-clipping algorithm.
- 6-9. Carefully discuss the rationale behind the various tests and methods for calculating the intersection parameters u_1 and u_2 in the Liang–Barsky line-clipping algorithm.
- 6-10. Compare the number of arithmetic operations performed in the Cohen–Sutherland and the Liang–Barsky line-clipping algorithms for several different line orientations relative to a clipping window.
- 6-11. Write a procedure to implement the Liang–Barsky line-clipping algorithm.
- 6-12. Devise symmetry transformations for mapping the intersection calculations for the three regions in Fig. 6-10 to the other six regions of the xy plane.
- 6-13. Set up a detailed algorithm for the Nicholl–Lee–Nicholl approach to line clipping for any input pair of line endpoints.
- 6-14. Compare the number of arithmetic operations performed in NLN algorithm to both the Cohen–Sutherland and the Liang–Barsky line-clipping algorithms for several different line orientations relative to a clipping window.

- 6-15. Write a routine to identify concave polygons by calculating cross products of pairs of edge vectors.
- 6-16. Write a routine to split a concave polygon using the vector method.
- 6-17. Write a routine to split a concave polygon using the rotational method.
- 6-18. Adapt the Liang–Barsky line-clipping algorithm to polygon clipping.
- 6-19. Set up a detailed algorithm for Weiler–Atherton polygon clipping assuming that the clipping window is a rectangle in standard position.
- 6-20. Devise an algorithm for Weiler–Atherton polygon clipping, where the clipping window can be any specified polygon.
- 6-21. Write a routine to clip an ellipse against a rectangular window.
- 6-22. Assuming that all characters in a text string have the same width, develop a text-clipping algorithm that clips a string according to the “all-or-none character-clipping” strategy.
- 6-23. Develop a text-clipping algorithm that clips individual characters assuming that the characters are defined in a pixel grid of a specified size.
- 6-24. Write a routine to implement exterior clipping on any part of a defined picture using any specified window.
- 6-25. Write a routine to perform both interior and exterior clipping, given a particular window-system display. Input to the routine is a set of window positions on the screen, the objects to be displayed in each window, and the window priorities. The individual objects are to be clipped to fit into their respective windows, then clipped to remove parts with overlapping windows of higher display priority.

Exercises

Three-Dimensional Concepts



When we model and display a three-dimensional scene, there are many more considerations we must take into account besides just including coordinate values for the third dimension. Object boundaries can be constructed with various combinations of plane and curved surfaces, and we sometimes need to specify information about object interiors. Graphics packages often provide routines for displaying internal components or cross-sectional views of solid objects. Also, some geometric transformations are more involved in three-dimensional space than in two dimensions. For example, we can rotate an object about an axis with any spatial orientation in three-dimensional space. Two-dimensional rotations, on the other hand, are always around an axis that is perpendicular to the xy plane. Viewing transformations in three dimensions are much more complicated because we have many more parameters to select when specifying how a three-dimensional scene is to be mapped to a display device. The scene description must be processed through viewing-coordinate transformations and projection routines that transform three-dimensional viewing coordinates onto two-dimensional device coordinates. Visible parts of a scene, for a selected view, must be identified; and surface-rendering algorithms must be applied if a realistic rendering of the scene is required.

9-1

THREE-DIMENSIONAL DISPLAY METHODS

To obtain a display of a three-dimensional scene that has been modeled in world coordinates, we must first set up a coordinate reference for the "camera". This coordinate reference defines the position and orientation for the plane of the camera film (Fig. 9-1), which is the plane we want to use to display a view of the objects in the scene. Object descriptions are then transferred to the camera reference coordinates and projected onto the selected display plane. We can then display

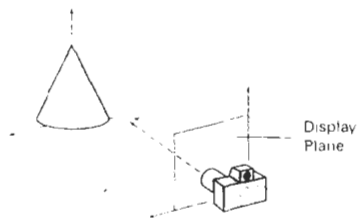


Figure 9-1
Coordinate reference for obtaining
a particular view of a
three-dimensional scene.

the objects in wireframe (outline) form, as in Fig. 9-2, or we can apply lighting and surface-rendering techniques to shade the visible surfaces.

Parallel Projection

One method for generating a view of a solid object is to project points on the object surface along parallel lines onto the display plane. By selecting different viewing positions, we can project visible points on the object onto the display plane to obtain different two-dimensional views of the object, as in Fig. 9-3. In a *parallel projection*, parallel lines in the world-coordinate scene project into parallel lines on the two-dimensional display plane. This technique is used in engineering and architectural drawings to represent an object with a set of views that maintain relative proportions of the object. The appearance of the solid object can then be reconstructed from the major views.

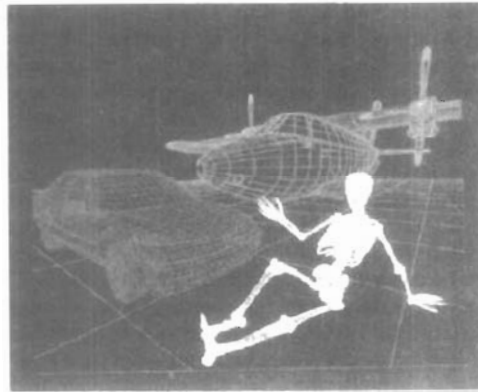


Figure 9-2
Wireframe display of three objects, with back lines removed, from a commercial database of object shapes. Each object in the database is defined as a grid of coordinate points, which can then be viewed in wireframe form or in a surface-rendered form. (Courtesy of Viewpoint DataLabs.)

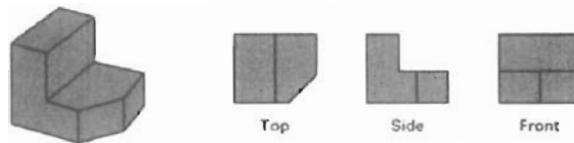


Figure 9-3
Three parallel-projection views of an object, showing relative proportions from different viewing positions.

Perspective Projection

Another method for generating a view of a three-dimensional scene is to project points to the display plane along converging paths. This causes objects farther from the viewing position to be displayed smaller than objects of the same size that are nearer to the viewing position. In a *perspective projection*, parallel lines in a scene that are not parallel to the display plane are projected into converging lines. Scenes displayed using perspective projections appear more realistic, since this is the way that our eyes and a camera lens form images. In the perspective-projection view shown in Fig. 9-4, parallel lines appear to converge to a distant point in the background, and distant objects appear smaller than objects closer to the viewing position.

Depth Cueing

With few exceptions, depth information is important so that we can easily identify, for a particular viewing direction, which is the front and which is the back of displayed objects. Figure 9-5 illustrates the ambiguity that can result when a wireframe object is displayed without depth information. There are several ways in which we can include depth information in the two-dimensional representation of solid objects.

A simple method for indicating depth with wireframe displays is to vary the intensity of objects according to their distance from the viewing position. Figure 9-6 shows a wireframe object displayed with *depth cueing*. The lines closest to

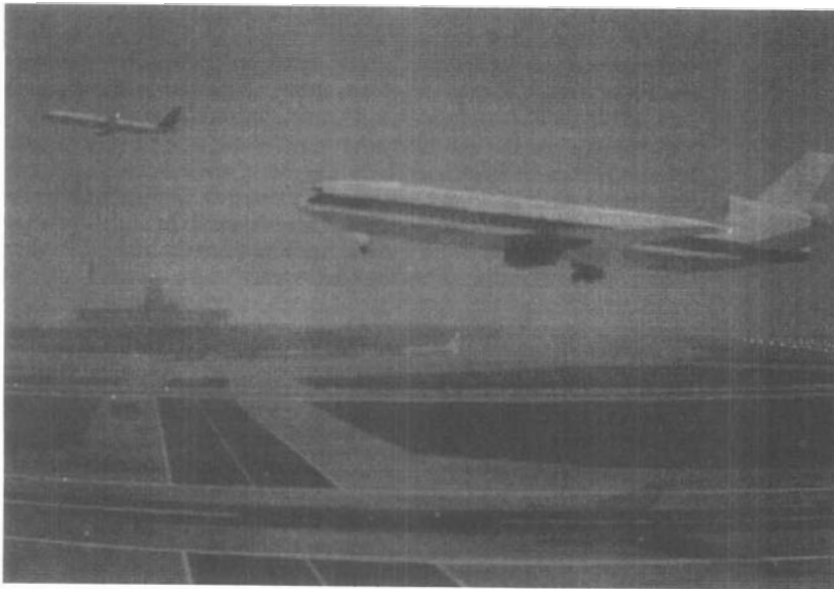
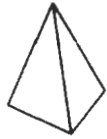


Figure 9-4

A perspective-projection view of an airport scene. (Courtesy of Evans & Sutherland.)



(a)



(b)



(c)

Figure 9-5
The wireframe representation of the pyramid in (a) contains no depth information to indicate whether the viewing direction is (b) downward from a position above the apex or (c) upward from a position below the base.

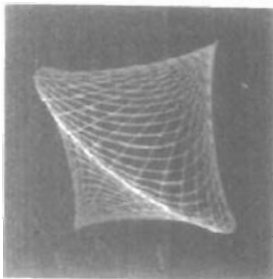


Figure 9-6
A wireframe object displayed with depth cueing, so that the intensity of lines decreases from the front to the back of the object.

the viewing position are displayed with the highest intensities, and lines farther away are displayed with decreasing intensities. Depth cueing is applied by choosing maximum and minimum intensity (or color) values and a range of distances over which the intensities are to vary.

Another application of depth cueing is modeling the effect of the atmosphere on the perceived intensity of objects. More distant objects appear dimmer to us than nearer objects due to light scattering by dust particles, haze, and smoke. Some atmospheric effects can change the perceived color of an object, and we can model these effects with depth cueing.

Visible Line and Surface Identification

We can also clarify depth relationships in a wireframe display by identifying visible lines in some way. The simplest method is to highlight the visible lines or to display them in a different color. Another technique, commonly used for engineering drawings, is to display the nonvisible lines as dashed lines. Another approach is to simply remove the nonvisible lines, as in Figs. 9-5(b) and 9-5(c). But removing the hidden lines also removes information about the shape of the back surfaces of an object. These visible-line methods also identify the visible surfaces of objects.

When objects are to be displayed with color or shaded surfaces, we apply surface-rendering procedures to the visible surfaces so that the hidden surfaces are obscured. Some visible-surface algorithms establish visibility pixel by pixel across the viewing plane; other algorithms determine visibility for object surfaces as a whole.

Surface Rendering

Added realism is attained in displays by setting the surface intensity of objects according to the lighting conditions in the scene and according to assigned surface characteristics. Lighting specifications include the intensity and positions of light sources and the general background illumination required for a scene. Surface properties of objects include degree of transparency and how rough or smooth the surfaces are to be. Procedures can then be applied to generate the correct illumination and shadow regions for the scene. In Fig. 9-7, surface-rendering methods are combined with perspective and visible-surface identification to generate a degree of realism in a displayed scene.

Exploded and Cutaway Views

Many graphics packages allow objects to be defined as hierarchical structures, so that internal details can be stored. Exploded and cutaway views of such objects can then be used to show the internal structure and relationship of the object parts. Figure 9-8 shows several kinds of exploded displays for a mechanical design. An alternative to exploding an object into its component parts is the cutaway view (Fig. 9-9), which removes part of the visible surfaces to show internal structure.

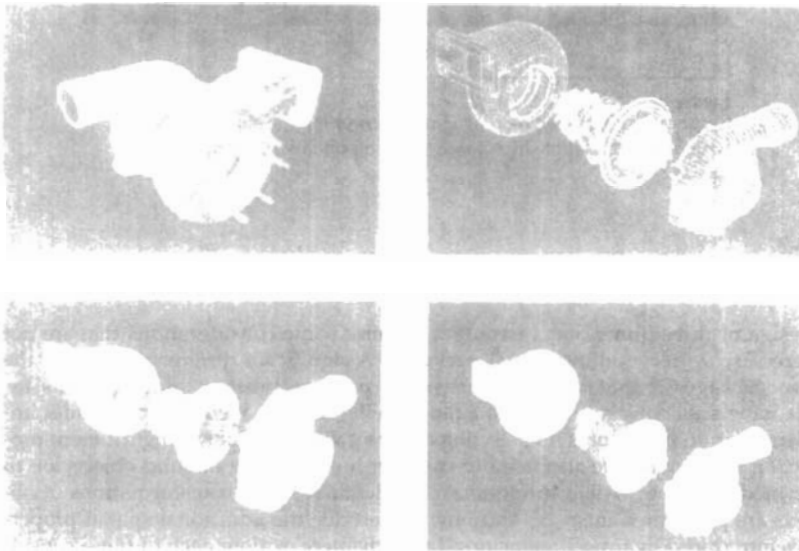
Three-Dimensional and Stereoscopic Views

Another method for adding a sense of realism to a computer-generated scene is to display objects using either three-dimensional or stereoscopic views. As we have seen in Chapter 2, three-dimensional views can be obtained by reflecting a

**Figure 9-7**

A realistic room display achieved with stochastic ray-tracing methods that apply a perspective projection, surface-texture mapping, and illumination models.

(Courtesy of John Snyder, Jed Lengyel, Devendra Kalra, and Al Barr, California Institute of Technology. Copyright © 1992 Caltech.)

**Figure 9-8**

A fully rendered and assembled turbine display (a) can also be viewed as (b) an exploded wireframe display, (c) a surface-rendered exploded display, or (d) a surface-rendered, color-coded exploded display.

(Courtesy of Autodesk, Inc.)

raster image from a vibrating flexible mirror. The vibrations of the mirror are synchronized with the display of the scene on the CRT. As the mirror vibrates, the focal length varies so that each point in the scene is projected to a position corresponding to its depth.

Stereoscopic devices present two views of a scene: one for the left eye and the other for the right eye. The two views are generated by selecting viewing positions that correspond to the two eye positions of a single viewer. These two views then can be displayed on alternate refresh cycles of a raster monitor, and viewed through glasses that alternately darken first one lens then the other in synchronization with the monitor refresh cycles.

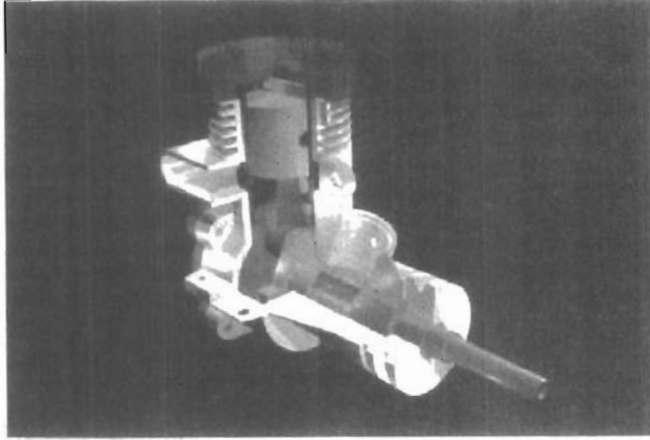


Figure 9-9
Color-coded cutaway view of a lawn mower engine showing the structure and relationship of internal components. (Courtesy of Autodesk, Inc.)

9-2

THREE-DIMENSIONAL GRAPHICS PACKAGES

Design of three-dimensional packages requires some considerations that are not necessary with two-dimensional packages. A significant difference between the two packages is that a three-dimensional package must include methods for mapping scene descriptions onto a flat viewing surface. We need to consider implementation procedures for selecting different views and for using different projection techniques. We also need to consider how surfaces of solid objects are to be modeled, how visible surfaces can be identified, how transformations of objects are performed in space, and how to describe the additional spatial properties introduced by three dimensions. Later chapters explore each of these considerations in detail.

Other considerations for three-dimensional packages are straightforward extensions from two-dimensional methods. World-coordinate descriptions are extended to three dimensions, and users are provided with output and input routines accessed with specifications such as

```
polyline3 (n, wcPoints)
fillarea3 (n, wcPoints)
text3 (wcPoint, string)
getLocator3 (wcPoint)
translate3(translateVector, matrixTranslate)
```

where points and vectors are specified with three components, and transformation matrices have four rows and four columns.

