

# Transformation Matrices

The machinery of linear algebra can be used to express many of the operations required to arrange objects in a 3D scene, view them with cameras, and get them onto the screen. *Geometric transformations* like rotation, translation, scaling, and projection can be accomplished with matrix multiplication, and the *transformation matrices* used to do this are the subject of this chapter.

We will show how a set of points transforms if the points are represented as offset vectors from the origin, and we will use the clock shown in Figure 6.1 as an example of a point set. So think of the clock as a bunch of points that are the ends of vectors whose tails are at the origin. We also discuss how these transforms operate differently on locations (points), displacement vectors, and surface normal vectors.

## 6.1 2D Linear Transformations

We can use a  $2 \times 2$  matrix to change, or transform, a 2D vector:

$$\begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} a_{11}x + a_{12}y \\ a_{21}x + a_{22}y \end{bmatrix}.$$

This kind of operation, which takes in a 2-vector and produces another 2-vector by a simple matrix multiplication, is a *linear transformation*.

By this simple formula we can achieve a variety of useful transformations, depending on what we put in the entries of the matrix, as will be discussed in

the following sections. For our purposes, consider moving along the  $x$ -axis a horizontal move and along the  $y$ -axis, a vertical move.

### 6.1.1 Scaling

The most basic transform is a *scale* along the coordinate axes. This transform can change length and possibly direction:

$$\text{scale}(s_x, s_y) = \begin{bmatrix} s_x & 0 \\ 0 & s_y \end{bmatrix}.$$

Note what this matrix does to a vector with Cartesian components  $(x, y)$ :

$$\begin{bmatrix} s_x & 0 \\ 0 & s_y \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} s_x x \\ s_y y \end{bmatrix}.$$

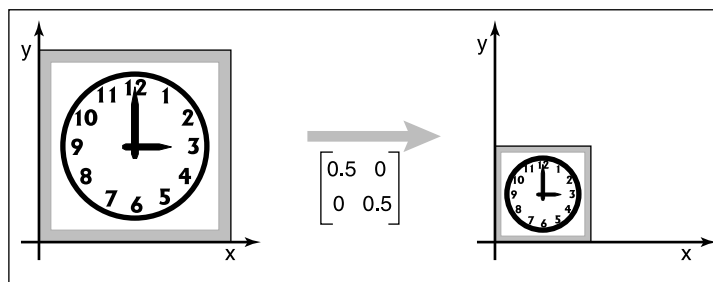
So just by looking at the matrix of an axis-aligned scale we can read off the two scale factors.

Example. The matrix that shrinks  $x$  and  $y$  uniformly by a factor of two is (Figure 6.1)

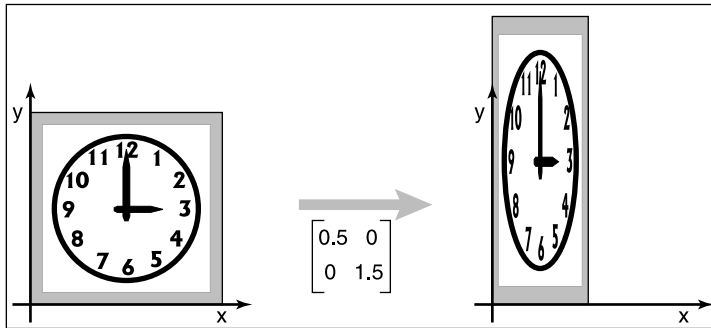
$$\text{scale}(0.5, 0.5) = \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix}.$$

A matrix which halves in the horizontal and increases by three-halves in the vertical is (see Figure 6.2)

$$\text{scale}(0.5, 1.5) = \begin{bmatrix} 0.5 & 0 \\ 0 & 1.5 \end{bmatrix}.$$



**Figure 6.1.** Scaling uniformly by half for each axis: The axis-aligned scale matrix has the proportion of change in each of the diagonal elements and zeroes in the off-diagonal elements.



**Figure 6.2.** Scaling non-uniformly in  $x$  and  $y$ : The scaling matrix is diagonal with non-equal elements. Note that the square outline of the clock becomes a rectangle and the circular face becomes an ellipse.

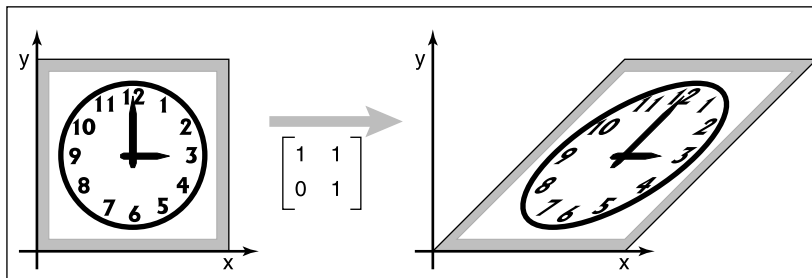
### 6.1.2 Shearing

A shear is something that pushes things sideways, producing something like a deck of cards across which you push your hand; the bottom card stays put and cards move more the closer they are to the top of the deck. The horizontal and vertical shear matrices are

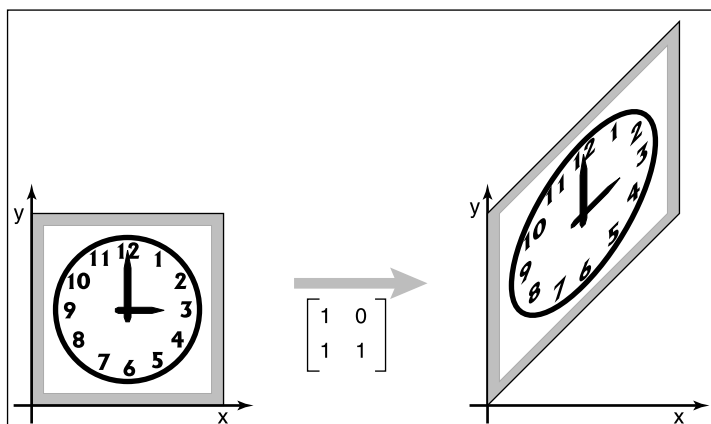
$$\text{shear-}x(s) = \begin{bmatrix} 1 & s \\ 0 & 1 \end{bmatrix}, \quad \text{shear-}y(s) = \begin{bmatrix} 1 & 0 \\ s & 1 \end{bmatrix}.$$

**Example.** The transform that shears horizontally so that vertical lines become  $45^\circ$  lines leaning towards the right is (see Figure 6.3)

$$\text{shear-}x(1) = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}.$$



**Figure 6.3.** An  $x$ -shear matrix moves points to the right in proportion to their  $y$ -coordinate. Now the square outline of the clock becomes a parallelogram and, as with scaling, the circular face of the clock becomes an ellipse.



**Figure 6.4.** A  $y$ -shear matrix moves points up in proportion to their  $x$ -coordinate.

An analogous transform vertically is (see Figure 6.4)

$$\text{shear-}y(1) = \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix}.$$

In fact, the image of a circle under any matrix transformation is an ellipse.

In both cases the square outline of the sheared clock becomes a parallelogram, and the circular face of the sheared clock becomes an ellipse.

Another way to think of a shear is in terms of rotation of only the vertical (or horizontal) axes. The shear transform that takes a vertical axis and tilts it clockwise by an angle  $\phi$  is

$$\begin{bmatrix} 1 & \tan \phi \\ 0 & 1 \end{bmatrix}.$$

Similarly, the shear matrix which rotates the horizontal axis counterclockwise by angle  $\phi$  is

$$\begin{bmatrix} 1 & 0 \\ \tan \phi & 1 \end{bmatrix}.$$

### 6.1.3 Rotation

Suppose we want to rotate a vector  $\mathbf{a}$  by an angle  $\phi$  counterclockwise to get vector  $\mathbf{b}$  (Figure 6.5). If  $\mathbf{a}$  makes an angle  $\alpha$  with the  $x$ -axis, and its length is  $r = x_a^2 + y_a^2$ , then we know that

$$x_a = r \cos \alpha,$$

$$y_a = r \sin \alpha.$$

Because  $\mathbf{b}$  is a rotation of  $\mathbf{a}$ , it also has length  $r$ . Because it is rotated an angle  $\phi$  from  $\mathbf{a}$ ,  $\mathbf{b}$  makes an angle  $(\alpha + \phi)$  with the  $x$ -axis. Using the trigonometric addition identities (Section 2.3.3):

$$\begin{aligned}x_b &= r \cos(\alpha + \phi) = r \cos \alpha \cos \phi - r \sin \alpha \sin \phi, \\y_b &= r \sin(\alpha + \phi) = r \sin \alpha \cos \phi + r \cos \alpha \sin \phi.\end{aligned}\quad (6.1)$$

Substituting  $x_a = r \cos \alpha$  and  $y_a = r \sin \alpha$  gives

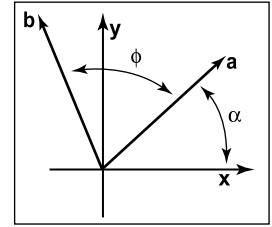
$$\begin{aligned}x_b &= x_a \cos \phi - y_a \sin \phi, \\y_b &= y_a \cos \phi + x_a \sin \phi.\end{aligned}$$

In matrix form, the transformation that takes  $\mathbf{a}$  to  $\mathbf{b}$  is then

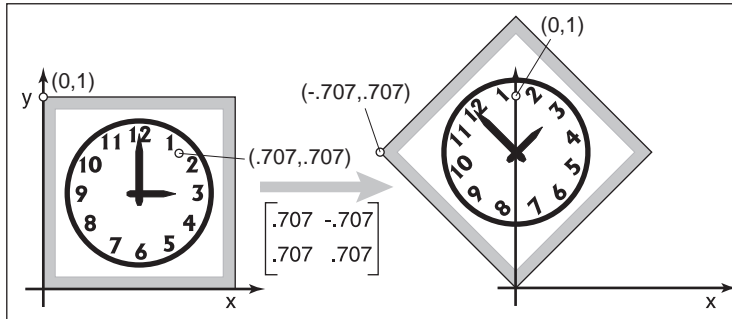
$$\text{rotate}(\phi) = \begin{bmatrix} \cos \phi & -\sin \phi \\ \sin \phi & \cos \phi \end{bmatrix}.$$

