# CHAPTER I

## CONTENTS:

## OBJECTIVE:

In this chapter we are concerned with basic architecture and the different operations related to explain the proper functioning of the computer. Also how we can specify the operations with the help of different instructions.

# CHAPTER II

## CONTENTS:

## OBJECTIVE:

Here the concept of digital hardware modules is discussed. Size and complexity of the system can be varied as per the requirement of today. The interconnection of various modules is explained in the text. The way by which data is transferred from one register to another is called micro operation. Different micro operations are explained in the chapter.

# CHAPTER III

## CONTENTS:

## OBJECTIVE:

There are various instructions with the help of which we can transfer the data from one place to another and manipulate the data as per our requirement. In this chapter we have included all the instructions, how they are being identified by the computer, what are their formats and many more details regarding the instructions.

# CHAPTER IV

## CONTENTS:

## OBJECTIVE:

Various examples of micro programs are discussed in this chapter. Complexity of the system is explained with the help of micro operations that are performed. Way out for mapping a instruction is discussed. Symbolic micro programs and micro operations are expressed in detail.

# CHAPTER V

## CONTENTS:

**5.1 INTRODUCTION**

**5.2 GENERAL REGISTER ORGANIZATION**

**5.3 STACK ORGANIZATION**

**5.4 INSTRUCTION FORMATS**

**5.5 ADDRESSING MODES**

**SUMMARY**

**SELF ASSESSMENT**

## OBJECTIVE:

**CPU is the heart of any computer system. Chapter contains the detail regarding the internal architectural details of CPU. Various functions performed by CPU are explained in the text. How registers are organized ,How they are being used during addressing schemes is presented in this chapter.**

# CHAPTER VI

## CONTENTS:

## OBJECTIVE:

**Without Input-Output devices any computer system is never a complete system. With the help of I/O devices we can instruct the computer to perform any task with help of input instructions, print out the results in the text form or get displays, with the help of printers. Mechanical movements of D.C motors or robots etc are again considered to be output effects. Concept of DMA is discussed in the chapter.**

# CHAPTER VII

## CONTENTS:

## OBJECTIVE:

One of the most important aspects of our computer system is Memory. Many concerns are associated with its organization. All the facts are encapsulated in this chapter. Cost factor, speed of storage, transfer speed of data are discussed here. Assessment

## REFERENCES:

- Computer Organization and Architecture, Rajaram & Radhakrishan, PHI.

- Computer Organization & Architecture: Designing for Performance, Stalling, PHI.

- Computer Organization and Design, Pal Choudhary, PHI.

- Computer Systems Organization & Architecture, Carpenelli, Pearson Education.

- Computer Organization and Architecture, Stalling, Pearson Education.

- Computer System Architecture, Morris Mano, PHI.

- Computer Architecture and Organization, McGraw Hill Company, New Delhi. J.P. Hayes.

# CHAPTER-I

# Instruction Codes

**Author: Dr. Manoj Duhan**                                   **Vetter: Mr. Sandeep Arya**

## 1.1   INTRODUCTION

In this chapter we introduce concept of a basic computer and show how its operation can be specified with register transfer statements. The organization of the computer is defined by its internal registers, the timing and control structure and the set of instructions that it uses. The design of the computer is then carried out in detail. Although the basic computer presented in this chapter is very small compared to commercial computers, it has the advantage of being simple enough so we can demonstrate the design process without too many complications.

The internal organization of a digital system is defined by the sequence of micro operations it performs on data stored in its registers. The general-purpose digital computer is capable of executing various micro operations and, in addition, can be instructed as to what specific sequence of operations it must perform. The user of a computer can control the process by means of a program. A program is a set of instructions that specify the operations, operations operands, and the sequence by which processing has to occur. The data processing task may be altered by specifying a new program with different instructions or specifying the same instructions with different data.

A computer instruction is a binary code that specifies a sequence of micro operations for the computer. Instruction codes together with data are stored in memory. The computer reads each instruction from memory and places it in a control register. The control then interprets the binary code of the instructions and proceeds to execute it by issuing a sequence of micro operations. Every computer has its own unique instruction set. The ability to store and execute instructions, the stored program concept, is the most important property of a general-purpose computer.

An instruction code is a group of bits that instruct the computer to perform a specific operation. It is usually divided into parts, each having its own particular interpretation. The most basic part of an instruction code is its operation part. The operation code of an instruction is a group of bits that define such operations as add, subtract, multiply, shift, and complement. The number of bits required for the operation code of an instruction depends on the total number of operations available in the computer. The operation code must consists of at least n bits for a given 2" (or less) distinct operations. As an illustration, consider a computer with 64 distinct operations, one of them being an ADD operation. The operation code consists of six bits, with a bit configuration 110010 assigned to the ADD operation. When this operation code is decoded in the control unit, the computer issues control signals to read an operand from memory and add the operand to a processor register.

At this point we must recognize the relationship between a computer operation and a micro operation. An operation is part of an instruction stored in computer memory. It is a binary code tells the computer to perform a specific operation. The control unit receives the instruction from memory and interprets the operation code bits. It then issues a sequence of control signals to initiate microoperations in internal computer registers. For every operation come, the control issues a sequence of microoperations needed for the hardware implementation of the specified operation. For this reason, an operation code is sometimes called a microoperation because it specifies a set of microoperations.
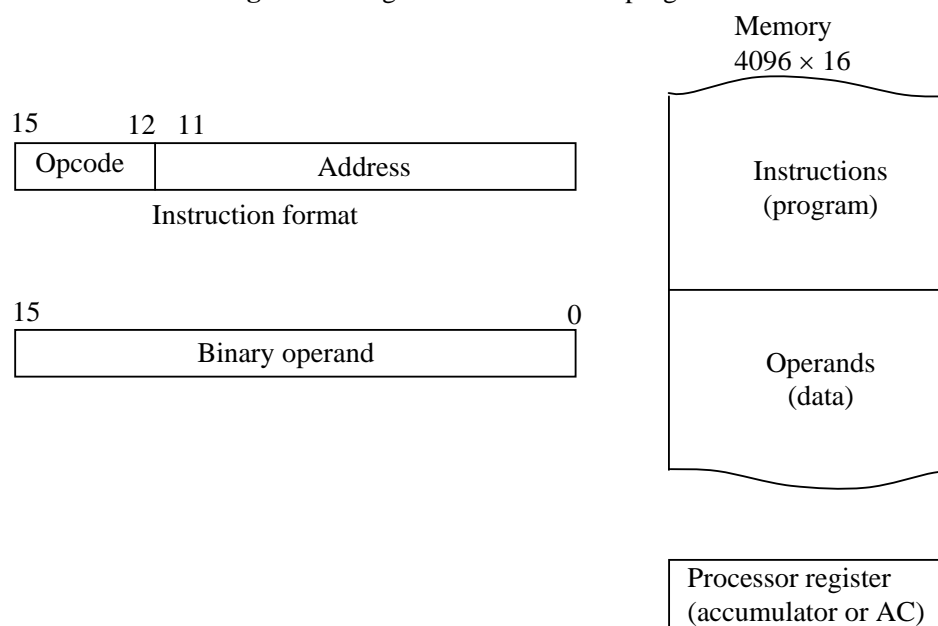
The operation part of an instruction code specifies the operation to be performed. This operation must be performed on some data stored in processor registers or in memory. An instruction code must therefore specify not only the operation but also the registers or the memory words where the operands are to be found, as well as the register or memory word where the result is to be stored. Memory words can be specified in instruction codes by their address. Processor registers can be specified by assigning to the instruction another binary code of k bits that specifies one of $2^k$ registers. There are many variations for arranging the binary code of instructions, and each computer has its own particular instruction code format. Instruction code formats are conceived computer designers who specify the architecture of the computer. In this chapter we choose a particular instruction code to explain the basic organization and design of digital computers.

## 1.2    STORED PROGRAM ORGANIZATION

The simplest way to organize a computer is to have one processor register and instruction code format with two parts. The first part specifies the operation to be performed and the second specifies an address. The memory address tells the control where to find an operand in memory. This operand is read from memory and used as the data to be operated on together with the data stored in the processor register.

Figure 1.1 depicts this type of organization. Instructions are stored in one section of memory and data in another. For a memory unit with 4096 words we need 12 bits to specify an address since $2^{12} = 4096$. If we store each instruction code in one 16-bit memory word, we have available four bits for the operation code (abbreviated op code) to specify one out of 16 possible operations, and 12 bits to specify the address of an operand. The control reads a 16-bit instruction from the program portion of memory. It uses the 12-bit address part of the instruction to read a 16-bit operand from the data portion of memory. It then executes the operation specified by the operation code.

**Figure 1.1** Organization for stored program

Computers that have a single- processor register usually assign to it the name accumulator and label it AC. The operation is performed with the memory operand and the content of AC.

If an operation in an instruction code does not need an operand from memory, the rest of the bits in the instruction can be used for other purposes. For example, operations such as clear AC, complement AC, and increment AC operate on data stored in the AC register. They do not need an operand from memory. For these types of operations, the second part of the instruction code (bits 0 through 11) is not needed for specifying a memory address and can be used to specify other operations for the computer.