Two-dimensional attribute functions that are independent of geometric considerations can be applied in both two-dimensional and three-dimensional applications. No new attribute functions need be defined for colors, line styles, marker

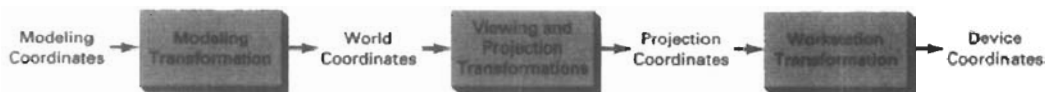
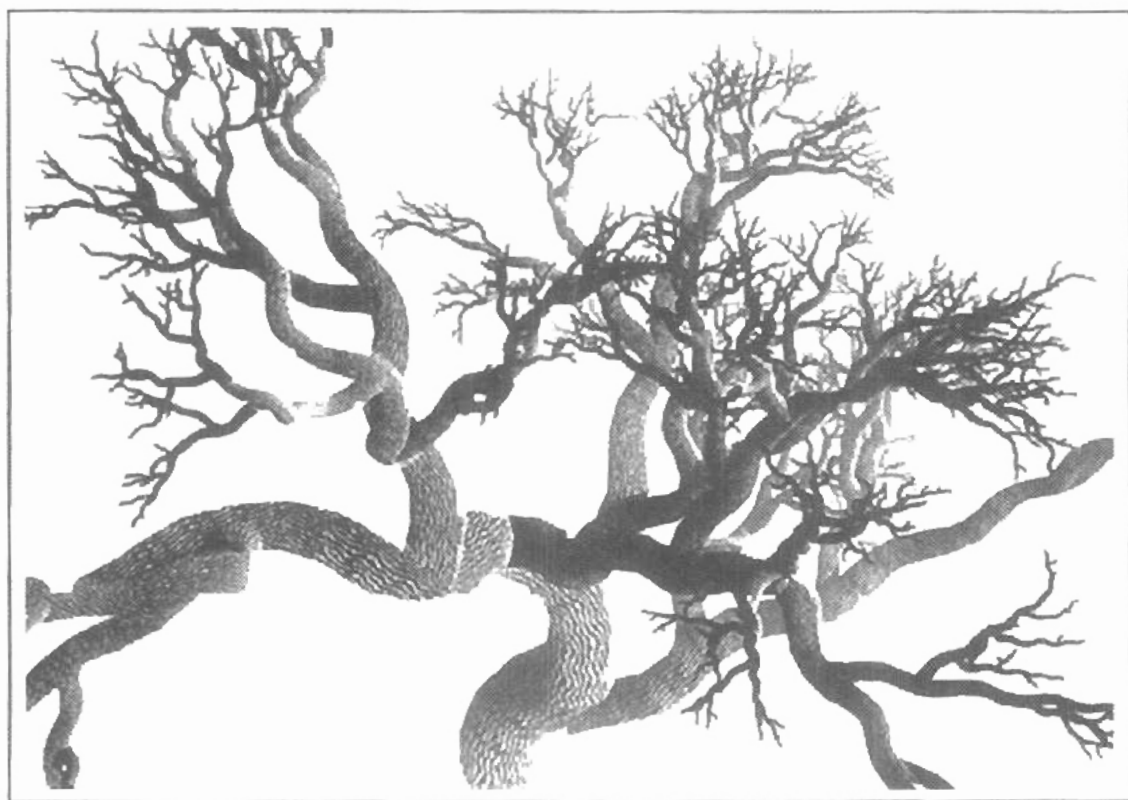


Figure 9-10

Pipeline for transforming a view of a world-coordinate scene to device coordinates.

attributes, or text fonts. Attribute procedures for orienting character strings, however, need to be extended to accommodate arbitrary spatial orientations. Text-attributed routines associated with the up vector require expansion to include z-coordinate data so that strings can be given any spatial orientation. Area-filling routines, such as those for positioning the pattern reference point and for mapping patterns onto a fill area, need to be expanded to accommodate various orientations of the fill-area plane and the pattern plane. Also, most of the two-dimensional structure operations discussed in earlier chapters can be carried over to a three-dimensional package.

Figure 9-10 shows the general stages in a three-dimensional transformation pipeline for displaying a world-coordinate scene. After object definitions have been converted to viewing coordinates and projected to the display plane, scan-conversion algorithms are applied to store the raster image.



Graphics scenes can contain many different kinds of objects: trees, flowers, clouds, rocks, water, bricks, wood paneling, rubber, paper, marble, steel, glass, plastic, and cloth, just to mention a few. So it is probably not too surprising that there is no one method that we can use to describe objects that will include all characteristics of these different materials. And to produce realistic displays of scenes, we need to use representations that accurately model object characteristics.

Polygon and quadric surfaces provide precise descriptions for simple Euclidean objects such as polyhedrons and ellipsoids; spline surfaces and construction techniques are useful for designing aircraft wings, gears, and other engineering structures with curved surfaces; procedural methods, such as fractal constructions and particle systems, allow us to give accurate representations for clouds, clumps of grass, and other natural objects; physically based modeling methods using systems of interacting forces can be used to describe the nonrigid behavior of a piece of cloth or a glob of jello; octree encodings are used to represent internal features of objects, such as those obtained from medical CT images; and isosurface displays, volume renderings, and other visualization techniques are applied to three-dimensional discrete data sets to obtain visual representations of the data.

Representation schemes for solid objects are often divided into two broad categories, although not all representations fall neatly into one or the other of these two categories. **Boundary representations (B-reps)** describe a three-dimensional object as a set of surfaces that separate the object interior from the environment. Typical examples of boundary representations are polygon facets and spline patches. **Space-partitioning representations** are used to describe interior properties, by partitioning the spatial region containing an object into a set of small, nonoverlapping, contiguous solids (usually cubes). A common space-partitioning description for a three-dimensional object is an octree representation. In this chapter, we consider the features of the various representation schemes and how they are used in applications.

10-1

POLYGON SURFACES

The most commonly used boundary representation for a three-dimensional graphics object is a set of surface polygons that enclose the object interior. Many graphics systems store all object descriptions as sets of surface polygons. This simplifies and speeds up the surface rendering and display of objects, since all surfaces are described with linear equations. For this reason, polygon descrip-

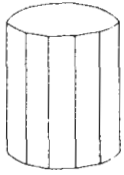


Figure 10-1
Wireframe representation of a cylinder with back (hidden) lines removed.

tions are often referred to as "standard graphics objects." In some cases, a polygonal representation is the only one available, but many packages allow objects to be described with other schemes, such as spline surfaces, that are then converted to polygonal representations for processing.

A polygon representation for a polyhedron precisely defines the surface features of the object. But for other objects, surfaces are *tessellated* (or *tiled*) to produce the polygon-mesh approximation. In Fig. 10-1, the surface of a cylinder is represented as a polygon mesh. Such representations are common in design and solid-modeling applications, since the wireframe outline can be displayed quickly to give a general indication of the surface structure. Realistic renderings are produced by interpolating shading patterns across the polygon surfaces to eliminate or reduce the presence of polygon edge boundaries. And the polygon-mesh approximation to a curved surface can be improved by dividing the surface into smaller polygon facets.

Polygon Tables

We specify a polygon surface with a set of vertex coordinates and associated attribute parameters. As information for each polygon is input, the data are placed into tables that are to be used in the subsequent processing, display, and manipulation of the objects in a scene. Polygon data tables can be organized into two groups: geometric tables and attribute tables. Geometric data tables contain vertex coordinates and parameters to identify the spatial orientation of the polygon surfaces. Attribute information for an object includes parameters specifying the degree of transparency of the object and its surface reflectivity and texture characteristics.

A convenient organization for storing geometric data is to create three lists: a vertex table, an edge table, and a polygon table. Coordinate values for each vertex in the object are stored in the vertex table. The edge table contains pointers back into the vertex table to identify the vertices for each polygon edge. And the polygon table contains pointers back into the edge table to identify the edges for each polygon. This scheme is illustrated in Fig. 10-2 for two adjacent polygons on an object surface. In addition, individual objects and their component polygon faces can be assigned object and facet identifiers for easy reference.

Listing the geometric data in three tables, as in Fig. 10-2, provides a convenient reference to the individual components (vertices, edges, and polygons) of each object. Also, the object can be displayed efficiently by using data from the edge table to draw the component lines. An alternative arrangement is to use just two tables: a vertex table and a polygon table. But this scheme is less convenient, and some edges could get drawn twice. Another possibility is to use only a polygon table, but this duplicates coordinate information, since explicit coordinate values are listed for each vertex in each polygon. Also edge information would have to be reconstructed from the vertex listings in the polygon table.

We can add extra information to the data tables of Fig. 10-2 for faster information extraction. For instance, we could expand the edge table to include forward pointers into the polygon table so that common edges between polygons could be identified more rapidly (Fig. 10-3). This is particularly useful for the rendering procedures that must vary surface shading smoothly across the edges from one polygon to the next. Similarly, the vertex table could be expanded so that vertices are cross-referenced to corresponding edges.

Additional geometric information that is usually stored in the data tables includes the slope for each edge and the coordinate extents for each polygon. As vertices are input, we can calculate edge slopes, and we can scan the coordinate

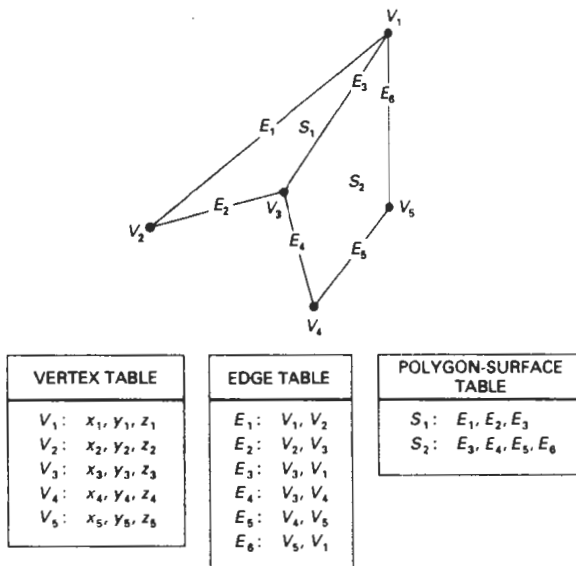


Figure 10-2
Geometric data table representation for two adjacent polygon surfaces, formed with six edges and five vertices.

values to identify the minimum and maximum x , y , and z values for individual polygons. Edge slopes and bounding-box information for the polygons are needed in subsequent processing, for example, surface rendering. Coordinate extents are also used in some visible-surface determination algorithms.

Since the geometric data tables may contain extensive listings of vertices and edges for complex objects, it is important that the data be checked for consistency and completeness. When vertex, edge, and polygon definitions are specified, it is possible, particularly in interactive applications, that certain input errors could be made that would distort the display of the object. The more information included in the data tables, the easier it is to check for errors. Therefore, error checking is easier when three data tables (vertex, edge, and polygon) are used, since this scheme provides the most information. Some of the tests that could be performed by a graphics package are (1) that every vertex is listed as an endpoint for at least two edges, (2) that every edge is part of at least one polygon, (3) that every polygon is closed, (4) that each polygon has at least one shared edge, and (5) that if the edge table contains pointers to polygons, every edge referenced by a polygon pointer has a reciprocal pointer back to the polygon.

E_1 :	V_1, V_2, S_1
E_2 :	V_2, V_3, S_1
E_3 :	V_3, V_1, S_1, S_2
E_4 :	V_3, V_4, S_2
E_5 :	V_4, V_5, S_2
E_6 :	V_5, V_1, S_2

Figure 10-3
Edge table for the surfaces of Fig. 10-2 expanded to include pointers to the polygon table.

Plane Equations

To produce a display of a three-dimensional object, we must process the input data representation for the object through several procedures. These processing steps include transformation of the modeling and world-coordinate descriptions to viewing coordinates, then to device coordinates; identification of visible surfaces; and the application of surface-rendering procedures. For some of these processes, we need information about the spatial orientation of the individual

surface components of the object. This information is obtained from the vertex-coordinate values and the equations that describe the polygon planes.

The equation for a plane surface can be expressed in the form

$$Ax + By + Cz + D = 0 \quad (10-1)$$

where (x, y, z) is any point on the plane, and the coefficients A, B, C , and D are constants describing the spatial properties of the plane. We can obtain the values of A, B, C , and D by solving a set of three plane equations using the coordinate values for three noncollinear points in the plane. For this purpose, we can select three successive polygon vertices, (x_1, y_1, z_1) , (x_2, y_2, z_2) , and (x_3, y_3, z_3) , and solve the following set of simultaneous linear plane equations for the ratios $A/D, B/D$, and C/D :

$$(A/D)x_k + (B/D)y_k + (C/D)z_k = -1, \quad k = 1, 2, 3 \quad (10-2)$$

The solution for this set of equations can be obtained in determinant form, using Cramer's rule, as

$$\begin{aligned} A &= \begin{vmatrix} 1 & y_1 & z_1 \\ 1 & y_2 & z_2 \\ 1 & y_3 & z_3 \end{vmatrix} & B &= \begin{vmatrix} x_1 & 1 & z_1 \\ x_2 & 1 & z_2 \\ x_3 & 1 & z_3 \end{vmatrix} \\ C &= \begin{vmatrix} x_1 & y_1 & 1 \\ x_2 & y_2 & 1 \\ x_3 & y_3 & 1 \end{vmatrix} & D &= - \begin{vmatrix} x_1 & y_1 & z_1 \\ x_2 & y_2 & z_2 \\ x_3 & y_3 & z_3 \end{vmatrix} \end{aligned} \quad (10-3)$$

Expanding the determinants, we can write the calculations for the plane coefficients in the form

$$\begin{aligned} A &= y_1(z_2 - z_3) + y_2(z_3 - z_1) + y_3(z_1 - z_2) \\ B &= z_1(x_2 - x_3) + z_2(x_3 - x_1) + z_3(x_1 - x_2) \\ C &= x_1(y_2 - y_3) + x_2(y_3 - y_1) + x_3(y_1 - y_2) \\ D &= -x_1(y_2z_3 - y_3z_2) - x_2(y_3z_1 - y_1z_3) - x_3(y_1z_2 - y_2z_1) \end{aligned} \quad (10-4)$$

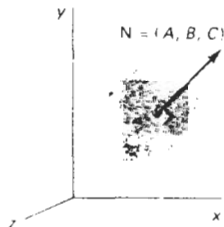


Figure 10-4
The vector N , normal to the surface of a plane described by the equation $Ax + By + Cz + D = 0$, has Cartesian components (A, B, C)

As vertex values and other information are entered into the polygon data structure, values for A, B, C , and D are computed for each polygon and stored with the other polygon data.

Orientation of a plane surface in space can be described with the normal vector to the plane, as shown in Fig. 10-4. This surface normal vector has Cartesian components (A, B, C) , where parameters A, B , and C are the plane coefficients calculated in Eqs. 10-4.

Since we are usually dealing with polygon surfaces that enclose an object interior, we need to distinguish between the two sides of the surface. The side of the plane that faces the object interior is called the "inside" face, and the visible or outward side is the "outside" face. If polygon vertices are specified in a counterclockwise direction when viewing the outer side of the plane in a right-handed coordinate system, the direction of the normal vector will be from inside to outside. This is demonstrated for one plane of a unit cube in Fig. 10-5.

To determine the components of the normal vector for the shaded surface shown in Fig. 10-5, we select three of the four vertices along the boundary of the polygon. These points are selected in a counterclockwise direction as we view from outside the cube toward the origin. Coordinates for these vertices, in the order selected, can be used in Eqs. 10-4 to obtain the plane coefficients: $A = 1$, $B = 0$, $C = 0$, $D = -1$. Thus, the normal vector for this plane is in the direction of the positive x axis.

The elements of the plane normal can also be obtained using a vector cross-product calculation. We again select three vertex positions, V_1 , V_2 , and V_3 , taken in counterclockwise order when viewing the surface from outside to inside in a right-handed Cartesian system. Forming two vectors, one from V_1 to V_2 and the other from V_1 to V_3 , we calculate N as the vector cross product:

$$N = (V_2 - V_1) \times (V_3 - V_1) \quad (10-5)$$

This generates values for the plane parameters A , B , and C . We can then obtain the value for parameter D by substituting these values and the coordinates for one of the polygon vertices in plane equation 10-1 and solving for D . The plane equation can be expressed in vector form using the normal N and the position P of any point in the plane as

$$N \cdot P = -D \quad (10-6)$$

Plane equations are used also to identify the position of spatial points relative to the plane surfaces of an object. For any point (x, y, z) not on a plane with parameters A, B, C, D , we have

$$Ax + By + Cz + D \neq 0$$

We can identify the point as either inside or outside the plane surface according to the sign (negative or positive) of $Ax + By + Cz + D$:

if $Ax + By + Cz + D < 0$, the point (x, y, z) is inside the surface

if $Ax + By + Cz + D > 0$, the point (x, y, z) is outside the surface

These inequality tests are valid in a right-handed Cartesian system, provided the plane parameters A, B, C , and D were calculated using vertices selected in a counterclockwise order when viewing the surface in an outside-to-inside direction. For example, in Fig. 10-5, any point outside the shaded plane satisfies the inequality $x - 1 > 0$, while any point inside the plane has an x -coordinate value less than 1.

Polygon Meshes

Some graphics packages (for example, PHIGS) provide several polygon functions for modeling objects. A single plane surface can be specified with a function such as `fillArea`. But when object surfaces are to be tiled, it is more convenient to specify the surface facets with a mesh function. One type of polygon mesh is the *triangle strip*. This function produces $n - 2$ connected triangles, as shown in Fig. 10-6, given the coordinates for n vertices. Another similar function is the *quadrilateral mesh*, which generates a mesh of $(n - 1)$ by $(m - 1)$ quadrilaterals, given

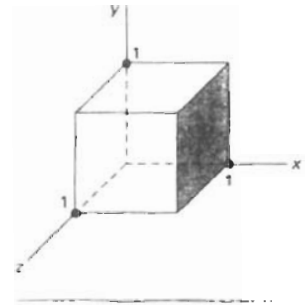


Figure 10-5
The shaded polygon surface of the unit cube has plane equation $x - 1 = 0$ and normal vector $N = (1, 0, 0)$.



