# 18

# Video Compression

## 18.1 Overview

Video compression can be viewed as image compression with a temporal component since video consists of a time sequence of images. From this point of view, the only "new" technique introduced in this chapter is a strategy to take advantage of this temporal correlation. However, there are different situations in which video compression becomes necessary, each requiring a solution specific to its peculiar conditions. In this chapter we briefly look at video compression algorithms and standards developed for different video communications applications.

## 18.2 Introduction

Of all the different sources of data, perhaps the one that produces the largest amount of data is video. Consider a video sequence generated using the CCIR 601 format (Section 18.4). Each image frame is made up of more than a quarter million pixels. At the rate of 30 frames per second and 16 bits per pixel, this corresponds to a data rate of about 21 Mbytes or 168 Mbits per second. This is certainly a change from the data rates of 2.4, 4.8, and 16 kbits per second that are the targets for speech coding systems discussed in Chapter 17.

Video compression can be viewed as the compression of a sequence of images; in other words, image compression with a temporal component. This is essentially the approach we will take in this chapter. However, there are limitations to this approach. We do not perceive motion video in the same manner as we perceive still images. Motion video may mask coding artifacts that would be visible in still images. On the other hand, artifacts that may not be visible in reconstructed still images can be very annoying in reconstructed motion video sequences. For example, consider a compression scheme that introduces a modest random amount of change in the average intensity of the pixels in the image. Unless a reconstructed

still image was being compared side by side with the original image, this artifact may go totally unnoticed. However, in a motion video sequence, especially one with low activity, random intensity changes can be quite annoying. As another example, poor reproduction of edges can be a serious problem in the compression of still images. However, if there is some temporal activity in the video sequence, errors in the reconstruction of edges may go unnoticed.

Although a more holistic approach might lead to better compression schemes, it is more convenient to view video as a sequence of correlated images. Most of the video compression algorithms make use of the temporal correlation to remove redundancy. The previous reconstructed frame is used to generate a prediction for the current frame. The difference between the prediction and the current frame, the prediction error or residual, is encoded and transmitted to the receiver. The previous reconstructed frame is also available at the receiver. Therefore, if the receiver knows the manner in which the prediction was performed, it can use this information to generate the prediction values and add them to the prediction error to generate the reconstruction. The prediction operation in video coding has to take into account motion of the objects in the frame, which is known as motion compensation (described in the next section).

We will also describe a number of different video compression algorithms. For the most part, we restrict ourselves to discussions of techniques that have found their way into international standards. Because there are a significant number of products that use proprietary video compression algorithms, it is difficult to find or include descriptions of them.

We can classify the algorithms based on the application area. While attempts have been made to develop standards that are "generic," the application requirements can play a large part in determining the features to be used and the values of parameters. When the compression algorithm is being designed for two-way communication, it is necessary for the coding delay to be minimal. Furthermore, compression and decompression should have about the same level of complexity. The complexity can be unbalanced in a broadcast application, where there is one transmitter and many receivers, and the communication is essentially one-way. In this case, the encoder can be much more complex than the receiver. There is also more tolerance for encoding delays. In applications where the video is to be decoded on workstations and personal computers, the decoding complexity has to be extremely low in order for the decoder to decode a sufficient number of images to give the illusion of motion. However, as the encoding is generally not done in real time, the encoder can be quite complex. When the video is to be transmitted over packet networks, the effects of packet loss have to be taken into account when designing the compression algorithm. Thus, each application will present its own unique requirements and demand a solution that fits those requirements.
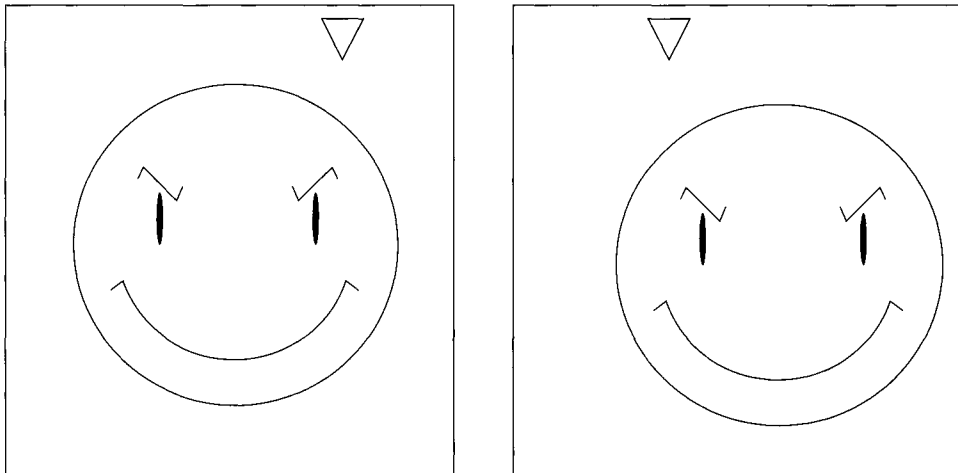
We will assume that you are familiar with the particular image compression technique being used. For example, when discussing transform-based video compression techniques, we assume that you have reviewed Chapter 13 and are familiar with the descriptions of transforms and the JPEG algorithm contained in that chapter.
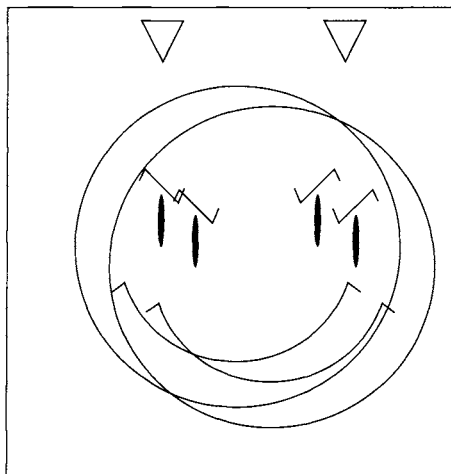
## 18.3  Motion Compensation

In most video sequences there is little change in the contents of the image from one frame to the next. Even in sequences that depict a great deal of activity, there are significant portions of the image that do not change from one frame to the next. Most video compression schemes take advantage of this redundancy by using the previous frame to generate a prediction for the current frame. We have used prediction previously when we studied differential encoding schemes. If we try to apply those techniques blindly to video compression by predicting the value of each pixel by the value of the pixel at the same location in the previous frame, we will run into trouble because we would not be taking into account the fact that objects tend to move between frames. Thus, the object in one frame that was providing the pixel at a certain location $(i_0, j_0)$ with its intensity value might be providing the same intensity value in the next frame to a pixel at location $(i_1, j_1)$. If we don't take this into account, we can actually increase the amount of information that needs to be transmitted.

### Example 18.3.1:

Consider the two frames of a motion video sequence shown in Figure 18.1. The only differences between the two frames are that the devious looking individual has moved slightly downward and to the right of the frame, while the triangular object has moved to the left. The differences between the two frames are so slight, you would think that if the first frame was available to both the transmitter and receiver, not much information would need to be transmitted to the receiver in order to reconstruct the second frame. However, if we simply



**FIGURE 18. 1    Two frames of a video sequence.**

**FIGURE 18. 2     Difference between the two frames.**

take the difference between the two frames, as shown in Figure 18.2, the displacement of the objects in the frame results in an image that contains more detail than the original image. In other words, instead of the differencing operation reducing the information, there is actually more information that needs to be transmitted.                                          ♦
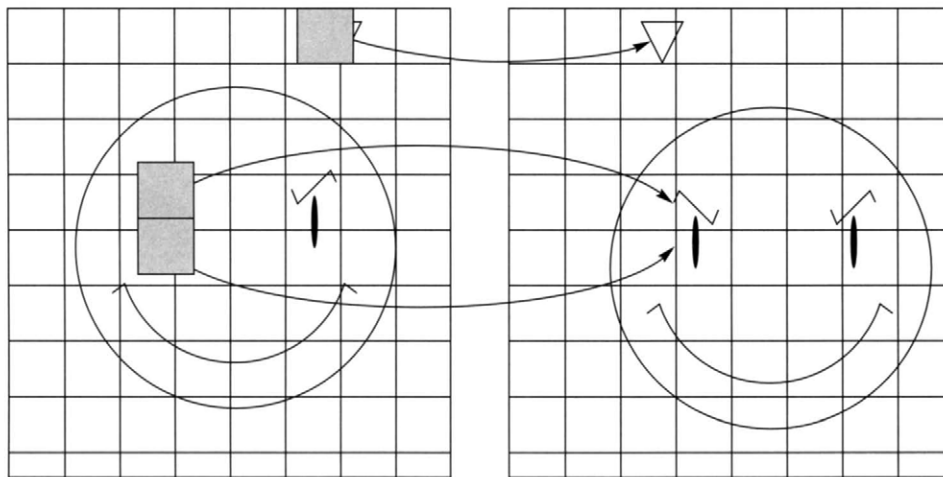
In order to use a previous frame to predict the pixel values in the frame being encoded, we have to take the motion of objects in the image into account. Although a number of approaches have been investigated, the method that has worked best in practice is a simple approach called *block-based motion compensation*. In this approach, the frame being encoded is divided into blocks of size $M \times M$. For each block, we search the previous reconstructed frame for the block of size $M \times M$ that most closely matches the block being encoded. We can measure the closeness of a match, or distance, between two blocks by the sum of absolute differences between corresponding pixels in the two blocks. We would obtain the same results if we used the sum of squared differences between the corresponding pixels as a measure of distance. Generally, if the distance from the block being encoded to the closest block in the previous reconstructed frame is greater than some prespecified threshold, the block is declared uncompensable and is encoded without the benefit of prediction. This decision is also transmitted to the receiver. If the distance is below the threshold, then a *motion vector* is transmitted to the receiver. The motion vector is the relative location of the block to be used for prediction obtained by subtracting the coordinates of the upper-left corner pixel of the block being encoded from the coordinates of the upper-left corner pixel of the block being used for prediction.

Suppose the block being encoded is an $8 \times 8$ block between pixel locations (24, 40) and (31, 47); that is, the upper-left corner pixel of the $8 \times 8$ block is at location (24, 40). If the block that best matches it in the previous frame is located between pixels at location (21, 43) and (28, 50), then the motion vector would be $(-3, 3)$. The motion vector was

obtained by subtracting the location of the upper-left corner of the block being encoded from the location of the upper-left corner of the best matching block. Note that the blocks are numbered starting from the top-left corner. Therefore, a positive $x$ component means that the best matching block in the previous frame is to the right of the location of the block being encoded. Similarly, a positive $y$ component means that the best matching block is at a location below that of the location of the block being encoded.

### Example 18.3.2:

Let us again try to predict the second frame of Example 18.3.1 using motion compensation. We divide the image into blocks and then predict the second frame from the first in the manner described above. Figure 18.3 shows the blocks in the previous frame that were used to predict some of the blocks in the current frame.



**FIGURE 18. 3** **Motion-compensated prediction.**

Notice that in this case all that needs to be transmitted to the receiver are the motion vectors. The current frame is completely predicted by the previous frame. ◆

We have been describing motion compensation where the displacement between the block being encoded and the best matching block is an integer number of pixels in the horizontal and vertical directions. There are algorithms in which the displacement is measured in half pixels. In order to do this, pixels of the coded frame being searched are interpolated to obtain twice as many pixels as in the original frame. This "doubled" image is then searched for the best matching block.

**TABLE 18.1**     *"Doubled" image.*

| A | $h_1$ | B |
|---|---|---|
| $v_1$ | $c$ | $v_2$ |
| C | $h_2$ | D |

The doubled image is obtained as follows: Consider Table 18.1. In this image $A$, $B$, $C$, and $D$ are the pixels of the original frame. The pixels $h_1, h_2, v_1,$ and $v_2$ are obtained by interpolating between the two neighboring pixels:

$$h_1 = \left\lfloor \frac{A+B}{2} + 0.5 \right\rfloor$$

$$h_2 = \left\lfloor \frac{C+D}{2} + 0.5 \right\rfloor$$

$$v_1 = \left\lfloor \frac{A+C}{2} + 0.5 \right\rfloor$$

$$v_2 = \left\lfloor \frac{B+D}{2} + 0.5 \right\rfloor \tag{18.1}$$

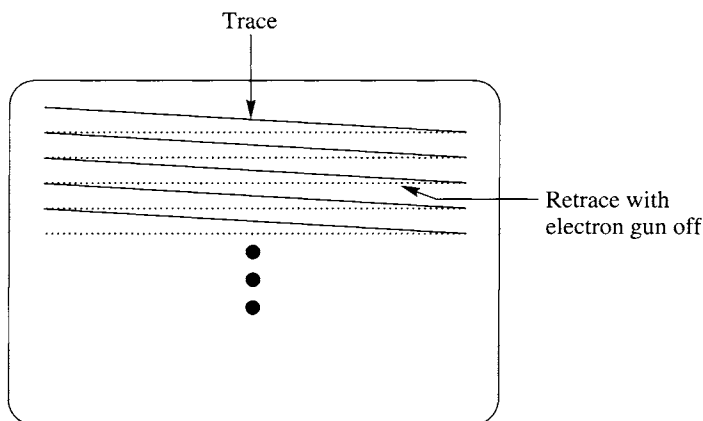while the pixel $c$ is obtained as the average of the four neighboring pixels from the coded original:

$$c = \left\lfloor \frac{A+B+C+D}{4} + 0.5 \right\rfloor .$$

We have described motion compensation in very general terms in this section. The various schemes in this chapter use specific motion compensation schemes that differ from each other. The differences generally involve the region of search for the matching block and the search procedure. We will look at the details with the study of the compression schemes. But before we begin our study of compression schemes, we briefly discuss how video signals are represented in the next section.