**Example.** A matrix that rotates vectors by  $\pi/4$  radians (45 degrees) is (see Figure 6.6)

$$\begin{bmatrix} \cos \frac{\pi}{4} & -\sin \frac{\pi}{4} \\ \sin \frac{\pi}{4} & \cos \frac{\pi}{4} \end{bmatrix} = \begin{bmatrix} 0.707 & -0.707 \\ 0.707 & 0.707 \end{bmatrix}.$$



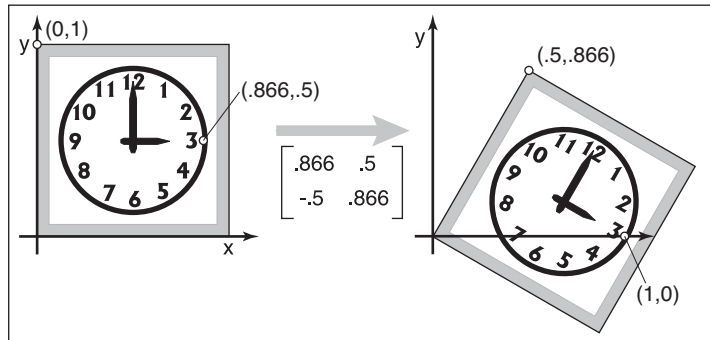
**Figure 6.5.** The geometry for Equation (6.1).



**Figure 6.6.** A rotation by 45 degrees. Note that the rotation is counterclockwise and that  $\cos(45^\circ) = \sin(45^\circ) \approx .707$ .

A matrix that rotates by  $\pi/6$  radians (30 degrees) in the *clockwise* direction is a rotation by  $-\pi/6$  radians in our framework (see Figure 6.7):

$$\begin{bmatrix} \cos \frac{-\pi}{6} & -\sin \frac{-\pi}{6} \\ \sin \frac{-\pi}{6} & \cos \frac{-\pi}{6} \end{bmatrix} = \begin{bmatrix} 0.866 & 0.5 \\ -0.5 & 0.866 \end{bmatrix}.$$



**Figure 6.7.** A rotation by minus thirty degrees. Note that the rotation is clockwise and that  $\cos(-30^\circ) \approx .866$  and  $\sin(-30^\circ) = -.5$ .

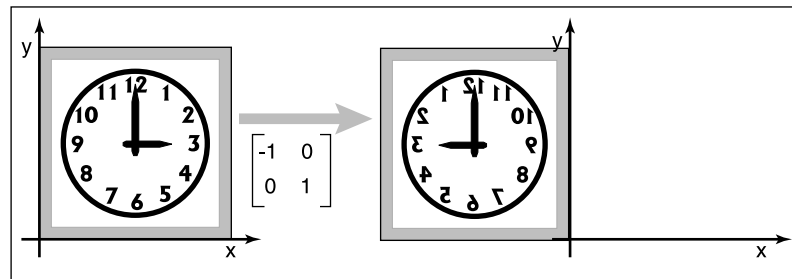
Because the norm of each row of a rotation matrix is one ( $\sin^2 \phi + \cos^2 \phi = 1$ ), and the rows are orthogonal ( $\cos \phi (-\sin \phi) + \sin \phi \cos \phi = 0$ ), we see that rotation matrices are orthogonal matrices (Section 5.2.4). By looking at the matrix we can read off two pairs of orthonormal vectors: the two columns, which are the vectors to which the transformation sends the canonical basis vectors  $(1, 0)$  and  $(0, 1)$ ; and the rows, which are the vectors that the transformations sends to the canonical basis vectors.

Said briefly,  $\mathbf{R}\mathbf{e}_i = \mathbf{u}_i$  and  $\mathbf{R}\mathbf{v}_i = \mathbf{u}_i$ , for a rotation with columns  $\mathbf{u}_i$  and rows  $\mathbf{v}_i$ .

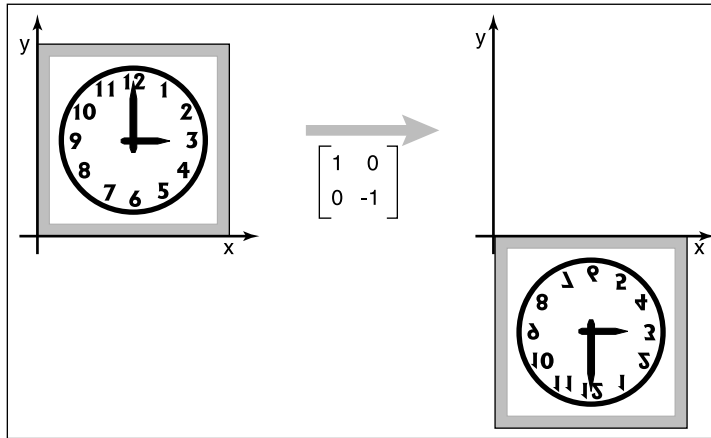
#### 6.1.4 Reflection

We can reflect a vector across either of the coordinate axes by using a scale with one negative scale factor (see Figures 6.8 and 6.9):

$$\text{reflect-}y = \begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix}, \quad \text{reflect-}x = \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix}.$$



**Figure 6.8.** A reflection about the  $y$ -axis is achieved by multiplying all  $x$ -coordinates by  $-1$ .



**Figure 6.9.** A reflection about the  $x$ -axis is achieved by multiplying all  $y$ -coordinates by  $-1$ .

While one might expect that the matrix with  $-1$  in both elements of the diagonal is also a reflection, in fact it is just a rotation by  $\pi$  radians.

This rotation can also be called a “reflection through the origin.”

### 6.1.5 Composition and Decomposition of Transformations

It is common for graphics programs to apply more than one transformation to an object. For example, we might want to first apply a scale  $\mathbf{S}$ , and then a rotation  $\mathbf{R}$ . This would be done in two steps on a 2D vector  $\mathbf{v}_1$ :

$$\text{first, } \mathbf{v}_2 = \mathbf{S}\mathbf{v}_1, \text{ then, } \mathbf{v}_3 = \mathbf{R}\mathbf{v}_2.$$

Another way to write this is

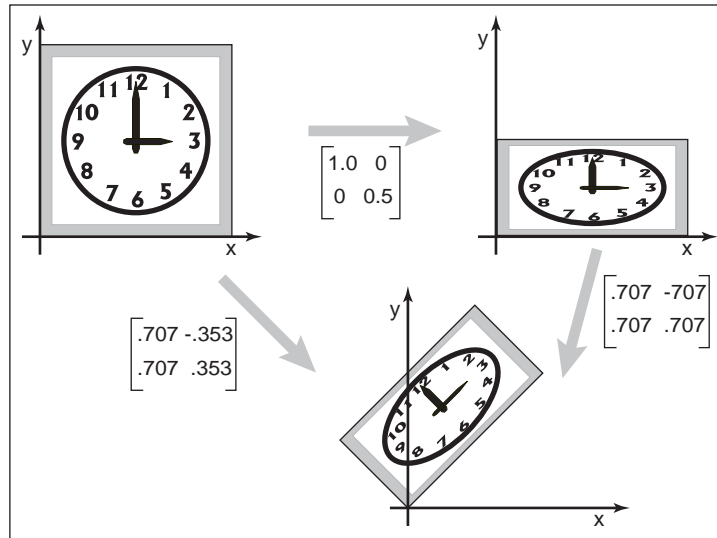
$$\mathbf{v}_3 = \mathbf{R}(\mathbf{S}\mathbf{v}_1).$$

Because matrix multiplication is associative, we can also write

$$\mathbf{v}_3 = (\mathbf{RS})\mathbf{v}_1.$$

In other words, we can represent the effects of transforming a vector by two matrices in sequence using a single matrix of the same size, which we can compute by multiplying the two matrices:  $\mathbf{M} = \mathbf{RS}$  (Figure 6.10).

It is *very important* to remember that these transforms are applied from the *right side first*. So the matrix  $\mathbf{M} = \mathbf{RS}$  first applies  $\mathbf{S}$  and then  $\mathbf{R}$ .



**Figure 6.10.** Applying the two transform matrices in sequence is the same as applying the product of those matrices once. This is a key concept that underlies most graphics hardware and software.

**Example.** Suppose we want to scale by one-half in the vertical direction and then rotate by  $\pi/4$  radians (45 degrees). The resulting matrix is

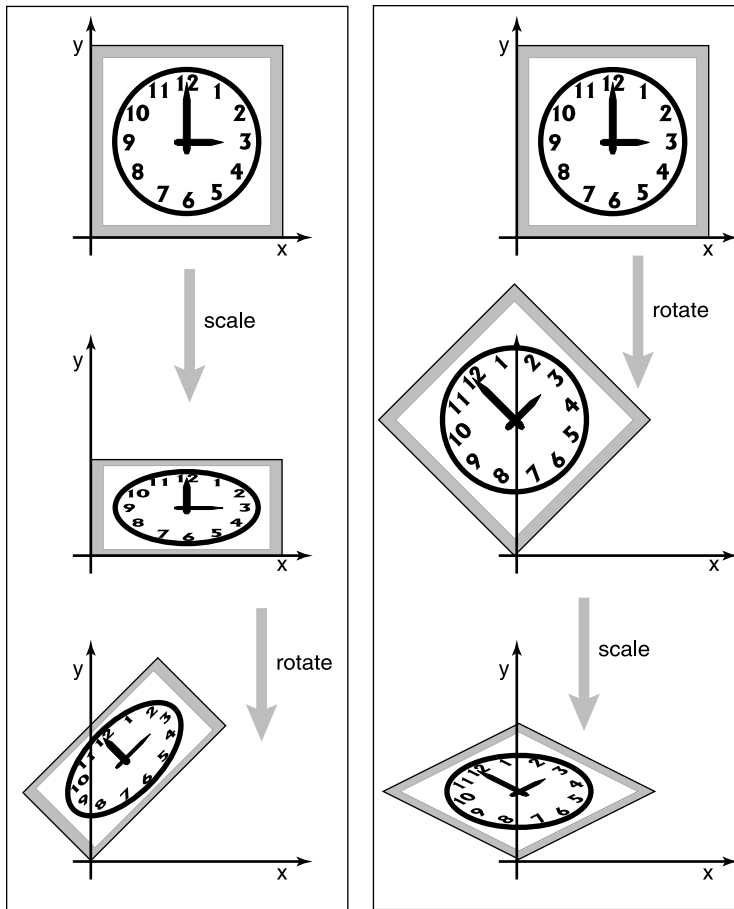
$$\begin{bmatrix} 0.707 & -0.707 \\ 0.707 & 0.707 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 0.5 \end{bmatrix} = \begin{bmatrix} 0.707 & -0.353 \\ 0.707 & 0.353 \end{bmatrix}.$$