Figure 10-6
A triangle strip formed with 11 triangles connecting 13 vertices.

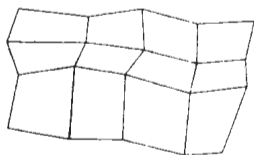


Figure 10-7
A quadrilateral mesh containing 12 quadrilaterals constructed from a 5 by 4 input vertex array.

the coordinates for an n by m array of vertices. Figure 10-7 shows 20 vertices forming a mesh of 12 quadrilaterals.

When polygons are specified with more than three vertices, it is possible that the vertices may not all lie in one plane. This can be due to numerical errors or errors in selecting coordinate positions for the vertices. One way to handle this situation is simply to divide the polygons into triangles. Another approach that is sometimes taken is to approximate the plane parameters A , B , and C . We can do this with averaging methods or we can project the polygon onto the coordinate planes. Using the projection method, we take A proportional to the area of the polygon projection on the yz plane, B proportional to the projection area on the xz plane, and C proportional to the projection area on the xy plane.

High-quality graphics systems typically model objects with polygon meshes and set up a database of geometric and attribute information to facilitate processing of the polygon facets. Fast hardware-implemented polygon renderers are incorporated into such systems with the capability for displaying hundreds of thousands to one million or more shaded polygons per second (usually triangles), including the application of surface texture and special lighting effects.

10-2

CURVED LINES AND SURFACES

Displays of three-dimensional curved lines and surfaces can be generated from an input set of mathematical functions defining the objects or from a set of user-specified data points. When functions are specified, a package can project the defining equations for a curve to the display plane and plot pixel positions along the path of the projected function. For surfaces, a functional description is often tessellated to produce a polygon-mesh approximation to the surface. Usually, this is done with triangular polygon patches to ensure that all vertices of any polygon are in one plane. Polygons specified with four or more vertices may not have all vertices in a single plane. Examples of display surfaces generated from functional descriptions include the quadrics and the superquadrics.

When a set of discrete coordinate points is used to specify an object shape, a functional description is obtained that best fits the designated points according to the constraints of the application. Spline representations are examples of this class of curves and surfaces. These methods are commonly used to design new object shapes, to digitize drawings, and to describe animation paths. Curve-fitting methods are also used to display graphs of data values by fitting specified curve functions to the discrete data set, using regression techniques such as the least-squares method.

Curve and surface equations can be expressed in either a parametric or a nonparametric form. Appendix A gives a summary and comparison of parametric and nonparametric equations. For computer graphics applications, parametric representations are generally more convenient.

10-3

QUADRIC SURFACES

A frequently used class of objects are the *quadric surfaces*, which are described with second-degree equations (quadratics). They include spheres, ellipsoids, tori,

paraboloids, and hyperboloids. Quadric surfaces, particularly spheres and ellipsoids, are common elements of graphics scenes, and they are often available in graphics packages as primitives from which more complex objects can be constructed.

Sphere

In Cartesian coordinates, a spherical surface with radius r centered on the coordinate origin is defined as the set of points (x, y, z) that satisfy the equation

$$x^2 + y^2 + z^2 = r^2 \quad (10-7)$$

We can also describe the spherical surface in parametric form, using latitude and longitude angles (Fig. 10-8):

$$\begin{aligned} x &= r \cos \phi \cos \theta, & -\pi/2 \leq \phi \leq \pi/2 \\ y &= r \cos \phi \sin \theta, & -\pi \leq \theta \leq \pi \\ z &= r \sin \phi \end{aligned} \quad (10-8)$$

The parametric representation in Eqs. 10-8 provides a symmetric range for the angular parameters θ and ϕ . Alternatively, we could write the parametric equations using standard spherical coordinates, where angle ϕ is specified as the colatitude (Fig. 10-9). Then, ϕ is defined over the range $0 \leq \phi \leq \pi$, and θ is often taken in the range $0 \leq \theta \leq 2\pi$. We could also set up the representation using parameters u and v defined over the range from 0 to 1 by substituting $\phi = \pi u$ and $\theta = 2\pi v$.

Ellipsoid

An ellipsoidal surface can be described as an extension of a spherical surface, where the radii in three mutually perpendicular directions can have different values (Fig. 10-10). The Cartesian representation for points over the surface of an ellipsoid centered on the origin is

$$\left(\frac{x}{r_x}\right)^2 + \left(\frac{y}{r_y}\right)^2 + \left(\frac{z}{r_z}\right)^2 = 1 \quad (10-9)$$

And a parametric representation for the ellipsoid in terms of the latitude angle ϕ and the longitude angle θ in Fig. 10-8 is

$$\begin{aligned} x &= r_x \cos \phi \cos \theta, & -\pi/2 \leq \phi \leq \pi/2 \\ y &= r_y \cos \phi \sin \theta, & -\pi \leq \theta \leq \pi \\ z &= r_z \sin \phi \end{aligned} \quad (10-10)$$

Torus

A torus is a doughnut-shaped object, as shown in Fig. 10-11. It can be generated by rotating a circle or other conic about a specified axis. The Cartesian represen-

Section 10-3

Quadric Surfaces

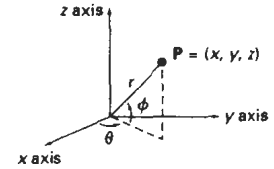


Figure 10-8
Parametric coordinate position (r, θ, ϕ) on the surface of a sphere with radius r .

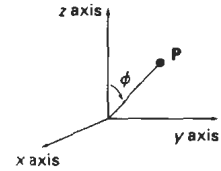


Figure 10-9
Spherical coordinate parameters (r, θ, ϕ) , using colatitude for angle ϕ .

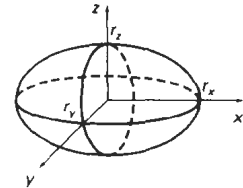


Figure 10-10
An ellipsoid with radii r_x , r_y , and r_z , centered on the coordinate origin.

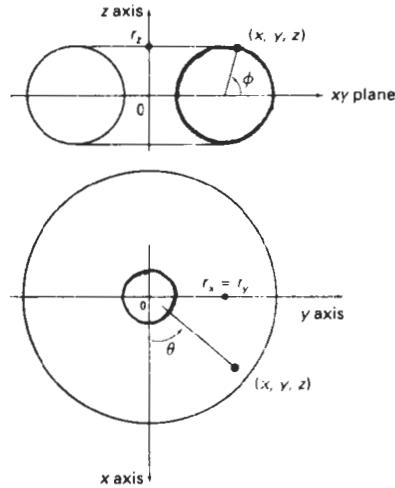


Figure 10-11
A torus with a circular cross section centered on the coordinate origin.

tation for points over the surface of a torus can be written in the form

$$\left[r - \sqrt{\left(\frac{x}{r_x}\right)^2 + \left(\frac{y}{r_y}\right)^2} \right]^2 + \left(\frac{z}{r_z}\right)^2 = 1 \quad (10-11)$$

where r is any given offset value. Parametric representations for a torus are similar to those for an ellipse, except that angle ϕ extends over 360° . Using latitude and longitude angles ϕ and θ , we can describe the torus surface as the set of points that satisfy

$$\begin{aligned} x &= r_x(r + \cos \phi)\cos \theta, & -\pi \leq \phi \leq \pi \\ y &= r_y(r + \cos \phi)\sin \theta, & -\pi \leq \theta \leq \pi \\ z &= r_z \sin \phi \end{aligned} \quad (10-12)$$

10-4

SUPERQUADRICS

This class of objects is a generalization of the quadric representations. **Superquadrics** are formed by incorporating additional parameters into the quadric equations to provide increased flexibility for adjusting object shapes. The number of additional parameters used is equal to the dimension of the object: one parameter for curves and two parameters for surfaces.

Superellipse

We obtain a Cartesian representation for a superellipse from the corresponding equation for an ellipse by allowing the exponent on the x and y terms to be vari-

able. One way to do this is to write the Cartesian superellipse equation in the form

$$\left(\frac{x}{r_x}\right)^{2/s} + \left(\frac{y}{r_y}\right)^{2/s} = 1 \quad (10-13)$$

where parameter s can be assigned any real value. When $s = 1$, we get an ordinary ellipse.

Corresponding parametric equations for the superellipse of Eq. 10-13 can be expressed as

$$\begin{aligned} x &= r_x \cos^s \theta, & -\pi \leq \theta \leq \pi \\ y &= r_y \sin^s \theta \end{aligned} \quad (10-14)$$

Figure 10-12 illustrates supercircle shapes that can be generated using various values for parameter s .

Superellipsoid

A Cartesian representation for a superellipsoid is obtained from the equation for an ellipsoid by incorporating two exponent parameters:

$$\left[\left(\frac{x}{r_x}\right)^{2/s_2} + \left(\frac{y}{r_y}\right)^{2/s_2}\right]^{s_2/s_1} + \left(\frac{z}{r_z}\right)^{2/s_1} = 1 \quad (10-15)$$

For $s_1 = s_2 = 1$, we have an ordinary ellipsoid.

We can then write the corresponding parametric representation for the superellipsoid of Eq. 10-15 as

$$\begin{aligned} x &= r_x \cos^{s_1} \phi \cos^{s_2} \theta, & -\pi/2 \leq \phi \leq \pi/2 \\ y &= r_y \cos^{s_1} \phi \sin^{s_2} \theta, & -\pi \leq \theta \leq \pi \\ z &= r_z \sin^{s_1} \phi \end{aligned} \quad (10-16)$$

Figure 10-13 illustrates supersphere shapes that can be generated using various values for parameters s_1 and s_2 . These and other superquadric shapes can be combined to create more complex structures, such as furniture, threaded bolts, and other hardware.

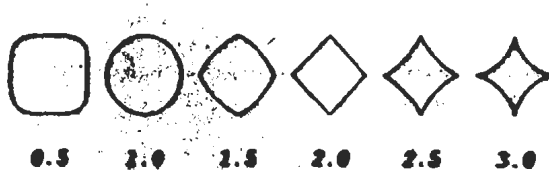


Figure 10-12

Superellipses plotted with different values for parameter s and with $r_x = r_y$.

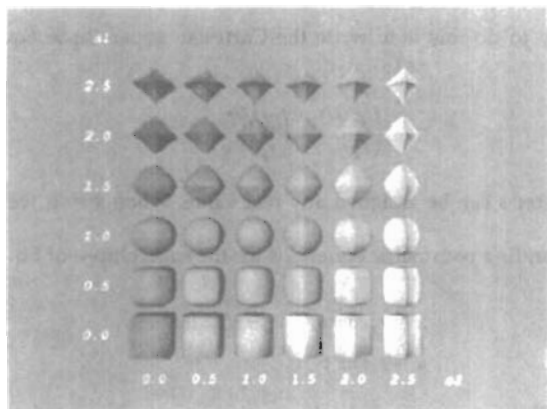


Figure 10-13
Superellipsoids plotted with different values for parameters s_1 and s_2 , and with $r_x = r_y = r_z$.

10-5 BLOBBY OBJECTS

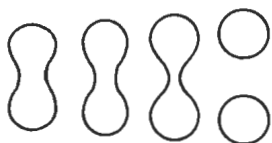


Figure 10-14
Molecular bonding. As two molecules move away from each other, the surface shapes stretch, snap, and finally contract into spheres.

Some objects do not maintain a fixed shape, but change their surface characteristics in certain motions or when in proximity to other objects. Examples in this class of objects include molecular structures, water droplets and other liquid effects, melting objects, and muscle shapes in the human body. These objects can be described as exhibiting “blobbiness” and are often simply referred to as **blobby** objects, since their shapes show a certain degree of fluidity.

A molecular shape, for example, can be described as spherical in isolation, but this shape changes when the molecule approaches another molecule. This distortion of the shape of the electron density cloud is due to the “bonding” that occurs between the two molecules. Figure 10-14 illustrates the stretching, snapping, and contracting effects on molecular shapes when two molecules move apart. These characteristics cannot be adequately described simply with spherical or elliptical shapes. Similarly, Fig. 10-15 shows muscle shapes in a human arm, which exhibit similar characteristics. In this case, we want to model surface shapes so that the total volume remains constant.

Several models have been developed for representing blobby objects as distribution functions over a region of space. One way to do this is to model objects as combinations of Gaussian density functions, or “bumps” (Fig. 10-16). A surface function is then defined as

$$f(x, y, z) = \sum_k b_k e^{-a_k r_k^2} - T = 0 \quad (10-17)$$

where $r_k^2 = \sqrt{x_k^2 + y_k^2 + z_k^2}$, parameter T is some specified threshold, and parameters a and b are used to adjust the amount of blobbiness of the individual objects. Negative values for parameter b can be used to produce dents instead of bumps. Figure 10-17 illustrates the surface structure of a composite object modeled with four Gaussian density functions. At the threshold level, numerical root-finding

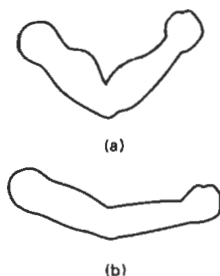


Figure 10-15
Blobby muscle shapes in a human arm.

techniques are used to locate the coordinate intersection values. The cross sections of the individual objects are then modeled as circles or ellipses. If two cross sections are near to each other, they are merged to form one blobby shape, as in Figure 10-14, whose structure depends on the separation of the two objects.

Other methods for generating blobby objects use density functions that fall off to 0 in a finite interval, rather than exponentially. The “metaball” model describes composite objects as combinations of quadratic density functions of the form

$$f(r) = \begin{cases} b(1 - 3r^2/d^2), & \text{if } 0 < r \leq d/3 \\ \frac{3}{2}b(1 - r/d)^2, & \text{if } d/3 < r \leq d \\ 0, & \text{if } r > d \end{cases} \quad (10-18)$$

And the “soft object” model uses the function

$$f(r) = \begin{cases} 1 - \frac{22r^2}{9d^2} + \frac{17r^4}{9d^4} - \frac{4r^6}{9d^6}, & \text{if } 0 < r \leq d \\ 0, & \text{if } r > d \end{cases} \quad (10-19)$$

Some design and painting packages now provide blobby function modeling for handling applications that cannot be adequately modeled with polygon or spline functions alone. Figure 10-18 shows a user interface for a blobby object modeler using metaballs.

Section 10-6

Spline Representations

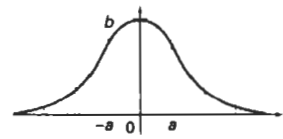


Figure 10-16

A three-dimensional Gaussian bump centered at position 0, with height b and standard deviation a .

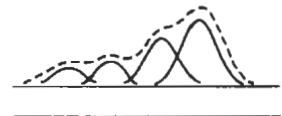


Figure 10-17

A composite blobby object formed with four Gaussian bumps.

10-6

SPLINE REPRESENTATIONS

In drafting terminology, a spline is a flexible strip used to produce a smooth curve through a designated set of points. Several small weights are distributed along the length of the strip to hold it in position on the drafting table as the curve is drawn. The term *spline curve* originally referred to a curve drawn in this manner. We can mathematically describe such a curve with a piecewise cubic

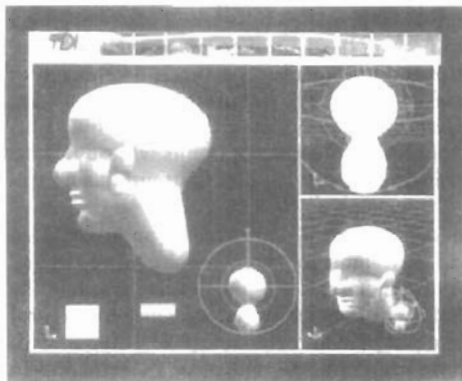


Figure 10-18

A screen layout, used in the Blob Modeler and the Blob Animator packages, for modeling objects with metaballs. (Courtesy of Thomson Digital Image.)

polynomial function whose first and second derivatives are continuous across the various curve sections. In computer graphics, the term **spline curve** now refers to any composite curve formed with polynomial sections satisfying specified continuity conditions at the boundary of the pieces. A **spline surface** can be described with two sets of orthogonal spline curves. There are several different kinds of spline specifications that are used in graphics applications. Each individual specification simply refers to a particular type of polynomial with certain specified boundary conditions.

Splines are used in graphics applications to design curve and surface shapes, to digitize drawings for computer storage, and to specify animation paths for the objects or the camera in a scene. Typical CAD applications for splines include the design of automobile bodies, aircraft and spacecraft surfaces, and ship hulls.



Figure 10-19
A set of six control points
interpolated with piecewise
continuous polynomial
sections.