## 1.3    INDIRECT ADDRESS

It is sometimes convenient to use the address bits of an instruction code not as an address but as the actual operand. When the second part of an instruction code specifies an operand, the instruction is said to have an immediate operand. When the second part specifies the address of an operand, the instruction is said to have a direct address. This is in contrast to a third possibility called indirect address, where the bits in the second part of the instruction designate an address of a memory word in which the address of the operand is found. One bit of the instruction code can be used to distinguish between a direct and an indirect address.

As an illustration of this configuration, consider the instruction code format shown in Fig. 1.2(a). It consists of a 3-bit operation code, a 12-bit address, and an indirect address mode bit designated by 1. The mode bit is 0 for a direct address and 1 for an indirect address. A direct address instruction is shown in Fig. 1.2(b). It is placed in address 22 in memory. The 1 bit is 0, so the instruction is recognized as a direct address instruction. The op code specifies an ADD instruction, and the address part is the binary equivalent of 457. The control finds the operand in memory at address 457 and adds it to the content of AC. The instruction in address 35 shown in Fig. 1-2(c) has mode bit I = 1. Therefore, it is recognized as an indirect address instruction. The address part is the binary equivalent of 300. The control gives to address 300 to find the address of the operand. The address of the operand in this case is 1350. The operand found in address 1350 is then added to the content of AC. The indirect address instruction needs two references to memory to fetch an operand. The first reference is needed to read the address of the operand; the second is for the operand itself.  We define the effective address to be the address of the operand in a computation-type instruction or the target address in a branch-type instruction. Thus the effective address  in the instruction of Fig. 1-2(b) is 457 and in the instruction of Fig 1-2(c) is 1350.

The direct and indirect addressing modes are used in the computer presented in this chapter. The memory word that holds the address of the operand in an indirect address instruction is used as a pointer to an array of data.  The pointer could be placed in processor register instead of memory as done in commercial computers.

## 1-4 COMPUTER REGISTERS

Computer instructions are normally stored in consecutive memory locations and are executed sequentially one at a time. The control reads an instruction from a specific address in memory and executes it. It then continues by reading the nest instruction in sequence and executes it, and so on. This type of instruction sequencing needs a counter to calculate the address of the next instruction after execution of the current instruction is completed. It is also necessary to provide a register in the control unit for storing the instruction code after it is read from memory. The computer needs processor registers for manipulating data and a register for holding a memory address. These requirements dictate the register configuration shown in Fig. 1-3. The registers are

also listed in Table 1-1 together with a brief description of their function and the number of bits that they contain.
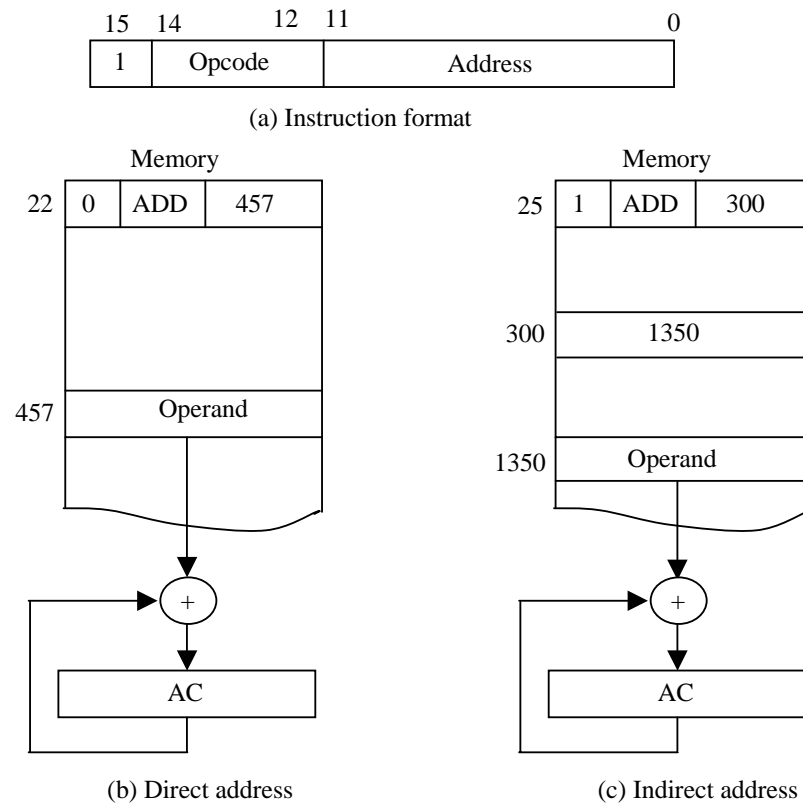


(a) Instruction format

(b) Direct address  (c) Indirect address

**Figure 1-2** Calculation of direct and indirect address.

The memory unit has a capacity of 4096 words and each word contains 16 bits. Twelve bits of an instruction word are needed to specify the address of an operand. This leaves three bits for the operation part of the instruction and a bit to specify a direct or indirect address. The data register (DR) holds the operand read from memory. The accumulator (AC) register is a general-purpose processing register. The instruction read form memory is placed in the instruction register (IR). The temporary register (TR) is used for holding temporary data during the processing.

**Table 1-1** List of Registers for Basic Computer

| Register symbol | Number of bits | Register name | Function |
|---|---|---|---|
| DR | 16 | Data register | Holds memory operand |
| AR | 12 | Address register | Holds address for memory |
| AC | 16 | Accumulator | Processor register |
| IR | 16 | Instruction register | Holds instruction code |

| PC | 12 | Program counter | Holds address of instruction |
|------|-----|-------------------|--------------------------------|
| TR | 16 | Temporary register | Holds temporary data |
| INPR | 8 | Input register | Holds input character |
| OUTR | 8 | Output register | Holds output character |

The memory address register (AR) has 12 bits since this is the width of a memory address. The program counter (PC) also has 12 bits and it holds address of the next instruction to be read from memory after the current the instruction is executed. The PC goes through a counting sequence and causes the computer to read sequential instructions previously stored in memory. Instruction words are and executed in sequence unless a branch instruction is encountered. A branch instruction calls for a transfer to a nonconsecutive instruction in the program. The address part of a branch instruction is transferred to PC to become the address of the next instruction. To read an instruction, the content of PC is taken as the address for memory and a memory read cycle is initiated. PC is then incremented by one, so it holds the address of the next instruction in sequence.

Two registers are used for input and output. The input register (INPR) receives an 8-bit character from an input device. The output register (OUTR) holds an 8-bit character for an output device.



**Figure 1-3** Basic computer registers and memory.

## 1.5     COMMON BUS SYSTEM

The basic computer has eight registers, a memory unit, and a control unit. Paths must be provided to transfer information from one register to another and between memory and registers. The number of wires will be excessive if connections are made between the outputs of each register and the inputs of the other registers. A more efficient scheme for transferring information in a system with many registers is to use a common bus. It is known that how to construct a bus

system using multiplexers or three-state buffer gates. The connection of the registers and memory of the basic computer to a common bus system is shown in Fig. 1-4.

The outputs of seven registers and memory are connected to the common bus. The specific output that is selected for the bus lines at any given time is determined from the binary value of the selection variables $S_2 S_1$, and $S_0$. The number along each output shows the decimal equivalent of the required binary selection. For example, the number along the output of DR is 3. the 16-bit outputs of DR are placed on the bus lines when $S_2 S_1 S_0 = 011$ since this is the binary value of decimal 3. The lines from the common bus are connected to the inputs of each register and the data input of each register and the data inputs of the memory. The particular register whose LD (load) input is enabled receives the data from the bus during the next clock pulse transition. The memory receives the contents of the bus when its write input is activated. The memory places its 16-bit output onto the bus when the read input is activated and $S_2 S_1 S_0 = 111$. Four registers, DR, AC, IR, and TR, have 16-bits each. Two registers, AR and PC, have 12 bits each since they hold a memory address. When the contents of AR or PC are applied to the 16-bit common bus, the four most significant bits are set to 0's. When AR or PC receives information from the bus, only the 12 least significant bits are transferred into the register.

The input register INPR and the output register OUTR have 8 bits each and communicate with the eight least significant bits (LSB) in the bus. INPR is connected to provide information to the bus but OUTR can only receive information from the bus. This is because INPR receives a character from an input device which is then transferred to AC. OUTR receives a character from AC and delivers it to an output device. There is no transfer from OUTR to any of the other registers.

The 16 lines of the common bus receive information from sex registers and the memory unit. The bus lines are connected to the inputs of six registers and the memory. Five registers have three control inputs: LD (load), INR (increment) and CLR (clear). This type of register is equivalent to a binary counter with parallel load and synchronous clear. The increment operation is achieved by enabling the count input of the counter. Two registers have only a LD input.

The input data and output data of the memory are connected to the common bus, but the memory address is connected to AR. Therefore, AR must always be used to specify a memory address. By using a single register for the address, we eliminate the need for an address bus that would have been needed otherwise. The content of any register can be specified for the memory data input during a write operation. Similarly, any register can receive the data from memory after a read operation except AC.