It is important to always remember that matrix multiplication is not commutative. So the order of transforms *does* matter. In this example, rotating first, and then scaling, results in a different matrix (see Figure 6.11):

$$\begin{bmatrix} 1 & 0 \\ 0 & 0.5 \end{bmatrix} \begin{bmatrix} 0.707 & -0.707 \\ 0.707 & 0.707 \end{bmatrix} = \begin{bmatrix} 0.707 & -0.707 \\ 0.353 & 0.353 \end{bmatrix}.$$

**Example.** Using the scale matrices we have presented, nonuniform scaling can only be done along the coordinate axes. If we wanted to stretch our clock by 50% along one of its diagonals, so that 8:00 through 1:00 move to the northwest and 2:00 through 7:00 move to the southeast, we can use rotation matrices in combination with an axis-aligned scaling matrix to get the result we want. The idea is to use a rotation to align the scaling axis with a coordinate axis, then scale along that axis, then rotate back. In our example, the scaling axis is the “backslash” diagonal of the square, and we can make it parallel to the  $x$ -axis with





**Figure 6.11.** The order in which two transforms are applied is usually important. In this example, we do a scale by one-half in  $y$  and then rotate by  $45^\circ$ . Reversing the order in which these two transforms are applied yields a different result.

a rotation by  $+45^\circ$ . Putting these operations together, the full transformation is

$$\text{rotate}(-45^\circ) \text{scale}(1.5, 1) \text{rotate}(45^\circ).$$

In mathematical notation, this can be written  $\mathbf{RSR}^T$ . The result of multiplying the three matrices together is

$$\begin{bmatrix} 1.25 & -0.25 \\ -0.25 & 1.25 \end{bmatrix}$$



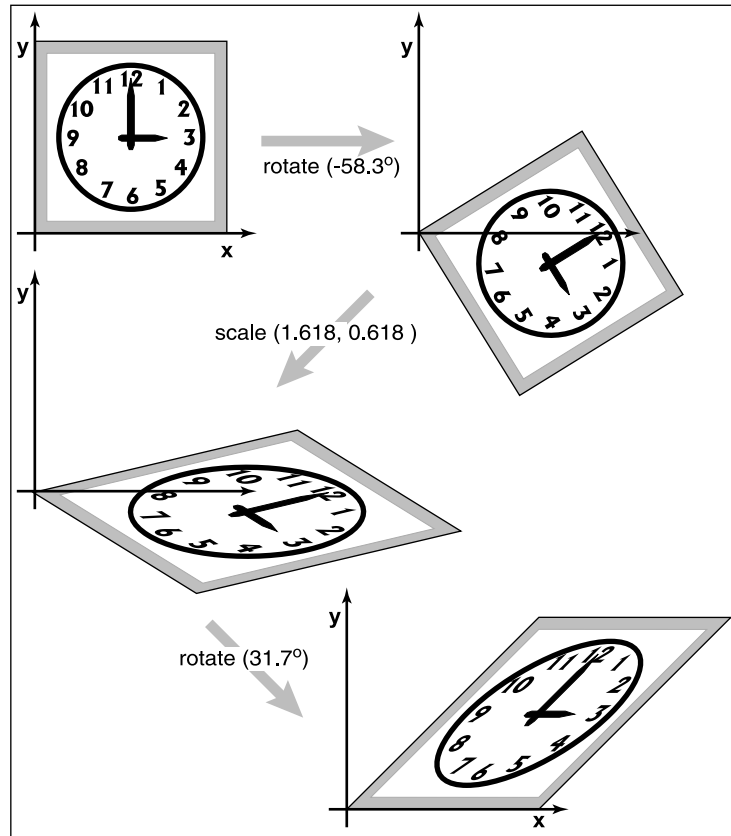
Remember to read the transformations from right to left.

It is no coincidence that this matrix is symmetric—try applying the transpose-of-product rule to the formula  $\mathbf{RSR}^T$ .

Building up a transformation from rotation and scaling transformations actually works for any linear transformation at all, and this fact leads to a powerful way of thinking about these transformations, as explored in the next section.

### 6.1.6 Decomposition of Transformations

Sometimes it's necessary to “undo” a composition of transformations, taking a transformation apart into simpler pieces. For instance, it's often useful to present a transformation to the user for manipulation in terms of separate rotations and scale factors, but a transformation might be represented internally simply as a



**Figure 6.12.** Singular Value Decomposition (SVD) for a shear matrix. Any 2D matrix can be decomposed into a product of rotation, scale, rotation. Note that the circular face of the clock must become an ellipse because it is just a rotated and scaled circle.



matrix, with the rotations and scales already mixed together. This kind of manipulation can be achieved if the matrix can be computationally disassembled into the desired pieces, the pieces adjusted, and the matrix reassembled by multiplying the pieces together again.

It turns out that this decomposition, or factorization, is possible, regardless of the entries in the matrix—and this fact provides a fruitful way of thinking about transformations and what they do to geometry that is transformed by them.

### Symmetric Eigenvalue Decomposition

Let's start with symmetric matrices. Recall from Section 5.4 that a symmetric matrix can always be taken apart using the eigenvalue decomposition into a product of the form

$$\mathbf{A} = \mathbf{R}\mathbf{S}\mathbf{R}^T$$

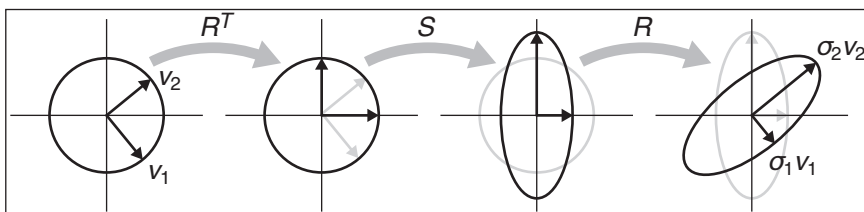
where  $\mathbf{R}$  is an orthogonal matrix and  $\mathbf{S}$  is a diagonal matrix; we will call the columns of  $\mathbf{R}$  (the eigenvectors) by the names  $\mathbf{v}_1$  and  $\mathbf{v}_2$ , and we'll call the diagonal entries of  $\mathbf{S}$  (the eigenvalues) by the names  $\lambda_1$  and  $\lambda_2$ .

In geometric terms we can now recognize  $\mathbf{R}$  as a rotation and  $\mathbf{S}$  as a scale, so this is just a multi-step geometric transformation (Figure 6.13):

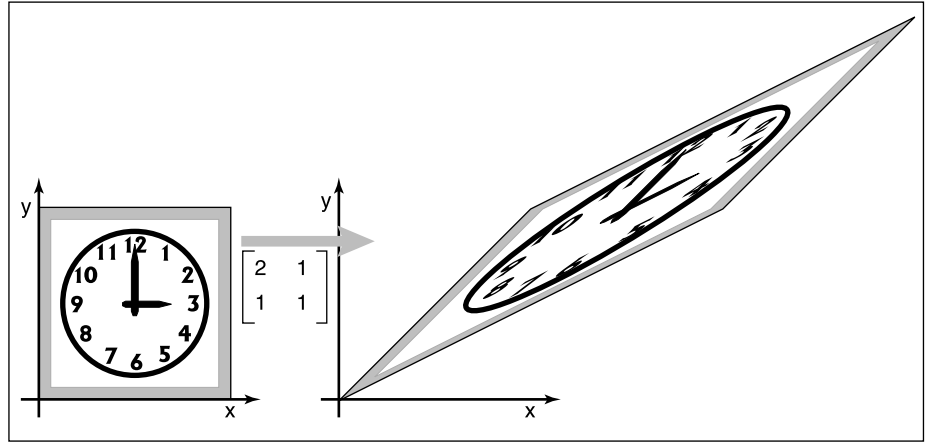
1. Rotate  $\mathbf{v}_1$  and  $\mathbf{v}_2$  to the  $x$ - and  $y$ -axes (the transform by  $\mathbf{R}^T$ ).
2. Scale in  $x$  and  $y$  by  $(\lambda_1, \lambda_2)$  (the transform by  $\mathbf{S}$ ).
3. Rotate the  $x$ - and  $y$ -axes back to  $\mathbf{v}_1$  and  $\mathbf{v}_2$  (the transform by  $\mathbf{R}$ ).

Looking at the effect of these three transforms together, we can see that they have the effect of a nonuniform scale along a pair of axes. As with an axis-aligned scale, the axes are perpendicular, but they aren't the coordinate axes; instead they

If you like to count dimensions: a symmetric  $2 \times 2$  matrix has 3 degrees of freedom, and the eigenvalue decomposition rewrites them as a rotation angle and two scale factors.



**Figure 6.13.** What happens when the unit circle is transformed by an arbitrary symmetric matrix  $\mathbf{A}$ , also known as a non-axis-aligned, nonuniform scale. The two perpendicular vectors  $\mathbf{v}_1$  and  $\mathbf{v}_2$ , which are the eigenvectors of  $\mathbf{A}$ , remain fixed in direction but get scaled. In terms of elementary transformations, this can be seen as first rotating the eigenvectors to the canonical basis, doing an axis-aligned scale, and then rotating the canonical basis back to the eigenvectors.



