

Figure 5.3 Typical 16 Megabit DRAM (4M x 4)

Figure 5.3 also indicates the inclusion of refresh circuitry. All DRAMs require a refresh operation. A simple technique for refreshing is, in effect, to disable the DRAM chip while all data cells are refreshed. The refresh counter steps through all of the row values. For each row, the output lines from the refresh counter are supplied to the row decoder and the RAS line is activated. The data are read out and written back into the same location. This causes each cell in the row to be refreshed.

### Chip Packaging

As was mentioned in Chapter 2, an integrated circuit is mounted on a package that contains pins for connection to the outside world.

Figure 5.4a shows an example EPROM package, which is an 8-Mbit chip organized as  $1M \times 8$ . In this case, the organization is treated as a one-word-per-chip package. The package includes 32 pins, which is one of the standard chip package sizes. The pins support the following signal lines:

- The address of the word being accessed. For 1M words, a total of 20 ( $2^{20} = 1M$ ) pins are needed (A0–A19).
- The data to be read out, consisting of 8 lines (D0–D7).
- The power supply to the chip ( $V_{cc}$ ).
- A ground pin ( $V_{ss}$ ).
- A chip enable (CE) pin. Because there may be more than one memory chip, each of which is connected to the same address bus, the CE pin is used to indicate whether or not the address is valid for this chip. The CE pin is activated by

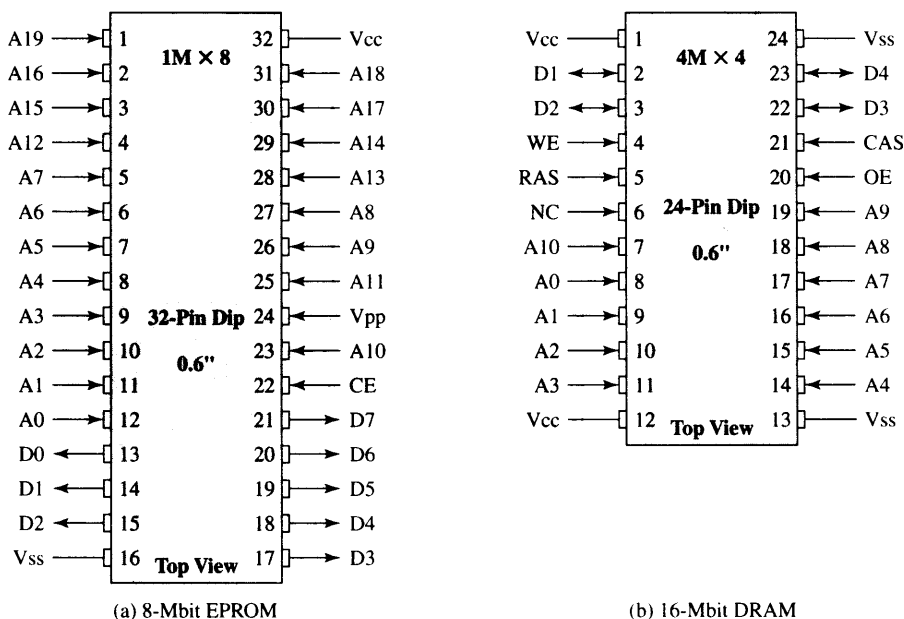


Figure 5.4 Typical Memory Package Pins and Signals

logic connected to the higher-order bits of the address bus (i.e., address bits above A19). The use of this signal is illustrated presently.

- A program voltage ( $V_{pp}$ ) that is supplied during programming (write operations).

A typical DRAM pin configuration is shown in Figure 5.4b, for a 16-Mbit chip organized as  $4M \times 4$ . There are several differences from a ROM chip. Because a RAM can be updated, the data pins are input/output. The write enable (WE) and output enable (OE) pins indicate whether this is a write or read operation. Because the DRAM is accessed by row and column, and the address is multiplexed, only 11 address pins are needed to specify the  $4M$  row/column combinations ( $2^{11} \times 2^{11} = 2^{22} = 4M$ ). The functions of the row address select (RAS) and column address select (CAS) pins were discussed previously. Finally, the no connect (NC) pin is provided so that there are an even number of pins.

### Module Organization

If a RAM chip contains only 1 bit per word, then clearly we will need at least a number of chips equal to the number of bits per word. As an example, Figure 5.5

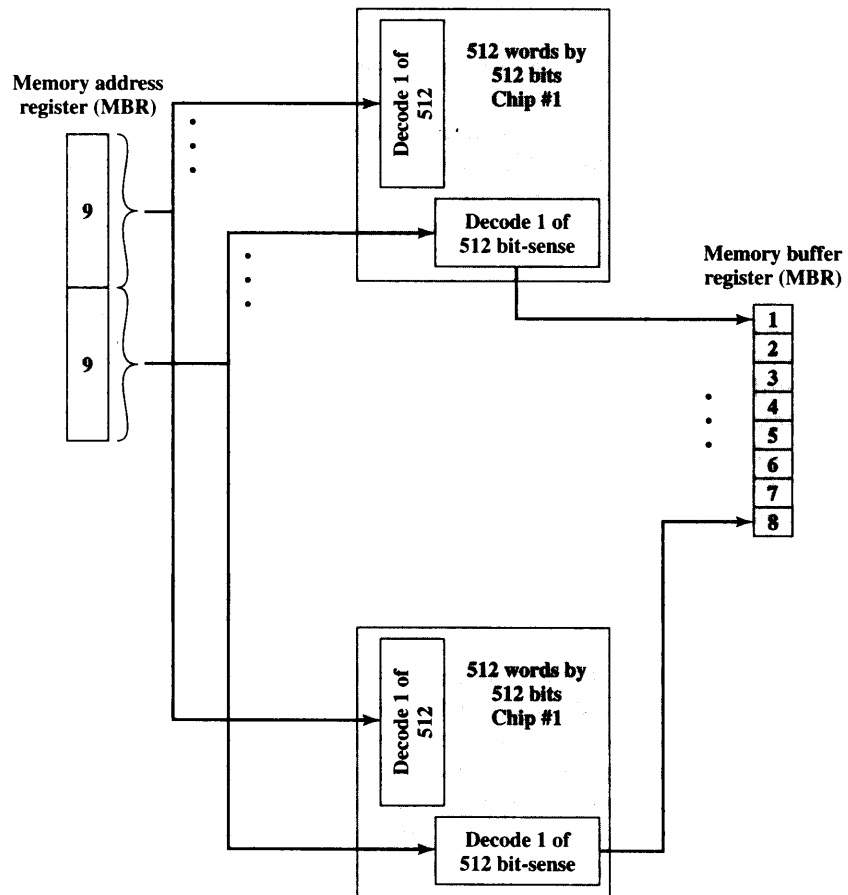


Figure 5.5 256-KByte Memory Organization

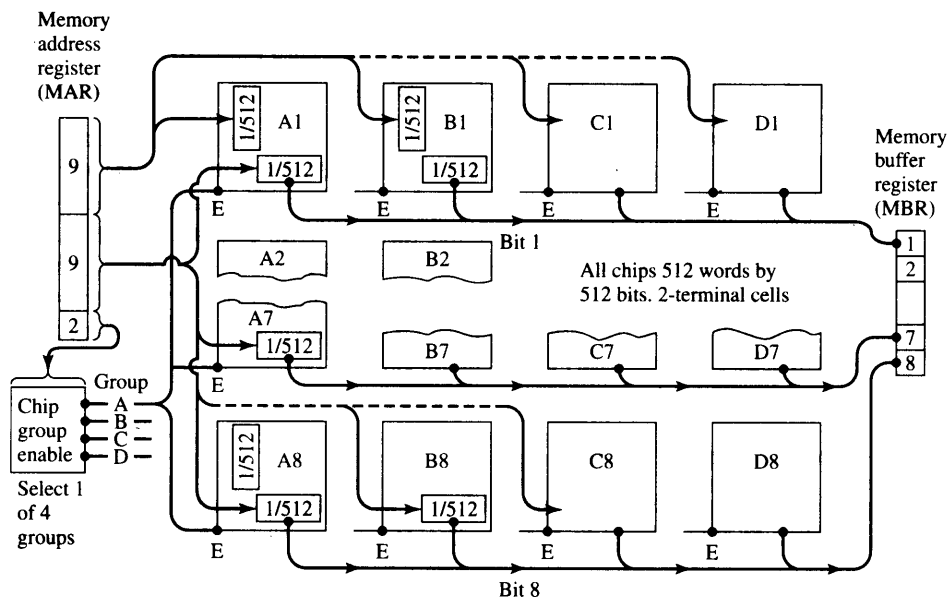


Figure 5.6 1-Mbyte Memory Organization

shows how a memory module consisting of 256K 8-bit words could be organized. For 256K words, an 18-bit address is needed and is supplied to the module from some external source (e.g., the address lines of a bus to which the module is attached). The address is presented to 8  $256K \times 1$ -bit chips, each of which provides the input/output of 1 bit.

This organization works as long as the size of memory equals the number of bits per chip. In the case in which larger memory is required, an array of chips is needed. Figure 5.6 shows the possible organization of a memory consisting of 1M word by 8 bits per word. In this case, we have four columns of chips, each column containing 256K words arranged as in Figure 5.5. For 1M word, 20 address lines are needed. The 18 least significant bits are routed to all 32 modules. The high-order 2 bits are input to a group select logic module that sends a chip enable signal to one of the four columns of modules.

## 5.2 ERROR CORRECTION

A semiconductor memory system is subject to errors. These can be categorized as hard failures and soft errors. A **hard failure** is a permanent physical defect so that the memory cell or cells affected cannot reliably store data, but become stuck at 0 or 1 or switch erratically between 0 and 1. Hard errors can be caused by harsh environmental abuse, manufacturing defects, and wear. A **soft error** is a random, nondestructive event that alters the contents of one or more memory cells, without damaging the memory. Soft errors can be caused by power supply problems or alpha particles. These particles result from radioactive decay and are distressingly common because radioactive nuclei are found in small quantities in nearly all

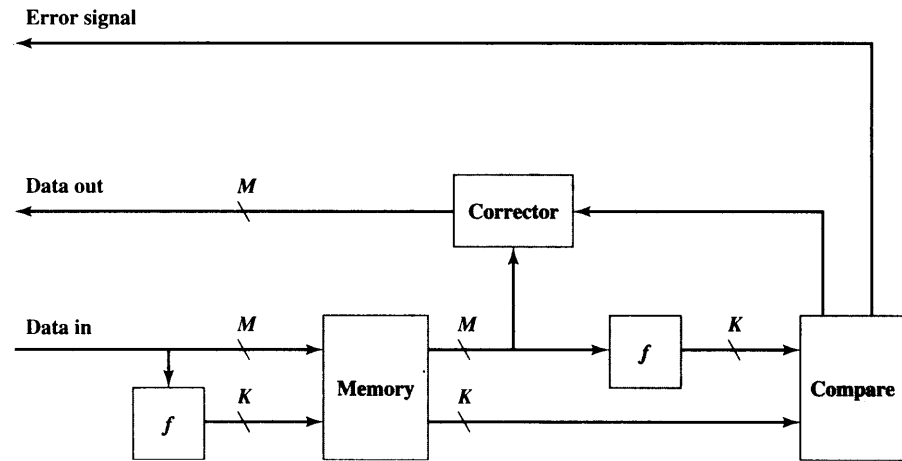


Figure 5.7 Error-Correcting Code Function

materials. Both hard and soft errors are clearly undesirable, and most modern main memory systems include logic for both detecting and correcting errors.

Figure 5.7 illustrates in general terms how the process is carried out. When data are to be read into memory, a calculation, depicted as a function  $f$ , is performed on the data to produce a code. Both the code and the data are stored. Thus, if an  $M$ -bit word of data is to be stored, and the code is of length  $K$  bits, then the actual size of the stored word is  $M + K$  bits.

When the previously stored word is read out, the code is used to detect and possibly correct errors. A new set of  $K$  code bits is generated from the  $M$  data bits and compared with the fetched code bits. The comparison yields one of three results:

- No errors are detected. The fetched data bits are sent out.
- An error is detected, and it is possible to correct the error. The data bits plus error correction bits are fed into a corrector, which produces a corrected set of  $M$  bits to be sent out.
- An error is detected, but it is not possible to correct it. This condition is reported.

Codes that operate in this fashion are referred to as *error-correcting codes*. A code is characterized by the number of bit errors in a word that it can correct and detect.

The simplest of the error-correcting codes is the *Hamming code* devised by Richard Hamming at Bell Laboratories. Figure 5.8 uses Venn diagrams to illustrate the use of this code on 4-bit words ( $M = 4$ ). With three intersecting circles, there are seven compartments. We assign the 4 data bits to the inner compartments (Figure 5.8a). The remaining compartments are filled with what are called *parity bits*. Each parity bit is chosen so that the total number of 1s in its circle is even (Figure 5.8b). Thus, because circle A includes three data 1s, the parity bit in that circle is set to 1. Now, if an error changes one of the data bits (Figure 5.8c), it is easily found. By checking the parity bits, discrepancies are found in circle A and circle C but not in circle B. Only one of the seven compartments is in A and C but not B. The error can therefore be corrected by changing that bit.

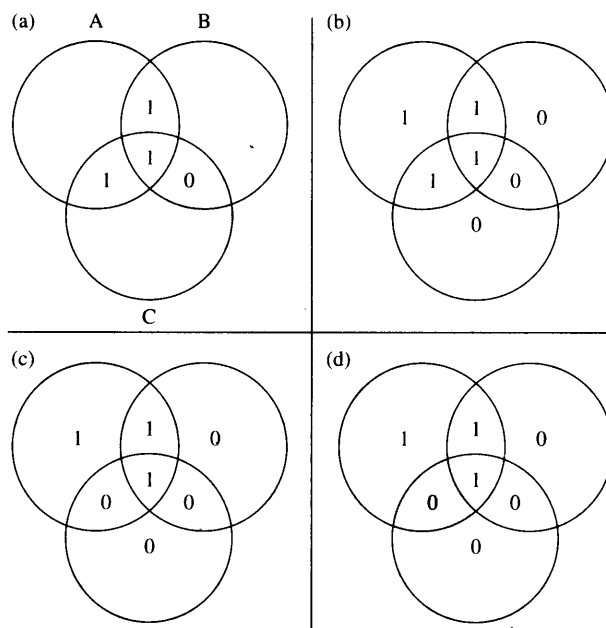


Figure 5.8 Hamming Error-Correcting Code

To clarify the concepts involved, we will develop a code that can detect and correct single-bit errors in 8-bit words.