## 18.4  Video Signal Representation

The development of different representations of video signals has depended a great deal on past history. We will also take a historical view, starting with black-and-white television proceeding to digital video formats. The history of the development of analog video signal formats for the United States has been different than for Europe. Although we will show the development using the formats used in the United States, the basic ideas are the same for all formats.
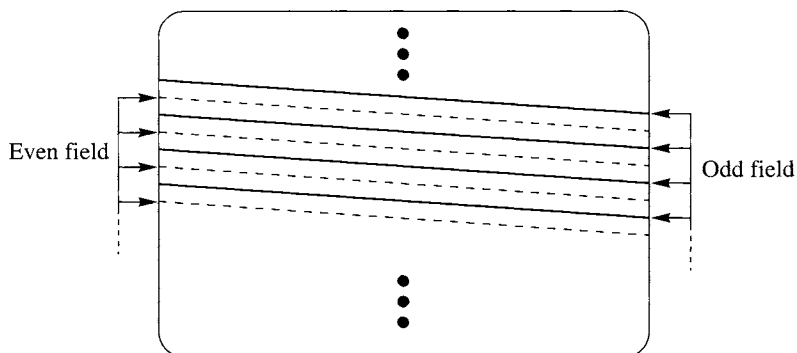
A black-and-white television picture is generated by exciting the phosphor on the television screen using an electron beam whose intensity is modulated to generate the image we see. The path that the modulated electron beam traces is shown in Figure 18.4. The line created by the horizontal traversal of the electron beam is called a line of the image. In order

Trace

Retrace with
electron gun off

**FIGURE 18. 4    The path traversed by the electron beam in a television.**

to trace a second line, the electron beam has to be deflected back to the left of the screen. During this period, the gun is turned off in order to prevent the retrace from becoming visible. The image generated by the traversal of the electron gun has to be updated rapidly enough for persistence of vision to make the image appear stable. However, higher rates of information transfer require higher bandwidths, which translate to higher costs.

In order to keep the cost of bandwidth low it was decided to send 525 lines 30 times a second. These 525 lines are said to constitute a *frame*. However, a thirtieth of a second between frames is long enough for the image to appear to flicker. To avoid the flicker, it was decided to divide the image into two interlaced fields. A field is sent once every sixtieth of a second. First, one field consisting of 262.5 lines is traced by the electron beam. Then, the second field consisting of the remaining 262.5 lines is traced *between* the lines of the first field. The situation is shown schematically in Figure 18.5. The first field is shown with

Even field

Odd field

**FIGURE 18. 5    A frame and its constituent fields.**

solid lines while the second field is shown with dashed lines. The first field begins on a full line and ends on a half line while the second field begins on a half line and ends on a full line. Not all 525 lines are displayed on the screen. Some are lost due to the time required for the electron gun to position the beam from the bottom to the top of the screen. We actually see about 486 lines per frame.

In a color television, instead of a single electron gun, we have three electron guns that act in unison. These guns excite red, green, and blue phosphor dots embedded in the screen. The beam from each gun strikes only one kind of phosphor, and the gun is named according to the color of the phosphor it excites. Thus, the red gun strikes only the red phosphor, the blue gun strikes only the blue phosphor, and the green gun strikes only the green phosphor. (Each gun is prevented from hitting a different type of phosphor by an aperture mask.)

In order to control the three guns we need three signals: a red signal, a blue signal, and a green signal. If we transmitted each of these separately, we would need three times the bandwidth. With the advent of color television, there was also the problem of backward compatibility. Most people had black-and-white television sets, and television stations did not want to broadcast using a format that most of the viewing audience could not see on their existing sets. Both issues were resolved with the creation of a composite color signal. In the United States, the specifications for the composite signal were created by the National Television Systems Committee, and the composite signal is often called an NTSC signal. The corresponding signals in Europe are PAL (Phase Alternating Lines), developed in Germany, and SECAM (Séquential Coleur avec Mémoire), developed in France. There is some (hopefully) good-natured rivalry between proponents of the different systems. Some problems with color reproduction in the NTSC signal have led to the name *Never Twice the Same Color*, while the idiosyncracies of the SECAM system have led to the name *Systéme Essentiallement Contre les Américains* (system essentially against the Americans).

The composite color signal consists of a *luminance* component, corresponding to the black-and-white television signal, and two *chrominance* components. The luminance component is denoted by $Y$:

$$Y = 0.299R + 0.587G + 0.114B \qquad (18.2)$$

where $R$ is the red component, $G$ is the green component, and $B$ is the blue component. The weighting of the three components was obtained through extensive testing with human observers. The two chrominance signals are obtained as

$$C_b = B - Y \qquad (18.3)$$

$$C_r = R - Y. \qquad (18.4)$$

These three signals can be used by the color television set to generate the red, blue, and green signals needed to control the electron guns. The luminance signal can be used directly by the black-and-white televisions.
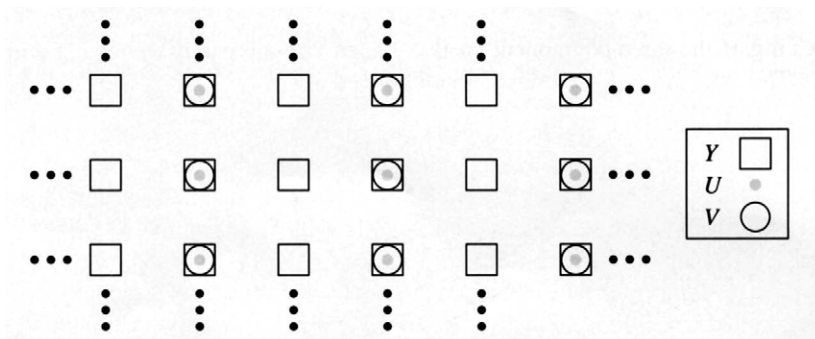
Because the eye is much less sensitive to changes of the chrominance in an image, the chrominance signal does not need to have higher frequency components. Thus, lower bandwidth of the chrominance signals along with a clever use of modulation techniques

permits all three signals to be encoded without need of any bandwidth expansion. (A simple and readable explanation of television systems can be found in [245].)

The early efforts toward digitization of the video signal were devoted to sampling the composite signal, and in the United States the Society of Motion Picture and Television Engineers developed a standard that required sampling the NTSC signal at a little more than 14 million times a second. In Europe, the efforts at standardization of video were centered around the characteristics of the PAL signal. Because of the differences between NTSC and PAL, this would have resulted in different "standards." In the late 1970s, this approach was dropped in favor of sampling the components and the development of a worldwide standard. This standard was developed under the auspices of the International Consultative Committee on Radio (CCIR) and was called CCIR recommendation 601-2. CCIR is now known as ITU-R, and the recommendation is officially known as ITU-R recommendation BT.601-2. However, the standard is generally referred to as recommendation 601 or CCIR 601.

The standard proposes a family of sampling rates based on the sampling frequency of 3.725 MHz (3.725 million samples per second). Multiples of this sampling frequency permit samples on each line to line up vertically, thus generating the rectangular array of pixels necessary for digital processing. Each component can be sampled at an integer multiple of 3.725 MHz, up to a maximum of four times this frequency. The sampling rate is represented as a triple of integers, with the first integer corresponding to the sampling of the luminance component and the remaining two corresponding to the chrominance components. Thus, 4:4:4 sampling means that all components were sampled at 13.5 MHz. The most popular sampling format is the 4:2:2 format, in which the luminance signal is sampled at 13.5 MHz, while the lower-bandwidth chrominance signals are sampled at 6.75 MHz. If we ignore the samples of the portion of the signal that do not correspond to active video, the sampling rate translates to 720 samples per line for the luminance signal and 360 samples per line for the chrominance signal. The sampling format is shown in Figure 18.6. The luminance component of the digital video signal is also denoted by $Y$, while the chrominance components are denoted by $U$ and $V$. The sampled analog values are converted to digital values as follows. The sampled values of $YC_bC_r$ are normalized so that the sampled $Y$ values, $Y_s$, take on values between 0 and 1, and the sampled chrominance values, $C_{rs}$ and $C_{bs}$, take on values



**FIGURE 18. 6      Recommendation 601 4:2:2 sampling format.**

between $\frac{-1}{2}$ and $\frac{1}{2}$. These normalized values are converted to 8-bit numbers according to the transformations

$$Y = 219Y_s + 16 \qquad (18.5)$$

$$U = 224C_{bs} + 128 \qquad (18.6)$$

$$V = 224C_{rs} + 128. \qquad (18.7)$$

Thus, the $Y$ component takes on values between 16 and 235, and the $U$ and $V$ components take on values between 16 and 240.
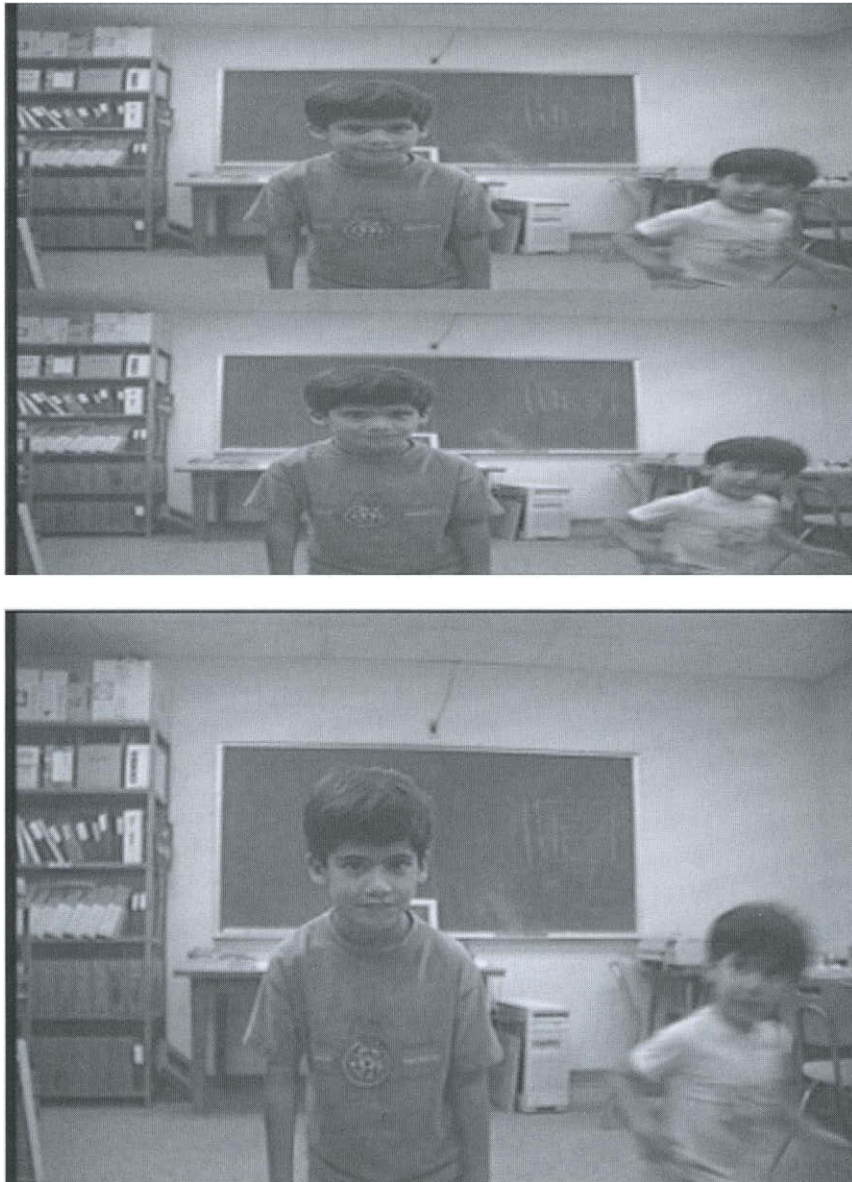
An example of the $Y$ component of a CCIR 601 frame is shown in Figure 18.7. In the top image we show the fields separately, while in the bottom image the fields have been interlaced. Notice that in the interlaced image the smaller figure looks blurred. This is because the individual moved in the sixtieth of a second between the two fields. (This is also proof—if any was needed—that a three-year-old cannot remain still, even for a sixtieth of a second!)

The $YUV$ data can also be arranged in other formats. In the Common Interchange Format (CIF), which is used for videoconferencing, the luminance of the image is represented by an array of $288 \times 352$ pixels, and the two chrominance signals are represented by two arrays consisting of $144 \times 176$ pixels. In the QCIF (Quarter CIF) format, we have half the number of pixels in both the rows and columns.
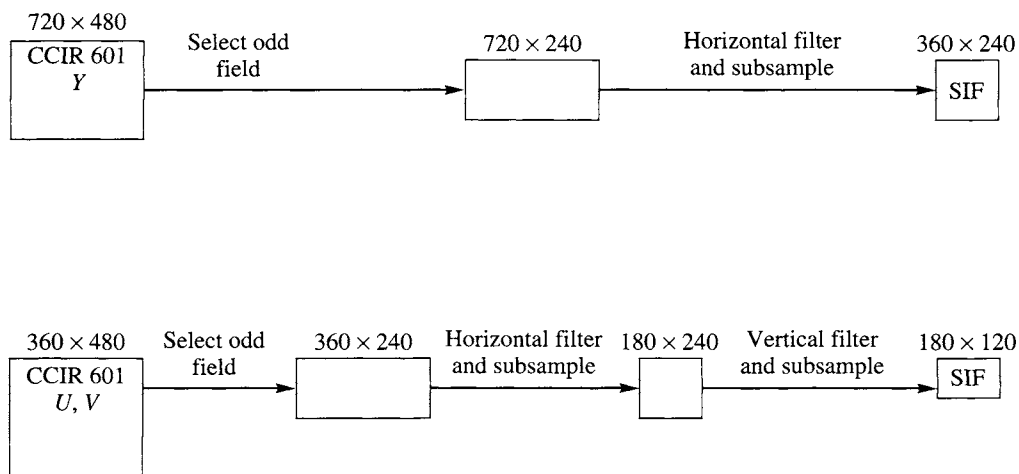
The MPEG-1 algorithm, which was developed for encoding video at rates up to 1.5 Mbits per second, uses a different subsampling of the CCIR 601 format to obtain the MPEG-SIF format. Starting from a 4:2:2, 480-line CCIR 601 format, the vertical resolution is first reduced by taking only the odd field for both the luminance and the chrominance components. The horizontal resolution is then reduced by filtering (to prevent aliasing) and then subsampling by a factor of two in the horizontal direction. This results in $360 \times 240$ samples of $Y$ and $180 \times 240$ samples each of $U$ and $V$. The vertical resolution of the chrominance samples is further reduced by filtering and subsampling in the vertical direction by a factor of two to obtain $180 \times 120$ samples for each of the chrominance signals. The process is shown in Figure 18.8, and the resulting format is shown in Figure 18.9.