**Figure 6.14.** A symmetric matrix is always a scale along some axis. In this case it is along the  $\phi = 31.7^\circ$  direction which means the real eigenvector for this matrix is in that direction.

are the eigenvectors of  $\mathbf{A}$ . This tells us something about what it means to be a symmetric matrix: symmetric matrices are just scaling operations—albeit potentially nonuniform and non-axis-aligned ones.

Example. Recall the example from Section 5.4:

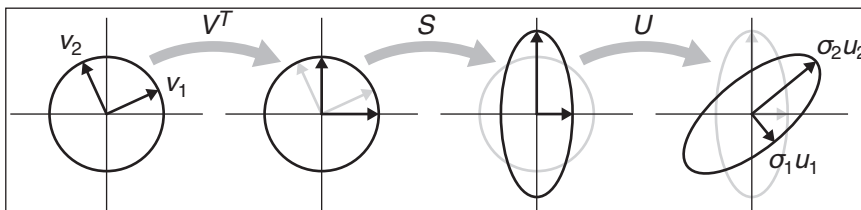
$$\begin{aligned}
 \begin{bmatrix} 2 & 1 \\ 1 & 1 \end{bmatrix} &= \mathbf{R} \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix} \mathbf{R}^T \\
 &= \begin{bmatrix} 0.8507 & -0.5257 \\ 0.5257 & 0.8507 \end{bmatrix} \begin{bmatrix} 2.618 & 0 \\ 0 & 0.382 \end{bmatrix} \begin{bmatrix} 0.8507 & 0.5257 \\ -0.5257 & 0.8507 \end{bmatrix} \\
 &= \text{rotate } (31.7^\circ) \text{ scale } (2.618, 0.382) \text{ rotate } (-31.7^\circ).
 \end{aligned}$$

The matrix above, then, according to its eigenvalue decomposition, scales in a direction  $31.7^\circ$  counterclockwise from three o'clock (the  $x$ -axis). This is a touch before 2 p.m. on the clockface as is confirmed by Figure 6.14.

We can also reverse the diagonalization process; to scale by  $(\lambda_1, \lambda_2)$  with the first scaling direction an angle  $\phi$  clockwise from the  $x$ -axis, we have

$$\begin{bmatrix} \cos \phi & \sin \phi \\ -\sin \phi & \cos \phi \end{bmatrix} \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix} \begin{bmatrix} \cos \phi & -\sin \phi \\ \sin \phi & \cos \phi \end{bmatrix} = \begin{bmatrix} \lambda_1 \cos^2 \phi + \lambda_2 \sin^2 \phi & (\lambda_2 - \lambda_1) \cos \phi \sin \phi \\ (\lambda_2 - \lambda_1) \cos \phi \sin \phi & \lambda_2 \cos^2 \phi + \lambda_1 \sin^2 \phi \end{bmatrix}.$$

We should take heart that this is a symmetric matrix as we know must be true since we constructed it from a symmetric eigenvalue decomposition. □



**Figure 6.15.** What happens when the unit circle is transformed by an arbitrary matrix  $\mathbf{A}$ . The two perpendicular vectors  $\mathbf{v}_1$  and  $\mathbf{v}_2$ , which are the right singular vectors of  $\mathbf{A}$ , get scaled and changed in direction to match the left singular vectors,  $\mathbf{u}_1$  and  $\mathbf{u}_2$ . In terms of elementary transformations, this can be seen as first rotating the right singular vectors to the canonical basis, doing an axis-aligned scale, and then rotating the canonical basis to the left singular vectors.

### Singular Value Decomposition

A very similar kind of decomposition can be done with non-symmetric matrices as well: it's the Singular Value Decomposition (SVD), also discussed in Section 5.4.1. The difference is that the matrices on either side of the diagonal matrix are no longer the same:

$$\mathbf{A} = \mathbf{U}\mathbf{S}\mathbf{V}^T$$

The two orthogonal matrices that replace the single rotation  $\mathbf{R}$  are called  $\mathbf{U}$  and  $\mathbf{V}$ , and their columns are called  $\mathbf{u}_i$  (the *left singular vectors*) and  $\mathbf{v}_i$  (the *right singular vectors*), respectively. In this context, the diagonal entries of  $\mathbf{S}$  are called *singular values* rather than eigenvalues. The geometric interpretation is very similar to that of the symmetric eigenvalue decomposition (Figure 6.15):

1. Rotate  $\mathbf{v}_1$  and  $\mathbf{v}_2$  to the  $x$ - and  $y$ -axes (the transform by  $\mathbf{V}^T$ ).
2. Scale in  $x$  and  $y$  by  $(\sigma_1, \sigma_2)$  (the transform by  $\mathbf{S}$ ).
3. Rotate the  $x$ - and  $y$ -axes to  $\mathbf{u}_1$  and  $\mathbf{u}_2$  (the transform by  $\mathbf{U}$ ).


The principal difference is between a single rotation and two different orthogonal matrices. This difference causes another, less important, difference. Because the SVD has different singular vectors on the two sides, there is no need for negative singular values: we can always flip the sign of a singular value, reverse the direction of one of the associated singular vectors, and end up with the same transformation again. For this reason, the SVD always produces a diagonal matrix with all positive entries, but the matrices  $\mathbf{U}$  and  $\mathbf{V}$  are not guaranteed to be rotations—they could include reflection as well. In geometric applications like graphics this is an inconvenience, but a minor one: it is easy to differentiate rotations from reflections by checking the determinant, which is  $+1$  for rotations

For dimension counters: a general  $2 \times 2$  matrix has 4 degrees of freedom, and the SVD rewrites them as two rotation angles and two scale factors. One more bit is needed to keep track of reflections, but that doesn't add a dimension.

and  $-1$  for reflections, and if rotations are desired, one of the singular values can be negated, resulting in a rotation–scale–rotation sequence where the reflection is rolled in with the scale, rather than with one of the rotations.

Example. The example used in Section 5.4.1 is in fact a shear matrix (Figure 6.12):

$$\begin{aligned} \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix} &= \mathbf{R}_2 \begin{bmatrix} \sigma_1 & 0 \\ 0 & \sigma_2 \end{bmatrix} \mathbf{R}_1 \\ &= \begin{bmatrix} 0.8507 & -0.5257 \\ 0.5257 & 0.8507 \end{bmatrix} \begin{bmatrix} 1.618 & 0 \\ 0 & 0.618 \end{bmatrix} \begin{bmatrix} 0.5257 & 0.8507 \\ -0.8507 & 0.5257 \end{bmatrix} \\ &= \text{rotate}(31.7^\circ) \text{ scale}(1.618, 0.618) \text{ rotate}(-58.3^\circ). \end{aligned}$$

An immediate consequence of the existence of SVD is that all the 2D transformation matrices we have seen can be made from rotation matrices and scale matrices. Shear matrices are a convenience, but they are not required for expressing transformations. 

In summary, every matrix can be decomposed via SVD into a rotation times a scale times another rotation. Only symmetric matrices can be decomposed via eigenvalue diagonalization into a rotation times a scale times the inverse-rotation, and such matrices are a simple scale in an arbitrary direction. The SVD of a symmetric matrix will yield the same triple product as eigenvalue decomposition via a slightly more complex algebraic manipulation.

### Paeth Decomposition of Rotations

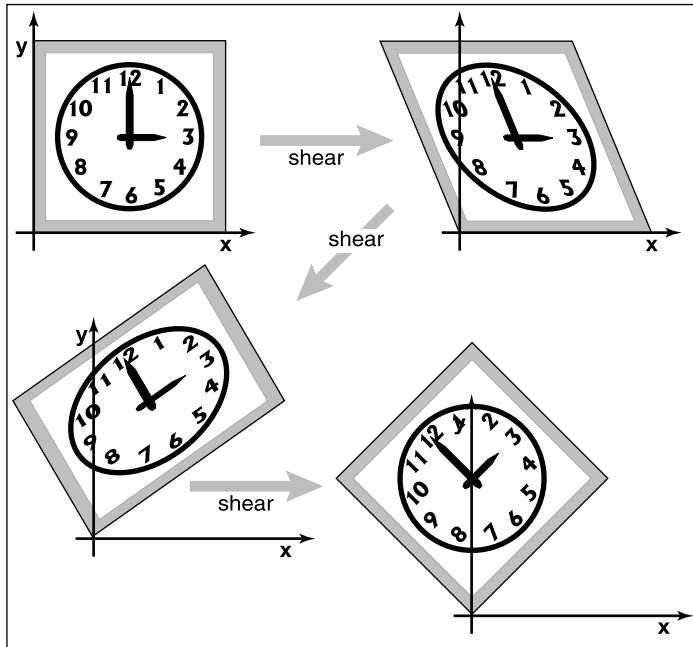
Another decomposition uses shears to represent non-zero rotations (Paeth, 1990). The following identity allows this:

$$\begin{bmatrix} \cos \phi & -\sin \phi \\ \sin \phi & \cos \phi \end{bmatrix} = \begin{bmatrix} 1 & \frac{\cos \phi - 1}{\sin \phi} \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ \sin \phi & 1 \end{bmatrix} \begin{bmatrix} 1 & \frac{\cos \phi - 1}{\sin \phi} \\ 0 & 1 \end{bmatrix}.$$

For example, a rotation by  $\pi/4$  (45 degrees) is (see Figure 6.16)

$$\text{rotate}\left(\frac{\pi}{4}\right) = \begin{bmatrix} 1 & 1 - \sqrt{2} \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ \frac{\sqrt{2}}{2} & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 - \sqrt{2} \\ 0 & 1 \end{bmatrix}. \quad (6.2)$$

This particular transform is useful for raster rotation because shearing is a very efficient raster operation for images; it introduces some jaggedness, but will



**Figure 6.16.** Any 2D rotation can be accomplished by three shears in sequence. In this case a rotation by  $45^\circ$  is decomposed as shown in Equation 6.2.

leave no holes. The key observation is that if we take a raster position  $(i, j)$  and apply a horizontal shear to it, we get

$$\begin{bmatrix} 1 & s \\ 0 & 1 \end{bmatrix} \begin{bmatrix} i \\ j \end{bmatrix} = \begin{bmatrix} i + sj \\ j \end{bmatrix}.$$

If we round  $sj$  to the nearest integer, this amounts to taking each row in the image and moving it sideways by some amount—a different amount for each row. Because it is the same displacement within a row, this allows us to rotate with no gaps in the resulting image. A similar action works for a vertical shear. Thus, we can implement a simple raster rotation easily.