Figure 10-20
A set of six control points
approximated with piecewise
continuous polynomial
sections

Interpolation and Approximation Splines

We specify a spline curve by giving a set of coordinate positions, called **control points**, which indicates the general shape of the curve. These control points are then fitted with piecewise continuous parametric polynomial functions in one of two ways. When polynomial sections are fitted so that the curve passes through each control point, as in Fig. 10-19, the resulting curve is said to **interpolate** the set of control points. On the other hand, when the polynomials are fitted to the general control-point path without necessarily passing through any control point, the resulting curve is said to **approximate** the set of control points (Fig. 10-20).

Interpolation curves are commonly used to digitize drawings or to specify animation paths. Approximation curves are primarily used as design tools to structure object surfaces. Figure 10-21 shows an approximation spline surface created for a design application. Straight lines connect the control-point positions above the surface.

A spline curve is defined, modified, and manipulated with operations on the control points. By interactively selecting spatial positions for the control points, a designer can set up an initial curve. After the polynomial fit is displayed for a given set of control points, the designer can then reposition some or all of the control points to restructure the shape of the curve. In addition, the curve can be translated, rotated, or scaled with transformations applied to the control points. CAD packages can also insert extra control points to aid a designer in adjusting the curve shapes.

The convex polygon boundary that encloses a set of control points is called the **convex hull**. One way to envision the shape of a convex hull is to imagine a rubber band stretched around the positions of the control points so that each control point is either on the perimeter of the hull or inside it (Fig. 10-22). Convex hulls provide a measure for the deviation of a curve or surface from the region bounding the control points. Some splines are bounded by the convex hull, thus ensuring that the polynomials smoothly follow the control points without erratic oscillations. Also, the polygon region inside the convex hull is useful in some algorithms as a clipping region.

A polyline connecting the sequence of control points for an approximation spline is usually displayed to remind a designer of the control-point ordering. This set of connected line segments is often referred to as the **control graph** of the curve. Other names for the series of straight-line sections connecting the control points in the order specified are **control polygon** and **characteristic polygon**. Figure 10-23 shows the shape of the control graph for the control-point sequences in Fig. 10-22.

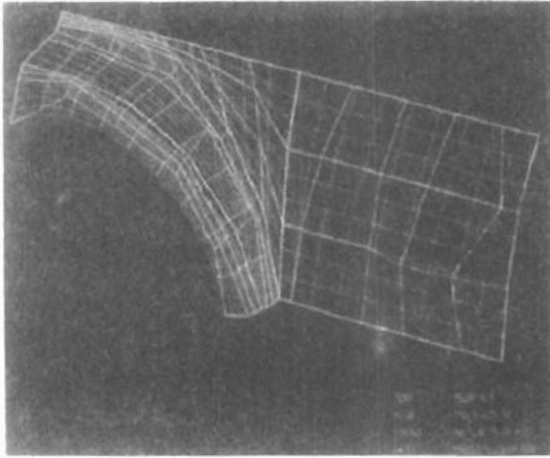


Figure 10-21
An approximation spline surface for a CAD application in automotive design. Surface contours are plotted with polynomial curve sections, and the surface control points are connected with straight-line segments. (Courtesy of Evans & Sutherland.)

Parametric Continuity Conditions

To ensure a smooth transition from one section of a piecewise parametric curve to the next, we can impose various **continuity conditions** at the connection points. If each section of a spline is described with a set of parametric coordinate functions of the form

$$x = x(u), \quad y = y(u), \quad z = z(u), \quad u_1 \leq u \leq u_2 \quad (10-20)$$

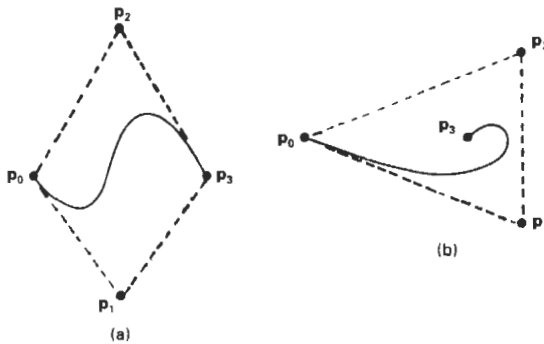


Figure 10-22
Convex-hull shapes (dashed lines) for two sets of control points.

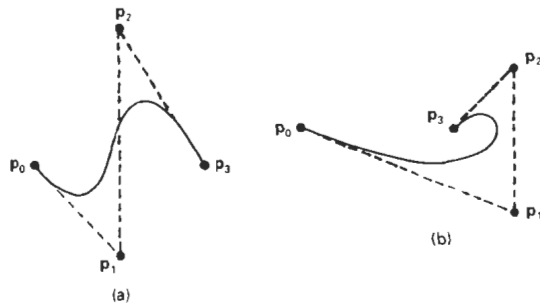


Figure 10-23
Control-graph shapes (dashed lines) for two different sets of control points.

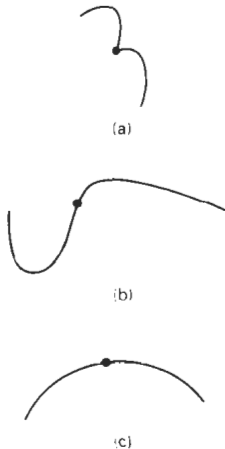


Figure 10-24
Piecewise construction of a curve by joining two curve segments using different orders of continuity: (a) zero-order continuity only, (b) first-order continuity, and (c) second-order continuity.

we set **parametric continuity** by matching the parametric derivatives of adjoining curve sections at their common boundary.

Zero-order parametric continuity, described as C^0 continuity, means simply that the curves meet. That is, the values of x , y , and z evaluated at u_2 for the first curve section are equal, respectively, to the values of x , y , and z evaluated at u_1 for the next curve section. **First-order parametric continuity**, referred to as C^1 continuity, means that the first parametric derivatives (tangent lines) of the coordinate functions in Eq. 10-20 for two successive curve sections are equal at their joining point. **Second-order parametric continuity**, or C^2 continuity, means that both the first and second parametric derivatives of the two curve sections are the same at the intersection. Higher-order parametric continuity conditions are defined similarly. Figure 10-24 shows examples of C^0 , C^1 , and C^2 continuity.

With second-order continuity, the rates of change of the tangent vectors for connecting sections are equal at their intersection. Thus, the tangent line transitions smoothly from one section of the curve to the next (Fig. 10-24(c)). But with first-order continuity, the rates of change of the tangent vectors for the two sections can be quite different (Fig. 10-24(b)), so that the general shapes of the two adjacent sections can change abruptly. First-order continuity is often sufficient for digitizing drawings and some design applications, while second-order continuity is useful for setting up animation paths for camera motion and for many precision CAD requirements. A camera traveling along the curve path in Fig. 10-24(b) with equal steps in parameter u would experience an abrupt change in acceleration at the boundary of the two sections, producing a discontinuity in the motion sequence. But if the camera were traveling along the path in Fig. 10-24(c), the frame sequence for the motion would smoothly transition across the boundary.

Geometric Continuity Conditions

An alternate method for joining two successive curve sections is to specify conditions for **geometric continuity**. In this case, we only require parametric derivatives of the two sections to be proportional to each other at their common boundary instead of equal to each other.

Zero-order geometric continuity, described as G^0 continuity, is the same as zero-order parametric continuity. That is, the two curves sections must have the

same coordinate position at the boundary point. **First-order geometric continuity**, or G^1 continuity, means that the parametric first derivatives are proportional at the intersection of two successive sections. If we denote the parametric position on the curve as $P(u)$, the direction of the tangent vector $P'(u)$, but not necessarily its magnitude, will be the same for two successive curve sections at their joining point under G^1 continuity. **Second-order geometric continuity**, or G^2 continuity, means that both the first and second parametric derivatives of the two curve sections are proportional at their boundary. Under G^2 continuity, curvatures of two curve sections will match at the joining position.

A curve generated with geometric continuity conditions is similar to one generated with parametric continuity, but with slight differences in curve shape. Figure 10-25 provides a comparison of geometric and parametric continuity. With geometric continuity, the curve is pulled toward the section with the greater tangent vector.

Spline Specifications

There are three equivalent methods for specifying a particular spline representation: (1) We can state the set of boundary conditions that are imposed on the spline; or (2) we can state the matrix that characterizes the spline; or (3) we can state the set of **blending functions** (or **basis functions**) that determine how specified geometric constraints on the curve are combined to calculate positions along the curve path.

To illustrate these three equivalent specifications, suppose we have the following parametric cubic polynomial representation for the x coordinate along the path of a spline section:

$$x(u) = a_x u^3 + b_x u^2 + c_x u + d_x, \quad 0 \leq u \leq 1 \quad (10-21)$$

Boundary conditions for this curve might be set, for example, on the endpoint coordinates $x(0)$ and $x(1)$ and on the parametric first derivatives at the endpoints $x'(0)$ and $x'(1)$. These four boundary conditions are sufficient to determine the values of the four coefficients a_x , b_x , c_x , and d_x .

From the boundary conditions, we can obtain the matrix that characterizes this spline curve by first rewriting Eq. 10-21 as the matrix product

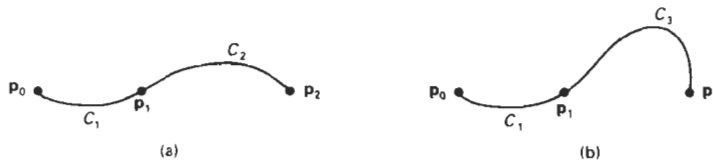


Figure 10-25
Three control points fitted with two curve sections joined with (a) parametric continuity and (b) geometric continuity, where the tangent vector of curve C_2 at point P_1 has a greater magnitude than the tangent vector of curve C_1 at P_1 .

$$\begin{aligned}
 x(u) &= [u^3 \ u^2 \ u \ 1] \begin{bmatrix} a_x \\ b_x \\ c_x \\ d_x \end{bmatrix} \\
 &= \mathbf{U} \cdot \mathbf{C}
 \end{aligned} \tag{10-22}$$

where \mathbf{U} is the row matrix of powers of parameter u , and \mathbf{C} is the coefficient column matrix. Using Eq. 10-22, we can write the boundary conditions in matrix form and solve for the coefficient matrix \mathbf{C} as

$$\mathbf{C} = \mathbf{M}_{\text{spline}} \cdot \mathbf{M}_{\text{geom}} \tag{10-23}$$

where \mathbf{M}_{geom} is a four-element column matrix containing the geometric constraint values (boundary conditions) on the spline; and $\mathbf{M}_{\text{spline}}$ is the 4-by-4 matrix that transforms the geometric constraint values to the polynomial coefficients and provides a characterization for the spline curve. Matrix \mathbf{M}_{geom} contains control-point coordinate values and other geometric constraints that have been specified. Thus, we can substitute the matrix representation for \mathbf{C} into Eq. 10-22 to obtain

$$x(u) = \mathbf{U} \cdot \mathbf{M}_{\text{spline}} \cdot \mathbf{M}_{\text{geom}} \tag{10-24}$$

The matrix, $\mathbf{M}_{\text{spline}}$, characterizing a spline representation, sometimes called the *basis matrix*, is particularly useful for transforming from one spline representation to another.

Finally, we can expand Eq. 10-24 to obtain a polynomial representation for coordinate x in terms of the geometric constraint parameters

$$x(u) = \sum_{k=0}^3 g_k \cdot BF_k(u) \tag{10-25}$$

where g_k are the constraint parameters, such as the control-point coordinates and slope of the curve at the control points, and $BF_k(u)$ are the polynomial blending functions. In the following sections, we discuss some commonly used splines and their matrix and blending-function specifications.

10-7

CUBIC SPLINE INTERPOLATION METHODS

This class of splines is most often used to set up paths for object motions or to provide a representation for an existing object or drawing, but interpolation splines are also used sometimes to design object shapes. Cubic polynomials offer a reasonable compromise between flexibility and speed of computation. Compared to higher-order polynomials, cubic splines require less calculations and memory and they are more stable. Compared to lower-order polynomials, cubic splines are more flexible for modeling arbitrary curve shapes.

Given a set of control points, cubic interpolation splines are obtained by fitting the input points with a piecewise cubic polynomial curve that passes through every control point. Suppose we have $n + 1$ control points specified with coordinates

$$\mathbf{p}_k = (x_k, y_k, z_k), \quad k = 0, 1, 2, \dots, n$$

A cubic interpolation fit of these points is illustrated in Fig. 10-26. We can describe the parametric cubic polynomial that is to be fitted between each pair of control points with the following set of equations:

$$\begin{aligned}x(u) &= a_x u^3 + b_x u^2 + c_x u + d_x \\y(u) &= a_y u^3 + b_y u^2 + c_y u + d_y, \quad (0 \leq u \leq 1) \\z(u) &= a_z u^3 + b_z u^2 + c_z u + d_z\end{aligned}\quad (10-26)$$

For each of these three equations, we need to determine the values of the four coefficients a , b , c , and d in the polynomial representation for each of the n curve sections between the $n + 1$ control points. We do this by setting enough boundary conditions at the "joints" between curve sections so that we can obtain numerical values for all the coefficients. In the following sections, we discuss common methods for setting the boundary conditions for cubic interpolation splines.

Natural Cubic Splines

One of the first spline curves to be developed for graphics applications is the **natural cubic spline**. This interpolation curve is a mathematical representation of the original drafting spline. We formulate a natural cubic spline by requiring that two adjacent curve sections have the same first and second parametric derivatives at their common boundary. Thus, natural cubic splines have C^2 continuity.

If we have $n + 1$ control points to fit, as in Fig. 10-26, then we have n curve sections with a total of $4n$ polynomial coefficients to be determined. At each of the $n - 1$ interior control points, we have four boundary conditions: The two curve sections on either side of a control point must have the same first and second parametric derivatives at that control point, and each curve must pass through that control point. This gives us $4n - 4$ equations to be satisfied by the $4n$ polynomial coefficients. We get an additional equation from the first control point p_0 , the position of the beginning of the curve, and another condition from control point p_n , which must be the last point on the curve. We still need two more conditions to be able to determine values for all coefficients. One method for obtaining the two additional conditions is to set the second derivatives at p_0 and p_n to 0. Another approach is to add two extra "dummy" control points, one at each end of the original control-point sequence. That is, we add a control point p_{-1} and a control point p_{n+1} . Then all of the original control points are interior points, and we have the necessary $4n$ boundary conditions.

Although natural cubic splines are a mathematical model for the drafting spline, they have a major disadvantage. If the position of any one control point is altered, the entire curve is affected. Thus, natural cubic splines allow for no "local control", so that we cannot restructure part of the curve without specifying an entirely new set of control points.

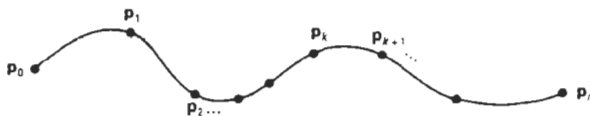


Figure 10-26
A piecewise continuous cubic-spline interpolation of $n + 1$ control points.

Hermite Interpolation

A **Hermite spline** (named after the French mathematician Charles Hermite) is an interpolating piecewise cubic polynomial with a specified tangent at each control point. Unlike the natural cubic splines, Hermite splines can be adjusted locally because each curve section is only dependent on its endpoint constraints.

If $\mathbf{P}(u)$ represents a parametric cubic point function for the curve section between control points \mathbf{p}_k and \mathbf{p}_{k+1} , as shown in Fig. 10-27, then the boundary conditions that define this Hermite curve section are

$$\begin{aligned}\mathbf{P}(0) &= \mathbf{p}_k \\ \mathbf{P}(1) &= \mathbf{p}_{k+1} \\ \mathbf{P}'(0) &= \mathbf{Dp}_k \\ \mathbf{P}'(1) &= \mathbf{Dp}_{k+1}\end{aligned}\quad (10-27)$$

with \mathbf{Dp}_k and \mathbf{Dp}_{k+1} specifying the values for the parametric derivatives (slope of the curve) at control points \mathbf{p}_k and \mathbf{p}_{k+1} , respectively.

We can write the vector equivalent of Eqs. 10-26 for this Hermite-curve section as

$$\mathbf{P}(u) = \mathbf{a}u^3 + \mathbf{b}u^2 + \mathbf{c}u + \mathbf{d}, \quad 0 \leq u \leq 1 \quad (10-28)$$

where the x component of \mathbf{P} is $x(u) = a_x u^3 + b_x u^2 + c_x u + d_x$, and similarly for the y and z components. The matrix equivalent of Eq. 10-28 is

$$\mathbf{P}(u) = [u^3 \ u^2 \ u \ 1] \cdot \begin{bmatrix} \mathbf{a} \\ \mathbf{b} \\ \mathbf{c} \\ \mathbf{d} \end{bmatrix} \quad (10-29)$$

and the derivative of the point function can be expressed as

$$\mathbf{P}'(u) = [3u^2 \ 2u \ 1 \ 0] \cdot \begin{bmatrix} \mathbf{a} \\ \mathbf{b} \\ \mathbf{c} \\ \mathbf{d} \end{bmatrix} \quad (10-30)$$

Substituting endpoint values 0 and 1 for parameter u into the previous two equations, we can express the Hermite boundary conditions 10-27 in the matrix form:

$$\begin{bmatrix} \mathbf{p}_k \\ \mathbf{p}_{k+1} \\ \mathbf{Dp}_k \\ \mathbf{Dp}_{k+1} \end{bmatrix} = \begin{bmatrix} 0 & 0 & 0 & 1 \\ 1 & 1 & 1 & 1 \\ 0 & 0 & 1 & 0 \\ 3 & 2 & 1 & 0 \end{bmatrix} \cdot \begin{bmatrix} \mathbf{a} \\ \mathbf{b} \\ \mathbf{c} \\ \mathbf{d} \end{bmatrix} \quad (10-31)$$

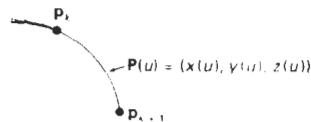


Figure 10-27
Parametric point function $\mathbf{P}(u)$ for a Hermite curve section between control points \mathbf{p}_k and \mathbf{p}_{k+1} .