To start, let us determine how long the code must be. Referring to Figure 5.7, the comparison logic receives as input two  $K$ -bit values. A bit-by-bit comparison is done by taking the exclusive-OR of the two inputs. The result is called the *syndrome word*. Thus, each bit of the syndrome is 0 or 1 according to if there is or is not a match in that bit position for the two inputs.

The syndrome word is therefore  $K$  bits wide and has a range between 0 and  $2^K - 1$ . The value 0 indicates that no error was detected, leaving  $2^K - 1$  values to indicate, if there is an error, which bit was in error. Now, because an error could occur on any of the  $M$  data bits or  $K$  check bits, we must have

$$2^K - 1 \geq M + K$$

This inequality gives the number of bits needed to correct a single bit error in a word containing  $M$  data bits. For example, for a word of 8 data bits ( $M = 8$ ), we have

- $K = 3: 2^3 - 1 < 8 + 3$
- $K = 4: 2^4 - 1 > 8 + 4$

Thus, eight data bits require four check bits. The first three columns of Table 5.2 lists the number of check bits required for various data word lengths.

For convenience, we would like to generate a 4-bit syndrome for an 8-bit data word with the following characteristics:

Table 5.2 Increase in Word Length with Error Correction

Data Bits	Single-Error Correction		Single-Error Correction/ Double-Error Detection	
	Check Bits	% Increase	Check Bits	% Increase
8	4	50	5	62.5
16	5	31.25	6	37.5
32	6	18.75	7	21.875
64	7	10.94	8	12.5
128	8	6.25	9	7.03
256	9	3.52	10	3.91

- If the syndrome contains all 0s, no error has been detected.
- If the syndrome contains one and only one bit set to 1, then an error has occurred in one of the 4 check bits. No correction is needed.
- If the syndrome contains more than one bit set to 1, then the numerical value of the syndrome indicates the position of the data bit in error. This data bit is inverted for correction.

To achieve these characteristics, the data and check bits are arranged into a 12-bit word as depicted in Figure 5.9. The bit positions are numbered from 1 to 12. Those bit positions whose position numbers are powers of 2 are designated as check bits. The check bits are calculated as follows, where the symbol  $\oplus$  designates the exclusive-OR operation:

$$\begin{aligned}
 C1 &= D1 \oplus D2 \oplus D4 \oplus D5 \oplus D7 \\
 C2 &= D1 \oplus D3 \oplus D4 \oplus D6 \oplus D7 \\
 C4 &= D2 \oplus D3 \oplus D4 \oplus D8 \\
 C8 &= D5 \oplus D6 \oplus D7 \oplus D8
 \end{aligned}$$

Each check bit operates on every data bit whose position number contains a 1 in the same bit position as the position number of that check bit. Thus, data bit positions 3, 5, 7, 9, and 11 ( $D1, D2, D4, D5, D7$ ) all contain a 1 in the least significant bit of their position number as does  $C1$ ; bit positions 3, 6, 7, 10, and 11 all contain a 1 in the second bit position, as does  $C2$ ; and so on. Looked at another way, bit position  $n$

Bit position	12	11	10	9	8	7	6	5	4	3	2	1
Position number	1100	1011	1010	1001	1000	0111	0110	0101	0100	0011	0010	0001
Data bit	D8	D7	D6	D5		D4	D3	D2		D1		
Check bit					C8				C4		C2	C1

Figure 5.9 Layout of Data Bits and Check Bits

is checked by those bits  $C_i$  such that  $\sum i = n$ . For example, position 7 is checked by bits in position 4, 2, and 1; and  $7 = 4 + 2 + 1$ .

Let us verify that this scheme works with an example. Assume that the 8-bit input word is 00111001, with data bit D1 in the rightmost position. The calculations are as follows:

$$\begin{aligned}
 C1 &= 1 \oplus 0 \oplus 1 \oplus 1 \oplus 0 = 1 \\
 C2 &= 1 \oplus 0 \oplus 1 \oplus 1 \oplus 0 = 1 \\
 C4 &= 0 \oplus 0 \oplus 1 \oplus 0 = 1 \\
 C8 &= 1 \oplus 1 \oplus 0 \oplus 0 = 0
 \end{aligned}$$

Suppose now that data bit 3 sustains an error and is changed from 0 to 1. When the check bits are recalculated, we have

$$\begin{aligned}
 C1 &= 1 \oplus 0 \oplus 1 \oplus 1 \oplus 0 = 1 \\
 C2 &= 1 \oplus 1 \oplus 1 \oplus 1 \oplus 0 = 0 \\
 C4 &= 0 \oplus 1 \oplus 1 \oplus 0 = 0 \\
 C8 &= 1 \oplus 1 \oplus 0 \oplus 0 = 0
 \end{aligned}$$

When the new check bits are compared with the old check bits, the syndrome word is formed:

$$\begin{array}{cccc}
 C8 & C4 & C2 & C1 \\
 0 & 1 & 1 & 1 \\
 \oplus & 0 & 0 & 0 \\
 \hline
 0 & 1 & 1 & 0
 \end{array}$$

The result is 0110, indicating that bit position 6, which contains data bit 3, is in error.

Figure 5.10 illustrates the preceding calculation. The data and check bits are positioned properly in the 12-bit word. Four of the data bits have a value 1 (shaded in the table), and their bit position values are XORed to produce the Hamming code 0111, which forms the four check digits. The entire block that is stored is 001101001111. Suppose now that data bit 3, in bit position 6, sustains an error and is

Bit position	12	11	10	9	8	7	6	5	4	3	2	1
Position number	1100	1011	1010	1001	1000	0111	0110	0101	0100	0011	0010	0001
Data bit	D8	D7	D6	D5		D4	D3	D2		D1		
Check bit					C8				C4		C2	C1
Word stored as	0	0	1	1	0	1	0	0	1	1	1	1
Word fetched as	0	0	1	1	0	1	1	0	1	1	1	1
Position number	1100	1011	1010	1001	1000	0111	0110	0101	0100	0011	0010	0001
Check bit					0				0		0	1

Figure 5.10 Check Bit Calculation



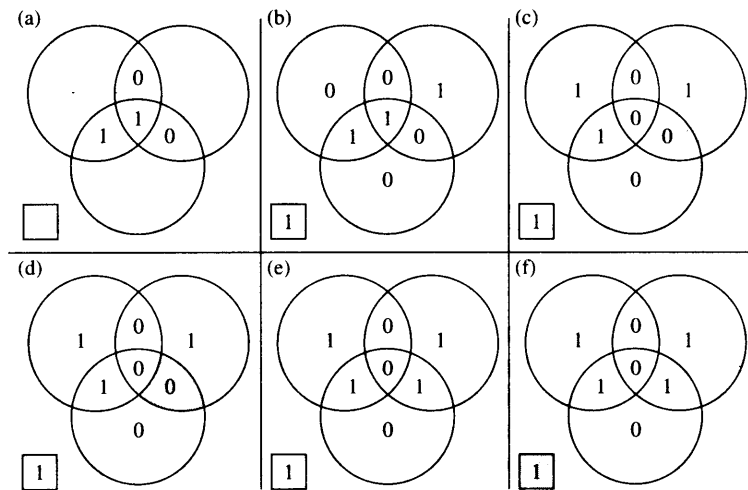


Figure 5.11 Hamming SEC-DEC Code

changed from 0 to 1. The resulting block is 001101101111, with a Hamming code of 0111. An XOR of the Hamming code and all of the bit position values for nonzero data bits results in 0110. The nonzero result detects an error and indicates that the error is in bit position 6.

The code just described is known as a *single-error-correcting* (SEC) code. More commonly, semiconductor memory is equipped with a single-error-correcting, double-error-detecting (SEC-DED) code. As Table 5.2 shows, such codes require one additional bit compared with SEC codes.

Figure 5.11 illustrates how such a code works, again with a 4-bit data word. The sequence shows that if two errors occur (Figure 5.11c), the checking procedure goes astray (d) and worsens the problem by creating a third error (e). To overcome the problem, an eighth bit is added that is set so that the total number of 1s in the diagram is even. The extra parity bit catches the error (f).

An error-correcting code enhances the reliability of the memory at the cost of added complexity. With a one-bit-per-chip organization, an SEC-DED code is generally considered adequate. For example, the IBM 30xx implementations used an 8-bit SEC-DED code for each 64 bits of data in main memory. Thus, the size of main memory is actually about 12% larger than is apparent to the user. The VAX computers used a 7-bit SEC-DED for each 32 bits of memory, for a 22% overhead. A number of contemporary DRAMs use 9 check bits for each 128 bits of data, for a 7% overhead [SHAR97].

## 5.3 ADVANCED DRAM ORGANIZATION

As was discussed in Chapter 2, one of the most critical system bottlenecks when using high-performance processors is the interface to main internal memory. This interface is the most important pathway in the entire computer system. The basic

Table 5.3 Performance Comparison of Some DRAM Alternatives

	<b>Clock frequency (MHz)</b>	<b>Transfer rate (GB/s)</b>	<b>Access time (ns)</b>	<b>Pin count</b>
<b>SDRAM</b>	166	1.3	18	168
<b>DDR</b>	200	3.2	12.5	184
<b>RDRAM</b>	600	4.8	12	162

building block of main memory remains the DRAM chip, as it has for decades; until recently, there had been no significant changes in DRAM architecture since the early 1970s. The traditional DRAM chip is constrained both by its internal architecture and by its interface to the processor's memory bus.

We have seen that one attack on the performance problem of DRAM main memory has been to insert one or more levels of high-speed SRAM cache between the DRAM main memory and the processor. But SRAM is much costlier than DRAM, and expanding cache size beyond a certain point yields diminishing returns.

In recent years, a number of enhancements to the basic DRAM architecture have been explored, and some of these are now on the market. The schemes that currently dominate the market are SDRAM, DDR-DRAM, and RDRAM. Table 5.3 provides a performance comparison. CDRAM has also received considerable attention. We examine each of these approaches in this section.

### Synchronous DRAM

One of the most widely used forms of DRAM is the synchronous DRAM (SDRAM) [VOGL94]. Unlike the traditional DRAM, which is asynchronous, the SDRAM exchanges data with the processor synchronized to an external clock signal and running at the full speed of the processor/memory bus without imposing wait states.

In a typical DRAM, the processor presents addresses and control levels to the memory, indicating that a set of data at a particular location in memory should be either read from or written into the DRAM. After a delay, the access time, the DRAM either writes or reads the data. During the access-time delay, the DRAM performs various internal functions, such as activating the high capacitance of the row and column lines, sensing the data, and routing the data out through the output buffers. The processor must simply wait through this delay, slowing system performance.

With synchronous access, the DRAM moves data in and out under control of the system clock. The processor or other master issues the instruction and address information, which is latched by the DRAM. The DRAM then responds after a set number of clock cycles. Meanwhile, the master can safely do other tasks while the SDRAM is processing the request.

Figure 5.12 shows the internal logic of IBM's 64-Mb SDRAM [IBM01], which is typical of SDRAM organization, and Table 5.4 defines the various pin assignments.

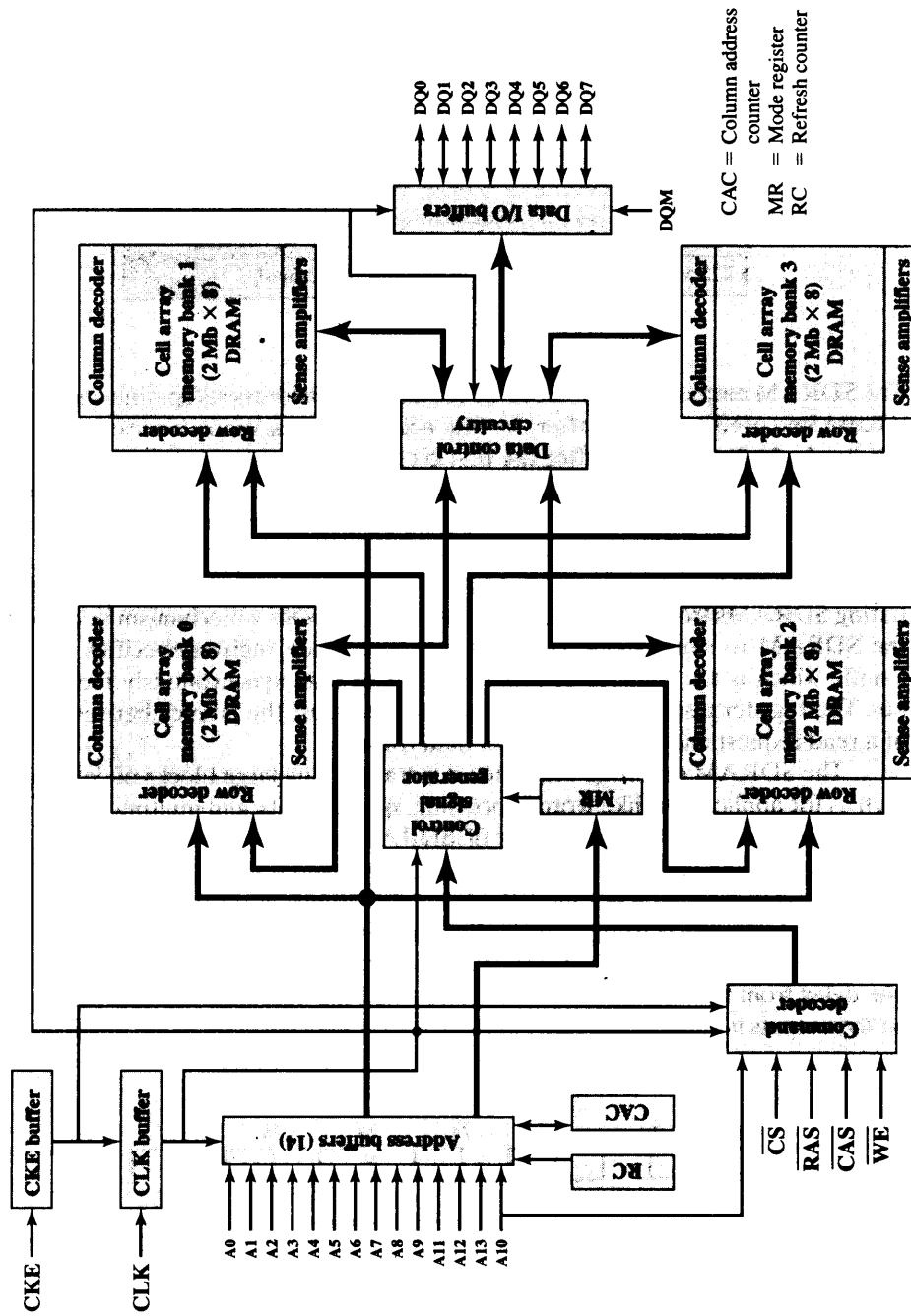


Figure 5.12 Synchronous Dynamic RAM (SDRAM)

Table 5.4 SDRAM Pin Assignments

<b>A0 to A13</b>	<b>Address inputs</b>
CLK	Clock input
CKE	Clock enable
$\overline{CS}$	Chip select
$\overline{RAS}$	Row address strobe
$\overline{CAS}$	Column address strobe
WE	Write enable
DQ0 to DQ7	Data input/output
DQM	Data mask

The SDRAM employs a burst mode to eliminate the address setup time and row and column line precharge time after the first access. In burst mode, a series of data bits can be clocked out rapidly after the first bit has been accessed. This mode is useful when all the bits to be accessed are in sequence and in the same row of the array as the initial access. In addition, the SDRAM has a multiple-bank internal architecture that improves opportunities for on-chip parallelism.