The 16 inputs of AC come from an adder and logic circuit. This circuit has three sets of inputs. One set of 16-bit inputs come from the outputs of AC. They are used to implement register micro-operations such as complement AC and shift AC. Another set of 16-bit inputs come from the data register DR. The inputs from DR and AC are used for arithmetic and logic micro-operations, such as add DR to AC or and DR to AC. The result of an addition is transferred to AC and the end carry-out of the addition is transferred to flip-flop E (extended AC bit). A third set of 8-bit inputs come from the input register INPR.

Note that the content of any register can be applied onto the bus and an operation can be performed in the adder and logic circuit during the same clock cycle. The clock transition at the end of the cycle transfers the content of the bus into the designated destination register and the output of the adder and logic circuit into AC. For example, the two micro-operations:

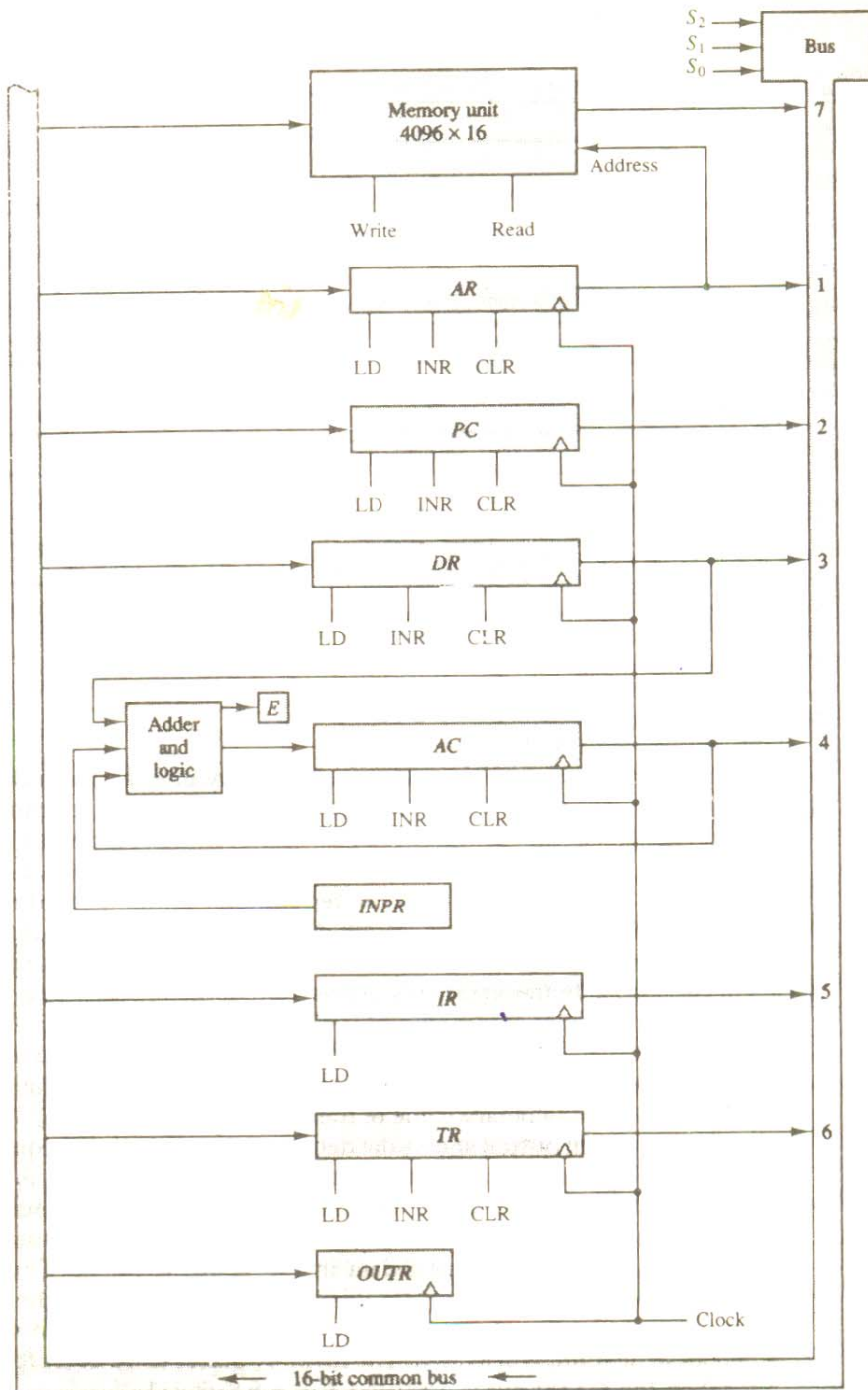$$DR \leftarrow AC \text{ and } AC \leftarrow DR$$

**Figure 1.4** Basic registers which are connected to a common bus

can be executed at the same time. This can be done by placing the content of AC on the bus (with $S_2S_1S_0 = 100$), enabling the LD(load) input of DR, transferring the content of DR through the

adder and logic circuit into AC, and enabling the LD (load) input of AC, all during the same clock cycle. The two transfers occur upon the arrival of the clock pulse transition at the end of the clock cycle.

## SUMMARY

1. The organization of the computer is defined by its internal registers, the timing and control structure, and the set of instructions that it uses.

2. The user of a computer can control the process by means of a program. A program is a set of instructions that specify the operations, operations operands, and the sequence by which processing has to occur.

3. The general-purpose digital computer is capable of executing various micro operations and, in addition, can be instructed as to what specific sequence of operations it must perform.

4. An operation is part of an instruction stored in computer memory. It is a binary code tells the computer to perform a specific operation.

5. The operation part of an instruction code specifies the operation to be performed.

6. Instruction code formats are conceived computer designers who specify the architecture of the computer.

7. The simplest way to organize a computer is to have one processor register and instruction code format with two parts. The first part specifies the operation to be performed and the second specifies an address.

8. Computer instructions are normally stored in consecutive memory locations and are executed sequentially one at a time.

9. The direct and indirect addressing modes are used in the computer.

10. The memory word that holds the address of the operand in an indirect address instruction is used as a pointer to an array of data. The pointer could be placed in processor register instead of memory as done in commercial computers.

11. The basic computer has eight registers, a memory unit, and a control unit. Paths should be provided to transfer information from one register to another and between memory and registers.

12. The output of seven registers and memory are connected to the common bus.

13. The lines from the common bus are connected to the inputs of each register and the data input of each register and the data inputs of the memory.

14. The 16 lines of the common bus receive information from sex registers and the memory unit. The bus lines are connected to the inputs of six registers and the memory.

15. The content of any register can be applied onto the bus and an operation can be performed in the adder and logic circuit during the same clock cycle.

16. The input data and output data of the memory are connected to the common bus, but the memory address is connected to AR.

17. Content of any register can be applied onto the bus and an operation can be performed in the adder and logic circuit during the same clock cycle.

**18.** The clock transition at the end of the cycle transfers the content of the bus into the designated destination register and the output of the adder and logic circuit into AC.

## SELF ASSESSMENT

1. What is a computer instruction and what are the ways for specifying them?

2. What is addressing? What are direct and indirect addressing?

3. Draw and explain the common bus system.

4. What are computer registers? Draw and explain the computer registers organization.

5. A computer uses a memory unit with 256K words of 32 bits each. A binary instruction code is stored in one word of memory. The instruction has four parts: an indirect bit, an operation code, a register code part to specify one of 64 registers, and an address part.

   a. How many bits are there in the operation code, the register code part, and the address part?

   b. Draw the instruction word format and indicate the number of bits in each part.

   c. How many bits are there in the data and address inputs of the memory?

6. What is the difference between a direct and an indirect address instruction? How many references to memory are needed for each type of instruction to bring an operand into a processor register?

7. The following control inputs are active in the bus system shown in Fig. 1-4. For each case, specify the register transfer that will be executed during the next clock transition.

|    | $S_2$ | $S_1$ | $S_0$ | LD of register | Memory | Adder |
|----|-------|-------|-------|----------------|--------|-------|
| a. | 1     | 1     | 1     | IR             | Read   | —     |
| b. | 1     | 1     | 0     | PC             | —      | —     |
| c. | 1     | 0     | 0     | DR             | Write  | —     |
| d. | 0     | 0     | 0     | AC             | —      | Add   |

8. The following register transfers are to be executed in the system of Fig. 1-4. For each transfer, specify: (1) the binary value that must be applied to bus select inputs $S_2$, $S_1$ and $S_0$; (2) the register whose LD control input must be active (if any); (3) a memory read or write operation (if needed); and (4) the operation in the adder and logic circuit (if any).

   a. $AR \leftarrow PC$

   b. $IR \leftarrow M[AR]$

    c.  M[AR] ← TR

    d.  AC ← DR, DR ← AC    (done simultaneously)

9.    What are the two instructions needed in the basic computer in order to set the E flip-flop to 1?

10.    The content of AC in the basic computer is hexadecimal A937 and the initial value of E is 1. Determine the contents of AC, E, PC, AR, and IR in hexadecimal after the execution of the CLA instruction. Repeat 11 more times, starting from each one of the register-reference instructions. The initial value of PC is hexadecimal 021.