Solving this equation for the polynomial coefficients, we have

$$\begin{aligned} \begin{bmatrix} a \\ b \\ c \\ d \end{bmatrix} &= \begin{bmatrix} 0 & 0 & 0 & 1 \\ 1 & 1 & 1 & 1 \\ 0 & 0 & 1 & 0 \\ 3 & 2 & 1 & 0 \end{bmatrix}^{-1} \cdot \begin{bmatrix} p_k \\ p_{k+1} \\ Dp_k \\ Dp_{k+1} \end{bmatrix} \\ &= \begin{bmatrix} 2 & -2 & 1 & 1 \\ -3 & 3 & -2 & -1 \\ 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 \end{bmatrix} \cdot \begin{bmatrix} p_k \\ p_{k+1} \\ Dp_k \\ Dp_{k+1} \end{bmatrix} \\ &= M_H \cdot \begin{bmatrix} p_k \\ p_{k+1} \\ Dp_k \\ Dp_{k+1} \end{bmatrix} \end{aligned} \quad (10-32)$$

where M_H , the Hermite matrix, is the inverse of the boundary constraint matrix. Equation 10-29 can thus be written in terms of the boundary conditions as

$$P(u) = [u^3 \ u^2 \ u \ 1] \cdot M_H \cdot \begin{bmatrix} p_k \\ p_{k+1} \\ Dp_k \\ Dp_{k+1} \end{bmatrix} \quad (10-33)$$

Finally, we can determine expressions for the Hermite blending functions by carrying out the matrix multiplications in Eq. 10-33 and collecting coefficients for the boundary constraints to obtain the polynomial form:

$$\begin{aligned} P(u) &= p_k(2u^3 - 3u^2 + 1) + p_{k+1}(-2u^3 + 3u^2) - Dp_k(u^3 - 2u^2 + u) \\ &\quad + Dp_{k+1}(u^3 - u^2) \\ &= p_k H_0(u) + p_{k+1} H_1(u) + Dp_k H_2(u) + Dp_{k+1} H_3(u) \end{aligned} \quad (10-34)$$

The polynomials $H_k(u)$ for $k = 0, 1, 2, 3$ are referred to as blending functions because they blend the boundary constraint values (endpoint coordinates and slopes) to obtain each coordinate position along the curve. Figure 10-28 shows the shape of the four Hermite blending functions.

Hermite polynomials can be useful for some digitizing applications where it may not be too difficult to specify or approximate the curve slopes. But for most problems in computer graphics, it is more useful to generate spline curves without requiring input values for curve slopes or other geometric information, in addition to control-point coordinates. Cardinal splines and Kochanek-Bartels splines, discussed in the following two sections, are variations on the Hermite splines that do not require input values for the curve derivatives at the control points. Procedures for these splines compute parametric derivatives from the coordinate positions of the control points.

Cardinal Splines

As with Hermite splines, **cardinal splines** are interpolating piecewise cubics with specified endpoint tangents at the boundary of each curve section. The difference

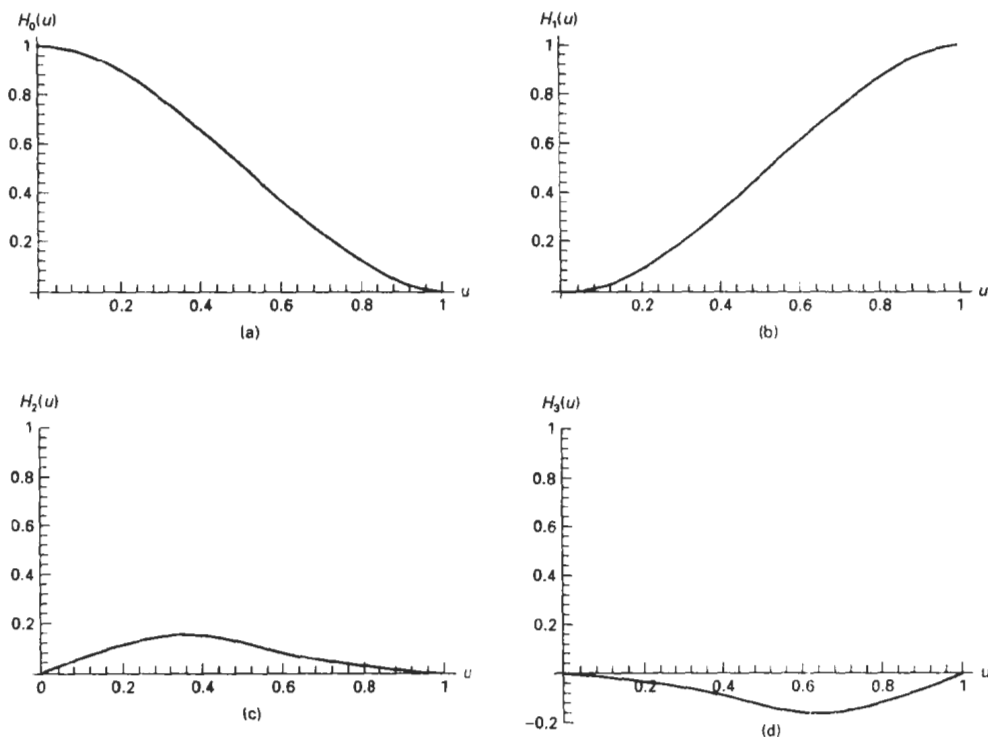


Figure 10-28
The Hermite blending functions.

is that we do not have to give the values for the endpoint tangents. For a cardinal spline, the value for the slope at a control point is calculated from the coordinates of the two adjacent control points.

A cardinal spline section is completely specified with four consecutive control points. The middle two control points are the section endpoints, and the other two points are used in the calculation of the endpoint slopes. If we take $P(u)$ as the representation for the parametric cubic point function for the curve section between control points p_k and p_{k+1} , as in Fig. 10-29, then the four control points from p_{k-1} to p_{k+1} are used to set the boundary conditions for the cardinal-spline section as

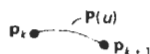


Figure 10-29
Parametric point function $P(u)$ for a cardinal-spline section between control points p_k and p_{k+1} .

$$\begin{aligned}
 P(0) &= p_k \\
 P(1) &= p_{k+1} \\
 P'(0) &= \frac{1}{2}(1-t)(p_{k+1} - p_{k-1}) \\
 P'(1) &= \frac{1}{2}(1-t)(p_{k+2} - p_k)
 \end{aligned}
 \tag{10-35}$$

Thus, the slopes at control points p_k and p_{k+1} are taken to be proportional, respectively, to the chords $\overline{p_{k-1}p_{k+1}}$ and $\overline{p_k p_{k+2}}$ (Fig. 10-30). Parameter t is called the **tension** parameter since it controls how loosely or tightly the cardinal spline fits

the input control points. Figure 10-31 illustrates the shape of a cardinal curve for very small and very large values of tension t . When $t = 0$, this class of curves is referred to as **Catmull-Rom splines**, or **Overhauser splines**.

Using methods similar to those for Hermite splines, we can convert the boundary conditions 10-35 into the matrix form

$$\mathbf{P}(u) = [u^3 \ u^2 \ u \ 1] \cdot \mathbf{M}_C \cdot \begin{bmatrix} \mathbf{p}_{k-1} \\ \mathbf{p}_k \\ \mathbf{p}_{k+1} \\ \mathbf{p}_{k+2} \end{bmatrix} \quad (10-36)$$

where the cardinal matrix is

$$\mathbf{M}_C = \begin{bmatrix} -s & 2-s & s-2 & s \\ 2s & s-3 & 3-2s & -s \\ -s & 0 & s & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix} \quad (10-37)$$

with $s = (1 - t)/2$.

Expanding matrix equation 10-36 into polynomial form, we have

$$\begin{aligned} \mathbf{P}(u) &= \mathbf{p}_{k-1}(-su^3 + 2su^2 - su) + \mathbf{p}_k[(2-s)u^3 + (s-3)u^2 + 1] \\ &\quad + \mathbf{p}_{k+1}[(s-2)u^3 + (3-2s)u^2 + su] + \mathbf{p}_{k+2}(su^3 - su^2) \quad (10-38) \\ &= \mathbf{p}_{k-1}\text{CAR}_0(u) + \mathbf{p}_k\text{CAR}_1(u) + \mathbf{p}_{k+1}\text{CAR}_2(u) + \mathbf{p}_{k+2}\text{CAR}_3(u) \end{aligned}$$

where the polynomials $\text{CAR}_k(u)$ for $k = 0, 1, 2, 3$ are the cardinal blending functions. Figure 10-32 gives a plot of the basis functions for cardinal splines with $t = 0$.

Kochanek-Bartels Splines

These interpolating cubic polynomials are extensions of the cardinal splines. Two additional parameters are introduced into the constraint equations defining **Kochanek-Bartels splines** to provide for further flexibility in adjusting the shape of curve sections.

Given four consecutive control points, labeled \mathbf{p}_{k-1} , \mathbf{p}_k , \mathbf{p}_{k+1} , and \mathbf{p}_{k+2} , we define the boundary conditions for a Kochanek-Bartels curve section between \mathbf{p}_k and \mathbf{p}_{k+1} as

$$\begin{aligned} \mathbf{P}(0) &= \mathbf{p}_k \\ \mathbf{P}(1) &= \mathbf{p}_{k+1} \\ \mathbf{P}'(0)_{\text{in}} &= \frac{1}{2}(1-t)[(1+b)(1-c)(\mathbf{p}_k - \mathbf{p}_{k-1}) \\ &\quad + (1-b)(1+c)(\mathbf{p}_{k+1} - \mathbf{p}_k)] \\ \mathbf{P}'(1)_{\text{out}} &= \frac{1}{2}(1-t)[(1+b)(1+c)(\mathbf{p}_{k+1} - \mathbf{p}_k) \\ &\quad + (1-b)(1-c)(\mathbf{p}_{k+2} - \mathbf{p}_{k+1})] \end{aligned} \quad (10-39)$$

where t is the **tension** parameter, b is the **bias** parameter, and c is the **continuity** parameter. In the Kochanek-Bartels formulation, parametric derivatives may not be continuous across section boundaries.

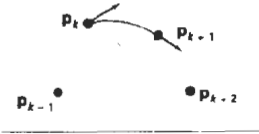


Figure 10-30
Tangent vectors at the endpoints of a cardinal-spline section are proportional to the chords formed with neighboring control points (dashed lines).

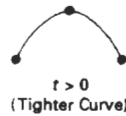
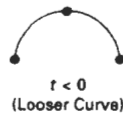


Figure 10-31
Effect of the tension parameter on the shape of a cardinal spline section.

Tension parameter t has the same interpretation as in the cardinal-spline formulation; that is, it controls the looseness or tightness of the curve sections. Bias (b) is used to adjust the amount that the curve bends at each end of a section, so that curve sections can be skewed toward one end or the other (Fig. 10-33). Parameter c controls the continuity of the tangent vector across the boundaries of sections. If c is assigned a nonzero value, there is a discontinuity in the slope of the curve across section boundaries.

Kochanek-Bartel splines were designed to model animation paths. In particular, abrupt changes in motion of a object can be simulated with nonzero values for parameter c .

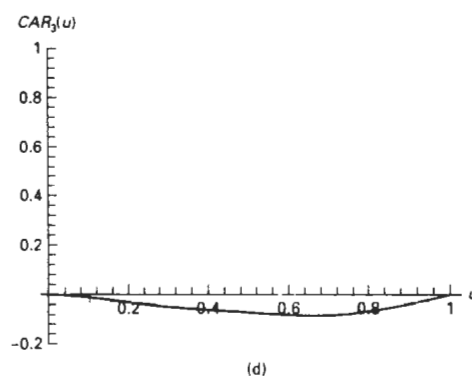
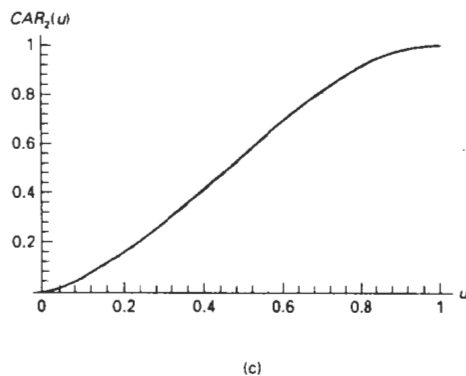
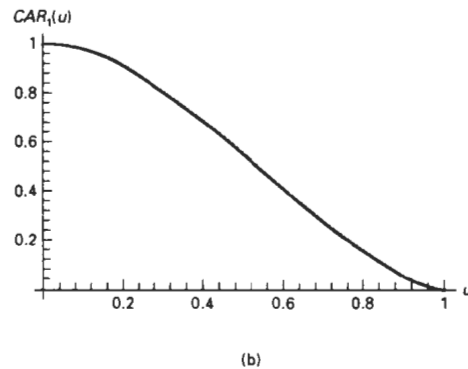
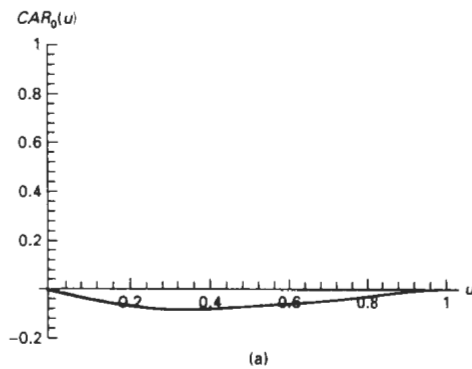


Figure 10-32
The cardinal blending functions for $t = 0$ and $s = 0.5$.

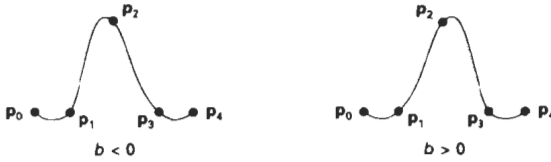


Figure 10-33
Effect of the bias parameter on the shape of a
Kochanek-Bartels spline section.

10-8

BÉZIER CURVES AND SURFACES

This spline approximation method was developed by the French engineer Pierre Bézier for use in the design of Renault automobile bodies. **Bézier splines** have a number of properties that make them highly useful and convenient for curve and surface design. They are also easy to implement. For these reasons, Bézier splines are widely available in various CAD systems, in general graphics packages (such as GL on Silicon Graphics systems), and in assorted drawing and painting packages (such as Aldus SuperPaint and Cricket Draw).

Bézier Curves

In general, a Bézier curve section can be fitted to any number of control points. The number of control points to be approximated and their relative position determine the degree of the Bézier polynomial. As with the interpolation splines, a Bézier curve can be specified with boundary conditions, with a characterizing matrix, or with blending functions. For general Bézier curves, the blending-function specification is the most convenient.

Suppose we are given $n + 1$ control-point positions: $\mathbf{p}_k = (x_k, y_k, z_k)$, with k varying from 0 to n . These coordinate points can be blended to produce the following position vector $\mathbf{P}(u)$, which describes the path of an approximating Bézier polynomial function between \mathbf{p}_0 and \mathbf{p}_n .

$$\mathbf{P}(u) = \sum_{k=0}^n \mathbf{p}_k \text{BEZ}_{k,n}(u), \quad 0 \leq u \leq 1 \quad (10-40)$$

The Bézier blending functions $\text{BEZ}_{k,n}(u)$ are the *Bernstein polynomials*:

$$\text{BEZ}_{k,n}(u) = C(n, k)u^k(1 - u)^{n-k} \quad (10-41)$$

where the $C(n, k)$ are the binomial coefficients:

$$C(n, k) = \frac{n!}{k!(n - k)!} \quad (10-42)$$

Equivalently, we can define Bézier blending functions with the recursive calculation

$$\text{BEZ}_{k,n}(u) = (1 - u) \text{BEZ}_{k,n-1}(u) + u \text{BEZ}_{k-1,n-1}(u), \quad n > k \geq 1 \quad (10-43)$$

with $BEZ_{k,k} = u^k$, and $BEZ_{0,k} = (1 - u)^k$. Vector equation 10-40 represents a set of three parametric equations for the individual curve coordinates:

$$\begin{aligned} x(u) &= \sum_{k=0}^n x_k BEZ_{k,n}(u) \\ y(u) &= \sum_{k=0}^n y_k BEZ_{k,n}(u) \\ z(u) &= \sum_{k=0}^n z_k BEZ_{k,n}(u) \end{aligned} \quad (10-44)$$

As a rule, a Bézier curve is a polynomial of degree one less than the number of control points used: Three points generate a parabola, four points a cubic curve, and so forth. Figure 10-34 demonstrates the appearance of some Bézier curves for various selections of control points in the xy plane ($z = 0$). With certain control-point placements, however, we obtain degenerate Bézier polynomials. For example, a Bézier curve generated with three collinear control points is a straight-line segment. And a set of control points that are all at the same coordinate position produces a Bézier "curve" that is a single point.