The mode register and associated control logic is another key feature differentiating SDRAMs from conventional DRAMs. It provides a mechanism to customize the SDRAM to suit specific system needs. The mode register specifies the burst length, which is the number of separate units of data synchronously fed onto the bus. The register also allows the programmer to adjust the latency between receipt of a read request and the beginning of data transfer.

The SDRAM performs best when it is transferring large blocks of data serially, such as for applications like word processing, spreadsheets, and multimedia.

Figure 5.13 shows an example of SDRAM operation. In this case, the burst length is 4 and the latency is 2. The burst read command is initiated by having  $\overline{CS}$  and  $\overline{CAS}$  low while holding  $\overline{RAS}$  and WE high at the rising edge of the clock. The address inputs determine the starting column address for the burst, and the mode register sets the type of burst (sequential or interleave) and the burst length (1, 2, 4, 8, full page). The delay from the start of the command to when the data from the first cell appears on the outputs is equal to the value of the  $\overline{CAS}$  latency that is set in the mode register.

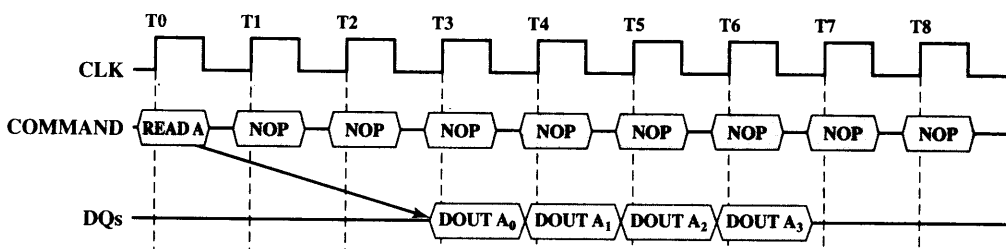


Figure 5.13 SDRAM Read Timing (burst length = 4,  $\overline{CAS}$  latency = 2)

There is now an enhanced version of SDRAM, known as double data rate SDRAM (DDR-SDRAM) that overcomes the once-per-cycle limitation. DDR-SDRAM can send data to the processor twice per clock cycle.

### Rambus DRAM

RDRAM, developed by Rambus [FARM92, CRIS97], has been adopted by Intel for its Pentium and Itanium processors. It has become the main competitor to SDRAM. RDRAM chips are vertical packages, with all pins on one side. The chip exchanges data with the processor over 28 wires no more than 12 centimeters long. The bus can address up to 320 RDRAM chips and is rated at 1.6 GBps.

The special RDRAM bus delivers address and control information using an asynchronous block-oriented protocol. After an initial 480 ns access time, this produces the 1.6 GBps data rate. What makes this speed possible is the bus itself, which defines impedances, clocking, and signals very precisely. Rather than being controlled by the explicit RAS, CAS, R/W, and CE signals used in conventional DRAMs, an RDRAM gets a memory request over the high-speed bus. This request contains the desired address, the type of operation, and the number of bytes in the operation.

Figure 5.14 illustrates the RDRAM layout. The configuration consists of a controller and a number of RDRAM modules connected together via a common bus. The controller is at one end of the configuration, and the far end of the bus is a parallel termination of the bus lines. The bus includes 18 data lines (16 actual data, two parity) cycling at twice the clock rate; that is, one bit is sent at the leading and following edge of each clock signal. This results in a signal rate on each data line of 800 Mbps. There is a separate set of 8 lines (RC) used for address and control signals. There is also a clock signal that starts at the far end from the controller propagates to the controller end and then loops back. A RDRAM module sends data to the controller synchronously to the clock to master, and the controller sends data to an RDRAM synchronously with the clock signal in the opposite direction. The remaining bus lines include a reference voltage, ground, and power source.

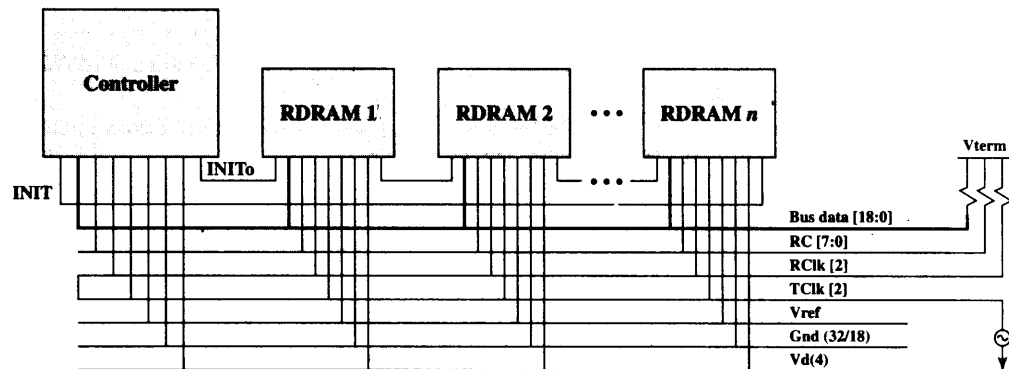


Figure 5.14 RDRAM Structure

### DDR SDRAM

SDRAM is limited by the fact that it can only send data to the processor once per bus clock cycle. A new version of SDRAM, referred to as double-data-rate SDRAM can send data twice per clock cycle, once on the rising edge of the clock pulse and once on the falling edge.

### Cache DRAM

Cache DRAM (CDRAM), developed by Mitsubishi [HIDA90, ZHAN01], integrates a small SRAM cache (16 Kb) onto a generic DRAM chip.

The SRAM on the CDRAM can be used in two ways. First, it can be used as a true cache, consisting of a number of 64-bit lines. The cache mode of the CDRAM is effective for ordinary random access to memory.

The SRAM on the CDRAM can also be used as a buffer to support the serial access of a block of data. For example, to refresh a bit-mapped screen, the CDRAM can prefetch the data from the DRAM into the SRAM buffer. Subsequent accesses to the chip result in accesses solely to the SRAM.

## 5.4 RECOMMENDED READING AND WEB SITES

[PRIN97] provides a comprehensive treatment of semiconductor memory technologies, including SRAM, DRAM, and flash memories. [SHAR97] covers the same material, with more emphasis on testing and reliability issues. [SHAR03] and [PRIN02] focus on advanced DRAM and SRAM architectures. For an in-depth look at DRAM, see [KEET01]. [CUPP01] provides an interesting performance comparison of various DRAM schemes. [BEZ03] is a comprehensive introduction to flash memory technology.

A good explanation of error-correcting codes is contained in [MCEL85]. For a deeper study, worthwhile book-length treatments are [ADAM91] and [BLAH83]. A quite readable theoretical and mathematical treatment of error-correcting codes is [ASH90]. [SHAR97] contains a good survey of codes used in contemporary main memories.

- ADAM91** Adamek, J. *Foundations of Coding*. New York: Wiley, 1991.
- ASH90** Ash, R. *Information Theory*. New York: Dover, 1990.
- BEZ03** Bez, R.; et al. Introduction to Flash Memory. *Proceedings of the IEEE*, April 2003.
- BLAH83** Blahut, R. *Theory and Practice of Error Control Codes*. Reading, MA: Addison-Wesley, 1983.
- CUPP01** Cuppu, V., et al. "High Performance DRAMS in Workstation Environments." *IEEE Transactions on Computers*, November 2001.
- KEET01** Keeth, B., and Baker, R. *DRAM Circuit Design: A Tutorial*. Piscataway, NJ: IEEE Press, 2001.
- MCEL85** McEliece, R. "The Reliability of Computer Memories." *Scientific American*, January 1985.
- PRIN97** Prince, B. *Semiconductor Memories*. New York: Wiley, 1997.
- PRIN02** Prince, B. *Emerging Memories: Technologies and Trends*. Norwell, MA: Kluwer, 2002.

**SHAR97** Sharma, A. *Semiconductor Memories: Technology, Testing, and Reliability*. New York: IEEE Press, 1997.

**SHAR03** Sharma, A. *Advanced Semiconductor Memories: Architectures, Designs, and Applications*. New York: IEEE Press, 2003.



### Recommended Web Sites:

- **The RAM Guide:** Good overview of RAM technology plus a number of useful links
- **RDRAM:** Another useful site for RDRAM information

## 5.5 KEY TERMS, REVIEW QUESTIONS, AND PROBLEMS

### Key Terms

cache DRAM (CDRAM) dynamic RAM (DRAM) electrically erasable program- mable ROM (EEPROM) erasable programmable ROM (EPROM) error-correcting code (ECC) error correction flash memory	Hamming code hard failure nonvolatile memory programmable ROM (PROM) RamBus DRAM (RDRAM) read-mostly memory read-only memory (ROM) semiconductor memory single-error-correcting (SEC) code	single-error-correcting, double-error-detecting (SEC-DED) code soft error static RAM (SRAM) synchronous DRAM (SDRAM) syndrome volatile memory
--	---	---

### Review Questions

- 5.1 What are the key properties of semiconductor memory?
- 5.2 What are two senses in which the term *random-access memory* is used?
- 5.3 What is the difference between DRAM and SRAM in terms of application?
- 5.4 What is the difference between DRAM and SRAM in terms of characteristics such as speed, size, and cost?
- 5.5 Explain why one type of RAM is considered to be analog and the other digital.
- 5.6 What are some applications for ROM?
- 5.7 What are the differences among EPROM, EEPROM, and flash memory?
- 5.8 Explain the function of each pin in Figure 5.4b.
- 5.9 What is a parity bit?
- 5.10 How is the syndrome for the Hamming code interpreted?
- 5.11 How does SDRAM differ from ordinary DRAM?

**Problems**

- 5.1 Suggest reasons why RAMs traditionally have been organized as only one bit per chip whereas ROMs are usually organized with multiple bits per chip.
- 5.2 Consider a dynamic RAM that must be given a refresh cycle 64 times per ms. Each refresh operation requires 150 ns; a memory cycle requires 250 ns. What percentage of the memory's total operating time must be given to refreshes?
- 5.3 Figure 5.15 shows a simplified timing diagram for a DRAM read operation over a bus. The access time is considered to last from  $t_1$  to  $t_2$ . Then there is a recharge time, lasting from  $t_2$  to  $t_3$ , during which the DRAM chips will have to recharge before the processor can access them again.
  - a. Assume that the access time is 60 ns and the recharge time is 40 ns. What is the memory cycle time? What is the maximum data rate this DRAM can sustain, assuming a 1-bit output?
  - b. Constructing a 32-bit wide memory system using these chips yields what data transfer rate?
- 5.4 Figure 5.6 indicates how to construct a module of chips that can store 1 Mbyte based on a group of four 256-Kbyte chips. Let's say this module of chips is packaged as a single 1-Mbyte chip, where the word size is 1 byte. Give a high-level chip diagram of how to construct a 8-Mbyte computer memory using eight 1-Mbyte chips. Be sure to show the address lines in your diagram and what the address lines are used for.
- 5.5 On a typical Intel 8086-based system, connected via system bus to DRAM memory, for a read operation,  $\overline{\text{RAS}}$  is activated by the trailing edge of the Address Enable signal (Figure 3.19). However, due to propagation and other delays,  $\overline{\text{RAS}}$  does not go active until 50 ns after Address Enable returns to a low. Assume the latter occurs in the middle of the second half of state  $T_2$  (somewhat earlier than in Figure 3.19). Data are read by the processor at the end of  $T_3$ . For timely presentation to the processor, however, data must be provided 60 ns earlier by memory. This interval accounts for propagation delays along the data paths (from memory to processor) and processor data hold time requirements. Assume a clocking rate of 10 MHz.
  - a. How fast (access time) should the DRAMs be if no wait states are to be inserted?
  - b. How many wait states do we have to insert per memory read operation if the access time of the DRAMs is 150 ns?