In the following we describe several of the video coding standards in existence today. Our order of description follows the historical development of the standards. As each standard has built upon features of previous standards this seems like a logical plan of attack. As in the case of image compression, most of the standards for video compression are based on the discrete cosine transform (DCT). The standard for teleconferencing applications, ITU-T recommendation H.261, is no exception. Most systems currently in use for videoconferencing use proprietary compression algorithms. However, in order for the equipment from different manufacturers to communicate with each other, these systems also offer the option of using H.261. We will describe the compression algorithm used in the H.261 standard in the next section. We will follow that with a description of the MPEG algorithms used in Video CDs, DVDs and HDTV, and a discussion of the latest joint offering from ITU and MPEG.
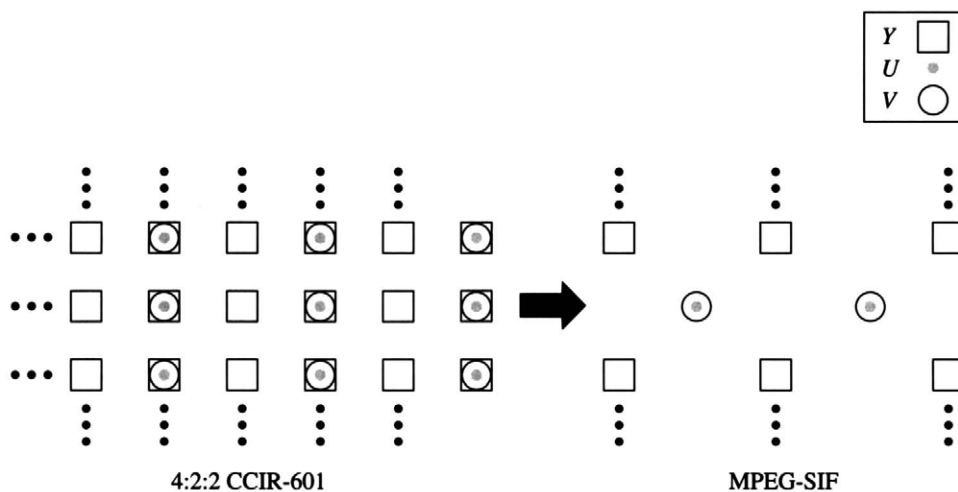
We will also describe a new approach towards compression of video for videophone applications called three-dimensional model-based coding. This approach is far from maturity, and our description will be rather cursory. The reason for including it here is the great promise it holds for the future.

**FIGURE 18. 7**     **Top: Fields of a CCIR 601 frame. Bottom: An interlaced CCIR 601 frame.**

**FIGURE 18. 8     Generation of an SIF frame.**



**FIGURE 18. 9     CCIR 601 to MPEG-SIF.**

## 18.5  ITU-T Recommendation H.261

The earliest DCT-based video coding standard is the ITU-T H.261 standard. This algorithm assumes one of two formats, CIF and QCIF. A block diagram of the H.261 video coder is shown in Figure 18.10. The basic idea is simple. An input image is divided into blocks of $8 \times 8$ pixels. For a given $8 \times 8$ block, we subtract the prediction generated using the previous frame. (If there is no previous frame or the previous frame is very different from the current frame, the prediction might be zero.) The difference between the block being encoded and

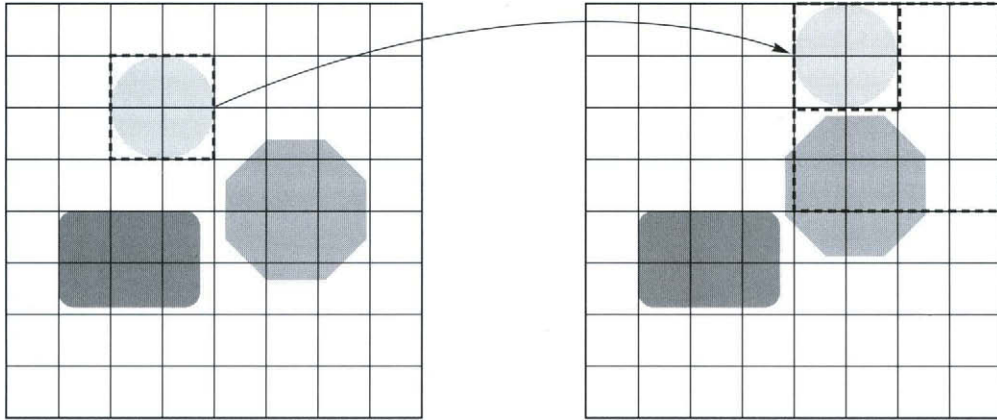**FIGURE 18. 10    Block diagram of the ITU-T H.261 encoder.**

the prediction is transformed using a DCT. The transform coefficients are quantized and the quantization label encoded using a variable-length code. In the following discussion, we will take a more detailed look at the various components of the compression algorithm.

## 18.5.1   Motion Compensation

Motion compensation requires a large amount of computation. Consider finding a matching block for an $8 \times 8$ block. Each comparison requires taking 64 differences and then computing the sum of the absolute value of the differences. If we assume that the closest block in the previous frame is located within 20 pixels in either the horizontal or vertical direction of the block to be encoded, we need to perform 1681 comparisons. There are several ways we can reduce the total number of computations.

One way is to increase the size of the block. Increasing the size of the block means more computations per comparison. However, it also means that we will have fewer blocks per frame, so the number of times we have to perform the motion compensation will decrease. However, different objects in a frame may be moving in different directions. The drawback to increasing the size of the block is that the probability that a block will contain objects moving in different directions increases with size. Consider the two images in Figure 18.11. If we use blocks that are made up of $2 \times 2$ squares, we can find a block that exactly matches the $2 \times 2$ block that contains the circle. However, if we increase the size of the block to $4 \times 4$ squares, the block that contains the circle also contains the upper part of the octagon. We cannot find a similar $4 \times 4$ block in the previous frame. Thus, there is a trade-off involved.

**FIGURE 18. 11    Effect of block size on motion compensation.**

Larger blocks reduce the amount of computation; however, they can also result in poor prediction, which in turn can lead to poor compression performance.

Another way we can reduce the number of computations is by reducing the search space. If we reduce the size of the region in which we search for a match, the number of computations will be reduced. However, reducing the search region also increases the probability of missing a match. Again, we have a trade-off between computation and the amount of compression.

The H.261 standard has balanced the trade-offs in the following manner. The $8 \times 8$ blocks of luminance and chrominance pixels are organized into *macroblocks*, which consist of four luminance blocks, and one each of the two types of chrominance blocks. The motion-compensated prediction (or motion compensation) operation is performed on the macroblock level. For each macroblock, we search the previous reconstructed frame for the macroblock that most closely matches the macroblock being encoded. In order to further reduce the amount of computations, only the luminance blocks are considered in this matching operation. The motion vector for the prediction of the chrominance blocks is obtained by halving the component values of the motion vector for the luminance macroblock. Therefore, if the motion vector for the luminance blocks was $(-3, 10)$, then the motion vector for the chrominance blocks would be $(-1, 5)$.

The search area is restricted to $\pm 15$ pixels of the macroblock being encoded in the horizontal and vertical directions. That is, if the upper-left corner pixel of the block being encoded is $(x_c, y_c)$, and the upper-left corner of the best matching macroblock is $(x_p, y_p)$, then $(x_c, y_c)$ and $(x_p, y_p)$ have to satisfy the constraints $|x_c - x_p| < 15$ and $|y_c - y_p| < 15$.

## 18.5.2 The Loop Filter

Sometimes sharp edges in the block used for prediction can result in the generation of sharp changes in the prediction error. This in turn can cause high values for the high-frequency coefficients in the transforms, which can increase the transmission rate. To avoid this, prior

to taking the difference, the prediction block can be smoothed by using a two-dimensional spatial filter. The filter is separable; it can be implemented as a one-dimensional filter that first operates on the rows, then on the columns. The filter coefficients are $\frac{1}{4}, \frac{1}{2}, \frac{1}{4}$, except at block boundaries where one of the filter taps would fall outside the block. To prevent this from happening, the block boundaries remain unchanged by the filtering operation.

### Example 18.5.1:

Let's filter the $4 \times 4$ block of pixel values shown in Table 18.2 using the filter specified for the H.261 algorithm. From the pixel values we can see that this is a gray square with a white $L$ in it. (Recall that small pixel values correspond to darker pixels and large pixel values correspond to lighter pixels, with 0 corresponding to black and 255 corresponding to white.)

**TABLE 18.2    Original block of pixels.**

| | | | |
|---|---|---|---|
| 110 | 218 | 116 | 112 |
| 108 | 210 | 110 | 114 |
| 110 | 218 | 210 | 112 |
| 112 | 108 | 110 | 116 |

Let's filter the first row. We leave the first pixel value the same. The second value becomes

$$\frac{1}{4} \times 110 + \frac{1}{2} \times 218 + \frac{1}{4} \times 116 = 165$$

where we have assumed integer division. The third filtered value becomes

$$\frac{1}{4} \times 218 + \frac{1}{2} \times 116 + \frac{1}{4} \times 112 = 140.$$

The final element in the first row of the filtered block remains unchanged. Continuing in this fashion with all four rows, we get the $4 \times 4$ block shown in Table 18.3.

**TABLE 18.3    After filtering the rows.**

| | | | |
|---|---|---|---|
| 110 | 165 | 140 | 112 |
| 108 | 159 | 135 | 114 |
| 110 | 188 | 187 | 112 |
| 112 | 109 | 111 | 116 |

Now repeat the filtering operation along the columns. The final $4 \times 4$ block is shown in Table 18.4. Notice how much more homogeneous this last block is compared to the original

block. This means that it will most likely not introduce any sharp variations in the difference block, and the high-frequency coefficients in the transform will be closer to zero, leading to compression.

**TABLE 18.4     Final block.**

| | | | |
|---|---|---|---|
| 110 | 165 | 140 | 112 |
| 108 | 167 | 148 | 113 |
| 110 | 161 | 154 | 113 |
| 112 | 109 | 111 | 116 |

◆

This filter is either switched on or off for each macroblock. The conditions for turning the filter on or off are not specified by the recommendations.
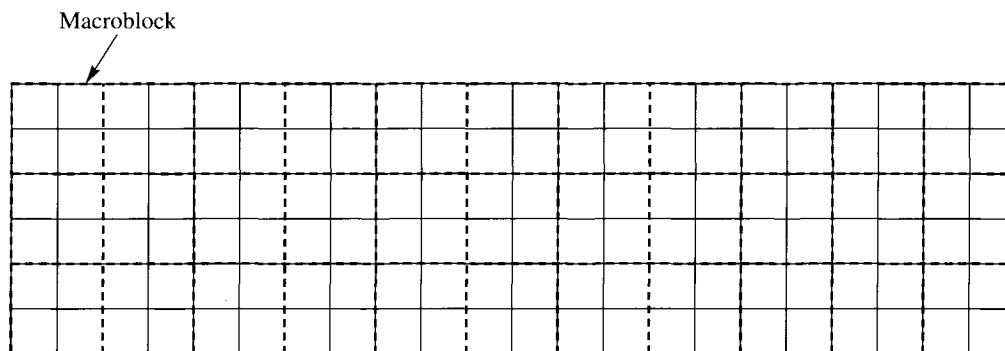
## 18.5.3   The Transform

The transform operation is performed with a DCT on an $8 \times 8$ block of pixels or pixel differences. If the motion compensation operation does not provide a close match, then the transform operation is performed on an $8 \times 8$ block of pixels. If the transform operation is performed on a block level, either a block or the difference between the block and its predicted value is quantized and transmitted to the receiver. The receiver performs the inverse operations to reconstruct the image. The receiver operation is also simulated at the transmitter, where the reconstructed images are obtained and stored in a frame store. The encoder is said to be in *intra* mode if it operates directly on the input image without the use of motion compensation. Otherwise, it is said to be in *inter* mode.

## 18.5.4   Quantization and Coding

Depending on how good or poor the prediction is, we can get a wide variation in the characteristics of the coefficients that are to be quantized. In the case of an intra block, the DC coefficients will take on much larger values than the other coefficients. Where there is little motion from frame to frame, the difference between the block being encoded and the prediction will be small, leading to small values for the coefficients.

In order to deal with this wide variation, we need a quantization strategy that can be rapidly adapted to the current situation. The H.261 algorithm does this by switching between 32 different quantizers, possibly from one macroblock to the next. One quantizer is reserved for the intra DC coefficient, while the remaining 31 quantizers are used for the other coefficients. The intra DC quantizer is a uniform midrise quantizer with a step size of 8. The other quantizers are midtread quantizers with a step size of an even value between 2 and 62. Given a particular block of coefficients, if we use a quantizer with smaller step size, we are likely to get a larger number of nonzero coefficients. Because of the manner in which the labels are encoded, the number of bits that will need to be transmitted will increase. Therefore, the availability of transmission resources will have a major impact on the quantizer selection. We will discuss this aspect further when we talk about the transmission

Macroblock



**FIGURE 18. 12     A GOB consisting of 33 macroblocks.**

buffer. Once a quantizer is selected, the receiver has to be informed about the selection. In H.261, this is done in one of two ways. Each macroblock is preceded by a header. The quantizer being used can be identified as part of this header. When the amount of activity or motion in the sequence is relatively constant, it is reasonable to expect that the same quantizer will be used for a large number of macroblocks. In this case, it would be wasteful to identify the quantizer being used with each macroblock. The macroblocks are organized into *groups of blocks* (GOBs), each of which consist of three rows of 11 macroblocks. This hierarchical arrangement is shown in Figure 18.12. Only the luminance blocks are shown. The header preceding each GOB contains a 5-bit field for identifying the quantizer. Once a quantizer has been identified in the GOB header, the receiver assumes that quantizer is being used, unless this choice is overridden using the macroblock header.