Bézier curves are commonly found in painting and drawing packages, as well as CAD systems, since they are easy to implement and they are reasonably powerful in curve design. Efficient methods for determining coordinate positions along a Bézier curve can be set up using recursive calculations. For example, successive binomial coefficients can be calculated as

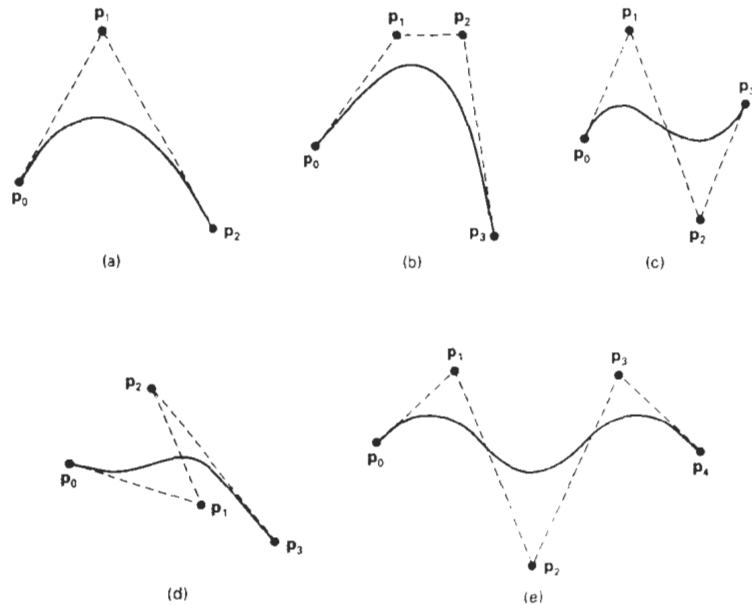


Figure 10-34
Examples of two-dimensional Bézier curves generated from three, four, and five control points. Dashed lines connect the control-point positions.

$$C(n, k) = \frac{n-k+1}{k} C(n, k-1) \quad (10-45)$$

for $n \geq k$. The following example program illustrates a method for generating Bézier curves.

```
#include <math.h>
#include "graphics.h"

void computeCoefficients (int n, int * c)
{
    int k, i;

    for (k=0; k<=n; k++) {
        /* Compute n!/(k!(n-k)!) */
        c[k] = 1;
        for (i=n; i>=k+1; i--)
            c[k] *= i;
        for (i=n-k; i>=2; i--)
            c[k] /= i;
    }
}

void computePoint
(float u, wcPt3 * pt, int nControls, wcPt3 * controls, int * c)
{
    int k, n = nControls - 1;
    float blend;

    pt->x = 0.0; pt->y = 0.0; pt->z = 0.0;

    /* Add in influence of each control point */
    for (k=0; k<nControls; k++) {
        blend = c[k] * powf (u,k) * powf (1-u,n-k);
        pt->x += controls[k].x * blend;
        pt->y += controls[k].y * blend;
        pt->z += controls[k].z * blend;
    }
}

void bezier (wcPt3 * controls, int nControls, int m, wcPt3 * curve)
{
    /* Allocate space for the coefficients */
    int * c = (int *) malloc (nControls * sizeof (int));
    int i;

    computeCoefficients (nControls-1, c);
    for (i=0; i<=m; i++)
        computePoint (i / (float) m, &curve[i], nControls, controls, c);
    free (c);
}
```

Properties of Bézier Curves

A very useful property of a Bézier curve is that it always passes through the first and last control points. That is, the boundary conditions at the two ends of the curve are

$$\begin{aligned} P(0) &= p_0 \\ P(1) &= p_n \end{aligned} \quad (10-46)$$

Values of the parametric first derivatives of a Bézier curve at the endpoints can be calculated from control-point coordinates as

$$\begin{aligned} \mathbf{P}'(0) &= -n\mathbf{p}_0 + n\mathbf{p}_1 \\ \mathbf{P}'(1) &= -n\mathbf{p}_{n-1} + n\mathbf{p}_n \end{aligned} \quad (10-47)$$

Thus, the slope at the beginning of the curve is along the line joining the first two control points, and the slope at the end of the curve is along the line joining the last two endpoints. Similarly, the parametric second derivatives of a Bézier curve at the endpoints are calculated as

$$\begin{aligned} \mathbf{P}''(0) &= n(n-1)[(\mathbf{p}_2 - \mathbf{p}_1) - (\mathbf{p}_1 - \mathbf{p}_0)] \\ \mathbf{P}''(1) &= n(n-1)[(\mathbf{p}_{n-2} - \mathbf{p}_{n-1}) - (\mathbf{p}_{n-1} - \mathbf{p}_n)] \end{aligned} \quad (10-48)$$

Another important property of any Bézier curve is that it lies within the convex hull (convex polygon boundary) of the control points. This follows from the properties of Bézier blending functions: They are all positive and their sum is always 1,

$$\sum_{k=0}^n \text{BEZ}_{k,n}(u) = 1 \quad (10-49)$$

so that any curve position is simply the weighted sum of the control-point positions. The convex-hull property for a Bézier curve ensures that the polynomial smoothly follows the control points without erratic oscillations.

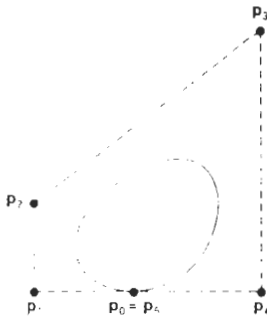


Figure 10-35
A closed Bézier curve generated by specifying the first and last control points at the same location

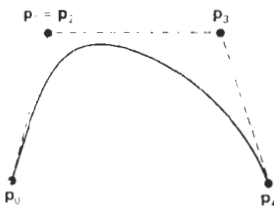


Figure 10-36
A Bézier curve can be made to pass closer to a given coordinate position by assigning multiple control points to that position

Design Techniques Using Bézier Curves

Closed Bézier curves are generated by specifying the first and last control points at the same position, as in the example shown in Fig. 10-35. Also, specifying multiple control points at a single coordinate position gives more weight to that position. In Fig. 10-36, a single coordinate position is input as two control points, and the resulting curve is pulled nearer to this position.

We can fit a Bézier curve to any number of control points, but this requires the calculation of polynomial functions of higher degree. When complicated curves are to be generated, they can be formed by piecing several Bézier sections of lower degree together. Piecing together smaller sections also gives us better control over the shape of the curve in small regions. Since Bézier curves pass through endpoints, it is easy to match curve sections (zero-order continuity). Also, Bézier curves have the important property that the tangent to the curve at an endpoint is along the line joining that endpoint to the adjacent control point. Therefore, to obtain first-order continuity between curve sections, we can pick control points \mathbf{p}'_0 and \mathbf{p}'_1 of a new section to be along the same straight line as control points \mathbf{p}_{n-1} and \mathbf{p}_1 of the previous section (Fig. 10-37). When the two curve sections have the same number of control points, we obtain C^1 continuity by choosing the first control point of the new section as the last control point of the previous section and by positioning the second control point of the new section at position

$$\mathbf{p}_{n-1} + (\mathbf{p}_n - \mathbf{p}_{n-1})$$

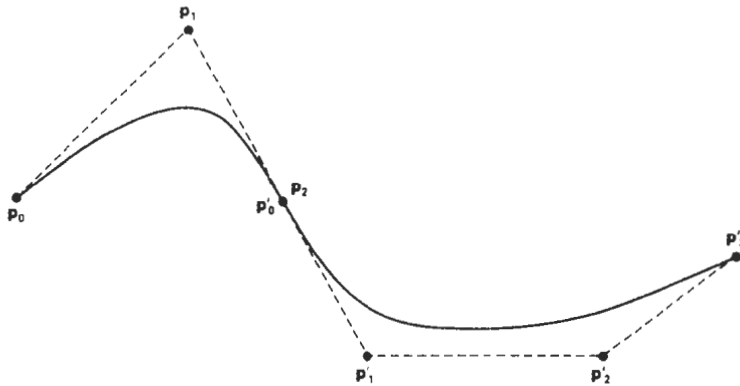


Figure 10-37
Piecewise approximation curve formed with two Bézier sections. Zero-order and first-order continuity are attained between curve sections by setting $p'_0 = p_2$ and by making points p_1 , p_2 , and p'_1 collinear.

Thus, the three control points are collinear and equally spaced.

We obtain C^2 continuity between two Bézier sections by calculating the position of the third control point of a new section in terms of the positions of the last three control points of the previous section as

$$p_{n-2} + 4(p_n - p_{n-1})$$

Requiring second-order continuity of Bézier curve sections can be unnecessarily restrictive. This is especially true with cubic curves, which have only four control points per section. In this case, second-order continuity fixes the position of the first three control points and leaves us only one point that we can use to adjust the shape of the curve segment.

Cubic Bézier Curves

Many graphics packages provide only cubic spline functions. This gives reasonable design flexibility while avoiding the increased calculations needed with higher-order polynomials. Cubic Bézier curves are generated with four control points. The four blending functions for cubic Bézier curves, obtained by substituting $n = 3$ into Eq. 10-41 are

$$\begin{aligned} BE_{0,3}(u) &= (1 - u)^3 \\ BE_{1,3}(u) &= 3u(1 - u)^2 \\ BE_{2,3}(u) &= 3u^2(1 - u) \\ BE_{3,3}(u) &= u^3 \end{aligned} \tag{10-50}$$

Plots of the four cubic Bézier blending functions are given in Fig. 10-38. The form of the blending functions determine how the control points influence the shape of the curve for values of parameter u over the range from 0 to 1. At $u = 0$,

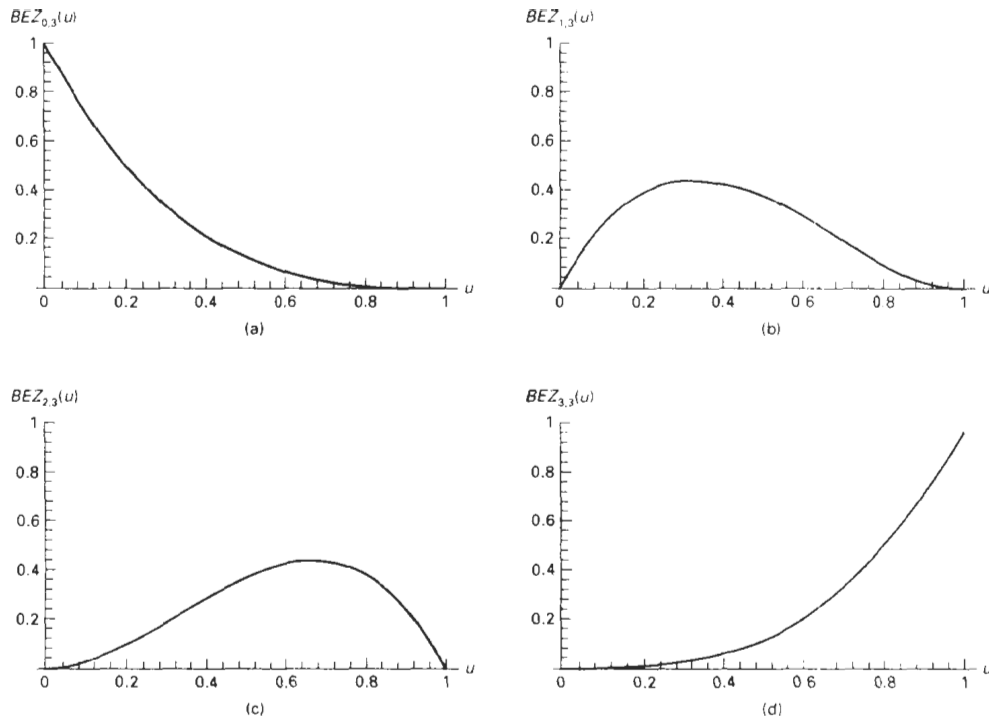


Figure 10-38
The four Bézier blending functions for cubic curves ($n = 3$)

the only nonzero blending function is $BEZ_{0,3}$, which has the value 1. At $u = 1$, the only nonzero function is $BEZ_{3,3}$, with a value of 1 at that point. Thus, the cubic Bézier curve will always pass through control points \mathbf{p}_0 and \mathbf{p}_3 . The other functions, $BEZ_{1,3}$ and $BEZ_{2,3}$, influence the shape of the curve at intermediate values of parameter u , so that the resulting curve tends toward points \mathbf{p}_1 and \mathbf{p}_2 . Blending function $BEZ_{1,3}$ is maximum at $u = 1/3$, and $BEZ_{2,3}$ is maximum at $u = 2/3$.

We note in Fig. 10-38 that each of the four blending functions is nonzero over the entire range of parameter u . Thus, Bézier curves do not allow for *local control* of the curve shape. If we decide to reposition any one of the control points, the entire curve will be affected.

At the end positions of the cubic Bézier curve, the parametric first derivatives (slopes) are

$$\mathbf{P}'(0) = 3(\mathbf{p}_1 - \mathbf{p}_0), \quad \mathbf{P}'(1) = 3(\mathbf{p}_3 - \mathbf{p}_2)$$

And the parametric second derivatives are

$$\mathbf{P}''(0) = 6(\mathbf{p}_0 - 2\mathbf{p}_1 + \mathbf{p}_2), \quad \mathbf{P}''(1) = 6(\mathbf{p}_1 - 2\mathbf{p}_2 + \mathbf{p}_3)$$

We can use these expressions for the parametric derivatives to construct piecewise curves with C^1 or C^2 continuity between sections.

By expanding the polynomial expressions for the blending functions, we can write the cubic Bézier point function in the matrix form

$$\mathbf{P}(u) = [u^3 \ u^2 \ u \ 1] \cdot \mathbf{M}_{\text{Bez}} \cdot \begin{bmatrix} \mathbf{p}_0 \\ \mathbf{p}_1 \\ \mathbf{p}_2 \\ \mathbf{p}_3 \end{bmatrix} \quad (10-51)$$

where the **Bézier matrix** is

$$\mathbf{M}_{\text{Bez}} = \begin{bmatrix} -1 & 3 & -3 & 1 \\ 3 & -6 & 3 & 0 \\ -3 & 3 & 0 & 0 \\ 1 & 0 & 0 & 0 \end{bmatrix} \quad (10-52)$$

We could also introduce additional parameters to allow adjustment of curve “tension” and “bias”, as we did with the interpolating splines. But the more useful B-splines, as well as β -splines, provide this capability.

Bézier Surfaces

Two sets of orthogonal Bézier curves can be used to design an object surface by specifying by an input mesh of control points. The parametric vector function for the Bézier surface is formed as the Cartesian product of Bézier blending functions:

$$\mathbf{P}(u, v) = \sum_{j=0}^m \sum_{k=0}^n \mathbf{p}_{j,k} \text{BEZ}_{j,m}(v) \text{BEZ}_{k,n}(u) \quad (10-53)$$

with $\mathbf{p}_{j,k}$ specifying the location of the $(m + 1)$ by $(n + 1)$ control points.

Figure 10-39 illustrates two Bézier surface plots. The control points are connected by dashed lines, and the solid lines show curves of constant u and constant v . Each curve of constant u is plotted by varying v over the interval from 0 to 1, with u fixed at one of the values in this unit interval. Curves of constant v are plotted similarly.

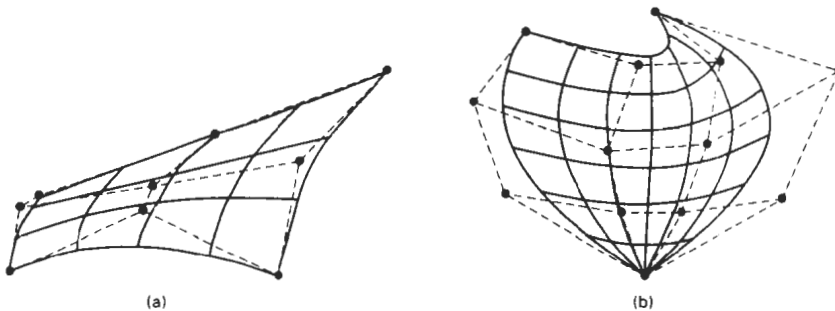


Figure 10-39

Bézier surfaces constructed for (a) $m = 3$, $n = 3$, and (b) $m = 4$, $n = 4$. Dashed lines connect the control points.