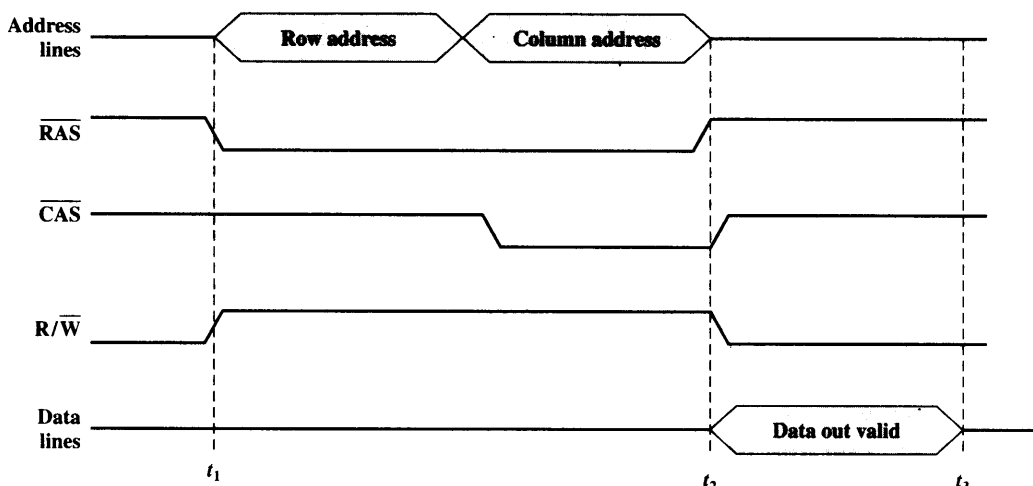
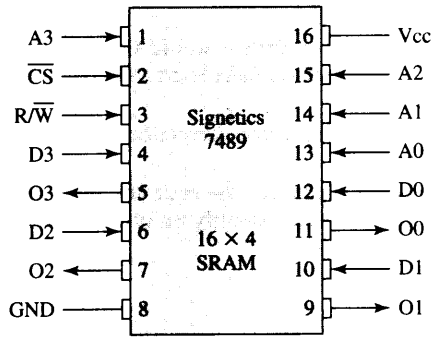


Figure 5.15 Simplified DRAM Read Timing



- 5.6 The memory of a particular microcomputer is built from  $64K \times 1$  DRAMs. According to the data sheet, the cell array of the DRAM is organized into 256 rows. Each row must be refreshed at least once every 4 ms. Suppose we refresh the memory on a strictly periodic basis.
- What is the time period between successive refresh requests?
  - How long a refresh address counter do we need?
- 5.7 Figure 5.16 shows one of the early SRAMs, the  $16 \times 4$  Signetics 7489 chip, which stores 16 4-bit words.
- List the mode of operation of the chip for each  $\overline{CS}$  input pulse shown in Figure 5.16c.
  - List the memory contents of word locations 0 through 6 after pulse n.
  - What is the state of the output data leads for the input pulses h through m?

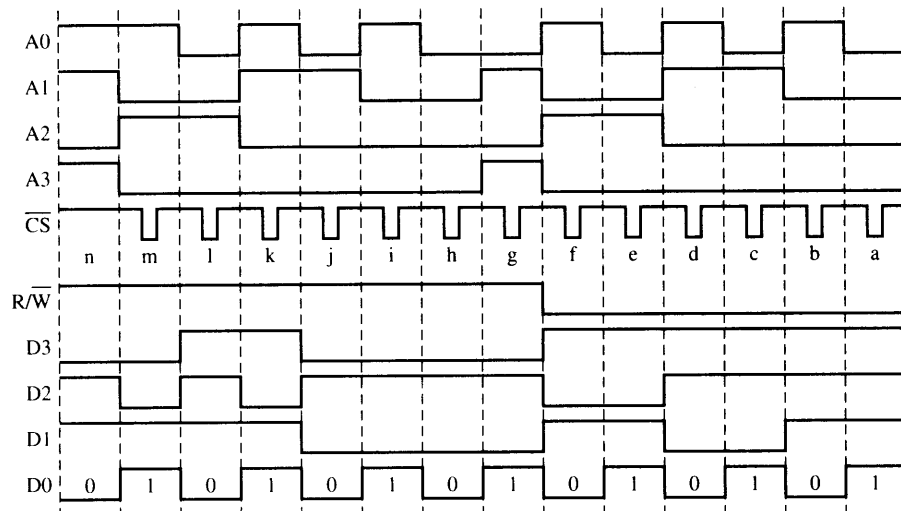


(a) Pin layout

Operating Mode	Inputs			Outputs
	$\overline{CS}$	$R/\overline{W}$	$D_n$	$O_n$
Write	L	L	L	L
	L	L	H	H
Read	L	H	X	Data
Inhibit writing	H	L	L	H
	H	L	H	L
Store - disable outputs	H	H	X	H

H = high voltage level  
 L = low voltage level  
 X = don't care

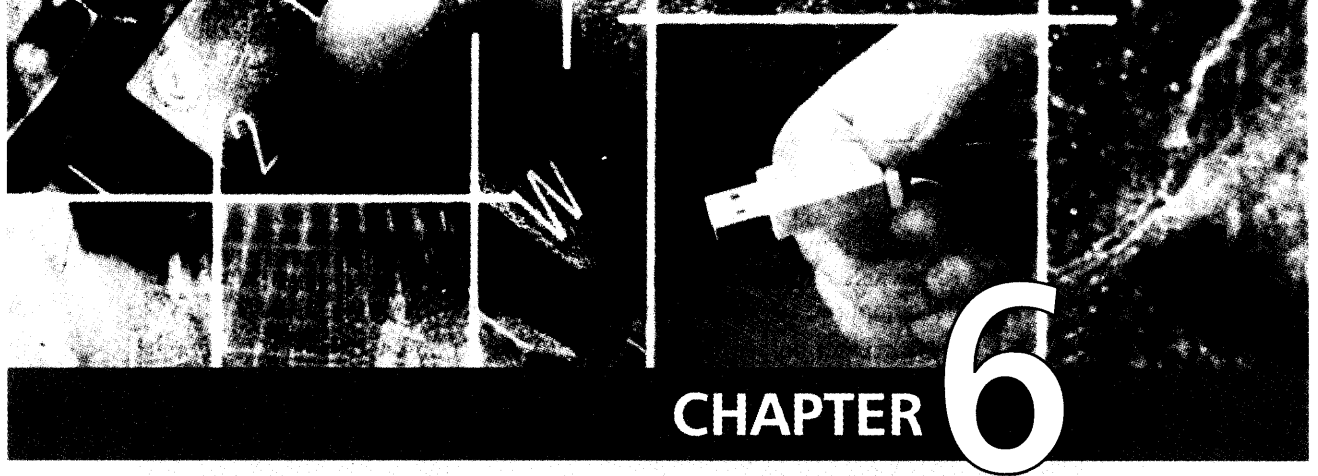
(b) Truth table



(c) Pulse train

Figure 5.16 The Signetics 7489 SRAM

- 5.8 Design a 16-bit memory of total capacity 8192 bits using SRAM chips of size  $64 \times 1$  bit. Give the array configuration of the chips on the memory board showing all required input and output signals for assigning this memory to the lowest address space. The design should allow for both byte and 16-bit word accesses.
- 5.9 A common unit of measure for failure rates of electronic components is the **Failure unit** (FIT), expressed as a rate of failures per billion device hours. Another well known but less used measure is **mean time between failures** (MTBF), which is the average time of operation of a particular component until it fails. Consider a 1 MB memory of a 16-bit microprocessor with  $256K \times 1$  DRAMs. Calculate its MTBF assuming 2000 FITS for each DRAM.
- 5.10 For the Hamming code shown in Figure 5.10, show what happens when a check bit rather than a data bit is in error?
- 5.11 Suppose an 8-bit data word stored in memory is 11000010. Using the Hamming algorithm, determine what check bits would be stored in memory with the data word. Show how you got your answer.
- 5.12 For the 8-bit word 00111001, the check bits stored with it would be 0111. Suppose when the word is read from memory, the check bits are calculated to be 1101. What is the data word that was read from memory?
- 5.13 How many check bits are needed if the Hamming error correction code is used to detect single bit errors in a 1024-bit data word?
- 5.14 Develop an SEC code for a 16-bit data word. Generate the code for the data word 0101000000111001. Show that the code will correctly identify an error in data bit 5.



# EXTERNAL MEMORY

## 6.1 Magnetic Disk

- Magnetic Read and Write Mechanisms
- Data Organization and Formatting
- Physical Characteristics
- Disk Performance Parameters

## 6.2 RAID

- RAID Level 0
- RAID Level 1
- RAID Level 2
- RAID Level 3
- RAID Level 4
- RAID Level 5
- RAID Level 6

## 6.3 Optical Memory

- Compact Disk
- Digital Versatile Disk

## 6.4 Magnetic Tape

## 6.5 Recommended Reading and Web Sites

## 6.6 Key Terms, Review Questions, and Problems

- Key Terms
- Review Questions
- Problems

---

### KEY POINTS

- ◆ **Magnetic disks remain the most important component of external memory. Both removable and fixed, or hard, disks are used in systems ranging from personal computers to mainframes and supercomputers.**
  - ◆ **To achieve greater performance and higher availability, servers and larger systems use RAID disk technology. RAID is a family of techniques for using multiple disks as a parallel array of data storage devices, with redundancy built in to compensate for disk failure.**
  - ◆ **Optical storage technology has become increasingly important in all types of computer systems. While CD-ROM has been widely used for many years, more recent technologies, such as writable CD and DVD, are becoming increasingly important.**
- 

This chapter examines a range of external memory devices and systems. We begin with the most important device, the magnetic disk. Magnetic disks are the foundation of external memory on virtually all computer systems. The next section examines the use of disk arrays to achieve greater performance, looking specifically at the family of systems known as RAID (Redundant Array of Independent Disks). An increasingly important component of many computer systems is external optical memory, and this is examined in the third section. Finally, magnetic tape is described.

## 6.1 MAGNETIC DISK

A disk is a circular platter constructed of nonmagnetic material, called the substrate, coated with a magnetizable material. Traditionally, the substrate has been an aluminum or aluminum alloy material. More recently, glass substrates have been introduced. The glass substrate has a number of benefits, including

- Improvement in the uniformity of the magnetic film surface to increase disk reliability
- A significant reduction in overall surface defects to help reduce read-write errors
- Ability to support lower fly heights (described subsequently)
- Better stiffness to reduce disk dynamics
- Greater ability to withstand shock and damage

### Magnetic Read and Write Mechanisms

Data are recorded on and later retrieved from the disk via a conducting coil named the **head**; in many systems, there are two heads, a read head and a write head.

During a read or write operation, the head is stationary while the platter rotates beneath it.

The write mechanism exploits the fact that electricity flowing through a coil produces a magnetic field. Electric pulses are sent to the write head, and the resulting magnetic patterns are recorded on the surface below, with different patterns for positive and negative currents. The write head itself is made of easily magnetizable material and is in the shape of a rectangular doughnut with a gap along one side and a few turns of conducting wire along the opposite side (Figure 6.1). An electric current in the wire induces a magnetic field across the gap, which in turn magnetizes a small area of the recording medium. Reversing the direction of the current reverses the direction of the magnetization on the recording medium.

The traditional read mechanism exploits the fact that a magnetic field moving relative to a coil produces an electrical current in the coil. When the surface of the disk passes under the head, it generates a current of the same polarity as the one already recorded. The structure of the head for reading is in this case essentially the same as for writing and therefore the same head can be used for both. Such single heads are used in floppy disk systems and in older rigid disk systems.

Contemporary rigid disk systems use a different read mechanism, requiring a separate read head, positioned for convenience close to the write head. The read head consists of a partially shielded magnetoresistive (MR) sensor. The MR material has an electrical resistance that depends on the direction of the magnetization of the medium moving under it. By passing a current through the MR sensor, resistance changes are detected as voltage signals. The MR design allows higher-frequency operation, which equates to greater storage densities and operating speeds.

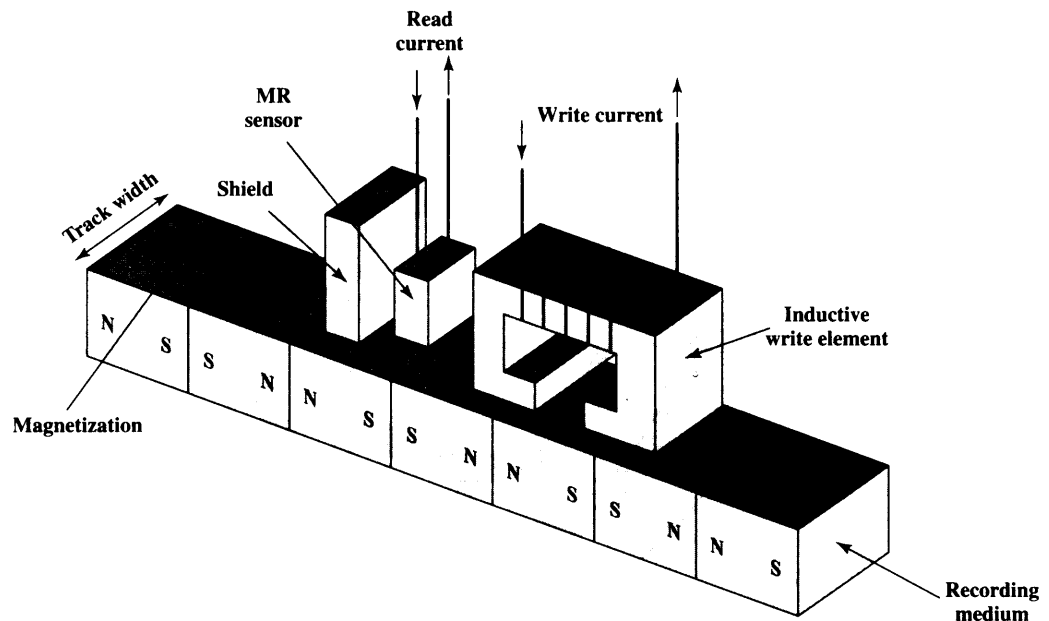


Figure 6.1 Inductive Write/Magnetoresistive Read Head

### Data Organization and Formatting

The head is a relatively small device capable of reading from or writing to a portion of the platter rotating beneath it. This gives rise to the organization of data on the platter in a concentric set of rings, called **tracks**. Each track is the same width as the head. There are thousands of tracks per surface.

Figure 6.2 depicts this data layout. Adjacent tracks are separated by **gaps**. This prevents, or at least minimizes, errors due to misalignment of the head or simply interference of magnetic fields.

Data are transferred to and from the disk in **sectors** (Figure 6.2). There are typically hundreds of sectors per track, and these may be of either fixed or variable length. In most contemporary systems, fixed-length sectors are used, with 512 bytes being the nearly universal sector size. To avoid imposing unreasonable precision requirements on the system, adjacent sectors are separated by intratrack (intersector) gaps.