The quantization labels are encoded in a manner similar to, but not exactly the same as, JPEG. The labels are scanned in a zigzag fashion like JPEG. The nonzero labels are coded along with the number, or run, of coefficients quantized to zero. The 20 most commonly occurring combinations of (run, label) are coded with a single variable-length codeword. All other combinations of (run, label) are coded with a 20-bit word, made up of a 6-bit escape sequence, a 6-bit code denoting the run, and an 8-bit code for the label.

In order to avoid transmitting blocks that have no nonzero quantized coefficient, the header preceding each macroblock can contain a variable-length code called the *coded block pattern* (CBP) that indicates which of the six blocks contain nonzero labels. The CBP can take on one of 64 different pattern numbers, which are then encoded by a variable-length code. The pattern number is given by

$$CBP = 32P_1 + 16P_2 + 8P_3 + 4P_4 + 2P_5 + P_6$$

where $P_1$ through $P_6$ correspond to the six different blocks in the macroblock, and is one if the corresponding block has a nonzero quantized coefficient and zero otherwise.

## 18.5.5   Rate Control

The binary codewords generated by the transform coder form the input to a transmission buffer. The function of the transmission buffer is to keep the output rate of the encoder fixed. If the buffer starts filling up faster than the transmission rate, it sends a message back to the transform coder to reduce the output from the quantization. If the buffer is in danger of becoming emptied because the transform coder is providing bits at a rate lower than the transmission rate, the transmission buffer can request a higher rate from the transform coder. This operation is called *rate control.*

The change in rate can be affected in two different ways. First, the quantizer being used will affect the rate. If a quantizer with a large step size is used, a larger number of coefficients will be quantized to zero. Also, there is a higher probability that those not quantized to zero will be one of the those values that have a shorter variable-length codeword. Therefore, if a higher rate is required, the transform coder selects a quantizer with a smaller step size, and if a lower rate is required, the transform coder selects a quantizer with a larger step size. The quantizer step size is set at the beginning of each GOB, but can be changed at the beginning of any macroblock. If the rate cannot be lowered enough and there is a danger of buffer overflow, the more drastic option of dropping frames from transmission is used.

The ITU-T H.261 algorithm was primarily designed for videophone and videoconferencing applications. Therefore, the algorithm had to operate with minimal coding delay (less than 150 milliseconds). Furthermore, for videophone applications, the algorithm had to operate at very low bit rates. In fact, the title for the recommendation is "Video Codec for Audiovisual Services at $p \times 64$ kbit/s," where $p$ takes on values from 1 to 30. A $p$ value of 2 corresponds to a total transmission rate of 128 kbps, which is the same as two voice-band telephone channels. These are very low rates for video, and the ITU-T H.261 recommendations perform relatively well at these rates.
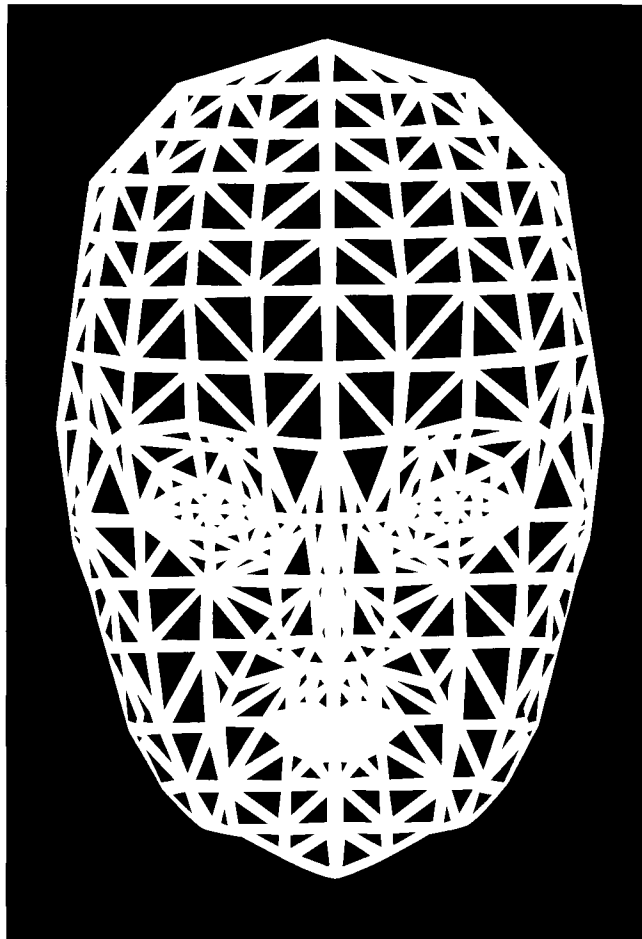
## 18.6   Model-Based Coding

In speech coding, a major decrease in rate is realized when we go from coding waveforms to an analysis/synthesis approach. An attempt at doing the same for video coding is described in the next section. A technique that has not yet reached maturity but shows great promise for use in videophone applications is an analysis/synthesis technique. The analysis/synthesis approach requires that the transmitter and receiver agree on a model for the information to be transmitted. The transmitter then analyzes the information to be transmitted and extracts the model parameters, which are transmitted to the receiver. The receiver uses these parameters to synthesize the source information. While this approach has been successfully used for speech compression for a long time (see Chapter 15), the same has not been true for images. In a delightful book, *Signals, Systems, and Noise—The Nature and Process of Communications,* published in 1961, J.R. Pierce [14] described his "dream" of an analysis/synthesis scheme for what we would now call a videoconferencing system:

> Imagine that we had at the receiver a sort of rubbery model of the human face. Or we might have a description of such a model stored in the memory of a huge electronic computer. . . . Then, as the person before the transmitter talked, the

transmitter would have to follow the movements of his eyes, lips, and jaws, and other muscular movements and transmit these so that the model at the receiver could do likewise.

Pierce's dream is a reasonably accurate description of a three-dimensional model-based approach to the compression of facial image sequences. In this approach, a generic wireframe model, such as the one shown in Figure 18.13, is constructed using triangles. When encoding the movements of a specific human face, the model is adjusted to the face by matching features and the outer contour of the face. The image textures are then mapped onto this wireframe model to synthesize the face. Once this model is available to both transmitter and receiver, only changes in the face are transmitted to the receiver. These changes can be



**FIGURE 18. 13     Generic wireframe model.**

classified as *global motion* or *local motion* [246]. Global motion involves movement of the head, while local motion involves changes in the features—in other words, changes in facial expressions. The global motion can be modeled in terms of movements of rigid bodies. The facial expressions can be represented in terms of relative movements of the vertices of the triangles in the wireframe model. In practice, separating a movement into global and local components can be difficult because most points on the face will be affected by both the changing position of the head and the movement due to changes in facial expression. Different approaches have been proposed to separate these effects [247, 246, 248].

The global movements can be described in terms of rotations and translations. The local motions, or facial expressions, can be described as a sum of *action units* (AU), which are a set of 44 descriptions of basic facial expressions [249]. For example, AU1 corresponds to the raising of the inner brow and AU2 corresponds to the raising of the outer brow; therefore, AU1 + AU2 would mean raising the brow.

Although the synthesis portion of this algorithm is relatively straightforward, the analysis portion is far from simple. Detecting changes in features, which tend to be rather subtle, is a very difficult task. There is a substantial amount of research in this area, and if this problem is resolved, this approach promises rates comparable to the rates of the analysis/synthesis voice coding schemes. A good starting point for exploring this fascinating area is [250].

## 18.7  Asymmetric Applications

There are a number of applications in which it is cost effective to shift more of the computational burden to the encoder. For example, in multimedia applications where a video sequence is stored on a CD-ROM, the decompression will be performed many times and has to be performed in real time. However, the compression is performed only once, and there is no need for it to be in real time. Thus, the encoding algorithms can be significantly more complex. A similar situation arises in broadcast applications, where for each transmitter there might be thousands of receivers. In this section we will look at the standards developed for such asymmetric applications. These standards have been developed by a joint committee of the International Standards Organization (ISO) and the International Electrotechnical Society (IEC), which is best known as MPEG (Moving Picture Experts Group). MPEG was initially set up in 1988 to develop a set of standard algorithms, at different rates, for applications that required storage of video and audio on digital storage media. Originally, the committee had three work items, nicknamed MPEG-1, MPEG-2, and MPEG-3, targeted at rates of 1.5, 10, and 40 Mbits per second, respectively. Later, it became clear that the algorithms developed for MPEG-2 would accommodate the MPEG-3 rates, and the third work item was dropped [251]. The MPEG-1 work item resulted in a set of standards, ISO/IEC IS 11172, "Information Technology—Coding of Moving Pictures and Associated Audio for Digital Storage Media Up to about 1.5 Mbit/s" [252]. During the development of the standard, the committee felt that the restriction to digital storage media was not necessary, and the set of standards developed under the second work item, ISO/IEC 13818 or MPEG-2, has been issued under the title "Information Technology—Generic Coding of Moving Pictures and Associated Audio Information" [253]. In July of 1993 the MPEG committee began working

on MPEG-4, the third and most ambitious of its standards. The goal of MPEG-4 was to provide an object-oriented framework for the encoding of multimedia. It took two years for the committee to arrive at a satisfactory definition of the scope of MPEG-4, and the call for proposals was finally issued in 1996. The standard ISO/IEC 14496 was finalized in 1998 and approved as an international standard in 1999. We have examined the audio standard in Chapter 16. In this section we briefly look at the video standards.

## 18.8 The MPEG-1 Video Standard

The basic structure of the compression algorithm proposed by MPEG is very similar to that of ITU-T H.261. Blocks ($8 \times 8$ in size) of either an original frame or the difference between a frame and the motion-compensated prediction are transformed using the DCT. The blocks are organized in macroblocks, which are defined in the same manner as in the H.261 algorithm, and the motion compensation is performed at the macroblock level. The transform coefficients are quantized and transmitted to the receiver. A buffer is used to smooth delivery of bits from the encoder and also for rate control.
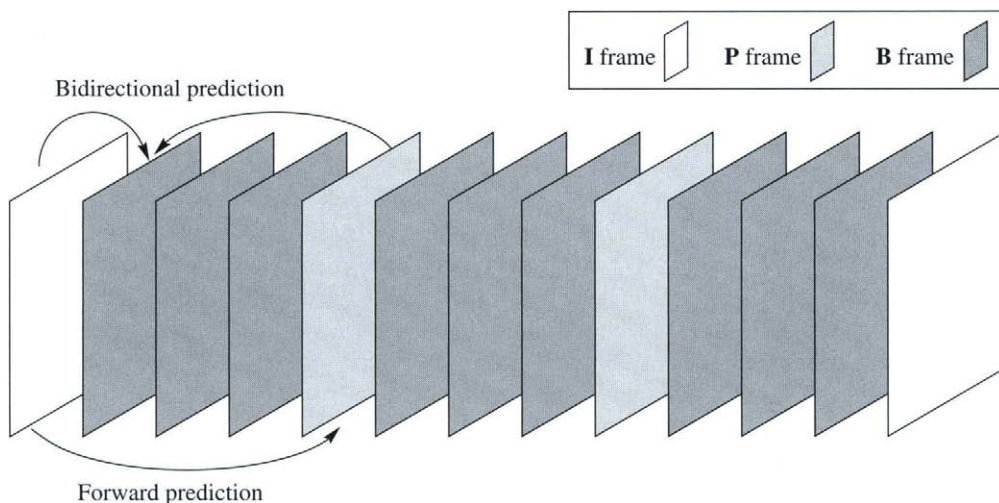
The basic structure of the MPEG-1 compression scheme may be viewed as very similar to that of the ITU-T H.261 video compression scheme; however, there are significant differences in the details of this structure. The H.261 standard has videophone and videoconferencing as the main application areas; the MPEG standard at least initially had applications that require digital storage and retrieval as a major focus. This does not mean that use of either algorithm is precluded in applications outside its focus, but simply that the features of the algorithm may be better understood if we keep in mind the target application areas. In videoconferencing a call is set up, conducted, and then terminated. This set of events always occurs together and in sequence. When accessing video from a storage medium, we do not always want to access the video sequence starting from the first frame. We want the ability to view the video sequence starting at, or close to, some arbitrary point in the sequence. A similar situation exists in broadcast situations. Viewers do not necessarily tune into a program at the beginning. They may do so at any random point in time. In H.261 each frame, after the first frame, may contain blocks that are coded using prediction from the previous frame. Therefore, to decode a particular frame in the sequence, it is possible that we may have to decode the sequence starting at the first frame. One of the major contributions of MPEG-1 was the provision of a random access capability. This capability is provided rather simply by requiring that there be frames periodically that are coded without any reference to past frames. These frames are referred to as **I** frames.

In order to avoid a long delay between the time a viewer switches on the TV to the time a reasonable picture appears on the screen, or between the frame that a user is looking for and the frame at which decoding starts, the **I** frames should occur quite frequently. However, because the **I** frames do not use temporal correlation, the compression rate is quite low compared to the frames that make use of the temporal correlations for prediction. Thus, the number of frames between two consecutive **I** frames is a trade-off between compression efficiency and convenience.

In order to improve compression efficiency, the MPEG-1 algorithm contains two other kinds of frames, the *predictive coded* (**P**) frames and the *bidirectionally predictive coded* (**B**) frames. The **P** frames are coded using motion-compensated prediction from the last **I** or **P** frame, whichever happens to be closest. Generally, the compression efficiency of **P** frames is substantially higher than **I** frames. The **I** and **P** frames are sometimes called *anchor* frames, for reasons that will become obvious.

To compensate for the reduction in the amount of compression due to the frequent use of **I** frames, the MPEG standard introduced **B** frames. The **B** frames achieve a high level of compression by using motion-compensated prediction from the most recent anchor frame and the closest future anchor frame. By using both past and future frames for prediction, generally we can get better compression than if we only used prediction based on the past. For example, consider a video sequence in which there is a sudden change between one frame and the next. This is a common occurrence in TV advertisements. In this situation, prediction based on the past frames may be useless. However, predictions based on future frames would have a high probability of being accurate. Note that a **B** frame can only be generated after the future anchor frame has been generated. Furthermore, the **B** frame is not used for predicting any other frame. This means that **B** frames can tolerate more error because this error will not be propagated by the prediction process.