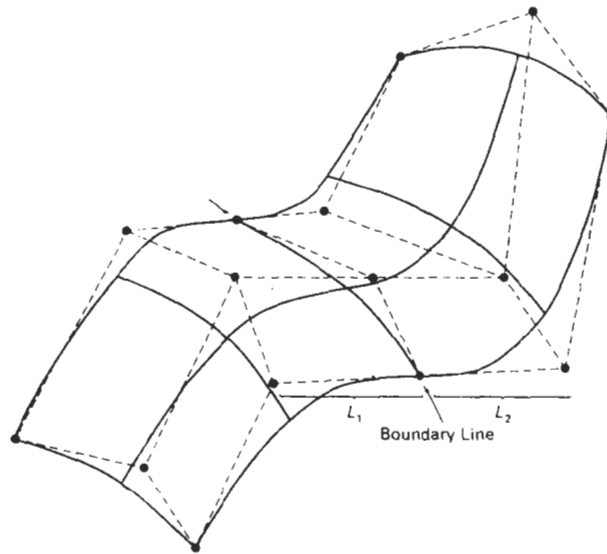


Figure 10-40

A composite Bézier surface constructed with two Bézier sections, joined at the indicated boundary line. The dashed lines connect specified control points. First-order continuity is established by making the ratio of length L_1 to length L_2 constant for each collinear line of control points across the boundary between the surface sections.

Bézier surfaces have the same properties as Bézier curves, and they provide a convenient method for interactive design applications. For each surface patch, we can select a mesh of control points in the xy “ground” plane, then we choose elevations above the ground plane for the z -coordinate values of the control points. Patches can then be pieced together using the boundary constraints.

Figure 10-40 illustrates a surface formed with two Bézier sections. As with curves, a smooth transition from one section to the other is assured by establishing both zero-order and first-order continuity at the boundary line. Zero-order continuity is obtained by matching control points at the boundary. First-order continuity is obtained by choosing control points along a straight line across the boundary and by maintaining a constant ratio of collinear line segments for each set of specified control points across section boundaries.

10-9

B-SPLINE CURVES AND SURFACES

These are the most widely used class of approximating splines. **B-splines** have two advantages over Bézier splines: (1) the degree of a B-spline polynomial can be set independently of the number of control points (with certain limitations), and (2) B-splines allow local control over the shape of a spline curve or surface. The trade-off is that B-splines are more complex than Bézier splines.

We can write a general expression for the calculation of coordinate positions along a B-spline curve in a blending-function formulation as

$$\mathbf{P}(u) = \sum_{k=0}^n \mathbf{p}_k B_{k,d}(u), \quad u_{\min} \leq u \leq u_{\max}, \quad 2 \leq d \leq n+1 \quad (10-54)$$

where the \mathbf{p}_k are an input set of $n+1$ control points. There are several differences between this B-spline formulation and that for Bézier splines. The range of parameter u now depends on how we choose the B-spline parameters. And the B-spline blending functions $B_{k,d}$ are polynomials of degree $d-1$, where parameter d can be chosen to be any integer value in the range from 2 up to the number of control points, $n+1$. (Actually, we can also set the value of d at 1, but then our “curve” is just a point plot of the control points.) Local control for B-splines is achieved by defining the blending functions over subintervals of the total range of u .

Blending functions for B-spline curves are defined by the Cox–deBoor recursion formulas:

$$B_{k,1}(u) = \begin{cases} 1, & \text{if } u_k \leq u < u_{k+1} \\ 0, & \text{otherwise} \end{cases} \quad (10-55)$$

$$B_{k,d}(u) = \frac{u - u_k}{u_{k+d-1} - u_k} B_{k,d-1}(u) + \frac{u_{k+d} - u}{u_{k+d} - u_{k+1}} B_{k+1,d-1}(u)$$

where each blending function is defined over d subintervals of the total range of u . The selected set of subinterval endpoints u_j is referred to as a **knot vector**. We can choose any values for the subinterval endpoints satisfying the relation $u_i \leq u_{j+1}$. Values for u_{\min} and u_{\max} then depend on the number of control points we select, the value we choose for parameter d , and how we set up the subintervals (knot vector). Since it is possible to choose the elements of the knot vector so that the denominators in the previous calculations can have a value of 0, this formulation assumes that any terms evaluated as 0/0 are to be assigned the value 0.

Figure 10-41 demonstrates the local-control characteristics of B-splines. In addition to local control, B-splines allow us to vary the number of control points used to design a curve without changing the degree of the polynomial. Also, any number of control points can be added or modified to manipulate curve shapes. Similarly, we can increase the number of values in the knot vector to aid in curve design. When we do this, however, we also need to add control points since the size of the knot vector depends on parameter n .

B-spline curves have the following properties:

- The polynomial curve has degree $d-1$ and C^{d-2} continuity over the range of u .
- For $n+1$ control points, the curve is described with $n+1$ blending functions.
- Each blending function $B_{k,d}$ is defined over d subintervals of the total range of u , starting at knot value u_k .
- The range of parameter u is divided into $n+d$ subintervals by the $n+d+1$ values specified in the knot vector.

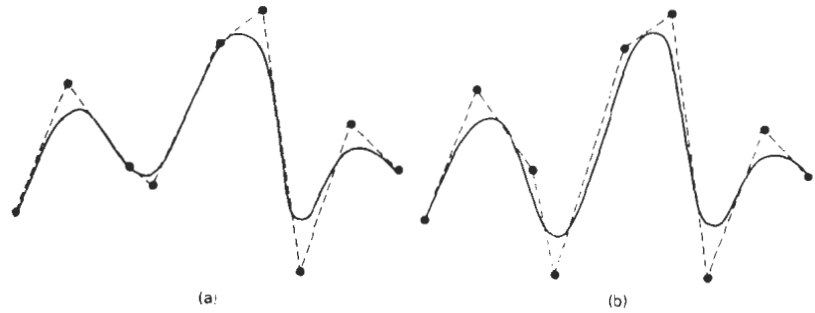


Figure 10-41

Local modification of a B-spline curve. Changing one of the control points in (a) produces curve (b), which is modified only in the neighborhood of the altered control point.

- With knot values labeled as $\{u_0, u_1, \dots, u_{n+d}\}$, the resulting B-spline curve is defined only in the interval from knot value u_{d-1} up to knot value u_{n+1} .
- Each section of the spline curve (between two successive knot values) is influenced by d control points.
- Any one control point can affect the shape of at most d curve sections.

In addition, a B-spline curve lies within the convex hull of at most $d + 1$ control points, so that B-splines are tightly bound to the input positions. For any value of u in the interval from knot value u_{d-1} to u_{n+1} , the sum over all basis functions is 1:

$$\sum_{k=0}^n B_{k,d}(u) = 1 \quad (10-56)$$

Given the control-point positions and the value of parameter d , we then need to specify the knot values to obtain the blending functions using the recurrence relations 10-55. There are three general classifications for knot vectors: uniform, open uniform, and nonuniform. B-splines are commonly described according to the selected knot-vector class.

Uniform, Periodic B-Splines

When the spacing between knot values is constant, the resulting curve is called a **uniform** B-spline. For example, we can set up a uniform knot vector as

$$\{-1.5, -1.0, -0.5, 0.0, 0.5, 1.0, 1.5, 2.0\}$$

Often knot values are normalized to the range between 0 and 1, as in

$$\{0.0, 0.2, 0.4, 0.6, 0.8, 1.0\}$$

It is convenient in many applications to set up uniform knot values with a separation of 1 and a starting value of 0. The following knot vector is an example of this specification scheme.

$$\{0, 1, 2, 3, 4, 5, 6, 7\}$$

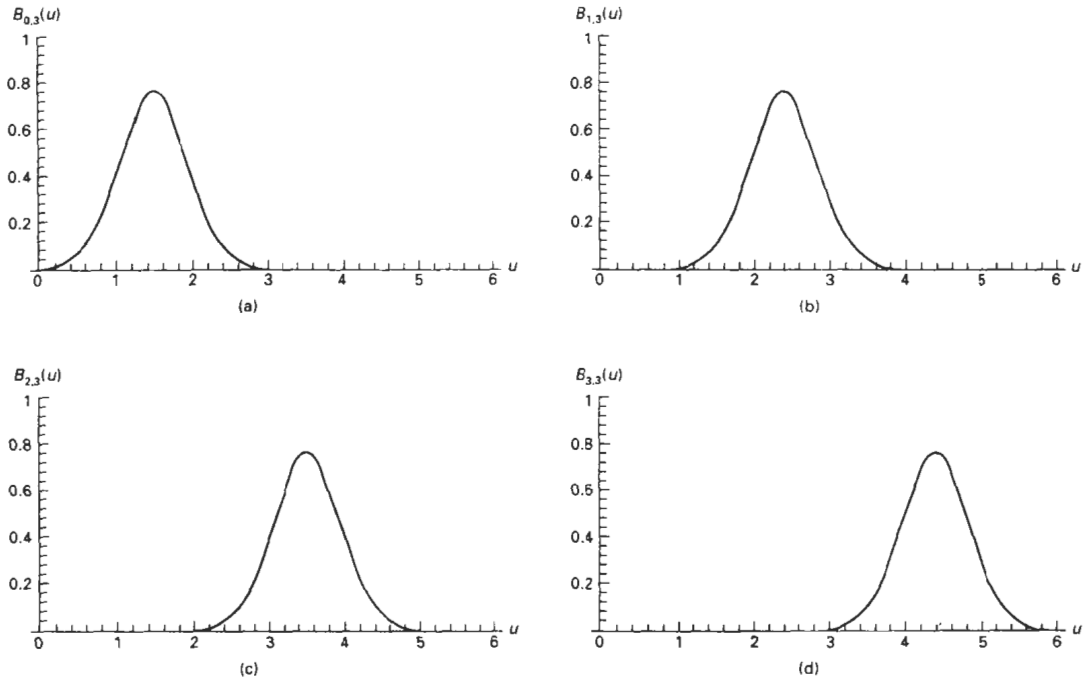


Figure 10-42
Periodic B-spline blending functions for $n = d = 3$ and a uniform, integer knot vector.

Uniform B-splines have **periodic** blending functions. That is, for given values of n and d , all blending functions have the same shape. Each successive blending function is simply a shifted version of the previous function:

$$B_{k,d}(u) = B_{k+1,d}(u + \Delta u) = B_{k-2,d}(u + 2 \Delta u) \quad (10-57)$$

where Δu is the interval between adjacent knot values. Figure 10-42 shows the quadratic, uniform B-spline blending functions generated in the following example for a curve with four control points.

Example 10-1 Uniform, Quadratic B-Splines

To illustrate the calculation of B-spline blending functions for a uniform, integer knot vector, we select parameter values $d = n = 3$. The knot vector must then contain $n + d + 1 = 7$ knot values:

$$\{0, 1, 2, 3, 4, 5, 6\}$$

and the range of parameter u is from 0 to 6, with $n + d = 6$ subintervals.

Each of the four blending functions spans $d = 3$ subintervals of the total range of u . Using the recurrence relations 10-55, we obtain the first blending function as

$$B_{0,3}(u) = \begin{cases} \frac{1}{2}u^2, & \text{for } 0 \leq u < 1 \\ \frac{1}{2}u(2-u) + \frac{1}{2}(u-1)(3-u), & \text{for } 1 \leq u < 2 \\ \frac{1}{2}(3-u)^2, & \text{for } 2 \leq u < 3 \end{cases}$$

We obtain the next periodic blending function using relationship 10-57, substituting $u - 1$ for u in $B_{0,3}$, and shifting the starting positions up by 1:

$$B_{1,3}(u) = \begin{cases} \frac{1}{2}(u-1)^2, & \text{for } 1 \leq u < 2 \\ \frac{1}{2}(u-1)(3-u) + \frac{1}{2}(u-2)(4-u), & \text{for } 2 \leq u < 3 \\ \frac{1}{2}(4-u)^2, & \text{for } 3 \leq u < 4 \end{cases}$$

Similarly, the remaining two periodic functions are obtained by successively shifting $B_{1,3}$ to the right:

$$B_{2,3}(u) = \begin{cases} \frac{1}{2}(u-2)^2, & \text{for } 2 \leq u < 3 \\ \frac{1}{2}(u-2)(4-u) + \frac{1}{2}(u-3)(5-u), & \text{for } 3 \leq u < 4 \\ \frac{1}{2}(5-u)^2, & \text{for } 4 \leq u < 5 \end{cases}$$

$$B_{3,3}(u) = \begin{cases} \frac{1}{2}(u-3)^2, & \text{for } 3 \leq u < 4 \\ \frac{1}{2}(u-3)(5-u) + \frac{1}{2}(u-4)(6-u), & \text{for } 4 \leq u < 5 \\ \frac{1}{2}(6-u)^2, & \text{for } 5 \leq u < 6 \end{cases}$$

A plot of the four periodic, quadratic blending functions is given in Fig. 10-42, which demonstrates the local feature of B-splines. The first control point is multiplied by blending function $B_{0,3}(u)$. Therefore, changing the position of the first control point only affects the shape of the curve up to $u = 3$. Similarly, the last control point influences the shape of the spline curve in the interval where $B_{3,3}$ is defined.

Figure 10-42 also illustrates the limits of the B-spline curve for this example. All blending functions are present in the interval from $u_{d-1} = 2$ to $u_{n+1} = 4$. Below 2 and above 4, not all blending functions are present. This is the range of the poly-



Figure 10-43
Quadratic, periodic B-spline fitted
to four control points in the xy
plane.

nomial curve, and the interval in which Eq. 10-56 is valid. Thus, the sum of all blending functions is 1 within this interval. Outside this interval, we cannot sum all blending functions, since they are not all defined below 2 and above 4.

Since the range of the resulting polynomial curve is from 2 to 4, we can determine the starting and ending positions of the curve by evaluating the blending functions at these points to obtain

$$\mathbf{P}_{\text{start}} = \frac{1}{2}(\mathbf{p}_0 + \mathbf{p}_1), \quad \mathbf{P}_{\text{end}} = \frac{1}{2}(\mathbf{p}_2 + \mathbf{p}_3)$$

Thus, the curve starts at the midposition between the first two control points and ends at the midposition between the last two control points.

We can also determine the parametric derivatives at the starting and ending positions of the curve. Taking the derivatives of the blending functions and substituting the endpoint values for parameter u , we find that

$$\mathbf{P}'_{\text{start}} = \mathbf{p}_1 - \mathbf{p}_0, \quad \mathbf{P}'_{\text{end}} = \mathbf{p}_3 - \mathbf{p}_2$$

The parametric slope of the curve at the start position is parallel to the line joining the first two control points, and the parametric slope at the end of the curve is parallel to the line joining the last two control points.

An example plot of the quadratic periodic B-spline curve is given in Figure 10-43 for four control points selected in the xy plane.

In the preceding example, we noted that the quadratic curve starts between the first two control points and ends at a position between the last two control points. This result is valid for a quadratic, periodic B-spline fitted to any number of distinct control points. In general, for higher-order polynomials, the start and end positions are each weighted averages of $d - 1$ control points. We can pull a spline curve closer to any control-point position by specifying that position multiple times.

General expressions for the boundary conditions for periodic B-splines can be obtained by reparameterizing the blending functions so that parameter u is mapped onto the unit interval from 0 to 1. Beginning and ending conditions are then obtained at $u = 0$ and $u = 1$.

Cubic, Periodic B-Splines

Since cubic, periodic B-splines are commonly used in graphics packages, we consider the formulation for this class of splines. Periodic splines are particularly useful for generating certain closed curves. For example, the closed curve in Fig. 10-44 can be generated in sections by cyclically specifying four of the six control

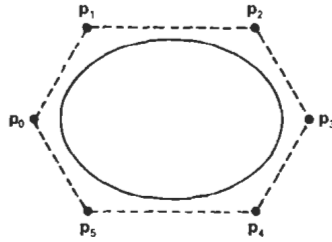


Figure 10-44
A closed, periodic, piecewise, cubic B-spline constructed with cyclic specification of the six control points.

points shown at each step. If any three consecutive control points are identical, the curve passes through that coordinate position.

For cubics, $d = 4$ and each blending function spans four subintervals of the total range of u . If we are to fit the cubic to four control points, then we could use the integer knot vector

$$\{0, 1, 2, 3, 4, 5, 6, 7\}$$

and recurrence relations 10-55 to obtain the periodic blending functions, as we did in the last section for quadratic periodic B-splines.

In this section, we consider an alternate formulation for periodic cubic B-splines. We start with the boundary conditions and obtain the blending functions normalized to the interval $0 \leq u \leq 1$. Using this formulation, we can also easily obtain the characteristic matrix. The boundary conditions for periodic cubic B-splines with four consecutive control points, labeled p_0 , p_1 , p_2 , and p_3 , are

$$\begin{aligned} \mathbf{P}(0) &= \frac{1}{6}(\mathbf{p}_0 + 4\mathbf{p}_1 + \mathbf{p}_2) \\ \mathbf{P}(1) &= \frac{1}{6}(\mathbf{p}_1 + 4\mathbf{p}_2 + \mathbf{p}_3) \\ \mathbf{P}'(0) &= \frac{1}{2}(\mathbf{p}_2 - \mathbf{p}_0) \\ \mathbf{P}'(1) &= \frac{1}{2}(\mathbf{p}_3 - \mathbf{p}_1) \end{aligned} \quad (10-58)$$

These boundary conditions are similar to those for cardinal splines: Curve sections are defined with four control points, and parametric derivatives (slopes) at the beginning and end of each curve section are parallel to the chords joining adjacent control points. The B-spline curve section starts at a position near p_1 and ends at a position near p_2 .