A bit near the center of a rotating disk travels past a fixed point (such as a read-write head) slower than a bit on the outside. Therefore, some way must be found to compensate for the variation in speed so that the head can read all the bits at the same rate. This can be done by increasing the spacing between bits of information recorded in segments of the disk. The information can then be scanned at the same rate by rotating the disk at a fixed speed, known as the **constant angular velocity (CAV)**. Figure 6.3a shows the layout of a disk using CAV. The disk is

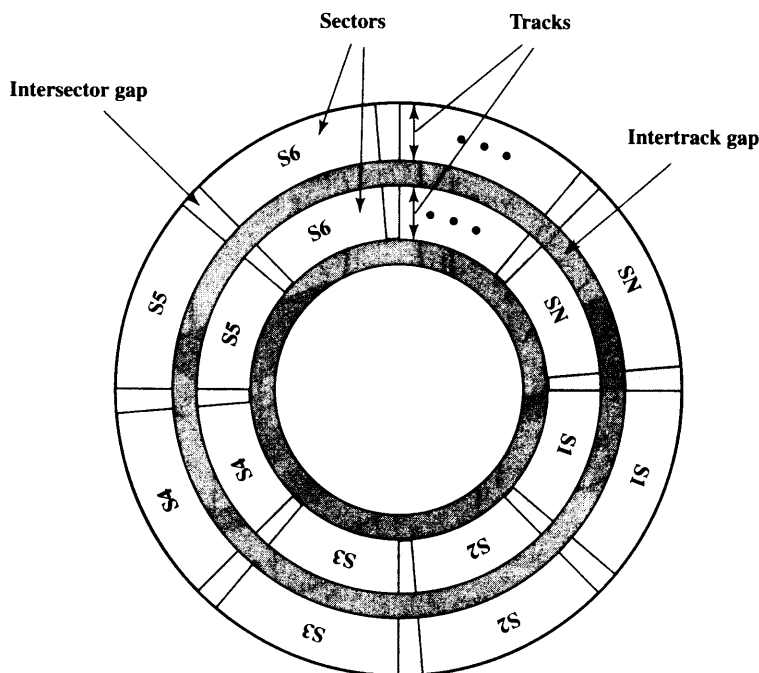


Figure 6.2 Disk Data Layout

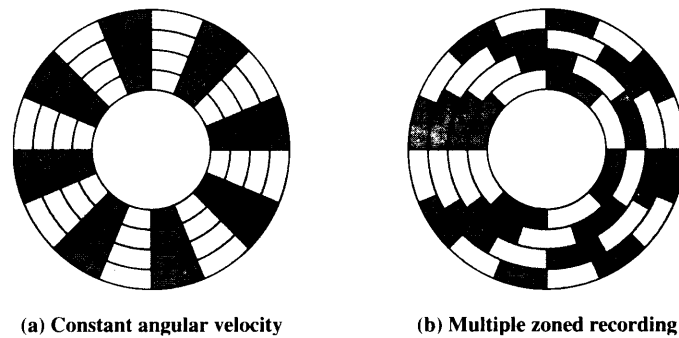


Figure 6.3 Comparison of Disk Layout Methods

divided into a number of pie-shaped sectors and into a series of concentric tracks. The advantage of using CAV is that individual blocks of data can be directly addressed by track and sector. To move the head from its current location to a specific address, it only takes a short movement of the head to a specific track and a short wait for the proper sector to spin under the head. The disadvantage of CAV is that the amount of data that can be stored on the long outer tracks is the only same as what can be stored on the short inner tracks.

Because the **density**, in bits per linear inch, increases in moving from the outermost track to the innermost track, disk storage capacity in a straightforward CAV system is limited by the maximum recording density that can be achieved on the innermost track. To increase density, modern hard disk systems use a technique known as **multiple zone recording**, in which the surface is divided into a number of concentric zones (16 is typical). Within a zone, the number of bits per track is constant. Zones farther from the center contain more bits (more sectors) than zones closer to the center. This allows for greater overall storage capacity at the expense of somewhat more complex circuitry. As the disk head moves from one zone to another, the length (along the track) of individual bits changes, causing a change in the timing for reads and writes. Figure 6.3b suggests the nature of multiple zone recording; in this illustration, each zone is only a single track wide.

Some means is needed to locate sector positions within a track. Clearly, there must be some starting point on the track and a way of identifying the start and end of each sector. These requirements are handled by means of control data recorded on the disk. Thus, the disk is formatted with some extra data used only by the disk drive and not accessible to the user.

An example of disk formatting is shown in Figure 6.4. In this case, each track contains 30 fixed-length sectors of 600 bytes each. Each sector holds 512 bytes of data plus control information useful to the disk controller. The ID field is a unique identifier or address used to locate a particular sector. The SYNCH byte is a special bit pattern that delimits the beginning of the field. The track number identifies a track on a surface. The head number identifies a head, because this disk has multiple surfaces (explained presently). The ID and data fields each contain an error-detecting code.

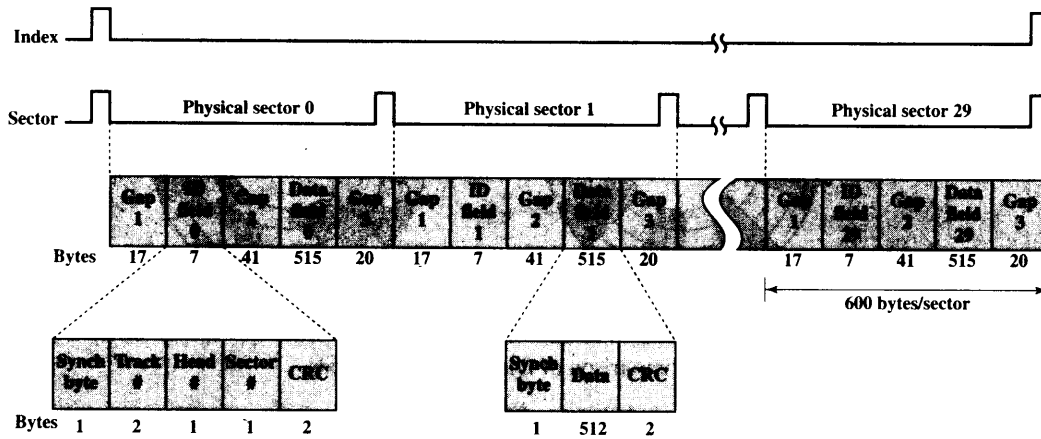


Figure 6.4 Winchester Disk Format (Seagate ST506)

### Physical Characteristics

Table 6.1 lists the major characteristics that differentiate among the various types of magnetic disks. First, the head may either be fixed or movable with respect to the radial direction of the platter. In a **fixed-head disk**, there is one read-write head per track. All of the heads are mounted on a rigid arm that extends across all tracks; such systems are rare today. In a **movable-head disk**, there is only one read-write head. Again, the head is mounted on an arm. Because the head must be able to be positioned above any track, the arm can be extended or retracted for this purpose.

The disk itself is mounted in a disk drive, which consists of the arm, a spindle that rotates the disk, and the electronics needed for input and output of binary data. A **nonremovable disk** is permanently mounted in the disk drive; the hard disk in a personal computer is a nonremovable disk. A **removable disk** can be removed and replaced with another disk. The advantage of the latter type is that unlimited amounts of data are available with a limited number of disk systems. Furthermore,

Table 6.1 Physical Characteristics of Disk Systems

Head location	Platters
Fixed head (one per track)	Single platter
Movable head (one per surface)	Multiple platters
Disk Portability	Head Mechanism
Nonremovable disk	Contact (floppy)
Removable disk	Fixed gap
Single sided	Aerodynamic gap (Winchester)
Double sided	



such a disk may be moved from one computer system to another. Floppy disks and ZIP cartridge disks are examples of removable disks.

For most disks, the magnetizable coating is applied to both sides of the platter, which is then referred to as **double sided**. Some less expensive disk systems use **single-sided** disks.

Some disk drives accommodate **multiple platters** stacked vertically a fraction of an inch apart. Multiple arms are provided (Figure 6.5). Multiple-platter disks employ a movable head, with one read-write head per platter surface. All of the heads are mechanically fixed so that all are at the same distance from the center of the disk and move together. Thus, at any time, all of the heads are positioned over tracks that are of equal distance from the center of the disk. The set of all the tracks in the same relative position on the platter is referred to as a **cylinder**. For example, all of the shaded tracks in Figure 6.6 are part of one cylinder.

Finally, the head mechanism provides a classification of disks into three types. Traditionally, the read-write head has been positioned a fixed distance above the platter, allowing an air gap. At the other extreme is a head mechanism that actually comes into physical contact with the medium during a read or write operation. This mechanism is used with the **floppy disk**, which is a small, flexible platter and the least expensive type of disk.

To understand the third type of disk, we need to comment on the relationship between data density and the size of the air gap. The head must generate or sense an electromagnetic field of sufficient magnitude to write and read properly. The narrower the head is, the closer it must be to the platter surface to function. A narrower head means narrower tracks and therefore greater data density, which is desirable. However, the closer the head is to the disk, the greater the risk of error from impurities or imperfections. To push the technology further, the Winchester disk was

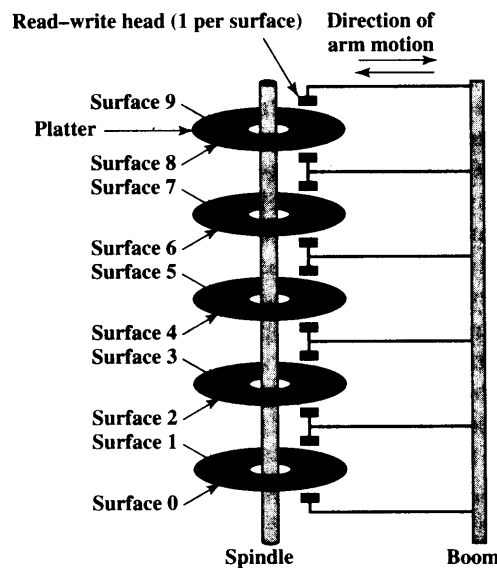


Figure 6.5 Components of a Disk Drive

is completed, the device determines when the data will rotate under the head. As that sector approaches the head, the device tries to reestablish the communication path back to the host. If either the control unit or the channel is busy with another I/O, then the reconnection attempt fails and the device must rotate one whole revolution before it can attempt to reconnect, which is called an RPS miss. This is an extra delay element that must be added to the timeline of Figure 6.7.

**Seek Time** Seek time is the time required to move the disk arm to the required track. It turns out that this is a difficult quantity to pin down. The seek time consists of two key components: the initial startup time, and the time taken to traverse the tracks that have to be crossed once the access arm is up to speed. Unfortunately, the traversal time is not a linear function of the number of tracks but includes a settling time (time after positioning the head over the target track until track identification is confirmed).

Much improvement comes from smaller and lighter disk components. Some years ago, a typical disk was 14 inches (36 cm) in diameter, whereas the most common size today is 3.5 inches (8.9 cm), reducing the distance that the arm has to travel. A typical average seek time on contemporary hard disks is under 10 ms.

**Rotational Delay** Disks, other than floppy disks, rotate at speeds ranging from 3600 rpm (for handheld devices such as digital cameras) up to, as of this writing, 15,000 rpm; at this latter speed, there is one revolution per 4 ms. Thus, on the average, the rotational delay will be 2 ms. Floppy disks typically rotate at between 300 and 600 rpm. Thus the average delay will be between 100 and 50 ms.

**Transfer Time** The transfer time to or from the disk depends on the rotation speed of the disk in the following fashion:

$$T = \frac{b}{rN}$$

where

$T$  = transfer time

$b$  = number of bytes to be transferred

$N$  = number of bytes on a track

$r$  = rotation speed, in revolutions per second

Thus the total average access time can be expressed as

$$T_a = T_s + \frac{1}{2r} + \frac{b}{rN}$$

where  $T_s$  is the average seek time. Note that on a zoned drive, the number of bytes per track is variable, complicating the calculation.

**A Timing Comparison** With the foregoing parameters defined, let us look at two different I/O operations that illustrate the danger of relying on average values. Consider a disk with an advertised average seek time of 4 ms, rotation speed of 15,000 rpm, and 512-byte sectors with 500 sectors per track. Suppose that we wish to read a file consisting of 2500 sectors for a total of 1.28 Mbytes. We would like to estimate the total time for the transfer.

First, let us assume that the file is stored as compactly as possible on the disk. That is, the file occupies all of the sectors on 5 adjacent tracks (5 tracks  $\times$  500 sectors/track = 2500 sectors): This is known as *sequential organization*. Now, the time to read the first track is as follows:

Average seek	4 ms
Average rotational delay	2 ms
Read 500 sectors	<u>4 ms</u>
	10 ms

Suppose that the remaining tracks can now be read with essentially no seek time. That is, the I/O operation can keep up with the flow from the disk. Then, at most, we need to deal with rotational delay for each succeeding track. Thus each successive track is read in  $2 + 4 = 6$  ms. To read the entire file,

$$\text{Total time} = 10 + (4 \times 6) = 34 \text{ ms} = 0.034 \text{ seconds}$$

Now let us calculate the time required to read the same data using random access rather than sequential access; that is, accesses to the sectors are distributed randomly over the disk. For each sector, we have

Average seek	4 ms
Rotational delay	2 ms
Read 1 sector	<u>0.008 ms</u>
	6.008 ms

$$\text{Total time} = 2500 \times 6.008 = 15020 \text{ ms} = 15.02 \text{ seconds}$$

It is clear that the order in which sectors are read from the disk has a tremendous effect on I/O performance. In the case of file access in which multiple sectors are read or written, we have some control over the way in which sectors of data are deployed. However, even in the case of a file access, in a multiprogramming environment, there will be I/O requests competing for the same disk. Thus, it is worthwhile to examine ways in which the performance of disk I/O can be improved over that achieved with purely random access to the disk. This leads to a consideration of disk scheduling algorithms, which is the province of the operating system and beyond the scope of this book (see [STAL05] for a discussion).

## 6.2 RAID

As discussed earlier, the rate in improvement in secondary storage performance has been considerably less than the rate for processors and main memory. This mismatch has made the disk storage system perhaps the main focus of concern in improving overall computer system performance.