The different frames are organized together in a *group of pictures* (GOP). A GOP is the smallest random access unit in the video sequence. The GOP structure is set up as a trade-off between the high compression efficiency of motion-compensated coding and the fast picture acquisition capability of periodic intra-only processing. As might be expected, a GOP has to contain at least one **I** frame. Furthermore, the first **I** frame in a GOP is either the first frame of the GOP, or is preceded by **B** frames that use motion-compensated prediction only from this **I** frame. A possible GOP is shown in Figure 18.14.



**FIGURE 18. 14     A possible arrangement for a group of pictures.**

**TABLE 18.5     A typical sequence of frames in display order.**

| I | B | B | P | B | B | P | B | B | P | B | B | I |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |

Because of the reliance of the **B** frame on future anchor frames, there are two different sequence orders. The *display order* is the sequence in which the video sequence is displayed to the user. A typical display order is shown in Table 18.5. Let us see how this sequence was generated. The first frame is an **I** frame, which is compressed without reference to any previous frame. The next frame to be compressed is the fourth frame. This frame is compressed using motion-compensated prediction from the first frame. Then we compress frame two, which is compressed using motion-compensated prediction from frame one and frame four. The third frame is also compressed using motion-compensated prediction from the first and fourth frames. The next frame to be compressed is frame seven, which uses motion-compensated prediction from frame four. This is followed by frames five and six, which are compressed using motion-compensated predictions from frames four and seven. Thus, there is a processing order that is quite different from the display order. The MPEG document calls this the *bitstream order*. The bitstream order for the sequence shown in Table 18.5 is given in Table 18.6. In terms of the bitstream order, the first frame in a GOP is always the **I** frame.

As we can see, unlike the ITU-T H.261 algorithm, the frame being predicted and the frame upon which the prediction is based are not necessarily adjacent. In fact, the number of frames between the frame being encoded and the frame upon which the prediction is based is variable. When searching for the best matching block in a neighboring frame, the region of search depends on assumptions about the amount of motion. More motion will lead to larger search areas than a small amount of motion. When the frame being predicted is always adjacent to the frame upon which the prediction is based, we can fix the search area based on our assumption about the amount of motion. When the number of frames between the frame being encoded and the prediction frame is variable, we make the search area a function of the distance between the two frames. While the MPEG standard does not specify the method used for motion compensation, it does recommend using a search area that grows with the distance between the frame being coded and the frame being used for prediction.

Once motion compensation has been performed, the block of prediction errors is transformed using the DCT and quantized, and the quantization labels are encoded. This procedure is the same as that recommended in the JPEG standard and is described in Chapter 12. The quantization tables used for the different frames are different and can be changed during the encoding process.

**TABLE 18.6     A typical sequence of frames in bitstream order.**

| I | P | B | B | P | B | B | P | B | B | I | B | B |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 4 | 2 | 3 | 7 | 5 | 6 | 10 | 8 | 9 | 13 | 11 | 12 |

Rate control in the MPEG standard can be performed at the sequence level or at the level of individual frames. At the sequence level, any reduction in bit rate first occurs with the **B** frames because they are not essential for the encoding of other frames. At the level of the individual frames, rate control takes place in two steps. First, as in the case of the H.261 algorithm, the quantizer step sizes are increased. If this is not sufficient, then the higher-order frequency coefficients are dropped until the need for rate reduction is past.

The format for MPEG is very flexible. However, the MPEG committee has provided some suggested values for the various parameters. For MPEG-1 these suggested values are called the *constrained parameter bitstream* (CPB). The horizontal picture size is constrained to be less than or equal to 768 pixels, and the vertical size is constrained to be less than or equal to 576 pixels. More importantly, the pixel rate is constrained to be less than 396 macroblocks per frame if the frame rate is 25 frames per second or less, and 330 macroblocks per frame if the frame rate is 30 frames per second or less. The definition of a macroblock is the same as in the ITU-T H.261 recommendations. Therefore, this corresponds to a frame size of $352 \times 288$ pixels at the 25-frames-per-second rate, or a frame size of $352 \times 240$ pixels at the 30-frames-per-second rate. Keeping the frame at this size allows the algorithm to achieve bit rates of between 1 and 1.5 Mbits per second. When referring to MPEG-1 parameters, most people are actually referring to the CPB.

The MPEG-1 algorithm provides reconstructed images of VHS quality for moderate-to low-motion video sequences, and worse than VHS quality for high-motion sequences at rates of around 1.2 Mbits per second. As the algorithm was targeted to applications such as CD-ROM, there is no consideration of interlaced video. In order to expand the applicability of the basic MPEG algorithm to interlaced video, the MPEG committee provided some additional recommendations, the MPEG-2 recommendations.

## 18.9   The MPEG-2 Video Standard—H.262

While MPEG-1 was specifically proposed for digital storage media, the idea behind MPEG-2 was to provide a generic, application-independent standard. To this end, MPEG-2 takes a "tool kit" approach, providing a number of subsets, each containing different options from the set of all possible options contained in the standard. For a particular application, the user can select from a set of *profiles* and *levels*. The profiles define the algorithms to be used, while the levels define the constraints on the parameters. There are five profiles: *simple, main, snr-scalable* (where *snr* stands for signal-to-noise ratio), *spatially scalable*, and *high*. There is an ordering of the profiles; each higher profile is capable of decoding video encoded using all profiles up to and including that profile. For example, a decoder designed for profile *snr-scalable* could decode video that was encoded using profiles *simple*, *main*, and *snr-scalable*. The *simple* profile eschews the use of **B** frames. Recall that the **B** frames require the most computation to generate (forward and backward prediction), require memory to store the coded frames needed for prediction, and increase the coding delay because of the need to wait for "future" frames for both generation and reconstruction. Therefore, removal of the **B** frames makes the requirements simpler. The *main* profile is very much the algorithm we have discussed in the previous section. The *snr-scalable, spatially scalable*, and *high* profiles may use more than one bitstream to encode the video. The base

bitstream is a lower-rate encoding of the video sequence. This bitstream could be decoded by itself to provide a reconstruction of the video sequence. The other bitstream is used to enhance the quality of the reconstruction. This layered approach is useful when transmitting video over a network, where some connections may only permit a lower rate. The base bitstream can be provided to these connections while providing the base and enhancement layers for a higher-quality reproduction over the links that can accommodate the higher bit rate. To understand the concept of layers, consider the following example.

## Example 18.9.1:

Suppose after the transform we obtain a set of coefficients, the first eight of which are

$$29.75 \quad 6.1 \quad -6.03 \quad 1.93 \quad -2.01 \quad 1.23 \quad -0.95 \quad 2.11$$

Let us suppose we quantize this set of coefficients using a step size of 4. For simplicity we will use the same step size for all coefficients. Recall that the quantizer label is given by

$$l_{ij} = \left\lfloor \frac{\theta_{ij}}{Q_{ij}^t} + 0.5 \right\rfloor \tag{18.8}$$

and the reconstructed value is given by

$$\hat{\theta}_{ij} = l_{ij} \times Q_{ij}^t. \tag{18.9}$$

Using these equations and the fact that $Q_{ij}^t = 4$, the reconstructed values of the coefficients are

$$28 \quad 8 \quad -8 \quad 0 \quad -4 \quad 0 \quad -0 \quad 4$$

The error in the reconstruction is

$$1.75 \quad -1.9 \quad 1.97 \quad 1.93 \quad 1.99 \quad 1.23 \quad -0.95 \quad -1.89$$

Now suppose we have some additional bandwidth made available to us. We can quantize the difference and send that to enhance the reconstruction. Suppose we used a step size of 2 to quantize the difference. The reconstructed values for this enhancement sequence would be

$$2 \quad -2 \quad 2 \quad 2 \quad 2 \quad 2 \quad 0 \quad -2$$

Adding this to the previous base-level reconstruction, we get an enhanced reconstruction of

$$30 \quad 6 \quad -6 \quad 2 \quad -2 \quad 2 \quad 0 \quad 2$$

which results in an error of

$$-0.25 \quad 0.1 \quad -0.03 \quad -0.07 \quad -0.01 \quad -0.77 \quad -0.95 \quad 0.11$$

The layered approach allows us to increase the accuracy of the reconstruction when bandwidth is available, while at the same time permitting a lower-quality reconstruction when there is not sufficient bandwidth for the enhancement. In other words, the quality is *scalable*. In this particular case, the error between the original and reconstruction decreases because of the enhancement. Because the signal-to-noise ratio is a measure of error, this can be called *snr-scalable*. If the enhancement layer contained a coded bitstream corresponding to frames that would occur between frames of the base layer, the system could be called *temporally scalable*. If the enhancement allowed an upsampling of the base layer, the system is *spatially scalable*.                                                                              ◆
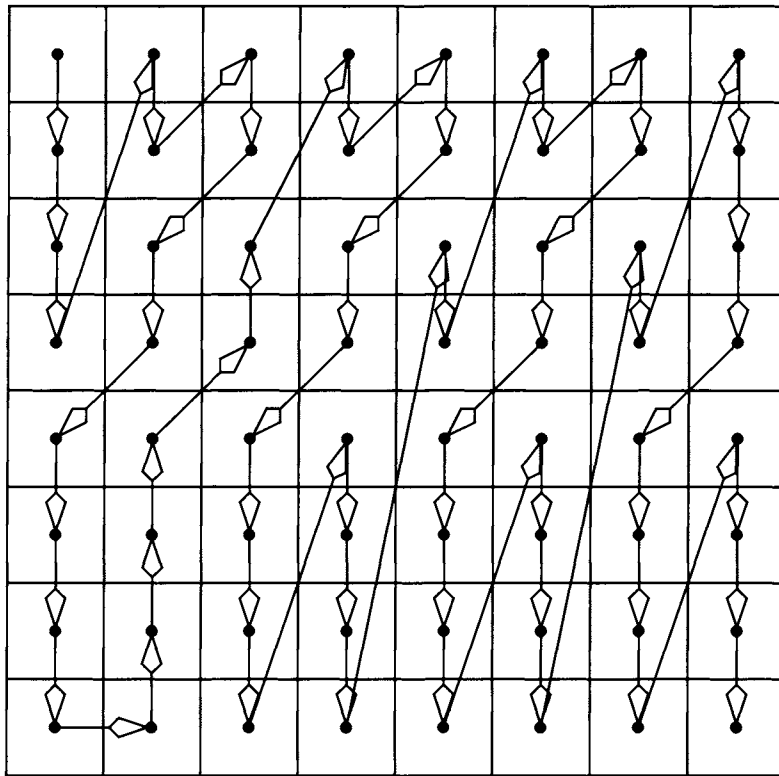
The levels are *low, main, high 1440*, and *high*. The *low* level corresponds to a frame size of 352 × 240, the *main* level corresponds to a frame size of 720 × 480, the *high 1440* level corresponds to a frame size of 1440 × 1152, and the *high* level corresponds to a frame size of 1920 × 1080. All levels are defined for a frame rate of 30 frames per second. There are many possible combinations of profiles and levels, not all of which are allowed in the MPEG-2 standard. Table 18.7 shows the allowable combinations [251]. A particular profile-level combination is denoted by *XX@YY* where *XX* is the two-letter abbreviation for the profile and *YY* is the two-letter abbreviation for the level. There are a large number of issues, such as bounds on parameters and decodability between different profile-level combinations, that we have not addressed here because they do not pertain to our main focus, compression (see the international standard [253] for these details).

Because MPEG-2 has been designed to handle interlaced video, there are field, based alternatives to the **I**, **P** and **B** frames. The **P** and **B** frames can be replaced by two **P** fields or two **B** fields. The **I** frame can be replaced by two **I** fields or an **I** field and a **P** field where the **P** field is obtained by predicting the bottom field by the top field. Because an 8 × 8 field block actually covers twice the spatial distance in the vertical direction as an 8 frame block, the zigzag scanning is adjusted to adapt to this imbalance. The scanning pattern for an 8 × 8 field block is shown in Figure 18.15

The most important addition from the point of view of compression in MPEG-2 is the addition of several new motion-compensated prediction modes: the field prediction and the dual prime prediction modes. MPEG-1 did not allow interlaced video. Therefore, there was no need for motion compensation algorithms based on fields. In the **P** frames, field predictions are obtained using one of the two most recently decoded fields. When the first field in a frame is being encoded, the prediction is based on the two fields from the

**TABLE 18.7      Allowable profile-level combinations in MPEG-2.**

|  | Simple Profile | Main Profile | SNR-Scalable Profile | Spatially Scalable Profile | High Profile |
|---|---|---|---|---|---|
| High Level |  | Allowed |  |  | Allowed |
| High 1440 |  | Allowed |  | Allowed | Allowed |
| Main Level | Allowed | Allowed | Allowed |  | Allowed |
| Low Level |  | Allowed | Allowed |  |  |

**FIGURE 18. 15**     *Scanning pattern for the DCT coefficients of a field block.*

previous frame. However, when the second field is being encoded, the prediction is based on the second field from the previous frame and the first field from the current frame. Information about which field is to be used for prediction is transmitted to the receiver. The field predictions are performed in a manner analogous to the motion-compensated prediction described earlier.

In addition to the regular frame and field prediction, MPEG-2 also contains two additional modes of prediction. One is the $16 \times 8$ motion compensation. In this mode, two predictions are generated for each macroblock, one for the top half and one for the bottom half. The other is called the dual prime motion compensation. In this technique, two predictions are formed for each field from the two recent fields. These predictions are averaged to obtain the final prediction.

## 18.9.1    The Grand Alliance HDTV Proposal

When the Federal Communications Commission (FCC) requested proposals for the HDTV standard, they received four proposals for digital HDTV from four consortia. After the

evaluation phase, the FCC declined to pick a winner among the four, and instead suggested that all these consortia join forces and produce a single proposal. The resulting partnership has the exalted title of the "Grand Alliance." Currently, the specifications for the digital HDTV system use MPEG-2 as the compression algorithm. The Grand Alliance system uses the *main* profile of the MPEG-2 standard implemented at the *high* level.