## 6.2 3D Linear Transformations

The linear 3D transforms are an extension of the 2D transforms. For example, a scale along Cartesian axes is

$$\text{scale}(s_x, s_y, s_z) = \begin{bmatrix} s_x & 0 & 0 \\ 0 & s_y & 0 \\ 0 & 0 & s_z \end{bmatrix}. \quad (6.3)$$

Rotation is considerably more complicated in 3D than in 2D, because there are more possible axes of rotation. However, if we simply want to rotate about the  $z$ -axis, which will only change  $x$ - and  $y$ -coordinates, we can use the 2D rotation matrix with no operation on  $z$ :

$$\text{rotate-}z(\phi) = \begin{bmatrix} \cos \phi & -\sin \phi & 0 \\ \sin \phi & \cos \phi & 0 \\ 0 & 0 & 1 \end{bmatrix}.$$

Similarly we can construct matrices to rotate about the  $x$ -axis and the  $y$ -axis:

$$\text{rotate-}x(\phi) = \begin{bmatrix} 1 & 0 & 0 \\ 0 & \cos \phi & -\sin \phi \\ 0 & \sin \phi & \cos \phi \end{bmatrix},$$

$$\text{rotate-}y(\phi) = \begin{bmatrix} \cos \phi & 0 & \sin \phi \\ 0 & 1 & 0 \\ -\sin \phi & 0 & \cos \phi \end{bmatrix}.$$

To understand why the minus sign is in the lower left for the  $y$ -axis rotation, think of the three axes in a circular sequence:  $y$  after  $x$ ;  $z$  after  $y$ ;  $x$  after  $z$ .

We will discuss rotations about arbitrary axes in the next section.

As in two dimensions, we can shear along a particular axis, for example,

$$\text{shear-}x(d_y, d_z) = \begin{bmatrix} 1 & d_y & d_z \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}.$$

As with 2D transforms, any 3D transformation matrix can be decomposed using SVD into a rotation, scale, and another rotation. Any symmetric 3D matrix has an eigenvalue decomposition into rotation, scale, and inverse-rotation. Finally, a 3D rotation can be decomposed into a product of 3D shear matrices.

### 6.2.1 Arbitrary 3D Rotations

As in 2D, 3D rotations are *orthogonal* matrices. Geometrically, this means that the three rows of the matrix are the Cartesian coordinates of three mutually-orthogonal unit vectors as discussed in Section 2.4.5. The columns are three, potentially different, mutually-orthogonal unit vectors. There are an infinite number of such rotation matrices. Let's write down such a matrix:

$$\mathbf{R}_{uvw} = \begin{bmatrix} x_u & y_u & z_u \\ x_v & y_v & z_v \\ x_w & y_w & z_w \end{bmatrix}.$$





Here,  $\mathbf{u} = x_u\mathbf{x} + y_u\mathbf{y} + z_u\mathbf{z}$  and so on for  $\mathbf{v}$  and  $\mathbf{w}$ . Since the three vectors are orthonormal we know that

$$\begin{aligned}\mathbf{u} \cdot \mathbf{u} &= \mathbf{v} \cdot \mathbf{v} = \mathbf{w} \cdot \mathbf{w} = 1, \\ \mathbf{u} \cdot \mathbf{v} &= \mathbf{v} \cdot \mathbf{w} = \mathbf{w} \cdot \mathbf{u} = 0.\end{aligned}$$

We can infer some of the behavior of the rotation matrix by applying it to the vectors  $\mathbf{u}$ ,  $\mathbf{v}$  and  $\mathbf{w}$ . For example,

$$\mathbf{R}_{uvw}\mathbf{u} = \begin{bmatrix} x_u & y_u & z_u \\ x_v & y_v & z_v \\ x_w & y_w & z_w \end{bmatrix} \begin{bmatrix} x_u \\ y_u \\ z_u \end{bmatrix} = \begin{bmatrix} x_u x_u + y_u y_u + z_u z_u \\ x_v x_u + y_v y_u + z_v z_u \\ x_w x_u + y_w y_u + z_w z_u \end{bmatrix}.$$

Note that those three rows of  $\mathbf{R}_{uvw}\mathbf{u}$  are all dot products:

$$\mathbf{R}_{uvw}\mathbf{u} = \begin{bmatrix} \mathbf{u} \cdot \mathbf{u} \\ \mathbf{v} \cdot \mathbf{u} \\ \mathbf{w} \cdot \mathbf{u} \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} = \mathbf{x}.$$

Similarly,  $\mathbf{R}_{uvw}\mathbf{v} = \mathbf{y}$ , and  $\mathbf{R}_{uvw}\mathbf{w} = \mathbf{z}$ . So  $\mathbf{R}_{uvw}$  takes the basis  $\mathbf{uvw}$  to the corresponding Cartesian axes via rotation.

If  $\mathbf{R}_{uvw}$  is a rotation matrix with orthonormal rows, then  $\mathbf{R}_{uvw}^T$  is also a rotation matrix with orthonormal columns, and in fact is the inverse of  $\mathbf{R}_{uvw}$  (the inverse of an orthogonal matrix is always its transpose). An important point is that for transformation matrices, the algebraic inverse is also the geometric inverse. So if  $\mathbf{R}_{uvw}$  takes  $\mathbf{u}$  to  $\mathbf{x}$ , then  $\mathbf{R}_{uvw}^T$  takes  $\mathbf{x}$  to  $\mathbf{u}$ . The same should be true of  $\mathbf{v}$  and  $\mathbf{y}$  as we can confirm:

$$\mathbf{R}_{uvw}^T\mathbf{y} = \begin{bmatrix} x_u & x_v & x_w \\ y_u & y_v & y_w \\ z_u & z_v & z_w \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} x_v \\ y_v \\ z_v \end{bmatrix} = \mathbf{v}.$$

So we can always create rotation matrices from orthonormal bases.

If we wish to rotate about an arbitrary vector  $\mathbf{a}$ , we can form an orthonormal basis with  $\mathbf{w} = \mathbf{a}$ , rotate that basis to the canonical basis  $\mathbf{xyz}$ , rotate about the  $z$ -axis, and then rotate the canonical basis back to the  $\mathbf{uvw}$  basis. In matrix form, to rotate about the  $w$ -axis by an angle  $\phi$ :

$$\begin{bmatrix} x_u & x_v & x_w \\ y_u & y_v & y_w \\ z_u & z_v & z_w \end{bmatrix} \begin{bmatrix} \cos \phi & -\sin \phi & 0 \\ \sin \phi & \cos \phi & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x_u & y_u & z_u \\ x_v & y_v & z_v \\ x_w & y_w & z_w \end{bmatrix}.$$

Here we have  $\mathbf{w}$  a unit vector in the direction of  $\mathbf{a}$  (i.e.  $\mathbf{a}$  divided by its own length). But what are  $\mathbf{u}$  and  $\mathbf{v}$ ? A method to find reasonable  $\mathbf{u}$  and  $\mathbf{v}$  is given in Section 2.4.6.

If we have a rotation matrix and we wish to have the rotation in axis-angle form, we can compute the one real eigenvalue (which will be  $\lambda = 1$ ), and the corresponding eigenvector is the axis of rotation. This is the one axis that is not changed by the rotation.

See Chapter 17 for a comparison of the few most-used ways to represent rotations, besides rotation matrices.

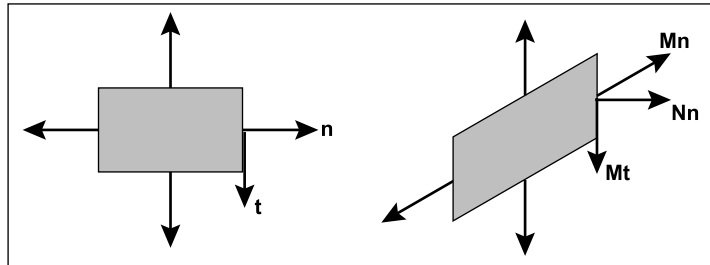
### 6.2.2 Transforming Normal Vectors

While most 3D vectors we use represent positions (offset vectors from the origin) or directions, such as where light comes from, some vectors represent *surface normals*. Surface normal vectors are perpendicular to the tangent plane of a surface. These normals do not transform the way we would like when the underlying surface is transformed. For example, if the points of a surface are transformed by a matrix  $M$ , a vector  $t$  that is tangent to the surface and is multiplied by  $M$  will be tangent to the transformed surface. However, a surface normal vector  $n$  that is transformed by  $M$  may not be normal to the transformed surface (Figure 6.17).