As in other areas of computer performance, disk storage designers recognize that if one component can only be pushed so far, additional gains in performance are to be had by using multiple parallel components. In the case of disk storage, this leads to the development of arrays of disks that operate independently and in

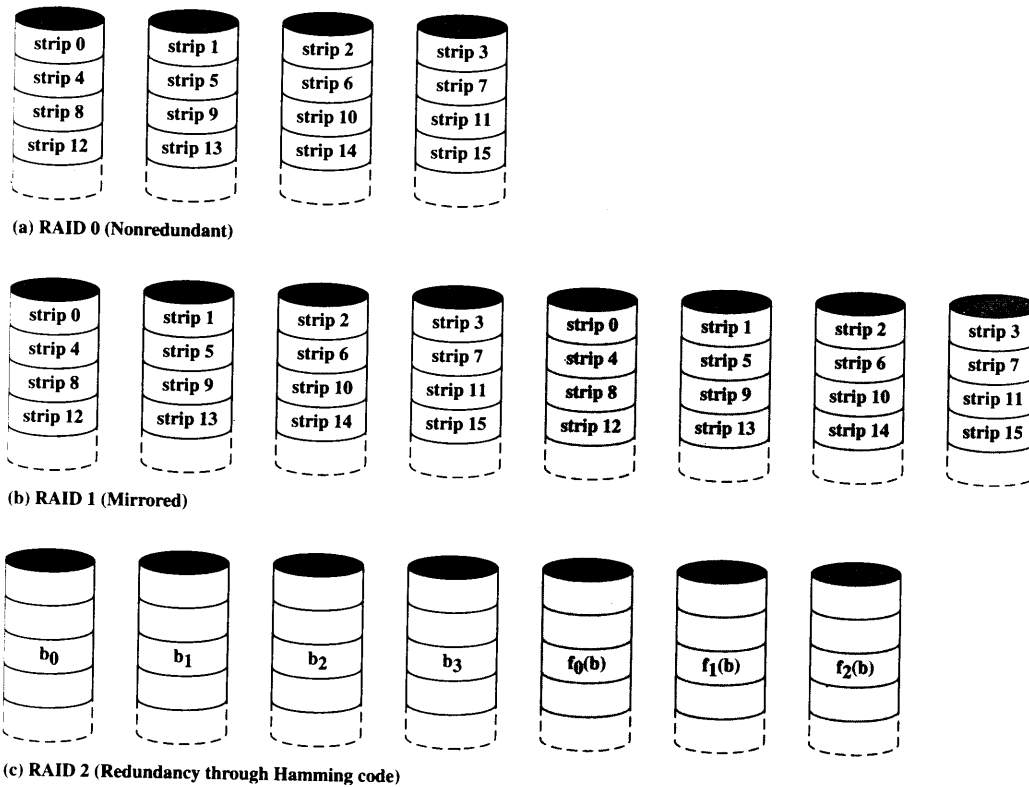


Figure 6.8 RAID Levels 0 through 3

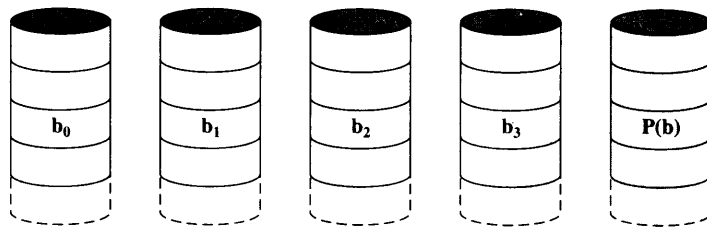
highlighted by shading. Figures 6.8 and 6.9 illustrate the use of the seven RAID schemes to support a data capacity requiring four disks with no redundancy. The figures highlight the layout of user data and redundant data and indicates the relative storage requirements of the various levels. We refer to these figures throughout the following discussion.

### RAID Level 0

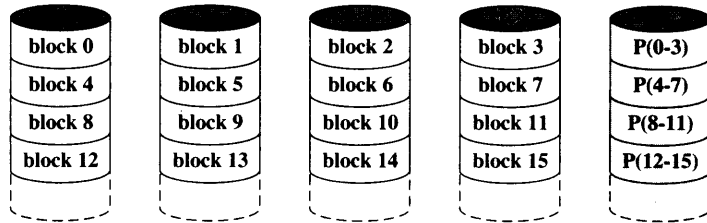
RAID level 0 is not a true member of the RAID family, because it does not include redundancy to improve performance. However, there are a few applications, such as some on supercomputers in which performance and capacity are primary concerns and low cost is more important than improved reliability.

For RAID 0, the user and system data are distributed across all of the disks in the array. This has a notable advantage over the use of a single large disk: If two different I/O requests are pending for two different blocks of data, then there is a good chance that the requested blocks are on different disks. Thus, the two requests can be issued in parallel, reducing the I/O queuing time.

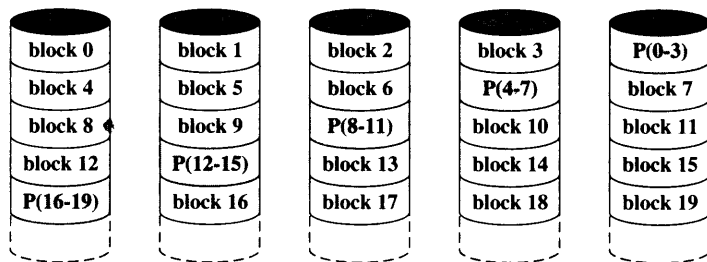
But RAID 0, as with all of the RAID levels, goes further than simply distributing the data across a disk array: The data are *striped* across the available disks. This is best understood by considering Figure 6.10. All of the user and system data are



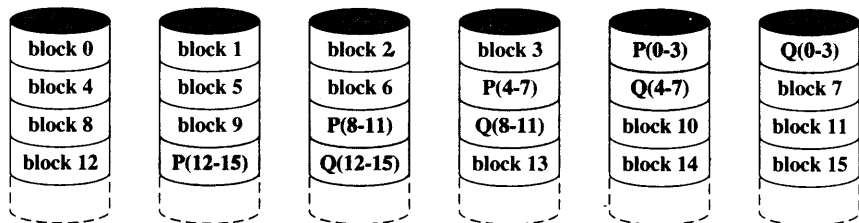
(a) RAID 3 (Bit-interleaved parity)



(b) RAID 4 (Block-level parity)



(c) RAID 5 (Block-level distributed parity)



(d) RAID 6 (Dual redundancy)

Figure 6.9 RAID Levels 3 through 6

viewed as being stored on a logical disk. The disk is divided into strips; these strips may be physical blocks, sectors, or some other unit. The strips are mapped round robin to consecutive array members. A set of logically consecutive strips that maps exactly one strip to each array member is referred to as a *stripe*. In an  $n$ -disk array, the first  $n$  logical strips are physically stored as the first strip on each of the  $n$  disks, forming the first stripe; the second  $n$  strips are distributed as the second strips on

RAID 0. RAID 1 may also provide improved performance over RAID 0 for data transfer intensive applications with a high percentage of reads. Improvement occurs if the application can split each read request so that both disk members participate.

### RAID Level 2

RAID levels 2 and 3 make use of a parallel access technique. In a parallel access array, all member disks participate in the execution of every I/O request. Typically, the spindles of the individual drives are synchronized so that each disk head is in the same position on each disk at any given time.

As in the other RAID schemes, data striping is used. In the case of RAID 2 and 3, the strips are very small, often as small as a single byte or word. With RAID 2, an error-correcting code is calculated across corresponding bits on each data disk, and the bits of the code are stored in the corresponding bit positions on multiple parity disks. Typically, a Hamming code is used, which is able to correct single-bit errors and detect double-bit errors.

Although RAID 2 requires fewer disks than RAID 1, it is still rather costly. The number of redundant disks is proportional to the log of the number of data disks. On a single read, all disks are simultaneously accessed. The requested data and the associated error-correcting code are delivered to the array controller. If there is a single-bit error, the controller can recognize and correct the error instantly, so that the read access time is not slowed. On a single write, all data disks and parity disks must be accessed for the write operation.

RAID 2 would only be an effective choice in an environment in which many disk errors occur. Given the high reliability of individual disks and disk drives, RAID 2 is overkill and is not implemented.

### RAID Level 3

RAID 3 is organized in a similar fashion to RAID 2. The difference is that RAID 3 requires only a single redundant disk, no matter how large the disk array. RAID 3 employs parallel access, with data distributed in small strips. Instead of an error-correcting code, a simple parity bit is computed for the set of individual bits in the same position on all of the data disks.

**Redundancy** In the event of a drive failure, the parity drive is accessed and data is reconstructed from the remaining devices. Once the failed drive is replaced, the missing data can be restored on the new drive and operation resumed.

Data reconstruction is simple. Consider an array of five drives in which X0 through X3 contain data and X4 is the parity disk. The parity for the  $i$ th bit is calculated as follows:

$$X4(i) = X3(i) \oplus X2(i) \oplus X1(i) \oplus X0(i)$$

where  $\oplus$  is exclusive-OR function.

Suppose that drive X1 has failed. If we add  $X4(i) \oplus X1(i)$  to both sides of the preceding equation, we get

$$X1(i) = X4(i) \oplus X3(i) \oplus X2(i) \oplus X0(i)$$

Thus, the contents of each strip of data on X1 can be regenerated from the contents of the corresponding strips on the remaining disks in the array. This principle is true for RAID levels 3 through 6.

In the event of a disk failure, all of the data are still available in what is referred to as reduced mode. In this mode, for reads, the missing data are regenerated on the fly using the exclusive-OR calculation. When data are written to a reduced RAID 3 array, consistency of the parity must be maintained for later regeneration. Return to full operation requires that the failed disk be replaced and the entire contents of the failed disk be regenerated on the new disk.

**Performance** Because data are striped in very small strips, RAID 3 can achieve very high data transfer rates. Any I/O request will involve the parallel transfer of data from all of the data disks. For large transfers, the performance improvement is especially noticeable. On the other hand, only one I/O request can be executed at a time. Thus, in a transaction-oriented environment, performance suffers.

#### RAID Level 4

RAID levels 4 through 6 make use of an independent access technique. In an independent access array, each member disk operates independently, so that separate I/O requests can be satisfied in parallel. Because of this, independent access arrays are more suitable for applications that require high I/O request rates and are relatively less suited for applications that require high data transfer rates.

As in the other RAID schemes, data striping is used. In the case of RAID 4 through 6, the strips are relatively large. With RAID 4, a bit-by-bit parity strip is calculated across corresponding strips on each data disk, and the parity bits are stored in the corresponding strip on the parity disk.

RAID 4 involves a write penalty when an I/O write request of small size is performed. Each time that a write occurs, the array management software must update not only the user data but also the corresponding parity bits. Consider an array of five drives in which X0 through X3 contain data and X4 is the parity disk. Suppose that a write is performed that only involves a strip on disk X1. Initially, for each bit  $i$ , we have the following relationship:

$$X4(i) = X3(i) \oplus X2(i) \oplus X1(i) \oplus X0(i) \quad (11.1)$$

After the update, with potentially altered bits indicated by a prime symbol:

$$\begin{aligned} X4'(i) &= X3(i) \oplus X2(i) \oplus X1'(i) \oplus X0(i) \\ &= X3(i) \oplus X2(i) \oplus X1'(i) \oplus X0(i) \oplus X1(i) \oplus X1(i) \\ &= X3(i) \oplus X2(i) \oplus X1(i) \oplus X0(i) \oplus X1(i) \oplus X1'(i) \\ &= X4(i) \oplus X1(i) \oplus X1'(i) \end{aligned}$$

The preceding set of equations is derived as follows. The first line shows that a change in X1 will also affect the parity disk X4. In the second line, we add the terms  $[\oplus X1(i) \oplus X1(i)]$ . Because the XOR of any quantity with itself is 0, this does not affect the equation. However, it is a convenience that is used to create the third line, by reordering. Finally, Equation (11.1) is used to replace the first four terms by X4(i).

To calculate the new parity, the array management software must read the old user strip and the old parity strip. Then it can update these two strips with the new data and the newly calculated parity. Thus, each strip write involves two reads and two writes.

In the case of a larger size I/O write that involves strips on all disk drives, parity is easily computed by calculation using only the new data bits. Thus, the parity drive can be updated in parallel with the data drives and there are no extra reads or writes.

In any case, every write operation must involve the parity disk, which therefore can become a bottleneck.

### RAID Level 5

RAID 5 is organized in a similar fashion to RAID 4. The difference is that RAID 5 distributes the parity strips across all disks. A typical allocation is a round-robin scheme, as illustrated in Figure 6.9c. For an  $n$ -disk array, the parity strip is on a different disk for the first  $n$  stripes, and the pattern then repeats.

The distribution of parity strips across all drives avoids the potential I/O bottleneck found in RAID 4.

### RAID Level 6

RAID 6 was introduced in a subsequent paper by the Berkeley researchers [KATZ89]. In the RAID 6 scheme, two different parity calculations are carried out and stored in separate blocks on different disks. Thus, a RAID 6 array whose user data require  $N$  disks consists of  $N + 2$  disks.

Figure 6.9d illustrates the scheme. P and Q are two different data check algorithms. One of the two is the exclusive-OR calculation used in RAID 4 and 5. But the other is an independent data check algorithm. This makes it possible to regenerate data even if two disks containing user data fail.

The advantage of RAID 6 is that it provides extremely high data availability. Three disks would have to fail within the MTTR (mean time to repair) interval to cause data to be lost. On the other hand, RAID 6 incurs a substantial write penalty, because each write affects two parity blocks.

Table 6.4 is a comparative summary of the seven levels.

## 6.3 OPTICAL MEMORY

In 1983, one of the most successful consumer products of all time was introduced: the compact disk (CD) digital audio system. The CD is a nonerasable disk that can store more than 60 minutes of audio information on one side. The huge commercial success of the CD enabled the development of low-cost optical-disk storage technology that has revolutionized computer data storage. A variety of optical-disk systems have been introduced (Table 6.5). We briefly review each of these.