# CHAPTER -II

# REGISTER TRANSFER AND MICROOPERATIONS

**Author: Dr. Manoj Duhan**                                      **Vetter: Mr. Sandeep Arya**

## 2.1 REGISTER TRANSFER LANGUAGE

A digital system is an interconnection of digital hardware modules that accomplish a specific information-processing task. Digital systems vary in size and complexity from a few integrated circuits to a complex of interconnected and interacting digital computers. Digital system design invariably uses a modular approach. The modules are constructed from such digital components as registers, decoders, arithmetic elements, and control logic. The various modules are interconnected with common data and control paths to form a digital computer system.

Digital modules are best defined by the registers they contain and the operations that are performed on the data stored in them. The operations executed on date stored in registers are called microoperations. A microoperation is an elementary operation performed on the on the information stored in one or more registers. The result of the operation may replace the previous binary information of a register or may be transferred to another register. Examples of microoperations are shift, count, clear, and load. Some of he digital components introduced here in this chapter are registers that implement microoperations. For example, a counter with parallel load is capable of performing the micro operations increment and load. A bidirectional shift register is capable of performing the shift right and shift left microoperations.

The internal hardware organization of a digital computer is best defined by specifying

1. The set of registers it contains and their function.

2. The sequence of microoperations performed on the binary information stored in the registers.

3. The control that initiates the sequence of microoperations.

It is possible to specify the sequence of microoperations in a computer by explaining every operation in words, but this procedure usually involves a lengthy descriptive explanation. It is more convenient to adopt a suitable symbolic representation to describe the sequence of transfers between registers and the various arithmetic and logic microoperations associated with the transfers. The use of symbols instead of a narrative explanation provides an organized and concise manner for listing the microoperation sequences in registers and the control functions that initiate them.

The symbolic notation used to describe the microoperation transfers among registers is called a register transfer language. The term "register transfer" implies the availability of hardware logic circuits that can perform a stated microoperation and transfer the result of the operation to the same or another register. The word "language" is borrowed from programmers, who apply this term to programming languages. A programming language is a procedure for writing symbols to specify a given computational process. Similarly, a natural language such as English is a system for writing symbols and combining them into words and sentences for the purpose of communication between people. A register transfer language is a system for expressing in symbolic form the microoperation sequences among the registers of a digital module. It is a convenient tool for describing the internal organization of digital computers in concise and precise manner. It can also be used to facilitate the design process of digital systems

The register transfer language adopted here is believed to be as simple as possible, so it should not take very long to memorize. We will proceed to define symbols for various types of microoperations, and at the same time, describe associated hardware that can implement the stated

microoperations. The symbolic designation introduced in this chapter will be utilized in subsequent chapters to specify the register transfers, the microoperations, and the control functions that describe the internal hardware organization of digital computers. Other symbology in use can easily be learned once this language has become familiar, for most of the differences between register transfer languages consist of variations in detail rather than in overall purpose.

## 2.2 REGISTER TRANSFER

Computer registers are designated by capital letters (sometimes followed by numerals) to denote the function of the register. For example, the register that holds an address for the memory unit is usually called a memory address register and is designated by the name MAR. Other designations for registers are PC (for program counter), IR (for instruction register, and R1 (for processor register). The individual flip-flops in an n-bit register are numbered in sequence from 0 through $n - 1$, starting from 0 in the rightmost position and increasing the numbers toward the left. Figure 2-1 shows the representation of registers in block diagram form. The most common way to represent a register is by a rectangular box with the name of the register inside, as in Fig. 2-1(a). The individual bits can be distinguished as in (b). The numbering of bits in a 16-bit register can be marked on top of the box as shown in (c). A 16-bit register is partitioned into two parts in (d). Bits 0 through 7 are assigned the symbol L (for low byte) and bits 8 through 15 are assigned the symbol H (for high byte). The name of the 16-bit register is PC. The symbol PC (0–7) or PC(L) refers to the low-order byte and PC(8–15) or PC(H) to the high-order byte.
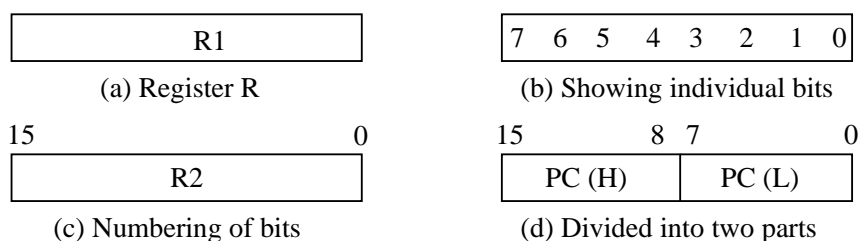
Information transfer from one register to another is designated in symbolic form by means of a replacement operator. The statement denotes a transfer of the content of

$$R2 \leftarrow R1$$

register R1 into register R2. It designates a replacement of the content of R2 by the content of R1. By definition, the content of the source register R1 does not change after the transfer.

A statement that specifies a register transfer implies that circuits are available from the outputs of the source register to the inputs of the destination register and that the destination register has a parallel load capability. Normally, we want the transfer to

**Figure 2.1** Block diagram of register.

| R1 |
|---|

(a) Register R

| 7 | 6 | 5 | 4 | 3 | 2 | 1 | 0 |
|---|---|---|---|---|---|---|---|

(b) Showing individual bits

15                                    0

| R2 |
|---|

(c) Numbering of bits

15              8 7              0

| PC (H) | PC (L) |
|---|---|

(d) Divided into two parts

Occur only under a predetermined control condition. This can be shown by means of an if-then statement.

$$\text{If } (P = 1) \text{ then } (R2 \leftarrow R1)$$

where P is a control signal generated in the control section. It is sometimes convenient to separate the control variables from the register transfer operation by specifying a control function. A control function is a Boolean variable that is equal to 1 or 0. The control function is included in the statement as follows

$$P : R2 \leftarrow R1$$

The control condition is terminated with a colon. It symbolizes the requirement that the transfer operation be executed by the hardware only if P = 1.

Every statement written in a register transfer notation implies a hardware construction for implementing the transfer. Figure 2.2 shows the block diagram that depicts the transfer from R1 to R2. The letter n will be used to indicate any number of bits for the register. It will be replaced by an actual number when the length of the register is known. Register R2 has a load input that is activated by the control variable P. It is assumed that the control variable is synchronized with the same clock as the one applied to the register. As shown in the timing diagram, P is activated in the control
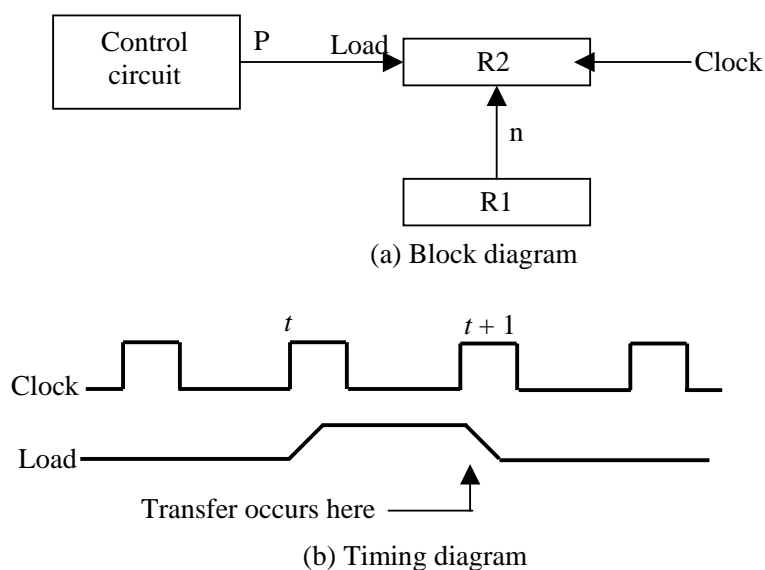


(a) Block diagram



(b) Timing diagram

Fig: 2.2

section by the rising edge of a clock pulse at time t. The next positive transition of the clock at time t + 1 finds the load input active and the data inputs of R2 are then loaded into the register in parallel. P may go back to 0 at time t + 1; otherwise, the transfer will occur with every clock pulse transition while P remains active.

Note that the clock is not included as a variable in the register transfer statements. It is assumed that all transfers occur during a clock edge transition. Even though the control condition such as P becomes active just after time t, the actual transfer dose not occur until the register is triggered by the next positive transition of the clock at time t + 1.

The basic symbols of the register transfer notation are listed in Table 2-1. Registers are denoted by capital letters, and numbers may follow the letters. Parentheses are used to denote a part of a register by specifying the range of bits or by giving a symbol name to a portion of a register. The arrow denotes a transfer of information and the direction of transfer. A comma is used to separate two or more operations that are executed at the same time. The statement

$$T: R2 \leftarrow R1, \quad R1 \leftarrow R2$$

denotes an operation that exchanges the contents of two registers during one common clock pulse provided that T = 1. This simultaneous operation is possible with registers that have edge-triggered flip-flops.