We can derive a transform matrix  $N$  which does take  $n$  to a vector perpendicular to the transformed surface. One way to attack this issue is to note that a surface normal vector and a tangent vector are perpendicular, so their dot product is zero, which is expressed in matrix form as

$$n^T t = 0. \quad (6.4)$$

If we denote the desired transformed vectors as  $t_M = Mt$  and  $n_N = Nn$ , our goal is to find  $N$  such that  $n_N^T t_M = 0$ . We can find  $N$  by some algebraic



**Figure 6.17.** When a normal vector is transformed using the same matrix that transforms the points on an object, the resulting vector may not be perpendicular to the surface as is shown here for the sheared rectangle. The tangent vector, however, does transform to a vector tangent to the transformed surface.



tricks. First, we can sneak an identity matrix into the dot product, and then take advantage of  $\mathbf{M}^{-1}\mathbf{M} = \mathbf{I}$ :

$$\mathbf{n}^T \mathbf{t} = \mathbf{n}^T \mathbf{I} \mathbf{t} = \mathbf{n}^T \mathbf{M}^{-1} \mathbf{M} \mathbf{t} = 0.$$

Although the manipulations above don't obviously get us anywhere, note that we can add parentheses that make the above expression more obviously a dot product:

$$(\mathbf{n}^T \mathbf{M}^{-1}) (\mathbf{M} \mathbf{t}) = (\mathbf{n}^T \mathbf{M}^{-1}) \mathbf{t}_M = 0.$$

This means that the row vector that is perpendicular to  $\mathbf{t}_M$  is the left part of the expression above. This expression holds for any of the tangent vectors in the tangent plane. Since there is only one direction in 3D (and its opposite) that is perpendicular to all such tangent vectors, we know that the left part of the expression above must be the row vector expression for  $\mathbf{n}_N$ , i.e., it is  $\mathbf{n}_N^T$ , so this allows us to infer  $\mathbf{N}$ :

$$\mathbf{n}_N^T = \mathbf{n}^T \mathbf{M}^{-1},$$

so we can take the transpose of that to get

$$\mathbf{n}_N = (\mathbf{n}^T \mathbf{M}^{-1})^T = (\mathbf{M}^{-1})^T \mathbf{n}. \quad (6.5)$$

Therefore, we can see that the matrix which correctly transforms normal vectors so they remain normal is  $\mathbf{N} = (\mathbf{M}^{-1})^T$ , i.e., the transpose of the inverse matrix. Since this matrix may change the length of  $\mathbf{n}$ , we can multiply it by an arbitrary scalar and it will still produce  $\mathbf{n}_N$  with the right direction. Recall from Section 5.3 that the inverse of a matrix is the transpose of the cofactor matrix divided by the determinant. Because we don't care about the length of a normal vector, we can skip the division and find that for a  $3 \times 3$  matrix,

$$\mathbf{N} = \begin{bmatrix} m_{11}^c & m_{12}^c & m_{13}^c \\ m_{21}^c & m_{22}^c & m_{23}^c \\ m_{31}^c & m_{32}^c & m_{33}^c \end{bmatrix}.$$

This assumes the element of  $\mathbf{M}$  in row  $i$  and column  $j$  is  $m_{ij}$ . So the full expression for  $\mathbf{N}$  is

$$\mathbf{N} = \begin{bmatrix} m_{22}m_{33} - m_{23}m_{32} & m_{23}m_{31} - m_{21}m_{33} & m_{21}m_{32} - m_{22}m_{31} \\ m_{13}m_{32} - m_{12}m_{33} & m_{11}m_{33} - m_{13}m_{31} & m_{12}m_{31} - m_{11}m_{32} \\ m_{12}m_{23} - m_{13}m_{22} & m_{13}m_{21} - m_{11}m_{23} & m_{11}m_{22} - m_{12}m_{21} \end{bmatrix}.$$

## 6.3 Translation and Affine Transformations

We have been looking at methods to change vectors using a matrix  $\mathbf{M}$ . In two dimensions, these transforms have the form,

$$\begin{aligned}x' &= m_{11}x + m_{12}y, \\y' &= m_{21}x + m_{22}y.\end{aligned}$$

We cannot use such transforms to *move* objects, only to scale and rotate them. In particular, the origin  $(0, 0)$  always remains fixed under a linear transformation. To move, or *translate*, an object by shifting all its points the same amount, we need a transform of the form,

$$\begin{aligned}x' &= x + x_t, \\y' &= y + y_t.\end{aligned}$$

There is just no way to do that by multiplying  $(x, y)$  by a  $2 \times 2$  matrix. One possibility for adding translation to our system of linear transformations is to simply associate a separate translation vector with each transformation matrix, letting the matrix take care of scaling and rotation and the vector take care of translation. This is perfectly feasible, but the bookkeeping is awkward and the rule for composing two transformations is not as simple and clean as with linear transformations.

Instead, we can use a clever trick to get a single matrix multiplication to do both operations together. The idea is simple: represent the point  $(x, y)$  by a 3D vector  $[x \ y \ 1]^T$ , and use  $3 \times 3$  matrices of the form

$$\begin{bmatrix} m_{11} & m_{12} & x_t \\ m_{21} & m_{22} & y_t \\ 0 & 0 & 1 \end{bmatrix}$$

The fixed third row serves to copy the 1 into the transformed vector, so that all vectors have a 1 in the last place, and the first two rows compute  $x'$  and  $y'$  as linear combinations of  $x$ ,  $y$ , and 1:

$$\begin{bmatrix} x' \\ y' \\ 1 \end{bmatrix} = \begin{bmatrix} m_{11} & m_{12} & x_t \\ m_{21} & m_{22} & y_t \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = \begin{bmatrix} m_{11}x + m_{12}y + x_t \\ m_{21}x + m_{22}y + y_t \\ 1 \end{bmatrix}.$$

The single matrix implements a linear transformation followed by a translation! This kind of transformation is called an *affine transformation*, and this way of implementing affine transformations by adding an extra dimension is called *homogeneous coordinates* (Roberts, 1965; Riesenfeld, 1981; Penna & Patterson, 1986). Homogeneous coordinates not only clean up the code for transformations,



but this scheme also makes it obvious how to compose two affine transformations: simply multiply the matrices.

A problem with this new formalism arises when we need to transform vectors that are not supposed to be positions—they represent directions, or offsets between positions. Vectors that represent directions or offsets should not change when we translate an object. Fortunately, we can arrange for this by setting the third coordinate to zero:

$$\begin{bmatrix} 1 & 0 & x_t \\ 0 & 1 & y_t \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ 0 \end{bmatrix} = \begin{bmatrix} x \\ y \\ 0 \end{bmatrix}.$$

If there is a scaling/rotation transformation in the upper-left  $2 \times 2$  entries of the matrix, it will apply to the vector, but the translation still multiplies with the zero and is ignored. Furthermore, the zero is copied into the transformed vector, so direction vectors remain direction vectors after they are transformed.

This is exactly the behavior we want for vectors, so they fit smoothly into the system: the extra (third) coordinate will be either 1 or 0 depending on whether we are encoding a position or a direction. We actually do need to store the homogeneous coordinate so we can distinguish between locations and other vectors. For example,

$$\begin{bmatrix} 3 \\ 2 \\ 1 \end{bmatrix} \text{ is a location} \quad \text{and} \quad \begin{bmatrix} 3 \\ 2 \\ 0 \end{bmatrix} \text{ is a displacement or direction.}$$

Later, when we do perspective viewing, we will see that it is useful to allow the homogeneous coordinate to take on values other than one or zero.

Homogeneous coordinates are used nearly universally to represent transformations in graphics systems. In particular, homogeneous coordinates underlie the design and operation of renderers implemented in graphics hardware. We will see in Chapter 7 that homogeneous coordinates also make it easy to draw scenes in perspective, another reason for their popularity.

Homogeneous coordinates can be considered just a clever way to handle the bookkeeping for translation, but there is also a different, geometric interpretation. The key observation is that when we do a 3D shear based on the  $z$ -coordinate we get this transform:

$$\begin{bmatrix} 1 & 0 & x_t \\ 0 & 1 & y_t \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \end{bmatrix} = \begin{bmatrix} x + x_t z \\ y + y_t z \\ z \end{bmatrix}.$$

Note that this almost has the form we want in  $x$  and  $y$  for a 2D translation, but has a  $z$  hanging around that doesn't have a meaning in 2D. Now comes the key

This gives an explanation for the name “homogeneous:” translation, rotation, and scaling of positions and directions all fit into a single system.

Homogeneous coordinates are also ubiquitous in computer vision.

decision: we will add a coordinate  $z = 1$  to all 2D locations. This gives us

$$\begin{bmatrix} 1 & 0 & x_t \\ 0 & 1 & y_t \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = \begin{bmatrix} x + x_t \\ y + y_t \\ 1 \end{bmatrix}.$$

By associating a ( $z = 1$ )-coordinate with all 2D points, we now can encode translations into matrix form. For example, to first translate in 2D by  $(t_x, t_y)$  and then rotate by angle  $\phi$  we would use the matrix

$$\mathbf{M} = \begin{bmatrix} \cos \phi & -\sin \phi & 0 \\ \sin \phi & \cos \phi & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & x_t \\ 0 & 1 & y_t \\ 0 & 0 & 1 \end{bmatrix}.$$

Note that the 2D rotation matrix is now  $3 \times 3$  with zeros in the “translation slots.” With this type of formalism, which uses shears along  $z = 1$  to encode translations, we can represent any number of 2D shears, 2D rotations, and 2D translations as one composite 3D matrix. The bottom row of that matrix will always be  $(0, 0, 1)$ , so we don’t really have to store it. We just need to remember it is there when we multiply two matrices together.

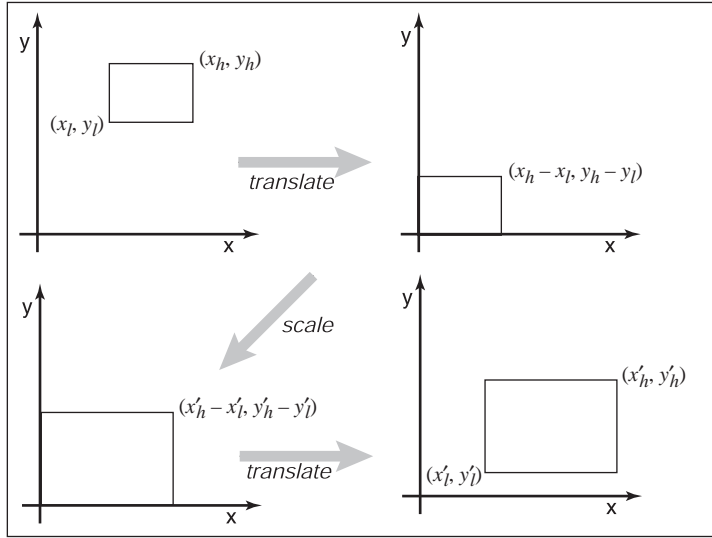
In 3D, the same technique works: we can add a fourth coordinate, a homogeneous coordinate, and then we have translations:

$$\begin{bmatrix} 1 & 0 & 0 & x_t \\ 0 & 1 & 0 & y_t \\ 0 & 0 & 1 & z_t \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \\ 1 \end{bmatrix} = \begin{bmatrix} x + x_t \\ y + y_t \\ z + z_t \\ 1 \end{bmatrix}.$$

Again, for a direction vector, the fourth coordinate is zero and the vector is thus unaffected by translations.