### Compact Disk

**CD-ROM** Both the audio CD and the CD-ROM (compact disk read-only memory) share a similar technology. The main difference is that CD-ROM players are more



Table 6.4 RAID Comparison

Level	Advantages	Disadvantages	Applications
0	<p>I/O performance is greatly improved by spreading the I/O load across many channels and drives</p> <p>No parity calculation overhead is involved</p> <p>Very simple design</p> <p>Easy to implement</p>	<p>The failure of just one drive will result in all data in an array being lost</p>	<p>Video production and editing</p> <p>Image editing</p> <p>Pre-press applications</p> <p>Any application requiring high bandwidth</p>
1	<p>100% redundancy of data means no rebuild is necessary in case of a disk failure, just a copy to the replacement disk</p> <p>Under certain circumstances, RAID 1 can sustain multiple simultaneous drive failures</p> <p>Simplest RAID storage subsystem design</p>	<p>Highest disk overhead of all RAID types (100%)—inefficient</p>	<p>Accounting</p> <p>Payroll</p> <p>Financial</p> <p>Any application requiring very high availability</p>
2	<p>Extremely high data transfer rates possible</p> <p>The higher the data transfer rate required, the better the ratio of data disks to ECC disks</p> <p>Relatively simple controller design compared to RAID levels 3,4 &amp; 5</p>	<p>Very high ratio of ECC disks to data disks with smaller word sizes—inefficient</p> <p>Entry level cost very high—requires very high transfer rate requirement to justify</p>	<p>No commercial implementations exist/ not commercially viable</p>
3	<p>Very high read data transfer rate</p> <p>Very high write data transfer rate</p> <p>Disk failure has an insignificant impact on throughput</p> <p>Low ratio of ECC (parity) disks to data disks means high efficiency</p>	<p>Transaction rate equal to that of a single disk drive at best (if spindles are synchronized)</p> <p>Controller design is fairly complex</p>	<p>Video production and live streaming</p> <p>Image editing</p> <p>Video editing</p> <p>Prepress applications</p> <p>Any application requiring high throughput</p>
4	<p>Very high read data transaction rate</p> <p>Low ratio of ECC (parity) disks to data disks means high efficiency</p>	<p>Quite complex controller design</p> <p>Worst write transaction rate and Write aggregate transfer rate</p> <p>Difficult and inefficient data rebuild in the event of disk failure</p>	<p>No commercial implementations exist/ not commercially viable</p>

(Continued)

Table 6.4 Continued

Level	Advantages	Disadvantages	Applications
5	<p>Highest read data transaction rate</p> <p>Low ratio of ECC (parity) disks to data disks means high efficiency</p> <p>Good aggregate transfer rate</p>	<p>Most complex controller design</p> <p>Difficult to rebuild in the event of a disk failure (as compared to RAID level 1)</p>	<p>File and application servers</p> <p>Database servers</p> <p>Web, e-mail, and news servers</p> <p>Intranet servers</p> <p>Most versatile RAID level</p>
6	<p>Provides for an extremely high data fault tolerance and can sustain multiple simultaneous drive failures</p>	<p>More complex controller design</p> <p>Controller overhead to compute parity addresses is extremely high</p>	<p>Perfect solution for mission critical applications</p>

rugged and have error correction devices to ensure that data are properly transferred from disk to computer. Both types of disk are made the same way. The disk is formed from a resin, such as polycarbonate. Digitally recorded information (either music or computer data) is imprinted as a series of microscopic pits on the surface of the polycarbonate. This is done, first of all, with a finely focused, high-intensity laser to create

Table 6.5 Optical Disk Products

<b>CD</b>	Compact Disk. A nonerasable disk that stores digitized audio information. The standard system uses 12-cm disks and can record more than 60 minutes of uninterrupted playing time.
<b>CD-ROM</b>	Compact Disk Read-Only Memory. A nonerasable disk used for storing computer data. The standard system uses 12-cm disks and can hold more than 650 Mbytes.
<b>CD-R</b>	CD Recordable. Similar to a CD-ROM. The user can write to the disk only once.
<b>CD-RW</b>	CD Rewritable. Similar to a CD-ROM. The user can erase and rewrite to the disk multiple times.
<b>DVD</b>	Digital Versatile Disk. A technology for producing digitized, compressed representation of video information, as well as large volumes of other digital data. Both 8 and 12 cm diameters are used, with a double-sided capacity of up to 17 Gbytes. The basic DVD is read-only (DVD-ROM).
<b>DVD-R</b>	DVD Recordable. Similar to a DVD-ROM. The user can write to the disk only once. Only one-sided disks can be used.
<b>DVD-RW</b>	DVD Rewritable. Similar to a DVD-ROM. The user can erase and rewrite to the disk multiple times. Only one-sided disks can be used.

a master disk. The master is used, in turn, to make a die to stamp out copies onto polycarbonate. The pitted surface is then coated with a highly reflective surface, usually aluminum or gold. This shiny surface is protected against dust and scratches by a top coat of clear acrylic. Finally, a label can be silkscreened onto the acrylic.

Information is retrieved from a CD or CD-ROM by a low-powered laser housed in an optical-disk player, or drive unit. The laser shines through the clear polycarbonate while a motor spins the disk past it (Figure 6.11). The intensity of the reflected light of the laser changes as it encounters a pit. Specifically, if the laser beam falls on a pit, which has a somewhat rough surface, the light scatters and a low intensity is reflected back to the source. The areas between pits are called *lands*. A land is a smooth surface, which reflects back at higher intensity. The change between pits and lands is detected by a photosensor and converted into a digital signal. The sensor tests the surface at regular intervals. The beginning or end of a pit represents a 1; when no change in elevation occurs between intervals, a 0 is recorded.

Recall that on a magnetic disk, information is recorded in concentric tracks. With the simplest constant angular velocity (CAV) system, the number of bits per track is constant. An increase in density is achieved with multiple zoned recording, in which the surface is divided into a number of zones, with zones farther from the center containing more bits than zones closer to the center. Although this technique increases capacity, it is still not optimal.

To achieve greater capacity, CDs and CD-ROMs do not organize information on concentric tracks. Instead, the disk contains a single spiral track, beginning near the center and spiraling out to the outer edge of the disk. Sectors near the outside of the disk are the same length as those near the inside. Thus, information is packed evenly across the disk in segments of the same size and these are scanned at the same rate by rotating the disk at a variable speed. The pits are then read by the laser at a **constant linear velocity (CLV)**. The disk rotates more slowly for accesses near the outer edge than for those near the center. Thus, the capacity of a track and the rotational delay both increase for positions nearer the outer edge of the disk. The data capacity for a CD-ROM is about 680 MB.

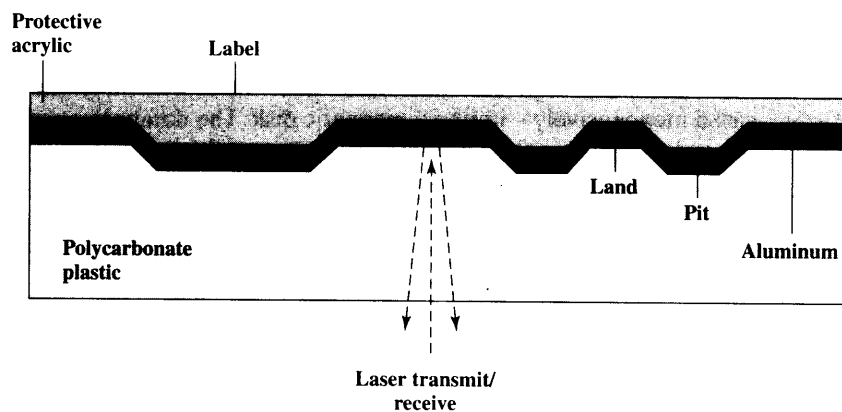


Figure 6.11 CD Operation

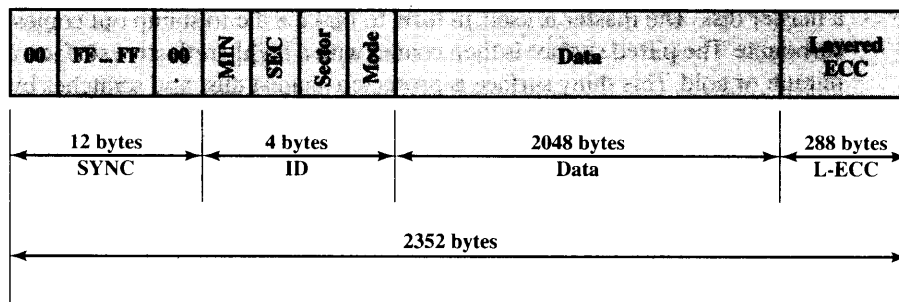


Figure 6.12 CD-ROM Block Format

Data on the CD-ROM are organized as a sequence of blocks. A typical block format is shown in Figure 6.12. It consists of the following fields:

- **Sync:** The sync field identifies the beginning of a block. It consists of a byte of all 0s, 10 bytes of all 1s, and a byte of all 0s.
- **Header:** The header contains the block address and the mode byte. Mode 0 specifies a blank data field; mode 1 specifies the use of an error-correcting code and 2048 bytes of data; mode 2 specifies 2336 bytes of user data with no error-correcting code.
- **Data:** User data.
- **Auxiliary:** Additional user data in mode 2. In mode 1, this is a 288-byte error-correcting code.

With the use of CLV, random access becomes more difficult. Locating a specific address involves moving the head to the general area, adjusting the rotation speed and reading the address, and then making minor adjustments to find and access the specific sector.

CD-ROM is appropriate for the distribution of large amounts of data to a large number of users. Because of the expense of the initial writing process, it is not appropriate for individualized applications. Compared with traditional hard disks, the CD-ROM has two advantages:

- The optical disk together with the information stored on it can be mass replicated inexpensively—unlike a magnetic disk. The database on a magnetic disk has to be reproduced by copying one disk at a time using two disk drives.
- The optical disk is removable, allowing the disk itself to be used for archival storage. Most magnetic disks are nonremovable. The information on nonremovable magnetic disks must first be copied to tape before the disk drive/disk can be used to store new information.

The disadvantages of CD-ROM are as follows:

- It is read-only and cannot be updated.
- It has an access time much longer than that of a magnetic disk drive, as much as half a second.

**CD Recordable** To accommodate applications in which only one or a small number of copies of a set of data is needed, the write-once read-many CD, known as the CD recordable (CD-R), has been developed. For CD-R, a disk is prepared in such a way that it can be subsequently written once with a laser beam of modest intensity. Thus, with a somewhat more expensive disk controller than for CD-ROM, the customer can write once as well as read the disk.

The CD-R medium is similar to but not identical to that of a CD or CD-ROM. For CDs and CD-ROMs, information is recorded by the pitting of the surface of the medium, which changes reflectivity. For a CD-R, the medium includes a dye layer. The dye is used to change reflectivity and is activated by a high-intensity laser. The resulting disk can be read on a CD-R drive or a CD-ROM drive.

The CD-R optical disk is attractive for archival storage of documents and files. It provides a permanent record of large volumes of user data.

**CD Rewritable** The CD-RW optical disk can be repeatedly written and overwritten, as with a magnetic disk. Although a number of approaches have been tried, the only pure optical approach that has proved attractive is called **phase change**. The phase change disk uses a material that has two significantly different reflectivities in two different phase states. There is an amorphous state, in which the molecules exhibit a random orientation and which reflects light poorly; and a crystalline state, which has a smooth surface that reflects light well. A beam of laser light can change the material from one phase to the other. The primary disadvantage of phase change optical disks is that the material eventually and permanently loses its desirable properties. Current materials can be used for between 500,000 and 1,000,000 erase cycles.

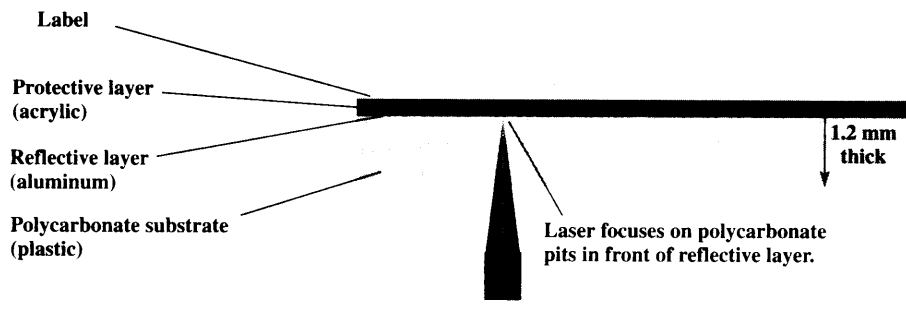
The CD-RW has the obvious advantage over CD-ROM and CD-R that it can be rewritten and thus used as a true secondary storage. As such, it competes with magnetic disk. A key advantage of the optical disk is that the engineering tolerances for optical disks are much less severe than for high-capacity magnetic disks. Thus, they exhibit higher reliability and longer life.

### Digital Versatile Disk

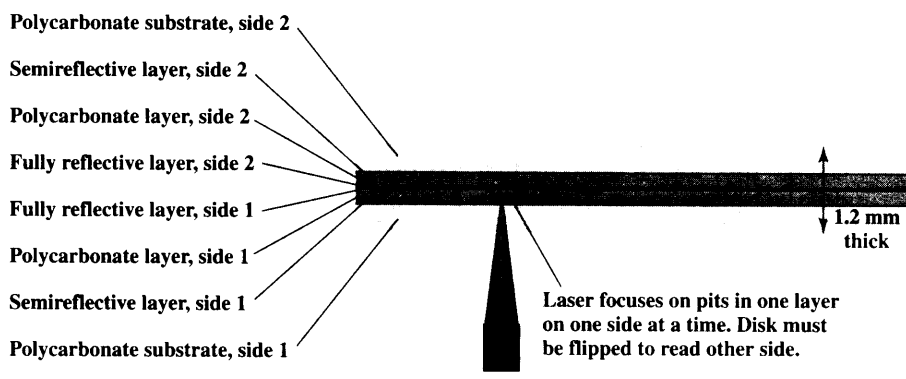
With the capacious digital versatile disk (DVD), the electronics industry has at last found an acceptable replacement for the analog VHS video tape. The DVD will replace the video tape used in video cassette recorders (VCRs) and, more important for this discussion, replace the CD-ROM in personal computers and servers. The DVD takes video into the digital age. It delivers movies with impressive picture quality, and it can be randomly accessed like audio CDs, which DVD machines can also play. Vast volumes of data can be crammed onto the disk, currently seven times as much as a CD-ROM. With DVD's huge storage capacity and vivid quality, PC games will become more realistic and educational software will incorporate more video. Following in the wake of these developments will be a new crest of traffic over the Internet and corporate intranets, as this material is incorporated into Web sites.

The DVD's greater capacity is due to three differences from CDs (Figure 6.13):

1. Bits are packed more closely on a DVD. The spacing between loops of a spiral on a CD is  $1.6 \mu\text{m}$  and the minimum distance between pits along the spiral is  $0.834 \mu\text{m}$ . The DVD uses a laser with shorter wavelength and achieves a loop spacing of  $0.74 \mu\text{m}$  and a minimum distance between pits of  $0.4 \mu\text{m}$ .