**TABLE** 2-1 Basic Symbols for Register Transfers

| Symbol | Description | Examples |
|---|---|---|
| Letters | Denotes a register | MAR, R2 |
| (and numerals) | | |
| Parentheses ( ) | Denotes a part of a register | R2(0–7), R2 (L) |
| Arrow ← | Denotes transfer of information | R2 ← R1 |
| Comma, | Separates two microoperations | R2 ← R1, R1 ← R2 |

## 2-3 BUS AND MEMORY TRANSFERS

A typical digital computer has many registers, and paths must be provided to transfer in formation form one register to another. The number of wires will be excessive if separate lines are used between each register and all other registers in the  system. A more efficient scheme for transferring information between registers in a multiple-register configuration is a common bus system. A bus structure consists of a set of common lines, one for each bit of a register, through which binary information is transferred one at a time. Control signals determine which register is selected by the onus during each particular register transfer.

One way of constructing a common bus system is with multiplexers. The multiplexers select the source register whose binary information is then placed on the bus. The construction of a  bus system for four registers is shown is Fig. 2-3. Each register has four bits, numbered 0  through 3. The bus consists of four  $4 \times 1$ multiplexers each having four data inputs, 0 through 3, and two selection inputs, $S_1$ and $S_0$. In order not to complicate the diagram with 16 lines crossing each other, we use labels to show the connections from the outputs of the registers to the inputs of the multiplexers. For example, output 1 of register A is connected to input 0 of MUX 1 because this input is labeled $A_1$. The diagram shows that the bits in the same significant position in each register are connected to the data inputs of one multiplexer to form one line of the bus. Thus MUX 0 multiplexes the four 0 bits of the registers, MUX 1 multiplexes the four 1 bits of the registers, and similarly for the other two bits.

The two selection lines $S_1$ and $S_0$ are connected to the selection inputs of all four multiplexers. The selection lines choose the four bits of one register and transfer them into the four-line common bus. When $S_1 S_0 = 00$, the 0 data inputs of all four multiplexers are selected and applied to the outputs that form the bus. This causes the bus lines to receive the content of register A since the outputs of this register are connected to the 0data inputs of the multiplexers. Similarly, register B is selected if $S_1 S_0 = 01$, and so on. Table 2-2 shows the register that is selected by the bus for each of the four possible binary value of the selection lines.

In general, a bus system will multiples k registers of n bits each to produce an n-line common bus. The number of multiplexers needed to construct the bus is equal to n, the number of bits in each register. The size of each multiplexer must be k ×1 since it multiplexes k data lines. For example, a common bus for eight registers of 16 bits each requires 16 multiplexers, one for each line in the bus. Each multiplexer must have eight data input lines and three selection lines to multiplex one significant bit in the eight registers.
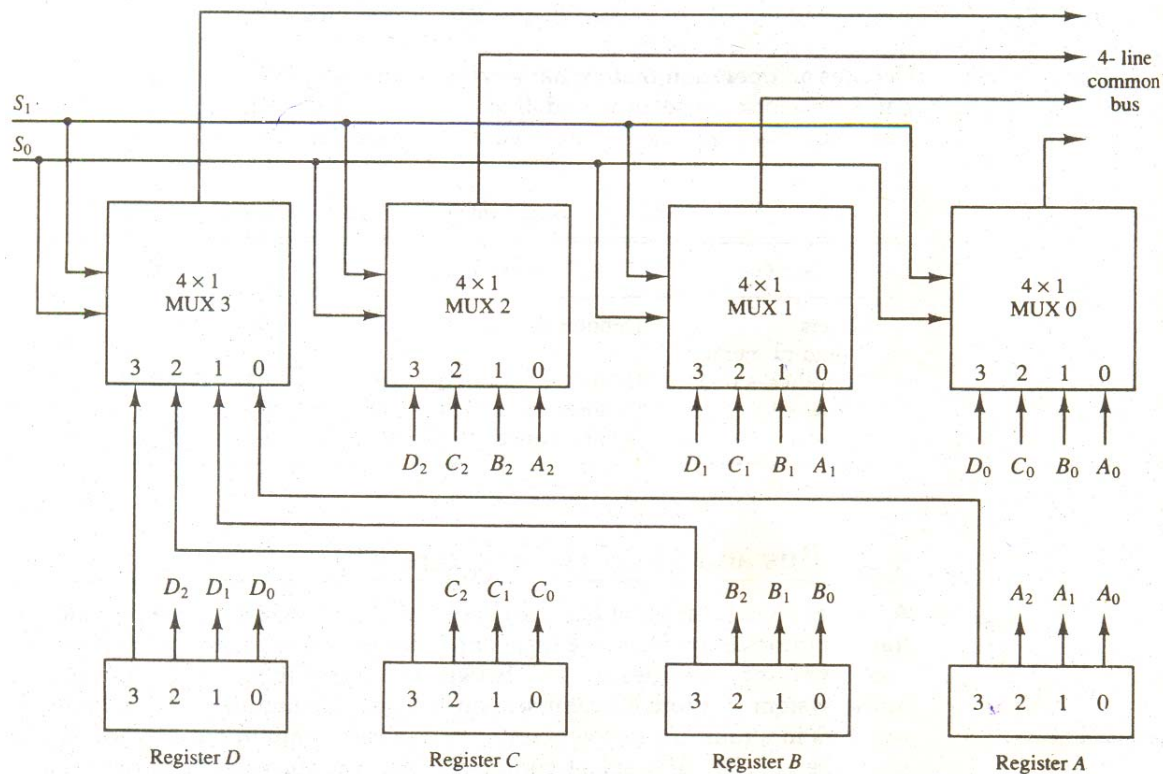
**Figure 2-3** Bus system for four registers.



TABLE 2-2

Function Table for Bus of Fig. 2-3

| $S_1$ | $S_0$ | Register selected |
|-------|-------|-------------------|
| 0 | 0 | A |
| 0 | 1 | B |
| 1 | 0 | C |
| 1 | 1 | D |

The transfer of information from a bus into one of many destination registers can be accomplished by connecting the bus lines to the inputs of all destination registers and activating the load control of the particular destination register selected. The symbolic statement for a bus transfer may mention the bus or its presence may be implied in the statement. When the bus is includes in the statement, the register transfer is symbolized as follows:

$$BUS \leftarrow C, \qquad R1 \leftarrow BUS$$

The content of register C is placed on the bus, and the content of the bus is loaded into register R1 by activating its load control input. If the bus is known to exist in the system, it may be convenient just to show the direct transfer.

$$R1 \leftarrow C$$

From this statement the designer knows which control signals must be activated to produce the transfer through the bus.

### 2.3.1 THREE-STATE BUS BUFFERS

A bus system can be constructed with three-state gates instead of multiplexes. A three-state gate is a digital circuit that exhibits three states. Two of the states are signals equivalent to logic 1 and 0 as in a conventional gate. The third state is a high-impedance state. The high-impedance state behaves like an open circuit which means that the output is disconnected and does not have a logic significance. Three-state gates may perform any conventional logic, such as AND or NAND. However, the one most commonly used in the design of a bus system is the buffer gate.
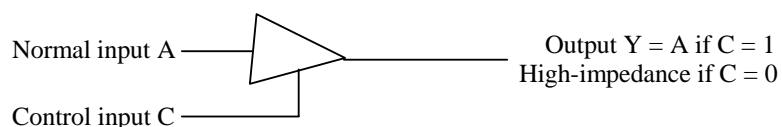
The graphic symbol of a three-state buffer gate is shown in Fig. 2-4. it is distinguished from a normal buffer by having both a normal input and a control input. The control input determines the output state. When the control input is equal to 1, the output is enabled and the gate behaves like any conventional buffer, with the output equal to the normal input. When the control input is 0, the output is disabled and the gate gives to a high-impedance state, regard less of the value in the normal input. The high-impedance state of a three-state gate provides a special feature not available in other gates. because of this feature, a large number of three-state gate outputs can be connected with wires to form a common bus line without endangering loading effects.

The construction of a bus system with three-state buffers is demonstrated in Fig. 2-5. The outputs of four buffers are connected together to form a single bus line. (It must be realized that this type of connection cannot be done with gates that do not have three-state outputs.) the control inputs to the buffers determine which of the four normal inputs will communicate with the us line. No more than one buffer may be in the active state at any given time. ;the connected buffers must be controlled so that only one three-state buffer bas access to the bus line while all other buffers are maintained in a high-impedance state.

One way to ensure that no more than one control input is active at any given time is to use a decoder, as shown in the diagram. When the enable input of the decoder is 0, all of tits four outputs are 0, and the bus line is in a high-impedance state because all four buffers are disabled. When the enable input is active, one of the three-state buffers will be active, depending on the binary value in the select inputs of the decoder. Careful investigation will reveal that Fig. 2-5 is another way of constructing a $4 \times 1$ multiplexer since the circuit can replace the multiplexer in Fig. 2-3.

To construct a common bus for four registers on n bits each using three- state buffers, we need n circuits with four buffers in each as shown in Fig. 2-5. Each group of four buffers receives one significant bit from the four registers. Each common output produces one of the lines for the common bus for a total of n lines. Only one decoder is necessary to select between the four register.