A matrix formulation for a cubic periodic B-splines with four control points can then be written as

$$\mathbf{P}(u) = [u^3 \ u^2 \ u \ 1] \cdot \mathbf{M}_B \cdot \begin{bmatrix} \mathbf{p}_0 \\ \mathbf{p}_1 \\ \mathbf{p}_2 \\ \mathbf{p}_3 \end{bmatrix} \quad (10-59)$$

where the B-spline matrix for periodic cubic polynomials is

$$\mathbf{M}_B = \frac{1}{6} \begin{bmatrix} -1 & 3 & -3 & 1 \\ 3 & -6 & 3 & 0 \\ -3 & 0 & 3 & 0 \\ 1 & 4 & 1 & 0 \end{bmatrix} \quad (10-60)$$

This matrix can be obtained by solving for the coefficients in a general cubic polynomial expression using the specified four boundary conditions.

We can also modify the B-spline equations to include a tension parameter t (as in cardinal splines). The periodic, cubic B-spline with tension matrix then has the form

$$\mathbf{M}_{Bt} = \frac{1}{6} \begin{bmatrix} -t & 12 - 9t & 9t - 12 & t \\ 3t & 12t - 18 & 18 - 15t & 0 \\ -3t & 0 & 3t & 0 \\ t & 6 - 2t & t & 0 \end{bmatrix} \quad (10-61)$$

which reduces to M_B when $t = 1$.

We obtain the periodic, cubic B-spline blending functions over the parameter range from 0 to 1 by expanding the matrix representation into polynomial form. For example, for the tension value $t = 1$, we have

$$\begin{aligned} B_{0,3}(u) &= \frac{1}{6}(1 - u)^3, & 0 \leq u \leq 1 \\ B_{1,3}(u) &= \frac{1}{6}(3u^3 - 6u^2 + 4) \\ B_{2,3}(u) &= \frac{1}{6}(-3u^3 + 3u^2 + 3u + 1) \\ B_{3,3}(u) &= \frac{1}{6}u^3 \end{aligned} \quad (10-62)$$

Open Uniform B-Splines

This class of B-splines is a cross between uniform B-splines and nonuniform B-splines. Sometimes it is treated as a special type of uniform B-spline, and sometimes it is considered to be in the nonuniform B-spline classification. For the **open uniform** B-splines, or simply **open** B-splines, the knot spacing is uniform except at the ends where knot values are repeated d times.

Following are two examples of open uniform, integer knot vectors, each with a starting value of 0:

$$\begin{aligned} \{0, 0, 1, 2, 3, 3\}, & \quad \text{for } d = 2 \text{ and } n = 3 \\ \{0, 0, 0, 0, 1, 2, 2, 2, 2\}, & \quad \text{for } d = 4 \text{ and } n = 4 \end{aligned}$$

We can normalize these knot vectors to the unit interval from 0 to 1:

$$\begin{aligned} \{0, 0, 0.33, 0.67, 1, 1\}, & \quad \text{for } d = 2 \text{ and } n = 3 \\ \{0, 0, 0, 0, 0.5, 1, 1, 1, 1\}, & \quad \text{for } d = 4 \text{ and } n = 4 \end{aligned}$$

For any values of parameters d and n , we can generate an open uniform knot vector with integer values using the calculations

$$u_j = \begin{cases} 0, & \text{for } 0 \leq j < d \\ j - d + 1, & \text{for } d \leq j \leq n \\ n - d + 2, & \text{for } j > n \end{cases} \quad (10-63)$$

for values of j ranging from 0 to $n + d$. With this assignment, the first d knots are assigned the value 0, and the last d knots have the value $n - d + 2$.

Open uniform B-splines have characteristics that are very similar to Bézier splines. In fact, when $d = n + 1$ (degree of the polynomial is n) open B-splines reduce to Bézier splines, and all knot values are either 0 or 1. For example, with a cubic, open B-spline ($d = 4$) and four control points, the knot vector is

$$\{0, 0, 0, 0, 1, 1, 1, 1\}$$

The polynomial curve for an open B-spline passes through the first and last control points. Also, the slope of the parametric curves at the first control point is parallel to the line connecting the first two control points. And the parametric slope at the last control point is parallel to the line connecting the last two control points. So geometric constraints for matching curve sections are the same as for Bézier curves.

As with Bézier curves, specifying multiple control points at the same coordinate position pulls any B-spline curve closer to that position. Since open B-splines start at the first control point and end at the last specified control point, closed curves are generated by specifying the first and last control points at the same position.

Example 10-2 Open Uniform, Quadratic B-Splines

From conditions 10-63 with $d = 3$ and $n = 4$ (five control points), we obtain the following eight values for the knot vector:

$$\{u_0, u_1, u_2, u_3, u_4, u_5, u_6, u_7\} = \{0, 0, 0, 1, 2, 3, 3, 3\}$$

The total range of u is divided into seven subintervals, and each of the five blending functions $B_{k,3}$ is defined over three subintervals, starting at knot position u_k . Thus, $B_{0,3}$ is defined from $u_0 = 0$ to $u_3 = 1$, $B_{1,3}$ is defined from $u_1 = 0$ to $u_4 = 2$, and $B_{4,3}$ is defined from $u_4 = 2$ to $u_7 = 3$. Explicit polynomial expressions are obtained for the blending functions from recurrence relations 10-55 as

$$B_{0,3}(u) = (1 - u)^2, \quad 0 \leq u < 1$$

$$B_{1,3}(u) = \begin{cases} \frac{1}{2}u(4 - 3u), & 0 \leq u < 1 \\ \frac{1}{2}(2 - u)^2, & 1 \leq u < 2 \end{cases}$$

$$B_{2,3}(u) = \begin{cases} \frac{1}{2}u^2, & 0 \leq u < 1 \\ \frac{1}{2}u(2-u) + \frac{1}{2}(u-1)(3-u), & 1 \leq u < 2 \\ \frac{1}{2}(3-u)^2, & 2 \leq u < 3 \end{cases}$$

$$B_{3,3}(u) = \begin{cases} \frac{1}{2}(u-1)^2, & 1 \leq u < 2 \\ \frac{1}{2}(3-u)(3u-5), & 2 \leq u < 3 \end{cases}$$

$$B_{4,3}(u) = (u-2)^2, \quad 2 \leq u < 3$$

Figure 10-45 shows the shape of these five blending functions. The local features of B-splines are again demonstrated. Blending function $B_{0,3}$ is nonzero only in the subinterval from 0 to 1, so the first control point influences the curve only in this interval. Similarly, function $B_{4,3}$ is zero outside the interval from 2 to 3, and the position of the last control point does not affect the shape of the beginning and middle parts of the curve.

Matrix formulations for open B-splines are not as conveniently generated as they are for periodic, uniform B-splines. This is due to the multiplicity of knot values at the beginning and end of the knot vector.

Nonuniform B-Splines

For this class of splines, we can specify any values and intervals for the knot vector. With **nonuniform** B-splines, we can choose multiple internal knot values and unequal spacing between the knot values. Some examples are

$$\begin{aligned} &\{0, 1, 2, 3, 3, 4\} \\ &\{0, 2, 2, 3, 3, 6\} \\ &\{0, 0, 0, 1, 1, 3, 3, 3\} \\ &\{0, 0.2, 0.6, 0.9, 1.0\} \end{aligned}$$

Nonuniform B-splines provide increased flexibility in controlling a curve shape. With unequally spaced intervals in the knot vector, we obtain different shapes for the blending functions in different intervals, which can be used to adjust spline shapes. By increasing knot multiplicity, we produce subtle variations in curve shape and even introduce discontinuities. Multiple knot values also reduce the continuity by 1 for each repeat of a particular value.

We obtain the blending functions for a nonuniform B-spline using methods similar to those discussed for uniform and open B-splines. Given a set of $n + 1$ control points, we set the degree of the polynomial and select the knot values. Then, using the recurrence relations, we could either obtain the set of blending functions or evaluate curve positions directly for the display of the curve. Graphics packages often restrict the knot intervals to be either 0 or 1 to reduce computations. A set of characteristic matrices then can be stored and used to compute

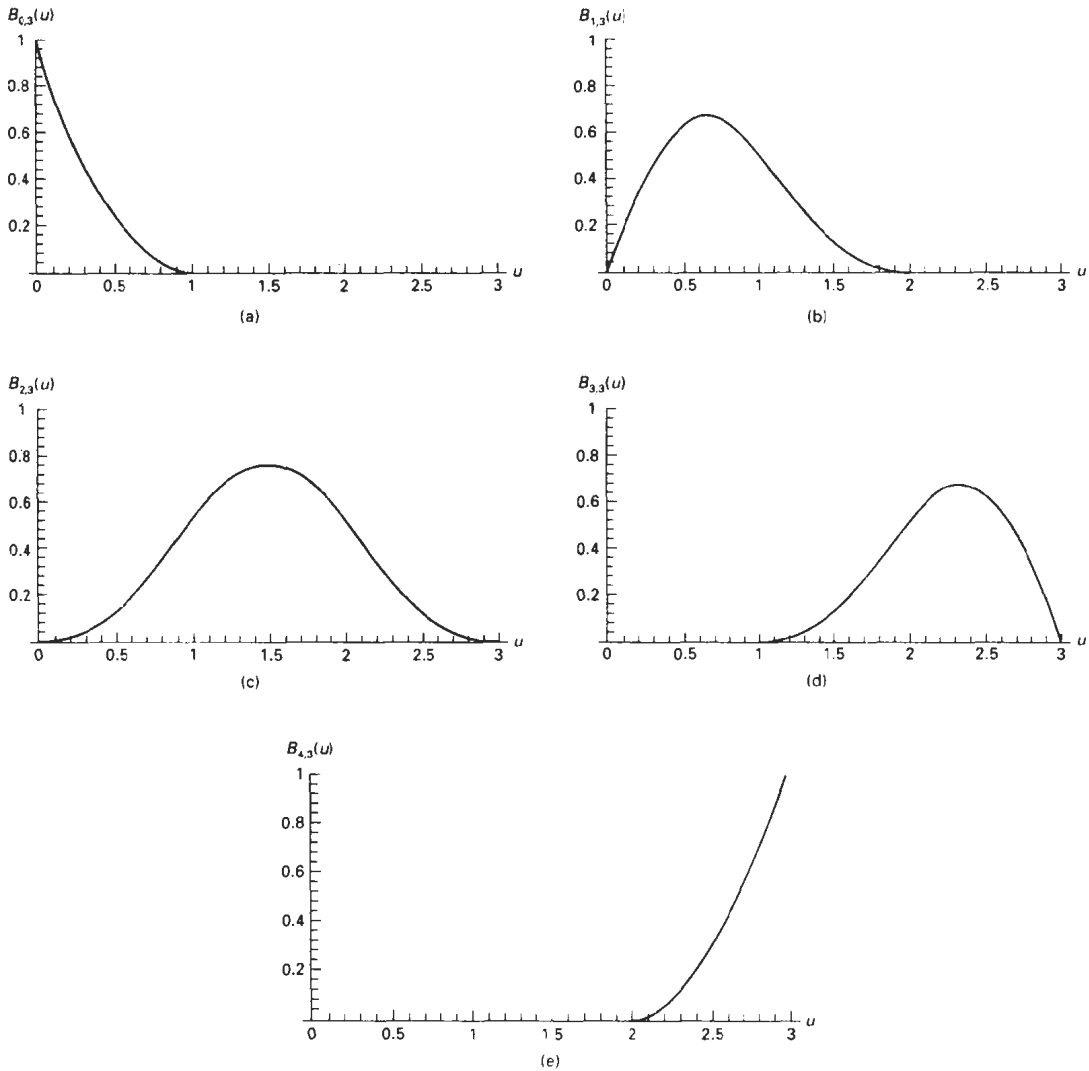


Figure 10-45

Open, uniform B-spline blending functions for $n = 4$ and $d = 3$.

values along the spline curve without evaluating the recurrence relations for each curve point to be plotted.

B-Spline Surfaces

Formulation of a B-spline surface is similar to that for Bézier splines. We can obtain a vector point function over a B-spline surface using the Cartesian product of B-spline blending functions in the form

**Figure 10-46**

A prototype helicopter, designed and modeled by Daniel Langlois of SOFTIMAGE, Inc., Montreal, using 180,000 B-spline surface patches. The scene was then rendered using ray tracing, bump mapping, and reflection mapping. (Courtesy of Silicon Graphics, Inc.)

$$\mathbf{P}(u, v) = \sum_{k_1=0}^{n_1} \sum_{k_2=0}^{n_2} \mathbf{P}_{k_1, k_2} B_{k_1, d_1}(u) B_{k_2, d_2}(v) \quad (10-64)$$

where the vector values for \mathbf{p}_{k_1, k_2} specify positions of the $(n_1 + 1)$ by $(n_2 + 1)$ control points.

B-spline surfaces exhibit the same properties as those of their component B-spline curves. A surface can be constructed from selected values for parameters d_1 and d_2 (which determine the polynomial degrees to be used) and from the specified knot vector. Figure 10-46 shows an object modeled with B-spline surfaces.

10-10

BETA-SPLINES

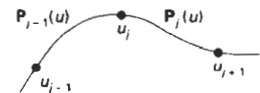
A generalization of B-splines are the **beta-splines**, also referred to as **β -splines**, that are formulated by imposing geometric continuity conditions on the first and second parametric derivatives. The continuity parameters for beta-splines are called **β parameters**.

Beta-Spline Continuity Conditions

For a specified knot vector, we can designate the spline sections to the left and right of a particular knot u_i with the position vectors $\mathbf{P}_{i-1}(u)$ and $\mathbf{P}_i(u)$ (Fig. 10-47). Zero-order continuity (*positional continuity*), G^0 , at u_i is obtained by requiring

$$\mathbf{P}_{i-1}(u_i) = \mathbf{P}_i(u_i) \quad (10-65)$$

First-order continuity (*unit tangent continuity*), G^1 , is obtained by requiring tangent vectors to be proportional:

**Figure 10-47**

Position vectors along curve sections to the left and right of knot u_i .

$$\beta_1 \mathbf{P}'_{j-1}(u_j) = \mathbf{P}'_j(u_j), \quad \beta_1 > 0 \quad (10-66)$$

Here, parametric first derivatives are proportional, and the unit tangent vectors are continuous across the knot.

Second-order continuity (*curvature vector continuity*), G^2 , is imposed with the condition

$$\beta_1 \mathbf{P}''_{j-1}(u_j) + \beta_2 \mathbf{P}'_{j-1}(u_j) = \mathbf{P}''_j(u_j) \quad (10-67)$$

where β_2 can be assigned any real number, and $\beta_1 > 0$. The curvature vector provides a measure of the amount of bending of the curve at position u_j . When $\beta_1 = 1$ and $\beta_2 = 0$, beta-splines reduce to B-splines.

Parameter β_1 is called the *bias parameter* since it controls the skewness of the curve. For $\beta_1 > 1$, the curve tends to flatten to the right in the direction of the unit tangent vector at the knots. For $0 < \beta_1 < 1$, the curve tends to flatten to the left. The effect of β_1 on the shape of the spline curve is shown in Fig. 10-48.

Parameter β_2 is called the *tension parameter* since it controls how tightly or loosely the spline fits the control graph. As β_2 increases, the curve approaches the shape of the control graph, as shown in Fig. 10-49.

Cubic, Periodic Beta-Spline Matrix Representation

Applying the beta-spline boundary conditions to a cubic polynomial with a uniform knot vector, we obtain the following matrix representation for a periodic beta-spline:

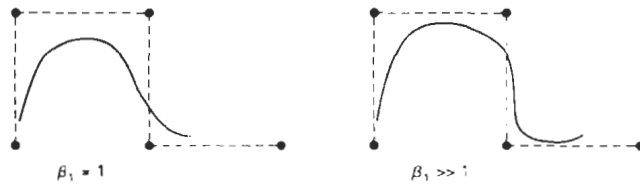


Figure 10-48
Effect of parameter β_1 on the shape of a beta-spline curve.

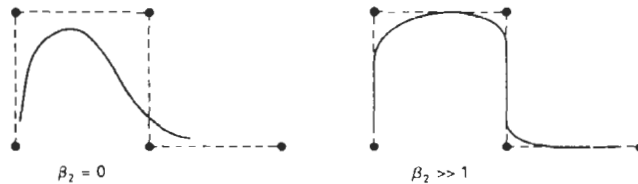


Figure 10-49
Effect of parameter β_2 on the shape of a beta-spline curve.