(a) CD-ROM—Capacity 682 MB



(b) DVD-ROM, double-sided, dual-layer—Capacity 17 GB

Figure 6.13 CD-ROM and DVD-ROM

The result of these two improvements is about a seven-fold increase in capacity, to about 4.7 GB.

2. The DVD employs a second layer of pits and lands on top of the first layer. A dual-layer DVD has a semireflective layer on top of the reflective layer, and by adjusting focus, the lasers in DVD drives can read each layer separately. This technique almost doubles the capacity of the disk, to about 8.5 GB. The lower reflectivity of the second layer limits its storage capacity so that a full doubling is not achieved.
3. The DVD-ROM can be two sided, whereas data is recorded on only one side of a CD. This brings total capacity up to 17 GB.

As with the CD, DVDs come in writeable as well as read-only versions (Table 6.5).

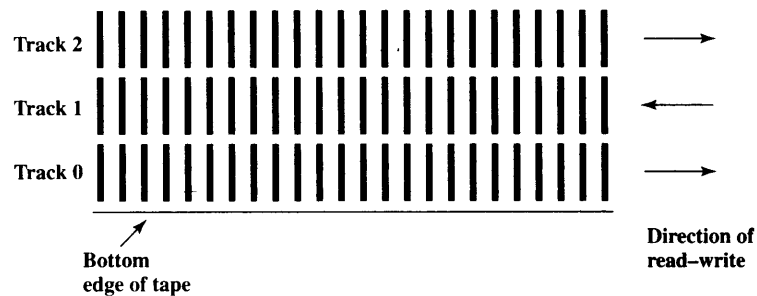
## 6.4 MAGNETIC TAPE

Tape systems use the same reading and recording techniques as disk systems. The medium is flexible polyester (similar to that used in some clothing) tape coated with magnetizable material. The coating may consist of particles of pure metal in special binders or vapor-plated metal films. The tape and the tape drive are analogous to a

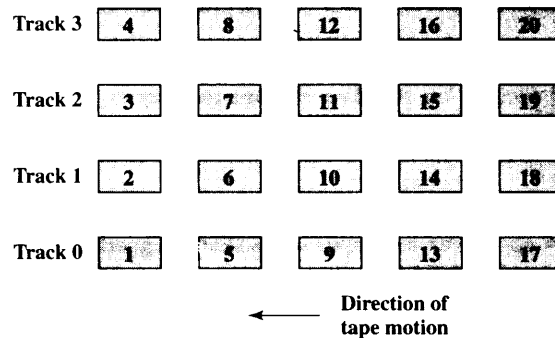
home tape recorder system. Tape widths vary from 0.38 cm (0.15 inch) to 1.27 cm (0.5 inch). Tapes used to be packaged as open reels that have to be threaded through a second spindle for use. Today, virtually all tapes are housed in cartridges.

Data on the tape are structured as a number of parallel tracks running lengthwise. Earlier tape systems typically used nine tracks. This made it possible to store data one byte at a time, with an additional parity bit as the ninth track. This was followed by tape systems using 18 or 36 tracks, corresponding to a digital word or double word. The recording of data in this form is referred to as **parallel recording**. Most modern systems instead use **serial recording**, in which data are laid out as a sequence of bits along each track, as is done with magnetic disks. As with the disk, data are read and written in contiguous blocks, called *physical records*, on a tape. Blocks on the tape are separated by gaps referred to as *interrecord gaps*. As with the disk, the tape is formatted to assist in locating physical records.

The typical recording technique used in serial tapes is referred to as **serpentine recording**. In this technique, when data are being recorded, the first set of bits is recorded along the whole length of the tape. When the end of the tape is reached, the heads are repositioned to record a new track, and the tape is again recorded on its whole length, this time in the opposite direction. That process continues, back and forth, until the tape is full (Figure 6.14a). To increase speed, the read-write head is



(a) Serpentine reading and writing



(b) Block layout for system that reads-writes four tracks simultaneously

Figure 6.14 Typical Magnetic Tape Features

Table 6.6 DLTape Drives

	DLT 4000	DLT 8000	SDLT 600
Capacity (GB)	20	40	300
Data rate (MB/s)	1.5	6.0	36.0
Bit density (Kb/cm)	32.3	38.6	92
Track density (t/cm)	101	164	387
Media length (m)	549	549	597
Media width (cm)	1.27	1.27	1.27
Number of tracks	128	208	448
Number of tracks read/write simultaneously	2	4	8

capable of reading and writing a number of adjacent tracks simultaneously (typically 2 to 8 tracks). Data are still recorded serially along individual tracks, but blocks in sequence are stored on adjacent tracks, as suggested by Figure 6.14b. Table 6.6 shows parameters for one system, known as DLTape.

A tape drive is a *sequential-access* device. If the tape head is positioned at record 1, then to read record  $N$ , it is necessary to read physical records 1 through  $N - 1$ , one at a time. If the head is currently positioned beyond the desired record, it is necessary to rewind the tape a certain distance and begin reading forward. Unlike the disk, the tape is in motion only during a read or write operation.

In contrast to the tape, the disk drive is referred to as a *direct-access* device. A disk drive need not read all the sectors on a disk sequentially to get to the desired one. It must only wait for the intervening sectors within one track and can make successive accesses to any track.

Magnetic tape was the first kind of secondary memory. It is still widely used as the lowest-cost, slowest-speed member of the memory hierarchy.

## 6.5 RECOMMENDED READING AND WEB SITES

[MEE96a] provides a good survey of the underlying recording technology of disk and tape systems. [MEE96b] focuses on the data storage techniques for disk and tape systems. [COME00] is a short but instructive article on current trends in magnetic disk storage technology. [ANDE03] provides a more recent discussion of magnetic disk storage technology.

An excellent survey of RAID technology, written by the inventors of the RAID concept, is [CHEN94]. A good overview paper is [FRIE96]. A good performance comparison of the RAID architectures is [CHEN96].

[MARC90] gives an excellent overview of the optical storage field. A good survey of the underlying recording and reading technology is [MANS97].

[ROSC03] provides a comprehensive overview of all types of external memory systems, with a modest amount of technical detail on each. [KHUR01] is another good survey.



- ANDE03** Anderson, D. "You Don't Know Jack About Disks." *ACM Queue*, June 2003.
- CHEN94** Chen, P.; Lee, E.; Gibson, G.; Katz, R.; and Patterson, D. "RAID: High-Performance, Reliable Secondary Storage." *ACM Computing Surveys*, June 1994.
- CHEN96** Chen, S., and Towsley, D. "A Performance Evaluation of RAID Architectures." *IEEE Transactions on Computers*, October 1996.
- COME00** Comerford, R. "Magnetic Storage: The Medium that Wouldn't Die." *IEEE Spectrum*, December 2000.
- FRIE96** Friedman, M. "RAID Keeps Going and Going and ..." *IEEE Spectrum*, April 1996.
- KHUR01** Khurshudov, A. *The Essential Guide to Computer Data Storage*. Upper Saddle River, NJ: Prentice Hall, 2001.
- MANS97** Mansuripur, M., and Sincerbox, G. "Principles and Techniques of Optical Data Storage." *Proceedings of the IEEE*, November 1997.
- MARC90** Marchant, A. *Optical Recording*. Reading, MA: Addison-Wesley, 1990.
- MEE96a** Mee, C., and Daniel, E. eds. *Magnetic Recording Technology*. New York: McGraw-Hill, 1996.
- MEE96b** Mee, C., and Daniel, E. eds. *Magnetic Storage Handbook*. New York: McGraw-Hill, 1996.
- ROSC03** Rosch, W. *Winn L. Rosch Hardware Bible*. Indianapolis, IN: Que Publishing, 2003.



#### Recommended Web Sites:

- **Optical Storage Technology Association:** Good source of information about optical storage technology and vendors, plus extensive list of relevant links
- **DLTtape:** Good collection of technical information and links to vendors

## 6.6 KEY TERMS, REVIEW QUESTIONS, AND PROBLEMS

### Key Terms

access time	fixed-head disk	platter
CD	floppy disk	RAID
CD-ROM	gap	removable disk
CD-R	head	rotational delay
CD-RW	land	sector
constant angular velocity (CAV)	magnetic disk	seek time
constant linear velocity (CLV)	magnetic tape	serpentine recording
cylinder	magnetoresistive	striped data
DVD	movable-head disk	substrate
DVD-ROM	multiple zoned recording	track
DVD-R	nonremovable disk	transfer time
DVD-RW	optical memory	
	pit	

### Review Questions

- 6.1 What are the advantages of using a glass substrate for a magnetic disk?
- 6.2 How are data written onto a magnetic disk?
- 6.3 How are data read from a magnetic disk?
- 6.4 Explain the difference between a simple CAV system and a multiple zoned recording system.
- 6.5 Define the terms *track*, *cylinder*, and *sector*.
- 6.6 What is the typical disk sector size?
- 6.7 Define the terms *seek time*, *rotational delay*, *access time*, and *transfer time*.
- 6.8 What common characteristics are shared by all RAID levels?
- 6.9 Briefly define the seven RAID levels.
- 6.10 Explain the term *striped data*.
- 6.11 How is redundancy achieved in a RAID system?
- 6.12 In the context of RAID, what is the distinction between parallel access and independent access?
- 6.13 What is the difference between CAV and CLV?
- 6.14 What differences between a CD and a DVD account for the larger capacity of the latter?
- 6.15 Explain serpentine recording.

### Problems

- 6.1 Consider a disk with  $N$  tracks numbered from 0 to  $(N - 1)$  and assume that requested sectors are distributed randomly and evenly over the disk. We want to calculate the average number of tracks traversed by a seek.
  - a. First, calculate the probability of a seek of length  $j$  when the head is currently positioned over track  $t$ . *Hint*: this is a matter of determining the total number of combinations, recognizing that all track positions for the destination of the seek are equally likely.
  - b. Next, calculate the probability of a seek of length  $K$ . *Hint*: This involves the summing over all possible combinations of movements of  $K$  tracks.
  - c. Calculate the average number of tracks traversed by a seek, using the formula for expected value:

$$E[x] = \sum_{i=0}^{N-1} i \times \Pr[x = i]$$

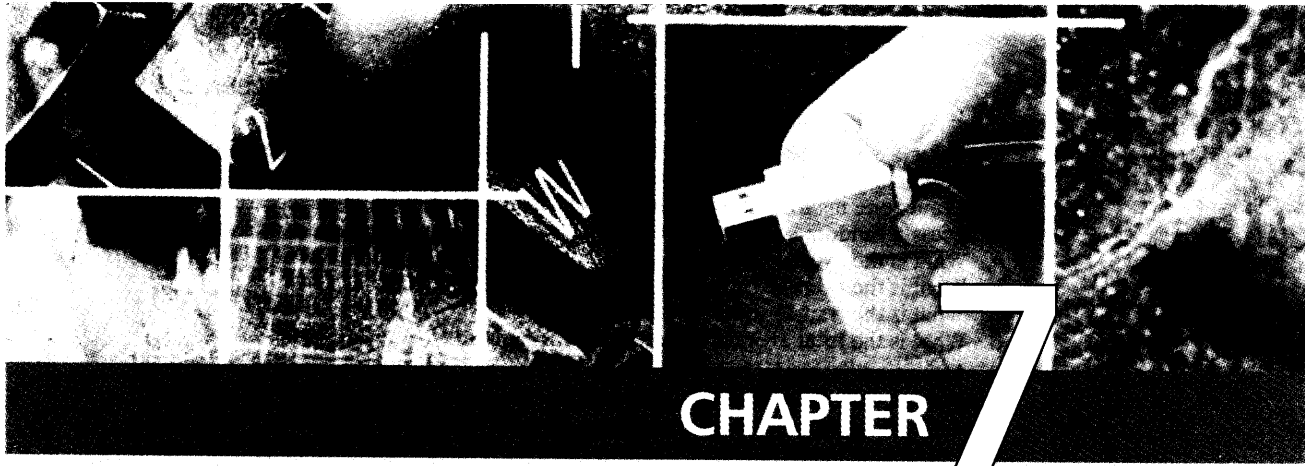
*Hint*: Use the equalities

$$\sum_{i=1}^n i = \frac{n(n+1)}{2}; \quad \sum_{i=1}^n i^2 = \frac{n(n+1)(2n+1)}{6}.$$

- d. Show that for large values of  $N$ , the average number of tracks traversed by a seek approaches  $N/3$ .
- 6.2 Define the following for a disk system:
    - $t_s$  = seek time; average time to position head over track
    - $r$  = rotation speed of the disk, in revolutions per second
    - $n$  = number of bits per sector
    - $N$  = capacity of a track, in bits
    - $t_A$  = time to access a sector

Develop a formula for  $t_A$  as a function of the other parameters.

- 6.3 Consider a single-platter disk with the following parameters: rotation speed: 7200 rpm; number of tracks on one side of platter: 30,000; number of sectors per track: 600; seek time: one ms for every hundred tracks traversed. Let the disk receive a request to access a random sector on a random track and assume the head starts at track 0.
- What is the average seek time?
  - What is the average rotational latency?
  - What is the transfer time for a sector?
  - What is the total average time to satisfy a request?
- 6.4 A distinction is made between physical records and logical records. A **logical record** is a collection of related data elements treated as a conceptual unit, independent of how or where the information is stored. A **physical record** is a contiguous area of storage space that is defined by the characteristics of the storage device and operating system. Assume a disk system in which each physical record contains thirty 120-byte logical records. Calculate how much disk space (in sectors, tracks, and surfaces) will be required to store 300,000 logical records if the disk is fixed-sector with 512 bytes/sector, with 96 sectors/track, 110 tracks per surface, and 8 usable surfaces. Ignore any file header record(s) and track indexes, and assume that records cannot span two sectors.
- 6.5 It should be clear that disk striping can improve data transfer rate when the strip size is small compared to the I/O request size. It should also be clear that RAID 0 provides improved performance relative to a single large disk, because multiple I/O requests can be handled in parallel. However, in this latter case, is disk striping necessary? That is, does disk striping improve I/O request rate performance compared to a comparable disk array without striping?



## INPUT/OUTPUT

- 7.1 External Devices
- 7.2 I/O Modules
- 7.3 Programmed I/O
- 7.4 Interrupt-Driven I/O
- 7.5 Direct Memory Access
- 7.6 I/O Channels and Processors
- 7.7 The External Interface: FireWire and InfiniBand
- 7.8 Recommended Reading and Web Sites
- 7.9 Key Terms, Review Questions, and Problems