**Example (Windowing transformations).** Often in graphics we need to create a transform matrix that takes points in the rectangle  $[x_l, x_h] \times [y_l, y_h]$  to the rectangle  $[x'_l, x'_h] \times [y'_l, y'_h]$ . This can be accomplished with a single scale and translate in sequence. However, it is more intuitive to create the transform from a sequence of three operations (Figure 6.18):

1. Move the point  $(x_l, y_l)$  to the origin.
2. Scale the rectangle to be the same size as the target rectangle.
3. Move the origin to point  $(x'_l, y'_l)$ .



**Figure 6.18.** To take one rectangle (window) to the other, we first shift the lower-left corner to the origin, then scale it to the new size, and then move the origin to the lower-left corner of the target rectangle.

Remembering that the right-hand matrix is applied first, we can write

$$\begin{aligned}
 \text{window} &= \text{translate}(x'_l, y'_l) \text{ scale} \left( \frac{x'_h - x'_l}{x_h - x_l}, \frac{y'_h - y'_l}{y_h - y_l} \right) \text{translate}(-x_l, -y_l) \\
 &= \begin{bmatrix} 1 & 0 & x'_l \\ 0 & 1 & y'_l \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} \frac{x'_h - x'_l}{x_h - x_l} & 0 & 0 \\ 0 & \frac{y'_h - y'_l}{y_h - y_l} & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & -x_l \\ 0 & 1 & -y_l \\ 0 & 0 & 1 \end{bmatrix} \\
 &= \begin{bmatrix} \frac{x'_h - x'_l}{x_h - x_l} & 0 & \frac{x'_l x_h - x'_h x_l}{x_h - x_l} \\ 0 & \frac{y'_h - y'_l}{y_h - y_l} & \frac{y'_l y_h - y'_h y_l}{y_h - y_l} \\ 0 & 0 & 1 \end{bmatrix}. \tag{6.6}
 \end{aligned}$$

It is perhaps not surprising to some readers that the resulting matrix has the form it does, but the constructive process with the three matrices leaves no doubt as to the correctness of the result.

An exactly analogous construction can be used to define a 3D windowing transformation, which maps the box  $[x_l, x_h] \times [y_l, y_h] \times [z_l, z_h]$  to the box

$$[x'_l, x'_h] \times [y'_l, y'_h] \times [z'_l, z'_h]:$$

$$\begin{bmatrix} \frac{x'_h - x'_l}{x_h - x_l} & 0 & 0 & \frac{x'_l x_h - x'_h x_l}{x_h - x_l} \\ 0 & \frac{y'_h - y'_l}{y_h - y_l} & 0 & \frac{y'_l y_h - y'_h y_l}{y_h - y_l} \\ 0 & 0 & \frac{z'_h - z'_l}{z_h - z_l} & \frac{z'_l z_h - z'_h z_l}{z_h - z_l} \\ 0 & 0 & 0 & 1 \end{bmatrix}. \quad (6.7)$$

It is interesting to note that if we multiply an arbitrary matrix composed of scales, shears, and rotations with a simple translation (translation comes second), we get

$$\begin{bmatrix} 1 & 0 & 0 & x_t \\ 0 & 1 & 0 & y_t \\ 0 & 0 & 1 & z_t \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} a_{11} & a_{12} & a_{13} & 0 \\ a_{21} & a_{22} & a_{23} & 0 \\ a_{31} & a_{32} & a_{33} & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} & a_{13} & x_t \\ a_{21} & a_{22} & a_{23} & y_t \\ a_{31} & a_{32} & a_{33} & z_t \\ 0 & 0 & 0 & 1 \end{bmatrix}.$$

Thus, we can look at any matrix and think of it as a scaling/rotation part and a translation part because the components are nicely separated from each other.

An important class of transforms are *rigid-body* transforms. These are composed only of translations and rotations, so they have no stretching or shrinking of the objects. Such transforms will have a pure rotation for the  $a_{ij}$  above.

## 6.4 Inverses of Transformation Matrices

While we can always invert a matrix algebraically, we can use geometry if we know what the transform does. For example, the inverse of  $\text{scale}(s_x, s_y, s_z)$  is  $\text{scale}(1/s_x, 1/s_y, 1/s_z)$ . The inverse of a rotation is the same rotation with the opposite sign on the angle. The inverse of a translation is a translation in the opposite direction. If we have a series of matrices  $\mathbf{M} = \mathbf{M}_1 \mathbf{M}_2 \cdots \mathbf{M}_n$  then  $\mathbf{M}^{-1} = \mathbf{M}_n^{-1} \cdots \mathbf{M}_2^{-1} \mathbf{M}_1^{-1}$ .

Also, certain types of transformation matrices are easy to invert. We've already mentioned scales, which are diagonal matrices; the second important example is rotations, which are orthogonal matrices. Recall (Section 5.2.4) that the inverse of an orthogonal matrix is its transpose. This makes it easy to invert rotations and rigid body transformations (see Exercise 6). Also, it's useful to know that a matrix with  $[0 \ 0 \ 0 \ 1]$  in the bottom row has an inverse that also has  $[0 \ 0 \ 0 \ 1]$  in the bottom row (see Exercise 7).

Interestingly, we can use SVD to invert a matrix as well. Since we know that any matrix can be decomposed into a rotation times a scale times a rotation,





inversion is straightforward. For example in 3D we have

$$\mathbf{M} = \mathbf{R}_1 \text{scale}(\sigma_1, \sigma_2, \sigma_3) \mathbf{R}_2,$$

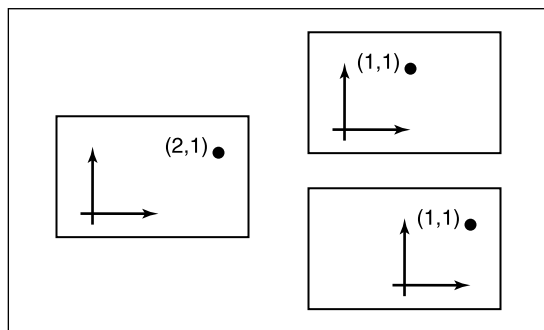
and from the rules above it follows easily that

$$\mathbf{M}^{-1} = \mathbf{R}_2^T \text{scale}(1/\sigma_1, 1/\sigma_2, 1/\sigma_3) \mathbf{R}_1^T.$$

## 6.5 Coordinate Transformations

All of the previous discussion has been in terms of using transformation matrices to move points around. We can also think of them as simply changing the coordinate system in which the point is represented. For example, in Figure 6.19, we see two ways to visualize a movement. In different contexts, either interpretation may be more suitable.

For example, a driving game may have a model of a city and a model of a car. If the player is presented with a view out the windshield, objects inside the car are always drawn in the same place on the screen, while the streets and buildings appear to move backward as the player drives. On each frame, we apply a transformation to these objects that moves them farther back than on the previous frame. One way to think of this operation is simply that it moves the buildings backward; another way to think of it is that the buildings are staying put but the coordinate system in which we want to draw them—which is attached to the car—is moving. In the second interpretation, the transformation is changing



**Figure 6.19.** The point (2,1) has a transform “translate by (-1,0)” applied to it. On the top right is our mental image if we view this transformation as a physical movement, and on the bottom right is our mental image if we view it as a change of coordinates (a movement of the origin in this case). The artificial boundary is just an artifice, and the relative position of the axes and the point are the same in either case.



the coordinates of the city geometry, expressing them as coordinates in the car's coordinate system. Both ways will lead to exactly the same matrix that is applied to the geometry outside the car.

If the game also supports an overhead view to show where the car is in the city, the buildings and streets need to be drawn in fixed positions while the car needs to move from frame to frame. The same two interpretations apply: we can think of the changing transformation as moving the car from its canonical position to its current location in the world; or we can think of the transformation as simply changing the coordinates of the car's geometry, which is originally expressed in terms of a coordinate system attached to the car, to express them instead in a coordinate system fixed relative to the city. The change-of-coordinates interpretation makes it clear that the matrices used in these two modes (city-to-car coordinate change vs. car-to-city coordinate change) are inverses of one another.

The idea of changing coordinate systems is much like the idea of type conversions in programming. Before we can add a floating-point number to an integer, we need to convert the integer to floating point or the floating-point number to an integer, depending on our needs, so that the types match. And before we can draw the city and the car together, we need to convert the city to car coordinates or the car to city coordinates, depending on our needs, so that the coordinates match.

When managing multiple coordinate systems, it's easy to get confused and wind up with objects in the wrong coordinates, causing them to show up in unexpected places. But with systematic thinking about transformations between coordinate systems, you can reliably get the transformations right.

In 2D, of course, there are two basis vectors.

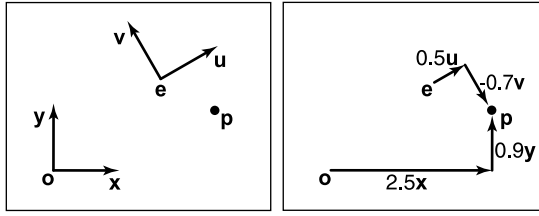
Geometrically, a coordinate system, or coordinate *frame*, consists of an origin and a basis—a set of three vectors. Orthonormal bases are so convenient that we'll normally assume frames are orthonormal unless otherwise specified. In a frame with origin  $\mathbf{p}$  and basis  $\{\mathbf{u}, \mathbf{v}, \mathbf{w}\}$ , the coordinates  $(u, v, w)$  describe the point

$$\mathbf{p} + u\mathbf{u} + v\mathbf{v} + w\mathbf{w}.$$

When we store these vectors in the computer, they need to be represented in terms of some coordinate system. To get things started, we have to designate some canonical coordinate system, often called “global” or “world” coordinates, which is used to describe all other systems. In the city example, we might adopt the street grid and use the convention that the  $x$ -axis points along Main Street, the  $y$ -axis points up, and the  $z$ -axis points along Central Avenue. Then when we write the origin and basis of the car frame in terms of these coordinates it is clear what we mean.