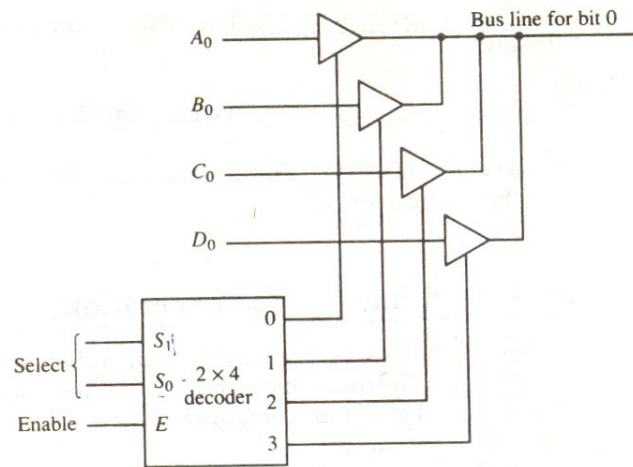**Figure 2-4** Graphic symbols for three-state buffer.

**Figure 2-5** Bus line with three state-buffers.

### 2.3.2 MEMORY TRANSFER

The transfer of information from a memory word to the outside environment is called a read operation. The transfer of new information to be stored into the memory is called a write operation. A memory word will be symbolized by the letter M. The particular memory word among the many available is selected by the memory address during he transfer. It is necessary to specify the address of M when writing memory transfer operations; this will be done by enclosing the address in square brackets following the letter M.

Consider a memory unit that receives the address form a register, called the address register, symbolized by AR. The data are transferred to another register, called the data register, symbolized by DR. the read operation can be stated as follows:

$$\text{Read:} \quad DR \leftarrow M\,[AR]$$

This causes a transfer to information into DR from the memory word M selected by the address in AR.

The write operation transfers the content of a data register to a memory word M selected by the address. Assume that the input data are in register R1 and the address is in AR. The write operation can be stated symbolically as follows:

$$\text{Write: } M\,[AR] \leftarrow R1$$

This causes transfer of information from R1 into the memory word M selected by the address in AR.

### 2-4 ARITHMETIC MICRO OPERATIONS

A micro operation is an elementary operation performed with the data stored in registers. The microoperations most often encountered in digital computers are classified into four categories:

1.      Register transfer microoperations transfer binary information from one register to another.

2.      Arithmetic micro operations perform arithmetic operation on numeric data stored in registers.

3.      Logic micro operations perform bit manipulation operations on non-numeric data stored in registers.

4.       Shift microoperations perform shift operations on data stored in registers.

The register transfer microoperation was introduced in Sec. 1-2. This type of microoperation dies not change the information content when the binary information moves from the source register to the destination register. The other three types of microoperation change the information content during the transfer. In this section we introduce a set of arithmetic microoperations.

The basic arithmetic microoperations are addition, subtraction, increment, decrement, and shift. Arithmetic shifts are explained later in conjunction with the shift microoperations. The arithmetic microoperation defined by the statement

$$R3 \leftarrow R1 + R2$$

specifies an add microoperation. It states that the contents of register R1 are added to the contents of register R2 and the sum transferred to register R3. To implement this statement with hardware we need three registers and the digital component that performs the addition operation. The other basic arithmetic microoperations are listed in Table 4-3. Subtraction is most often implemented through complementation and addition. Instead of using the minus operator, we can specify the subtraction by the following statement:

$$R3 \leftarrow R1 + \overline{R2} + 1$$

$\overline{R2}$ is the symbol for the 1's complement of R2. Adding 1 to the 1's complement produces the 2's complement. Adding the contents of R1 to the 2's complement of R2 is equivalent to R1 – R2.

TABLE 2-3 Arithmetic Microoperations

| Symbolic designation | Description |
| --- | --- |
| $R3 \leftarrow R1 + R2$ | Contents of R1 plus R2 transferred to R3 |
| $R3 \leftarrow R1 - R2$ | Contents of R1 minus R2 transferred to R3 |
| $R2 \leftarrow \overline{R2}$ | Complement the contents of R2 (1's complement) |
| $R2 \leftarrow \overline{R2} + 1$ | 2's complement the contents of R2 (negate) |
| $R3 \leftarrow R1 + \overline{R2} + 1$ | R1 plus the 2's complement of R2 (subtraction) |
| $R1 \leftarrow R1 + 1$ | Increment the contents of R1 by one |
| $R1 \leftarrow R1 - 1$ | Decrement the contents of R1 by one |

The increment and decrement microoperations are symbolized by plus-one  and minus-one operations, respectively. These microoperations are implemented with a combinational circuit or with a binary up-down counter.

The arithmetic operations of multiply and divide are not listed in Table 2-3. These two operations are valid arithmetic operations but are not included in the basic set of microoperations. The only place where these operations can be considered as microoperations is in a digital system, where they are implemented by means of a combinational circuit. In such a case, the signals that perform these operations propagate through gates, and the result of the operation can be transferred into a destination register by a clock pulse as soon as the output signal propagates through the combinational circuit. In most computers, the multiplication operation is implemented with a sequence of add and shift microoperations. Division is implemented with a sequence of subtract and

shift microoperations. To specify the hardware in such a case requires a list of statements that use the basic microoperations of add, subtract, and shift (see Chapter 10).

### 2.4.1 BINARY ADDER

To implement the add microoperation with hardware, we need the registers that hold the data and the digital component that performs the arithmetic addition. The digital circuit that forms the arithmetic sum of two bits and a previous carry is called a full-adder (see Fig. 2-6). The digital circuit that generates the arithmetic sum of two binary numbers of any length is called a binry adder. The binary adder is constructed with full-adder circuits
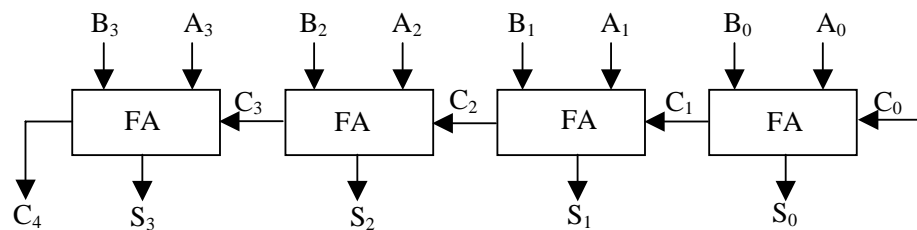


Figure 2-6 4-bits binary adder.

connected in cascade, with the output carry from one full-adder connected to the input carry of the next full-adder. Figure 2-6 shows the interconnections of four full-adders (FA) to provide a 4-bit binary adder. The augends bits of A and the addend bits of B are designated by subscript numbers from right to left, with subscript 0 denoting the low-order bit. The carries are connected in a chain through the full-adders. The input carry to the binary adder is $C_0$ and the output carry is $C_4$. This S outputs of the full-adders generate the required sum bits.

An n-bit binary adder requires n full –adders. The output carry from each full-adder is connected to the input carry of the next-high-order full-adder. The n data bits for the A inputs come from tone register (such as R1), and the n data bits for the B inputs come from another register (such as R2). The sum can be transferred to a third register or to one of the sourc3e registers (R1 or R2), replacing its previous content,

### 2.4.2 Binary Adder-Subtractor

The subtraction of binary numbers can be done most conveniently by means of complements. Remember that the subtraction A –B can be done by taking the 2's complement of B and adding it to A. The 2's complement can be obtained by taking the 1's complement and adding one to the least significant pair of bits. The 1's complement can be implemented with inverters and a one can be added to the sum through the input carry.

The addition and subtraction operations can be combined into one common circuit by including an exclusive-OR gate with each full-adder. A 4-bit adder-subtract or circuit is shown in Fig. 2-7. The mode input M controls the operation. When M = 0 the circuit is an adder and when M = 1 the circuit becomes a subtractor. Each exclusive-OR gate receives input M and one of the inputs of B. When M = 0, we have $B \oplus 0 = B$. the full-adders receive the value of B, the input carry is 0, and the circuit performs A plus B. When M = 1, we have $B \oplus 1 = B'$ and $C_0 = 1$. The B inputs are all complemented and a 1 is added through the input carry. The circuit performs the operation A plus the
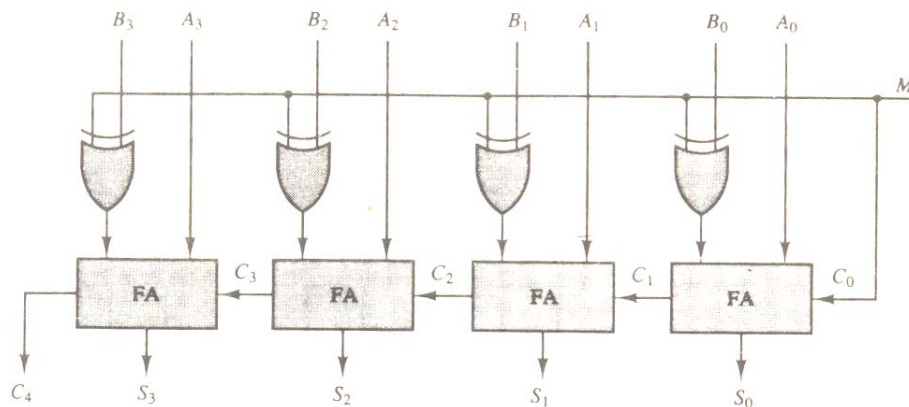
**Figure 2-7** 4-bit adder - subtractor.