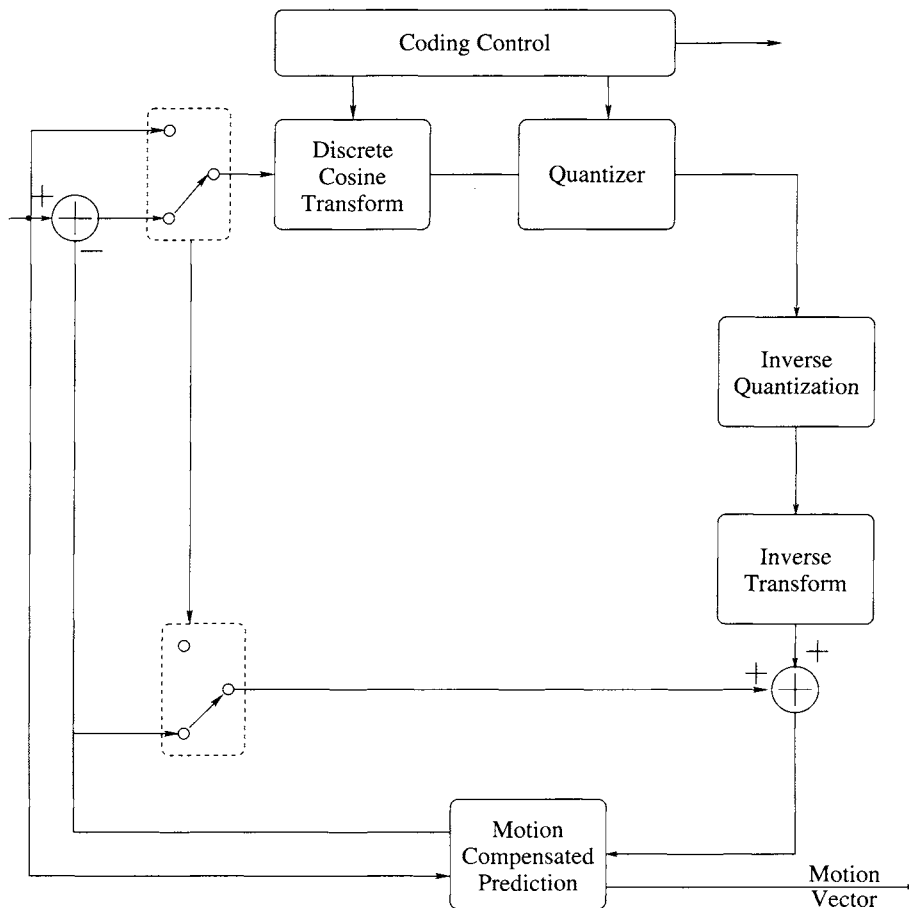## 18.10   ITU-T Recommendation H.263

The H.263 standard was developed to update the H.261 video conferencing standard with the experience acquired in the development of the MPEG and H.262 algorithms. The initial algorithm provided incremental improvement over H.261. After the development of the core algorithm, several optional updates were proposed, which significantly improved the compression performance. The standard with these optional components is sometimes referred to as H.263+ (or H.263 + +).

In the following sections we first describe the core algorithm and then describe some of the options. The standard focuses on noninterlaced video. The different picture formats addressed by the standard are shown in Table 18.8. The picture is divided into *Groups of Blocks* (GOBs) or slices. A Group of Blocks is a strip of pixels across the picture with a height that is a multiple of 16 lines. The number of multiples depends on the size of the picture, and the bottom-most GOB may have less than 16 lines. Each GOB is divided into macroblocks, which are defined as in the H.261 recommendation.

A block diagram of the baseline video coder is shown in Figure 18.16. It is very similar to Figure 18.10, the block diagram for the H.261 encoder. The only major difference is the ability to work with both predicted or **P** frames and intra or **I** frames. As in the case of H.261, the motion-compensated prediction is performed on a macroblock basis. The vertical and horizontal components of the motion vector are restricted to the range $[-16, 15.5]$. The transform used for representing the prediction errors in the case of the **P** frame and the pixels in the case of the **I** frames is the discrete cosine transform. The transform coefficients are quantized using uniform midtread quantizers. The DC coefficient of the intra block is quantized using a uniform quantizer with a step size of 8. There are 31 quantizers available for the quantization of all other coefficients with stepsizes ranging from 2 to 62. Apart from the DC coefficient of the intra block, all coefficients in a macroblock are quantized using the same quantizer.

**TABLE 18.8     The standardized H.263 formats [254].**

| Picture format | Number of luminance pixels (columns) | Number of luminance lines (rows) | Number of chrominance pixels (columns) | Number of chrominance lines(rows) |
|---|---|---|---|---|
| sub-QCIF | 128 | 96 | 64 | 48 |
| QCIF | 176 | 144 | 88 | 72 |
| CIF | 352 | 288 | 176 | 144 |
| 4CIF | 704 | 576 | 352 | 288 |
| 16CIF | 1408 | 1152 | 704 | 576 |

**FIGURE 18. 16    A block diagram of the H.263 video compression algorithm.**

The motion vectors are differentially encoded. The prediction is the median of the motion vectors in the neighboring blocks. The H.263 recommendation allows half pixel motion compensation as opposed to only integer pixel compensation used in H.261. Notice that the sign of the component is encoded in the last bit of the variable length code, a "0" for positive values and a "1" for negative values. Two values that differ only in their sign differ only in the least significant bit.

The code for the quantized transform coefficients is indexed by three indicators. The first indicates whether the coefficient being encoded is the last nonzero coefficient in the zigzag scan. The second indicator is the number of zero coefficients preceding the coefficient being encoded, and the last indicates the absolute value of the quantized coefficient level. The sign bit is appended as the last bit of the variable length code.

Here we describe some of the optional modes of the H.263 recommendation. The first four options were part of the initial H.263 specification. The remaining options were added later and the resulting standard is sometimes referred to as the H.263+ standard.

### 18.10.1   Unrestricted Motion Vector Mode

In this mode the motion vector range is extended to $[-31.5, 31.5]$, which is particularly useful in improving the compression performance of the algorithm for larger picture sizes. The mode also allows motion vectors to point outside the picture. This is done by repeating the edge pixels to create the picture beyond its boundary.

### 18.10.2   Syntax-Based Arithmetic Coding Mode

In this mode the variable length codes are replaced with an arithmetic coder. The word length for the upper and lower limits is 16. The option specifies several different Cum Count tables that can be used for arithmetic coding. There are separate Cum Count tables for encoding motion vectors, intra DC component, and intra and inter coefficients.

### 18.10.3   Advanced Prediction Mode

In the baseline mode a single motion vector is sent for each macroblock. Recall that a macroblock consists of four $8 \times 8$ luminance blocks and two chrominance blocks. In the advanced prediction mode the encoder can send four motion vectors, one for each luminance block. The chrominance motion vectors are obtained by adding the four luminance motion vectors and dividing by 8. The resulting values are adjusted to the nearest half pixel position. This mode also allows for *Overlapped Block Motion Compensation (OBMC)*. In this mode the motion vector is obtained by taking a weighted sum of the motion vector of the current block and two of the four vertical and horizontal neighboring blocks.

### 18.10.4   PB-frames and Improved PB-frames Mode

The PB frame consists of a **P** picture and a **B** picture in the same frame. The blocks for the **P** frame and the **B** frame are interleaved so that a macroblock consists of six blocks of a **P** picture followed by six blocks of a **B** picture. The motion vector for the **B** picture is derived from the motion vector for the **P** picture by taking into account the time difference between the **P** picture and the **B** picture. If the motion cannot be properly derived, a delta correction is included. The improved PB-frame mode updates the PB-frame mode to include forward, backward, and bidirectional prediction.

### 18.10.5   Advanced Intra Coding Mode

The coefficients for the **I** frames are obtained directly by transforming the pixels of the picture. As a result, there can be significant correlation between some of the coefficients of neighboring blocks. For example, the DC coefficient represents the average value of a

block. It is very likely that the average value will not change significantly between blocks. The same may be true, albeit to a lesser degree, for the low-frequency horizontal and vertical coefficients. The advanced intra coding mode allows the use of this correlation by using coefficients from neighboring blocks for predicting the coefficients of the block being encoded. The prediction errors are then quantized and coded.

When this mode is used, the quantization approach and variable length codes have to be adjusted to adapt to the different statistical properties of the prediction errors. Furthermore, it might also become necessary to change the scan order. The recommendation provides alternate scanning patterns as well as alternate variable length codes and quantization strategies.

## 18.10.6  Deblocking Filter Mode

This mode is used to remove blocking effects from the $8 \times 8$ block edges. This smoothing of block boundaries allows for better prediction. This mode also permits the use of four motion vectors per macroblock and motion vectors that point beyond the edge of the picture.

## 18.10.7  Reference Picture Selection Mode

This mode is used to prevent error propagation by allowing the algorithm to use a picture other than the previous picture to perform prediction. The mode permits the use of a back-channel that the decoder uses to inform the encoder about the correct decoding of parts of the picture. If a part of the picture is not correctly decoded, it is not used for prediction. Instead, an alternate frame is selected as the reference frame. The information about which frame was selected as the reference frame is transmitted to the decoder. The number of possible reference frames is limited by the amount of frame memory available.

## 18.10.8  Temporal, SNR, and Spatial Scalability Mode

This is very similar to the scalability structures defined earlier for the MPEG-2 algorithm. Temporal scalability is achieved by using separate **B** frames, as opposed to the PB frames. SNR scalability is achieved using the kind of layered coding described earlier. Spatial scalability is achieved using upsampling.

## 18.10.9  Reference Picture Resampling

Reference picture resampling allows a reference picture to be "warped" in order to permit the generation of better prediction. It can be used to adaptively alter the resolution of pictures during encoding.

### 18.10.10   Reduced-Resolution Update Mode

This mode is used for encoding highly active scenes. The macroblock in this mode is assumed to cover an area twice the height and width of the regular macroblock. The motion vector is assumed to correspond to this larger area. Using this motion vector a predicted macroblock is created. The transform coefficients are decoded and then upsampled to create the expanded texture block. The predicted and texture blocks are then added to obtain the reconstruction.

### 18.10.11   Alternative Inter VLC Mode

The variable length codes for inter and intra frames are designed with different assumptions. In the case of the inter frames it is assumed that the values of the coefficients will be small and there can be large numbers of zero values between nonzero coefficients. This is a result of prediction which, if successfully employed, would reduce the magnitude of the differences, and hence the coefficients, and would also lead to large numbers of zero-valued coefficients. Therefore, coefficients indexed with large runs and small coefficient values are assigned shorter codes. In the case of the intra frames, the opposite is generally true. There is no prediction, therefore there is a much smaller probability of runs of zero-valued coefficients. Also, large-valued coefficients are quite possible. Therefore, coefficients indexed by small run values and larger coefficient values are assigned shorter codes. During periods of increased temporal activity, prediction is generally not as good and therefore the assumptions under which the variable length codes for the inter frames were created are violated. In these situations it is likely that the variable length codes designed for the intra frames are a better match. The alternative inter VLC mode allows for the use of the intra codes in these sitations, improving the compression performance. Note that the codewords used in intra and inter frame coding are the same. What is different is the interpretation. To detect the proper interpretation, the decoder first decodes the block assuming an inter frame codebook. If the decoding results in more than 64 coefficients it switches its interpretation.

### 18.10.12   Modified Quantization Mode

In this mode, along with changes in the signalling of changes in quantization parameters, the quantization process is improved in several ways. In the baseline mode, both the luminanace and chrominance components in a block are quantized using the same quantizer. In the modified quantization mode, the quantizer used for the luminance coefficients is different from the quantizer used for the chrominance component. This allows the quantizers to be more closely matched to the statistics of the input. The modified quantization mode also allows for the quantization of a wider range of coefficient values, preventing significant overload. If the coefficient exceeds the range of the baseline quantizer, the encoder sends an escape symbol followed by an 11-bit representation of the coefficient. This relaxation of the structured representation of the quantizer outputs makes it more likely that bit errors will be accepted as valid quantizer outputs. To reduce the chances of this happening, the mode prohibits "unreasonable" coefficient values.

## 18.10.13  Enhanced Reference Picture Selection Mode

Motion-compensated prediction is accomplished by searching the previous picture for a block similar to the block being encoded. The enhanced reference picture selection mode allows the encoder to search more than one picture to find the best match and then use the best-suited picture to perform motion-compensated prediction. Reference picture selection can be accomplished on a macroblock level. The selection of the pictures to be used for motion compensation can be performed in one of two ways. A sliding window of $M$ pictures can be used and the last $M$ decoded, with reconstructed pictures stored in a multipicture buffer. A more complex adaptive memory (not specified by the standard) can also be used in place of the simple sliding window. This mode significantly enhances the prediction, resulting in a reduction in the rate for equivalent quality. However, it also increases the computational and memory requirements on the encoder. This memory burden can be mitigated to some extent by assigning an unused label to pictures or portions of pictures. These pictures, or portions of pictures, then do not need to be stored in the buffer. This unused label can also be used as part of the adaptive memory control to manage the pictures that are stored in the buffer.

## 18.11  ITU-T Recommendation H.264, MPEG-4 Part 10, Advanced Video Coding

As described in the previous section, the H.263 recommendation started out as an incremental improvement over H.261 and ended up with a slew of optional features, which in fact make the improvement over H.261 more than incremental. In H.264 we have a standard that started out with a goal of significant improvement over the MPEG-1/2 standards and achieved those goals. The standard, while initiated by ITU-T's Video Coding Experts Group (VCEG), ended up being a collaboration of the VCEG and ISO/IEC's MPEG committees which joined to form the Joint Video Team (JVT) in December of 2001 [255]. The collaboration of various groups in the development of this standard has also resulted in the richness of names. It is variously known as ITU-T H.264, MPEG-4 Part 10, MPEG-4 Advanced Video Coding (AVC), as well as the name under which it started its life, H.26L. We will just refer to it as H.264.

The basic block diagram looks very similar to the previous schemes. There are intra and inter pictures. The inter pictures are obtained by subtracting a motion compensated prediction from the original picture. The residuals are transformed into the frequency domain. The transform coefficients are scanned, quantized, and encoded using variable length codes. A local decoder reconstructs the picture for use in future predictions. The intra picture is coded without reference to past pictures.