In 2D, right handed means  $\mathbf{y}$  is counter-clockwise from  $\mathbf{x}$ .

In 2D our convention is to use the point  $\mathbf{o}$  for the origin, and  $\mathbf{x}$  and  $\mathbf{y}$  for the right-handed orthonormal basis vectors  $\mathbf{x}$  and  $\mathbf{y}$  (Figure 6.20).



**Figure 6.20.** The point  $\mathbf{p}$  can be represented in terms of either coordinate system.

Another coordinate system might have an origin  $\mathbf{e}$  and right-handed orthonormal basis vectors  $\mathbf{u}$  and  $\mathbf{v}$ . Note that typically the canonical data  $\mathbf{o}$ ,  $\mathbf{x}$ , and  $\mathbf{y}$  are never stored explicitly. They are the frame-of-reference for all other coordinate systems. In that coordinate system, we often write down the location of  $\mathbf{p}$  as an ordered pair, which is shorthand for a full vector expression:

$$\mathbf{p} = (x_p, y_p) \equiv \mathbf{o} + x_p \mathbf{x} + y_p \mathbf{y}.$$

For example, in Figure 6.20,  $(x_p, y_p) = (2.5, 0.9)$ . Note that the pair  $(x_p, y_p)$  implicitly assumes the origin  $\mathbf{o}$ . Similarly, we can express  $\mathbf{p}$  in terms of another equation:

$$\mathbf{p} = (u_p, v_p) \equiv \mathbf{e} + u_p \mathbf{u} + v_p \mathbf{v}.$$

In Figure 6.20, this has  $(u_p, v_p) = (0.5, -0.7)$ . Again, the origin  $\mathbf{e}$  is left as an implicit part of the coordinate system associated with  $\mathbf{u}$  and  $\mathbf{v}$ .

We can express this same relationship using matrix machinery, like this:

$$\begin{bmatrix} x_p \\ y_p \\ 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 & x_e \\ 0 & 1 & y_e \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x_u & x_v & 0 \\ y_u & y_v & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} u_p \\ v_p \\ 1 \end{bmatrix} = \begin{bmatrix} x_u & x_v & x_e \\ y_u & y_v & y_e \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} u_p \\ v_p \\ 1 \end{bmatrix}.$$

Note that this assumes we have the point  $\mathbf{e}$  and vectors  $\mathbf{u}$  and  $\mathbf{v}$  stored in canonical coordinates; the  $(x, y)$ -coordinate system is the first among equals. In terms of the basic types of transformations we've discussed in this chapter, this is a rotation (involving  $\mathbf{u}$  and  $\mathbf{v}$ ) followed by a translation (involving  $\mathbf{e}$ ). Looking at the matrix for the rotation and translation together, you can see it's very easy to write down: we just put  $\mathbf{u}$ ,  $\mathbf{v}$ , and  $\mathbf{e}$  into the columns of a matrix, with the usual  $[0 \ 0 \ 1]$  in the third row. To make this even clearer we can write the matrix like this:

$$\mathbf{P}_{xy} = \begin{bmatrix} \mathbf{u} & \mathbf{v} & \mathbf{e} \\ 0 & 0 & 1 \end{bmatrix} \mathbf{P}_{uv}.$$

We call this matrix the *frame-to-canonical* matrix for the  $(u, v)$  frame. It takes points expressed in the  $(u, v)$  frame and converts them to the same points expressed in the canonical frame.

The name “frame-to-canonical” is based on thinking about changing the coordinates of a vector from one system to another. Thinking in terms of moving vectors around, the frame-to-canonical matrix maps the canonical frame to the  $(u, v)$  frame.

To go in the other direction we have

$$\begin{bmatrix} u_p \\ v_p \\ 1 \end{bmatrix} = \begin{bmatrix} x_u & y_u & 0 \\ x_v & y_v & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & -x_e \\ 0 & 1 & -y_e \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x_p \\ y_p \\ 1 \end{bmatrix}.$$

This is a translation followed by a rotation; they are the inverses of the rotation and translation we used to build the frame-to-canonical matrix, and when multiplied together they produce the inverse of the frame-to-canonical matrix, which is (not surprisingly) called the canonical-to-frame matrix:

$$\mathbf{p}_{uv} = \begin{bmatrix} \mathbf{u} & \mathbf{v} & \mathbf{e} \\ 0 & 0 & 1 \end{bmatrix}^{-1} \mathbf{p}_{xy}.$$

The canonical-to-frame matrix takes points expressed in the canonical frame and converts them to the same points expressed in the  $(u,v)$  frame. We have written this matrix as the inverse of the frame-to-canonical matrix because it can't immediately be written down using the canonical coordinates of  $\mathbf{e}$ ,  $\mathbf{u}$ , and  $\mathbf{v}$ . But remember that all coordinate systems are equivalent; it's only our convention of storing vectors in terms of  $x$ - and  $y$ -coordinates that creates this seeming asymmetry. The canonical-to-frame matrix *can* be expressed simply in terms of the  $(u, v)$  coordinates of  $\mathbf{o}$ ,  $\mathbf{x}$ , and  $\mathbf{y}$ :

$$\mathbf{p}_{uv} = \begin{bmatrix} \mathbf{x}_{uv} & \mathbf{y}_{uv} & \mathbf{o}_{uv} \\ 0 & 0 & 1 \end{bmatrix} \mathbf{p}_{xy}.$$

All these ideas work strictly analogously in 3D, where we have

$$\begin{aligned} \begin{bmatrix} x_p \\ y_p \\ z_p \\ 1 \end{bmatrix} &= \begin{bmatrix} 1 & 0 & 0 & x_e \\ 0 & 1 & 0 & y_e \\ 0 & 0 & 1 & z_e \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x_u & x_v & x_w & 0 \\ y_u & y_v & y_w & 0 \\ z_u & z_v & z_w & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} u_p \\ v_p \\ w_p \\ 1 \end{bmatrix} \\ \mathbf{p}_{xyz} &= \begin{bmatrix} \mathbf{u} & \mathbf{v} & \mathbf{w} & \mathbf{e} \\ 0 & 0 & 0 & 1 \end{bmatrix} \mathbf{p}_{uvw}, \end{aligned} \quad (6.8)$$

and

$$\begin{aligned} \begin{bmatrix} u_p \\ v_p \\ w_p \\ 1 \end{bmatrix} &= \begin{bmatrix} x_u & y_u & z_u & 0 \\ x_v & y_v & z_v & 0 \\ x_w & y_w & z_w & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & -x_e \\ 0 & 1 & 0 & -y_e \\ 0 & 0 & 1 & -z_e \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x_p \\ y_p \\ z_p \\ 1 \end{bmatrix} \\ \mathbf{p}_{uvw} &= \begin{bmatrix} \mathbf{u} & \mathbf{v} & \mathbf{w} & \mathbf{e} \\ 0 & 0 & 0 & 1 \end{bmatrix}^{-1} \mathbf{p}_{xyz}. \end{aligned} \quad (6.9)$$



## Frequently Asked Questions

- Can't I just hardcode transforms rather than use the matrix formalisms?

Yes, but in practice it is harder to derive, harder to debug, and not any more efficient. Also, all current graphics APIs use this matrix formalism so it must be understood even to use graphics libraries.

- The bottom row of the matrix is always (0,0,0,1). Do I have to store it?

You do not have to store it unless you include perspective transforms (Chapter 7).

## Notes

The derivation of the transformation properties of normals is based on *Properties of Surface Normal Transformations* (Turkowsky, 1990). In many treatments through the mid-1990s, vectors were represented as row vectors and premultiplied, e.g.,  $\mathbf{b} = \mathbf{aM}$ . In our notation this would be  $\mathbf{b}^T = \mathbf{a}^T \mathbf{M}^T$ . If you want to find a rotation matrix  $\mathbf{R}$  that takes one vector  $\mathbf{a}$  to a vector  $\mathbf{b}$  of the same length:  $\mathbf{b} = \mathbf{Ra}$  you could use two rotations constructed from orthonormal bases. A more efficient method is given in *Efficiently Building a Matrix to Rotate One Vector to Another* (Akenine-Möller et al., 2008).

## Exercises

1. Write down the  $4 \times 4$  3D matrix to move by  $(x_m, y_m, z_m)$ .
2. Write down the  $4 \times 4$  3D matrix to rotate by an angle  $\theta$  about the  $y$ -axis.
3. Write down the  $4 \times 4$  3D matrix to scale an object by 50% in all directions.
4. Write the 2D rotation matrix that rotates by 90 degrees clockwise.
5. Write the matrix from Exercise 4 as a product of three shear matrices.
6. Find the inverse of the rigid body transformation:

$$\begin{bmatrix} \mathbf{R} & \mathbf{t} \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

where  $\mathbf{R}$  is a  $3 \times 3$  rotation matrix and  $\mathbf{t}$  is a 3-vector.



7. Show that the inverse of the matrix for an affine transformation (one that has all zeros in the bottom row except for a one in the lower right entry) also has the same form.
8. Describe in words what this 2D transform matrix does:

$$\begin{bmatrix} 0 & -1 & 1 \\ 1 & 0 & 1 \\ 0 & 0 & 1 \end{bmatrix}.$$

9. Write down the  $3 \times 3$  matrix that rotates a 2D point by angle  $\theta$  about a point  $\mathbf{p} = (x_p, y_p)$ .
10. Write down the  $4 \times 4$  rotation matrix that takes the orthonormal 3D vectors  $\mathbf{u} = (x_u, y_u, z_u)$ ,  $\mathbf{v} = (x_v, y_v, z_v)$ , and  $\mathbf{w} = (x_w, y_w, z_w)$ , to orthonormal 3D vectors  $\mathbf{a} = (x_a, y_a, z_a)$ ,  $\mathbf{b} = (x_b, y_b, z_b)$ , and  $\mathbf{c} = (x_c, y_c, z_c)$ . So  $M\mathbf{u} = \mathbf{a}$ ,  $M\mathbf{v} = \mathbf{b}$ , and  $M\mathbf{w} = \mathbf{c}$ .
11. What is the inverse matrix for the answer to the previous problem?