2's complement of B. For unsigned, this gives A – B if A ≥ B or the 2's complement of (B−A) if A < B. For signed numbers, the result is A – B provided that there is no overflow.

### 2.4.3    BINARY INCREMENTER

The increment microoperation adds one to a number in a register. For example, if a 4-bit register ahs a binary value 0110, it will go to 0111 after it is incremented. This microoperation is easily implemented with a binary counter. Every time the count enable is active, the clock pulse transition increments the content of the register by one. There may be occasions when the increment microoperation must be done with a combinational circuit independent of a particular register. This can be accomplished by means of half-adders connected in cascade.

The diagram of a 4-bit combinational circuit incrementer is shown in Fig. 2-8. One of the inputs to the least significant half-adder (HA) is connected to logic-1 and the other input is connected to the least significant bit of the number to be incremented. The output carry from one half-adder is connected to one of the inputs of the next-higher-order half-adder. The circuit receives the four bits from $A_0$ through $A_3$, adds one to it, and generates the incremented output in $S_0$ through $S_3$. The output carry $C_4$ will be 1 only after incrementing binary 1111. This also causes outputs $S_0$ through $S_3$ to go to 0.

The circuit of Fig. 2-8 can be extended to an N-bit binary incremented by extending the
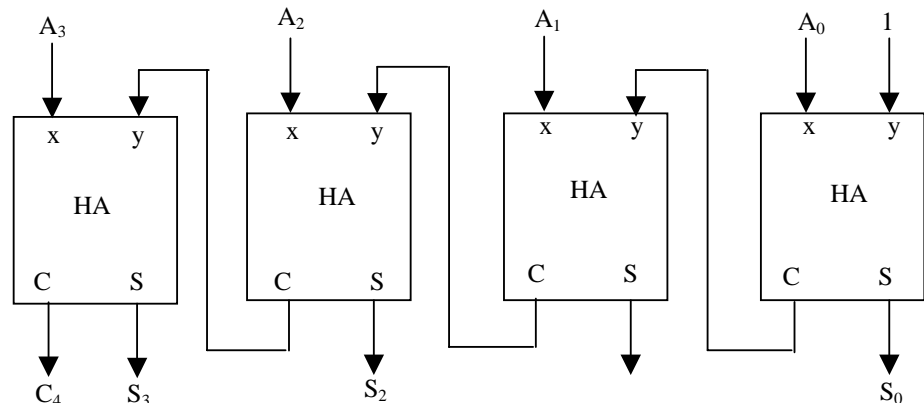


**Figure  2-8**  4-bit binary incrementer.

diagram to include n half-adders. The least significant bit must have one input connected to logic-. The other inputs receive the number to be incremented or the carry from the previous stage.

### 2.4.5   ARITHMETIC CIRCUIT

The arithmetic microoperations listed in Table 2-3 can be implemented in one composite arithmetic circuit. The basic component of an arithmetic circuit is the parallel adder. By controlling the data inputs to the adder, it is possible to obtain different types of arithmetic operations.

The diagram of a 4-bit arithmetic circuit is shown in Fig. 2-9. It has four full-adder circuits that constitute the 4-bit adder and four multiplexers for choosing different operations. There are two 4-bit inputs A and B and a 4-bit output D. The four inputs from A go directly to the X inputs of the binary adder. Each of the for inputs from B are connected to the data inputs of the multiplexers. The multiplexer's data inputs also receive the complement of B. The other two data inputs are connected to logic-0 ad logic -1. Logic-0 is fixed voltage value (0 volts for TTL integrated circuits) and the logic-1 signal can be generated through an inverter whose input is 0. The four multiplexers are controlled by two selection inputs, $S_1$ and $S_0$. The input carry $C_{in}$ goes to the carry input of the FA in the least significant position. The other carries are connected from one stage to the next.

The output of the binary adder is calculated from the following arithmetic sum:

$$D = A + Y + C_{in}$$

where A is the 4-bit binary number at the X inputs and Y is the 4-bit binary number at the Y inputs of the binary adder. $C_{in}$ is the input carry, which can be equal to 0 or 1. Note that the symbol + in the equation above denotes an arithmetic plus. By controlling the value of Y with the two selection inputs $S_1$ and $S_0$ ad making $C_{in}$ equal to 0 or 1, it is possible to generate the eight arithmetic microoperations listed in Table 2-4.

**TABLE 2-4**  Arithmetic Circuit Function Table

| Select | | | Input | Output | |
|---|---|---|---|---|---|
| $S_1$ | $S_0$ | $C_{in}$ | Y | $D = A + Y + C_{in}$ | Microoperation |
| 0 | 0 | 0 | B | $D = A + B$ | Add |
| 0 | 0 | 1 | B | $D = A + B + 1$ | Add with carry |
| 0 | 1 | 0 | $\overline{R}$ | $D = A + \overline{B}$ | Subtract with borrow |
| 0 | 1 | 1 | $\overline{R}$ | $D = A + \overline{B} + 1$ | Subtract |
| 1 | 0 | 0 | 0 | $D = A$ | Transfer A |
| 1 | 0 | 1 | 0 | $D = A + 1$ | Increment A |
| 1 | 1 | 0 | 1 | $D = A - 1$ | Decrement A |
| 1 | 1 | 1 | 1 | $D = A$ | Transfer A |

When $S_1 S_0 = 00$, the value of B is applied to the Y inputs of the adder. If $C_{in} = 0$, the output $D = A + B$. If $C_{in} = 1$, output $D = A + B + 1$. Both cases perform the add microoperation with or without adding the input carry.

When $S_1 S_0 = 01$, the complement of B is applied to the Y inputs of the adder. If $C_{in} = 1$, then $D = A + \overline{B} + 1$. This produces A plus the 2's complement of B, which is equivalent to a subtract with borrow, that is, $A - B - 1$.

**Figure 2-9** 4-bit arithmetic circuit

When $S_1 S_0 = 10$, the input from B are neglected, and instead, all 0's are inserted into the Y inputs. The output becomes $D = A + 0 + C_{in}$. This gives $D = A$ when $C_{in} = 0$ and $D = A + 1$ when $C_{in} = 1$. In the first case we have a direct transfer from input A to output D. In the second case, the value of A is incremented by 1.

When $S_1 S_0 = 11$, all 1's are inserted into the Y inputs of the adder to produce the decrement operation $D = A - 1$ when $C_{in}$. This is because a number with all 1's is equal to the 2's

complement of 1 (the 2's complement of binary 0001 is 1111). Adding a number A to the 2's complement of 1 produces F = A +2's complement of 1 = A − 1. When $C_{in}$ = 1, then D = A − 1 + 1 = A, which causes a direct transfer from input A to output D. Note that the microoperation D = A is generated twice, so there are only seven distinct microoperations in the arithmetic circuit.

## 2-5 LOGIC MICROOPERATIONS

Logic micro operations specify binary operations for strings of bits stored in registers. These operations consider each it of the register separately and treat them as binary variables. For example, the exclusive-OR microoperation with the contents of two registers R1 and R2 is symbolized by the statement

$$P: R1 \leftarrow R1 \oplus R2$$

It specifies a logic microoperation to be executed on the individual bits of the registers provided that the control variable P = 1. As a numerical example, assume that each register has four bits. Let the content of R1 be 1010 and the content of R2 be 1100. The exclusive-OR microoperation stated above symbolizes the following logic computation:

$$
\begin{array}{ll}
1010 & \text{Content of R1} \\
1100 & \text{Content of R2} \\
\overline{0110} & \text{Content of R1 after P = 1}
\end{array}
$$

The content of R1, after the execution of the microoperation, is equal to the bit-by-bit exclusive-OR operation on pairs of bits in R2 and previous values of R1. the logic microoperations are seldom used in scientific computations, but they are very useful for bit manipulation of binary data and for making logical decisions.

Special symbols will be adopted for the logic microoperations OR, AND, and complement, to distinguish them from the corresponding symbols used to express Boolean functions. The Symbol $\vee$ will be used to denote an OR microoperation and the symbol $\wedge$ to denote an AND microoperation. The complement microoperation is the same as the 1's complement and uses a bar on top of the symbol that denotes the register name. By using different symbols, it will be possible to differentiate between a logic microoperation and a control (or Boolean) function. Another reason for adopting two sets of symbols is to be able to distinguish the symbol +, when used to symbolize an arithmetic plus, from a logic OR operation. Although the + symbol has two meaning, it will be possible to distinguish between them by noting where the symbol occurs. When the symbol + occurs in a control (or Boolean) function, it will denote an OR operation. We will never use it to symbolize an OR microoperation. For example, in the statement

$$P + Q: R1 \leftarrow R2 + R3, R4 \leftarrow R5 \vee R6$$

the + between P and Q is an OR operation between two binary variables of a control function. The + between R2 and R3 specifies an add microoperation. The OR microoperation is designated by the symbol $\vee$ between register R5 and R6.