While the basic block diagram is very similar to the previous standards the details are quite different. We will look at these details in the following sections. We begin by looking at the basic structural elements, then look at the decorrelation of the inter frames. The decorrelation process includes motion-compensated prediction and transformation of the prediction error. We then look at the decorrelation of the intra frames. This includes intra prediction

modes and transforms used in this mode. We finally look at the different binary coding options.

The macroblock structure is the same as used in the other standards. Each macroblock consists of four $8 \times 8$ luminance blocks and two chrominance blocks. An integer number of sequential macroblocks can be put together to form a slice. In the previous standards the smallest subdivision of the macroblock was into its $8 \times 8$ component blocks. The H.264 standard allows $8 \times 8$ macroblock partitions to be further divided into sub-macroblocks of size $8 \times 4$, $4 \times 8$, and $4 \times 4$. These smaller blocks can be used for motion-compensated prediction, allowing for tracking of much finer details than is possible with the other standards. Along with the $8 \times 8$ partition, the macroblock can also be partitioned into two $8 \times 16$ or $16 \times 8$ blocks. In field mode the H.264 standard groups $16 \times 8$ blocks from each field to form a $16 \times 16$ macroblock.

### 18.11.1   Motion-Compensated Prediction

The H.264 standard uses its macroblock partitions to develop a tree-structured motion compensation algorithm. One of the problems with motion-compensated prediction has always been the selection of the size and shape of the block used for prediction. Different parts of a video scene will move at different rates in different directions or stay put. A smaller-size block allows tracking of diverse movement in the video frame, leading to better prediction and hence lower bit rates. However, more motion vectors need to be encoded and transmitted, using up valuable bit resources. In fact, in some video sequences the bits used to encode the motion vectors may make up most of the bits used. If we use small blocks, the number of motion vectors goes up, as does the bit rate. Because of the variety of sizes and shapes available to it, the H.264 algorithm provides a high level of accuracy and efficiency in its prediction. It uses small block sizes in regions of activity and larger block sizes in stationary regions. The availability of rectangular shapes allows the algorithm to focus more precisely on regions of activity.

The motion compensation is accomplished using quarter-pixel accuracy. To do this the reference picture is "expanded" by interpolating twice between neighboring pixels. This results in a much smoother residual. The prediction process is also enhanced by the use of filters on the 4 block edges. The standard allows for searching of up to 32 pictures to find the best matching block. The selection of the reference picture is done on the macroblock partion level, so all sub-macroblock partitions use the same reference picture.

As in H.263, the motion vectors are differentially encoded. The basic scheme is the same. The median values of the three neighboring motion vectors are used to predict the current motion vector. This basic strategy is modified if the block used for motion compensation is a $16 \times 16$, $16 \times 8$, or $8 \times 16$ block.

For **B** pictures, as in the case of the previous standards, two motion vectors are allowed for each macroblock or sub-macroblock partition. The prediction for each pixel is the weighted average of the two prediction pixels.

Finally, a $\mathbf{P}_{skip}$ type macroblock is defined for which $16 \times 16$ motion compensation is used and the prediction error is not transmitted. This type of macroblock is useful for regions of little change as well as for slow pans.

## 18.11.2  The Transform

Unlike the previous video coding schemes, the transform used is not an $8 \times 8$ DCT. For most blocks the transform used is a $4 \times 4$ integer DCT-like matrix. The transform matrix is given by

$$H = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 2 & 1 & -1 & 2 \\ 1 & -1 & -1 & 1 \\ 1 & -2 & 2 & -1 \end{bmatrix}$$

The inverse transform matrix is given by

$$H^I = \begin{bmatrix} 1 & 1 & 1 & \frac{1}{2} \\ 1 & \frac{1}{2} & -1 & -1 \\ 1 & -\frac{1}{2} & -1 & 1 \\ 1 & -1 & 1 & -\frac{1}{2} \end{bmatrix}$$
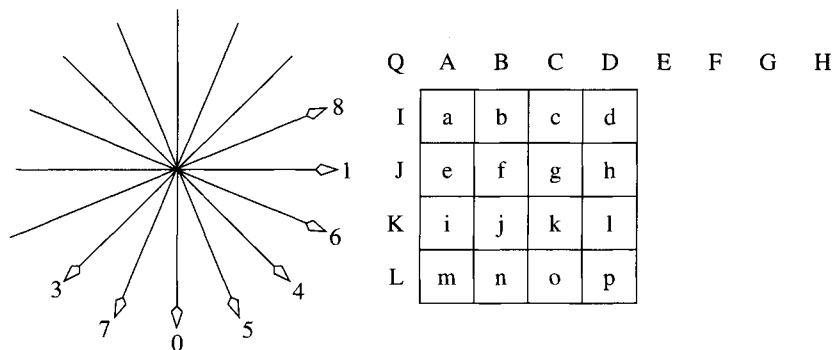
The transform operations can be implemented using addition and shifts. Multiplication by 2 is a single-bit left shift and division by 2 is a single-bit right shift. However, there is a price for the simplicity. Notice that the norm of the rows is not the same and the product of the forward and inverse transforms does not result in the identity matrix. This discrepancy is compensated for by the use of scale factors during quantization. There are several advantages to using a smaller integer transform. The integer nature makes the implementation simple and also avoids error accumulation in the transform process. The smaller size allows better representation of small stationary regions of the image. The smaller blocks are less likely to contain a wide range of values. Where there are sharp transitions in the blocks, any ringing effect is contained within a small number of pixels.

## 18.11.3  Intra Prediction

In the previous standards the **I** pictures were transform coded without any decorrelation. This meant that the number of bits required for the **I** frames is substantially higher than for the other pictures. When asked why he robbed banks, the notorious robber Willie Sutton is supposed to have said simply, "because that's where the money is." Because most of the bits in video coding are expended in encoding the **I** frame, it made a lot of sense for the JVT to look at improving the compression of the **I** frame in order to substantially reduce the bitrate.

The H.264 standard contains a number of spatial prediction modes. For $4 \times 4$ blocks there are nine prediction modes. Eight of these are summarized in Figure 18.17. The sixteen pixels in the block $a-p$ are predicted using the thirteen pixels on the boundary (and extending from it).[1] The arrows corresponding to the mode numbers show the direction of prediction. For example, mode 0 corresponds to the downward pointing arrow. In this case pixel $A$ is used to predict pixels $a, e, i, m$, pixel $B$ is used to predict pixels $b, f, j, n$, pixel $C$ is used

---

[1] The jump from pixel $L$ to $Q$ is a historical artifact. In an earlier version of the standard, pixels below $L$ were also used in some prediction modes.

**FIGURE 18.17    Prediction modes for 4 × 4 intra prediction.**

to predict pixels $c, g, k, o$, and pixel $D$ is used to predict pixels $d, h, l, p$. In mode 3, also called the diagonal down/left mode, $B$ is used to predict $a$, $C$ is used to predict $b, e$, pixel $D$ is used to predict pixels $c, f, i$, pixel $E$ is used to predict pixels $d, g, j, m$, pixel $F$ is used to predict pixels $h, k, n$, pixel $G$ is used to predict pixels $l, o$, and pixel $H$ is used to predict pixel $p$. If pixels $E, F, G$, and $H$ are not available, pixel $D$ is repeated four times. Notice that no direction is availble for mode 2. This is called the DC mode, in which the average of the left and top boundary pixels is used as a prediction for all sixteen pixels in the $4 \times 4$ block. In most cases the prediction modes used for all the $4 \times 4$ blocks in a macroblock are heavily correlated. The standard uses this correlation to efficiently encode the mode information.

In smooth regions of the picture it is more convenient to use prediction on a macroblock basis. In case of the full macroblock there are four prediction modes. Three of them correspond to modes 0, 1, and 2 (vertical, horizontal, and DC). The fourth prediction mode is called the planar prediction mode. In this mode a three-parameter plane is fitted to the pixel values of the macroblock.

## 18.11.4  Quantization

The H.264 standard uses a uniform scalar quantizer for quantizing the coefficients. There are 52 scalar quantizers indexed by $Q_{step}$. The step size doubles for every sixth $Q_{step}$. The quantization incorporates scaling necessitated by the approximations used to make the transform simple. If $\alpha_{i,j}(Q_{step})$ are the weights for the $(i, j)^{th}$ coefficient then

$$l_{i,j} = sign(\theta_{i,j}) \left\lfloor \frac{|\theta_{i,j}| \alpha_{i,j}(Q_{step})}{Q_{step}} \right\rfloor$$

In order to broaden the quantization interval around the origin we add a small value in the numerator.

$$l_{i,j} = sign(\theta_{i,j}) \left\lfloor \frac{|\theta_{i,j}| \alpha_{i,j}(Q_{step}) + f(Q_{step})}{Q_{step}} \right\rfloor$$

In actual implementation we do away with divisions and the quantization is implemented as [255]

$$l_{i,j} = sign(\theta_{i,j})[|\theta_{i,j}|M(Q_M, r) + f2^{17+Q_E}] >> 17 + Q_E$$

where

$$Q_M = Q_{step}(mod6)$$

$$Q_E = \left\lfloor \frac{Q_{step}}{6} \right\rfloor$$

$$r = \begin{cases} 0 & i, j \text{ even} \\ 1 & i, j \text{ odd} \\ 2 & \text{otherwise} \end{cases}$$

and $M$ is given in Table 18.9

The inverse quantization is given by

$$\hat{\theta}_{i,j} = l_{i,j}S(Q_M, r) << Q_E$$

where $S$ is given in Table 18.10

Prior to quantization, the transforms of the $16 \times 16$ luminance residuals and the $8 \times 8$ chrominance residuals of the macroblock-based intra prediction are processed to further remove redundancy. Recall that macroblock-based prediction is used in smooth regions of the **I** picture. Therefore, it is very likely that the DC coefficients of the $4 \times 4$ transforms

**TABLE 18.9**    **$M(Q_M, r)$ values in H.264.**

| $Q_M$ | $r = 0$ | $r = 1$ | $r = 2$ |
|---|---|---|---|
| 0 | 13107 | 5243 | 8066 |
| 1 | 11916 | 4660 | 7490 |
| 2 | 10082 | 4194 | 6554 |
| 3 | 9362 | 3647 | 5825 |
| 4 | 8192 | 3355 | 5243 |
| 5 | 7282 | 2893 | 4559 |

**TABLE 18.10**    **$S(Q_M, r)$ values in H.264.**

| $Q_M$ | $r = 0$ | $r = 1$ | $r = 2$ |
|---|---|---|---|
| 0 | 10 | 16 | 13 |
| 1 | 11 | 18 | 14 |
| 2 | 13 | 20 | 16 |
| 3 | 14 | 23 | 18 |
| 4 | 16 | 25 | 20 |
| 5 | 18 | 29 | 23 |

are heavily correlated. To remove this redundancy, a discrete Walsh-Hadamard transform is used on the DC coefficients in the macroblock. In the case of the luminance block, this is a $4 \times 4$ transform for the sixteen DC coefficients. The smaller chrominance block contains four DC coefficients, so we use a $2 \times 2$ discrete Walsh-Hadamard transform.

## 18.11.5   Coding

The H.264 standard contains two options for binary coding. The first uses exponential Golomb codes to encode the parameters and a context-adaptive variable length code (CAVLC) to encode the quantizer labels [255]. The second binarizes all the values and then uses a context-adaptive binary arithmetic code (CABAC) [256].

An exponential Golomb code for a positive number $x$ can be obtained as the unary code for $M = \lfloor \log_2(x+1) \rfloor$ concatenated with the $M$ bit natural binary code for $x + 1$. The unary code for a number $x$ is given as $x$ zeros followed by a 1. The exponential Golomb code for zero is 1.

The quantizer labels are first scanned in a zigzag fashion. In many cases the last nonzero labels in the zigzag scan have a magnitude of 1. The number $N$ of nonzero labels and the number $T$ of trailing ones are used as an index into a codebook that is selected based on the values of $N$ and $T$ for the neighboring blocks. The maximum allowed value of $T$ is 3. If the number of trailing labels with a magnitude of 1 is greater than 3, the remaining are encoded in the same manner as the other nonzero labels. The nonzero labels are then coded in reverse order. That is, the quantizer labels corresponding to the higher-frequency coefficients are encoded first. First the signs of the trailing 1s are encoded with 0s signifying positive values and 1s signifying negative values. Then the remaining quantizer labels are encoded in reverse scan order. After this, the total number of 0s in the scan between the beginning of the scan and the last nonzero label is encoded. This will be a number between 0 and $16 - N$. Then the run of zeros before each label, starting with the last nonzero label is encoded until we run out of zeros or coefficients. The number of bits used to code each zero run will depend on the number of zeros remaining to be assigned.
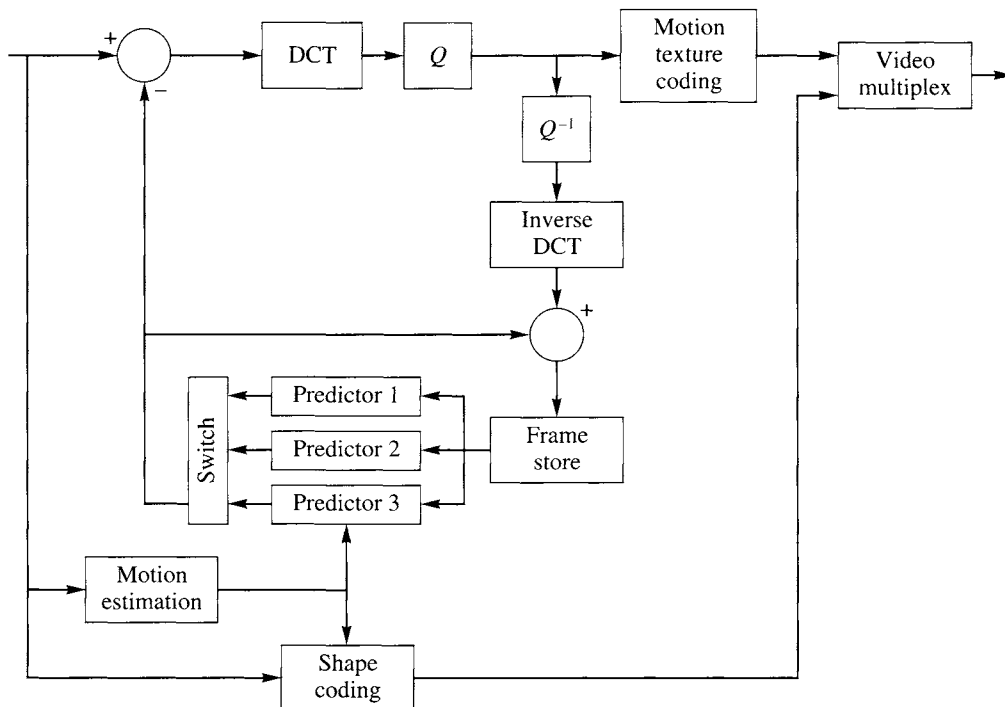
In the second technique, which provides higher compression, all values are first converted to binary strings. This binarization is performed, depending on the data type, using unary codes, truncated unary codes, exponential Golomb codes, and fixed-length codes, plus five specific binary trees for encoding macroblock and sub-macroblock types. The binary string is encoded in one of two ways. Redundant strings are encoded using a context-adaptive binary arithmetic code. Binary strings that are random, such as the suffixes of the exponential Golomb codes, bypass the arithmetic coder. The arithmetic coder has 399 contexts available to it, with 325 of these contexts used for encoding the quantizer labels. These numbers include contexts for both frame and field slices. In a pure frame or field slice only 277 of the 399 context models are used. These context models are simply Cum_Count tables for use with a binary arithmetic coder. The H.264 standard recommends a multiplication-free implementation of binary arithmetic coding.

The H.264 standard is substantially more flexible than previous standards, with a much broader range of applications. In terms of performance, it claims a 50% reduction in bit rate over previous standards for equivalent perceptual quality [255].

## 18.12  MPEG-4 Part 2

The MPEG-4 standard provides a more abstract approach to the coding of multimedia. The standard views a multimedia "scene" as a collection of objects. These objects can be visual, such as a still background or a talking head, or aural, such as speech, music, background noise, and so on. Each of these objects can be coded independently using different techniques to generate separate elementary bitstreams. These bitstreams are multiplexed along with a scene description. A language called the Binary Format for Scenes (BIFS) based on the Virtual Reality Modeling Language (VRML) has been developed by MPEG for scene descriptions. The decoder can use the scene description and additional input from the user to combine or compose the objects to reconstruct the original scene or create a variation on it. The protocol for managing the elementary streams and their multiplexed version, called the Delivery Multimedia Integration Framework (DMIF), is an important part of MPEG-4. However, as our focus in this book is on compression, we will not discuss the protocol (for details, see the standard [213]).

A block diagram for the basic video coding algorithm is shown in Figure 18.18. Although shape coding occupies a very small portion of the diagram, it is a major part of the algorithm. The different objects that make up the scene are coded and sent to the multiplexer. The information about the presence of these objects is also provided to the motion-compensated



**FIGURE 18. 18    A block diagram for video coding.**

predictor, which can use object-based motion compensation algorithms to improve the compression efficiency. What is left after the prediction can be transmitted using a DCT-based coder. The video coding algorithm can also use a background "sprite"—generally a large panoramic still image that forms the background for the video sequence. The sprite is transmitted once, and the moving foreground video objects are placed in front of different portions of the sprite based on the information about the background provided by the encoder.

The MPEG-4 standard also envisions the use of model-based coding, where a triangular mesh representing the moving object is transmitted followed by texture information for covering the mesh. Information about movement of the mesh nodes can then be transmitted to animate the video object. The texture coding technique suggested by the standard is the embedded zerotree wavelet (EZW) algorithm. In particular, the standard envisions the use of a facial animation object to render an animated face. The shape, texture, and expressions of the face are controlled using facial definition parameters (FDPs) and facial action parameters (FAPs). BIFS provides features to support custom models and specialized interpretation of FAPs.

The MPEG-2 standard allows for SNR and spatial scalability. The MPEG-4 standard also allows for object scalability, in which certain objects may not be sent in order to reduce the bandwidth requirement.

## 18.13   Packet Video

The increasing popularity of communication over networks has led to increased interest in the development of compression schemes for use over networks. In this section we look at some of the issues involved in developing video compression schemes for use over networks.

## 18.14   ATM Networks

With the explosion of information, we have also seen the development of new ways of transmitting the information. One of the most efficient ways of transferring information among a large number of users is the use of asynchronous transfer mode (ATM) technology. In the past, communication has usually taken place over dedicated channels; that is, in order to communicate between two points, a channel was dedicated only to transferring information between those two points. Even if there was no information transfer going on during a particular period, the channel could not be used by anyone else. Because of the inefficiency of this approach, there is an increasing movement away from it. In an ATM network, the users divide their information into packets, which are transmitted over channels that can be used by more than one user.

We could draw an analogy between the movement of packets over a communication network and the movement of automobiles over a road network. If we break up a message into packets, then the movement of the message over the network is like the movement of a number of cars on a highway system going from one point to the other. Although two cars may not occupy the same position at the same time, they can occupy the same road at the same time. Thus, more than one group of cars can use the road at any given time.

Furthermore, not all the cars in the group have to take the same route. Depending on the amount of traffic on the various roads that run between the origin of the traffic and the destination, different cars can take different routes. This is a more efficient utilization of the road than if the entire road was blocked off until the first group of cars completed its traversal of the road.

Using this analogy, we can see that the availability of transmission capacity, that is, the number of bits per second that we can transmit, is affected by factors that are outside our control. If at a given time there is very little traffic on the network, the available capacity will be high. On the other hand, if there is congestion on the network, the available capacity will be low. Furthermore, the ability to take alternate routes through the network also means that some of the packets may encounter congestion, leading to a variable amount of delay through the network. In order to prevent congestion from impeding the flow of vital traffic, networks will prioritize the traffic, with higher-priority traffic being permitted to move ahead of lower-priority traffic. Users can negotiate with the network for a fixed amount of guaranteed traffic. Of course, such guarantees tend to be expensive, so it is important that the user have some idea about how much high-priority traffic they will be transmitting over the network.

## 18.14.1 Compression Issues in ATM Networks

In video coding, this situation provides both opportunities and challenges. In the video compression algorithms discussed previously, there is a buffer that smooths the output of the compression algorithm. Thus, if we encounter a high-activity region of the video and generate more than the average number of bits per second, in order to prevent the buffer from overflowing, this period has to be followed by a period in which we generate fewer bits per second than the average. Sometimes this may happen naturally, with periods of low activity following periods of high activity. However, it is quite likely that this would not happen, in which case we have to reduce the quality by increasing the step size or dropping coefficients, or maybe even entire frames.

The ATM network, if it is not congested, will accommodate the variable rate generated by the compression algorithm. But if the network is congested, the compression algorithm will have to operate at a reduced rate. If the network is well designed, the latter situation will not happen too often, and the video coder can function in a manner that provides uniform quality. However, when the network is congested, it may remain so for a relatively long period. Therefore, the compression scheme should have the ability to operate for significant periods of time at a reduced rate. Furthermore, congestion might cause such long delays that some packets arrive after they can be of any use; that is, the frame they were supposed to be a part of might have already been reconstructed.
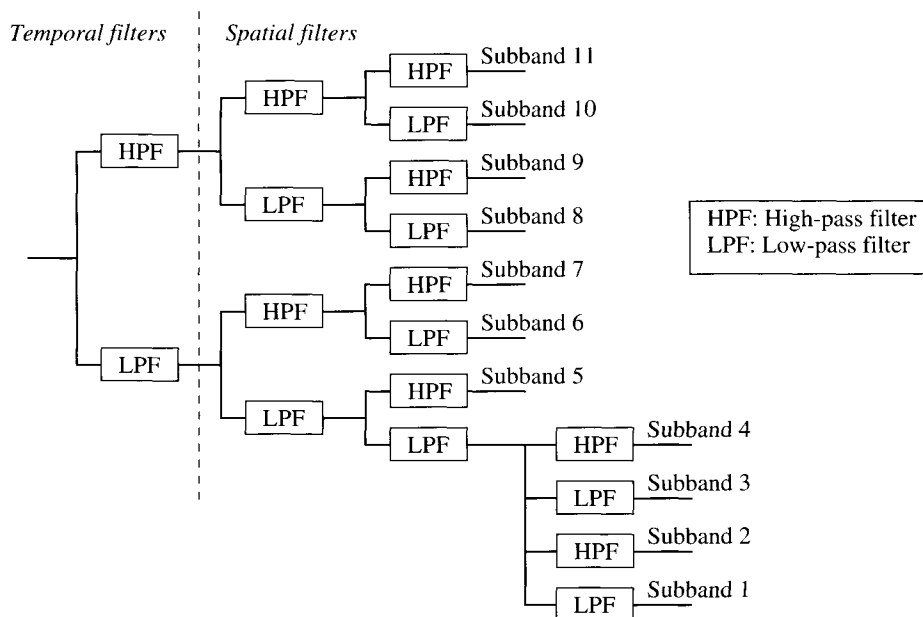
In order to deal with these problems, it is useful if the video compression algorithm provides information in a layered fashion, with a low-rate high-priority layer that can be used to reconstruct the video, even though the reconstruction may be poor, and low-priority enhancement layers that enhance the quality of the reconstruction. This is similar to the idea of progressive transmission, in which we first send a crude but low-rate representation of the image, followed by higher-rate enhancements. It is also useful if the bit rate required for the high-priority layer does not vary too much.

## 18.14.2   Compression Algorithms for Packet Video

Almost any compression algorithm can be modified to perform in the ATM environment, but some approaches seem more suited to this environment. We briefly present two approaches (see the original papers for more details).

One compression scheme that functions in an inherently layered manner is subband coding. In subband coding, the lower-frequency bands can be used to provide the basic reconstruction, with the higher-frequency bands providing the enhancement. As an example, consider the compression scheme proposed for packet video by Karlsson and Vetterli [257]. In their scheme, the video is divided into 11 bands. First, the video signal is divided into two temporal bands. Each band is then split into four spatial bands. The low-low band of the temporal low-frequency band is then split into four spatial bands. A graphical representation of this splitting is shown in Figure 18.19. The subband denoted 1 in the figure contains the basic information about the video sequence. Therefore, it is transmitted with the highest priority. If the data in all the other subbands are lost, it will still be possible to reconstruct the video using only the information in this subband. We can also prioritize the output of the other bands, and if the network starts getting congested and we are required to reduce our rate, we can do so by not transmitting the information in the lower-priority subbands. Subband 1 also generates the least variable data rate. This is very helpful when negotiating with the network for the amount of priority traffic.

Given the similarity of the ideas behind progressive transmission and subband coding, it should be possible to use progressive transmission algorithms as a starting point in the



**FIGURE 18. 19     Analysis filter bank.**

design of layered compression schemes for packet video. Chen, Sayood, and Nelson [258] use a DCT-based progressive transmission scheme [259] to develop a compression algorithm for packet video. In their scheme, they first encode the difference between the current frame and the prediction for the current frame using a $16 \times 16$ DCT. They only transmit the DC coefficient and the three lowest-order AC coefficients to the receiver. The coded coefficients make up the highest-priority layer.

The reconstructed frame is then subtracted from the original. The sum of squared errors is calculated for each $16 \times 16$ block. Blocks with squared error greater than a prescribed threshold are subdivided into four $8 \times 8$ blocks, and the coding process is repeated using an $8 \times 8$ DCT. The coded coefficients make up the next layer. Because only blocks that fail to meet the threshold test are subdivided, information about which blocks are subdivided is transmitted to the receiver as side information.

The process is repeated with $4 \times 4$ blocks, which make up the third layer, and $2 \times 2$ blocks, which make up the fourth layer. Although this algorithm is a variable-rate coding scheme, the rate for the first layer is constant. Therefore, the user can negotiate with the network for a fixed amount of high-priority traffic. In order to remove the effect of delayed packets from the prediction, only the reconstruction from the higher-priority layers is used for prediction.

This idea can be used with many different progressive transmission algorithms to make them suitable for use over ATM networks.

## 18.15  Summary

In this chapter we described a number of different video compression algorithms. The only new information in terms of compression algorithms was the description of motion-compensated prediction. While the compression algorithms themselves have already been studied in previous chapters, we looked at how these algorithms are used under different requirements. The three scenarios that we looked at are teleconferencing, asymmetric applications such as broadcast video, and video over packet networks. Each application has slightly different requirements, leading to different ways of using the compression algorithms. We have by no means attempted to cover the entire area of video compression. However, by now you should have sufficient background to explore the subject further using the following list as a starting point.

### Further Reading

1. An excellent source for information about the technical issues involved with digital video is the book *The Art of Digital Video*, by J. Watkinson [260].

2. The MPEG-1 standards document [252], "Information Technology—Coding of Moving Pictures and Associated Audio for Digital Storage Media Up to about 1.5 Mbit/s," has an excellent description of the video compression algorithm.

**3.** Detailed information about the MPEG-1 and MPEG-2 video standards can also be found in *MPEG Video Compression Standard*, by J.L. Mitchell, W.B. Pennebaker, C.E. Fogg, and D.J. LeGall [261].

**4.** To find more on model-based coding, see "Model Based Image Coding: Advanced Video Coding Techniques for Very Low Bit-Rate Applications," by K. Aizawa and T.S. Huang [250], in the February 1995 issue of the *Proceedings of the IEEE*.

**5.** A good place to begin exploring the various areas of research in packet video is the June 1989 issue of the *IEEE Journal on Selected Areas of Communication*.

**6.** The MPEG 1/2 and MPEG 4 standards are covered in an accesible manner in the book. *The MPEG Handbook* by J. Watkinson [261]. Focal press 2001.

**7.** A good source for information about H.264 and MPEG-4 is H.264 and MPEG-4 video compression, by I.E.G. Richardson. Wiley, 2003.

## 18.16   Projects and Problems

**1. (a)** Take the DCT of the Sinan image and plot the average squared value of each coefficient.

**(b)** Circularly shift each line of the image by eight pixels. That is, $new\_image[i, j] = old\_image[i, j + 8 \pmod{256}]$. Take the DCT of the difference and plot the average squared value of each coefficient.

**(c)** Compare the results in parts (a) and (b) above. Comment on the differences.