## 2.5.1   LIST OF LOGIC MICROOPERATIONS

There are 16 different logic operations that can be performed with two binary variables. They can be determined from all possible truth tables obtained with two binary variables as shown in Table 2-5. in this table, each of the 16 columns $F_0$ through $F_{15}$ represents a truth table of one possible Boolean function for the

**Table 2-5** Truth Tables for 16 Functions of Two Variables

| x | y | $F_0$ | $F_1$ | $F_2$ | $F_3$ | $F_4$ | $F_5$ | $F_6$ | $F_7$ | $F_8$ | $F_9$ | $F_{10}$ | $F_{11}$ | $F_{12}$ | $F_{13}$ | $F_{14}$ | $F_{15}$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 |
| 1 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 1 | 1 | 0 | 0 | 1 | 1 | 0 | 0 | 1 | 1 |
| 1 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 |

two variables x and y. Note that the functions are determined from the 16 binary combinations that can be assigned to F.

The 16 Boolean functions of two variables x and y are expressed in algebraic form in the first column of Table 2-6. The 16 logic microoperations are derived from these functions by replacing variable x by the binary content of register A and variable y by the binary content of register B. It is important to realize that the Boolean functions listed in the first column of Table 2-6 represent a relationship between two binary variables x and y. The logic microoperations listed in the second column represent a relationship between the binary content of two registers A and B. Each bit of the register is treated as a binary variable and the microoperation is performed on the string of bits stored in the registers.

**TABLE 2-6**     Sixteen Logic Micro operations

| Boolean function | Microoperation | Name |
|---|---|---|
| | $F \leftarrow 0$ | Clear |
| $F_0 = 0$ | $F \leftarrow A \wedge B$ | AND |
| $F_1 = xy$ | $F \leftarrow A \wedge \overline{B}$ | |
| $F_2 = xy'$ | $F \leftarrow A$ | Transfer A |
| $F_3 = x$ | $F \leftarrow \overline{A} \wedge B$ | |
| $F_4 = x'y$ | $F \leftarrow B$ | Transfer B |
| $F_5 = y$ | $F \leftarrow A \oplus B$ | Exclusive-OR |
| $F_6 = x \oplus y$ | $F \leftarrow A \vee B$ | OR |
| $F_7 = x + y$ | $F \leftarrow \overline{A \vee B}$ | NOR |
| $F_8 = (x + y)'$ | $F \leftarrow \overline{A \oplus B}$ | |
| $F_9 = (x \oplus y)'$ | $F \leftarrow \overline{B}$ | Exclusive-NOR |
| $F_{10} = y'$ | $F \leftarrow A \vee \overline{B}$ | Complement B |
| $F_{11} = x + y'$ | $F \leftarrow \overline{A}$ | |
| $F_{12} = x'$ | $F \leftarrow A \vee B$ | Complement A |
| $F_{13} = x' + y$ | | |

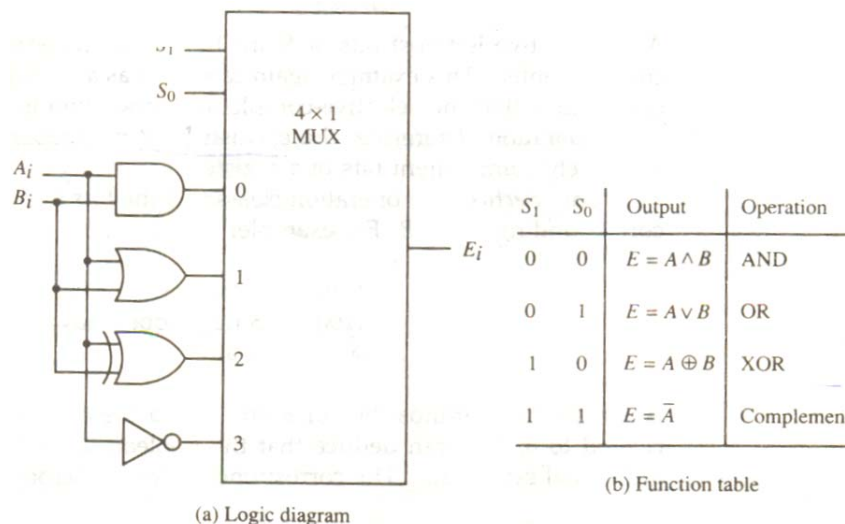| | | NAND |
|---|---|---|
| $F_{14} = (xy)'$ | $F \leftarrow \overline{A \wedge B}$ | |
| | | Set to all 1's |
| $F_{15} = 1$ | $F \leftarrow$ all 1's | |

## 2.5.2   HARDWARE IMPLEMENTATION

The hardware implementation of logic microoperations requires that logic gates be inserted for each bit or pair of bits in the registers to perform the required logic function. Although there are 16 logic microoperations, most computers use only four --- AND, OR, XOR (exclusive-OR), and complement by which all others can be derived.

Figure 4-10 shows one stage of a circuit that generates the four basic logic microoperations. It consists of four gates and a multiplexer. Each of the four logic operations is generated through a gate that performs the required logic. The outputs of the gates are applied to the data inputs of the multiplexer. The two selection inputs $S_1$ and $S_0$ choose one of the data inputs of the multiplexer and direct its value to the output. The diagram shows one typical stage with subscript i. For a logic circuit with n bits, the diagram must be repeated n times for $i = 0, 1, 2, \ldots, N-1$. The selection variables are applied to all stages. The function table in Fig. 2-10 (b) lists the logic microoperations obtained for each combination of the selection variables.

## 2.5.3   SOME APPLICATIONS

Logic micro operations are very useful for manipulating individual bits or a portion of a word stored in a register. They can be used to change bit values, delete a group of bits, or insert new bit values into a register. The following examples show how the bits of one register (designated by A) are manipulated

**Figure 2-10**  One stage of logic circuit



| $S_1$ | $S_0$ | Output | Operation |
|---|---|---|---|
| 0 | 0 | $E = A \wedge B$ | AND |
| 0 | 1 | $E = A \vee B$ | OR |
| 1 | 0 | $E = A \oplus B$ | XOR |
| 1 | 1 | $E = \overline{A}$ | Complement |

(b) Function table

(a) Logic diagram

by logic microoperations as a function of the bits of another register (designated by B). In a typical application, register A is a processor register and the bits of register B constitute a logic operand extracted from memory and placed in register B.

The selective-set operation sets to 1 the bits in register A where there are corresponding 1's in register B. It does not affect bit positions that have 0's in B. The following numerical example clarifies this operation:

$$
\begin{array}{ll}
1010 & \text{A before} \\
\underline{1100} & \text{B(logic operand)} \\
1110 & \text{A after}
\end{array}
$$

The two leftmost bits of B are 1's, so the corresponding bits of A are set to 1. One of these two bits was already set and the other has been changed from 0 to 1. The two bits of A with corresponding 0's in B remain unchanged. The example above serves as a truth table since it has all four possible combinations of two binary variables. From the truth table we note that the bits of A after the operation are obtained from the logic-OR operation of bits in B and previous values of A. therefore, the OR microoperation can be used to selectively set bits of a register.

The selective-complement operation complements bits in A where there are corresponding 1's in B. It does not affect bit positions that have 0's in B. For example:

$$
\begin{array}{ll}
1010 & \text{A before} \\
\underline{1100} & \text{b (logic operand)} \\
0110 & \text{A after}
\end{array}
$$

Again the two leftmost bits of B are 1's, so the corresponding bits of A are complemented. This example again can serve as a truth table from which one can deduce that the selective-complement operation is just an exclusive-OR microoperation. Therefore, the exclusive-OR microoperation can be used to selectively complement bits of a register.

The selective-clear operation clears to 0 the bits in A only where there are corresponding 1's in B. or example:

$$
\begin{array}{ll}
1010 & \text{A before} \\
\underline{1100} & \text{B (logic operand)} \\
0010 & \text{A after}
\end{array}
$$

Again the two leftmost bits of B are 1's, so the corresponding bits of A are cleared to 0. One can deduce that the Boolean operation performed on the individual bits is AB'. The corresponding logic microoperation is

$$ A \leftarrow A \wedge \overline{B} $$

The mask operation is similar to the selective-clear operation except that the bits of A are cleared only where there are corresponding 0's in B. The mask operation is an AND micro operation as seen from the following numerical example:

$$
\begin{array}{ll}
1010 & \text{A before} \\
\underline{1100} & \text{B (logic operand)} \\
1000 & \text{A after masking}
\end{array}
$$

The two rightmost bits of A are cleared because the corresponding bits of B are 0's. The two leftmost bits are left unchanged because the corresponding bits of B are 1's. The mask operation is more convenient to use than the selective-clear operation because most computers provide an AND instruction, and few provide an instruction that executes the microoperation for selective-clear.

The insert operation inserts a new value into a group of bits. This is done by first masking the bits and then ORing them with the required value. For example, suppose that an A register contains eight bits, 0110 10101. To replace the four leftmost bits by the value 1001 we first mask the four unwanted bits:

$$
\begin{array}{ll}
0110\ 1010 & \text{A before} \\
\underline{0000\ 1111} & \text{B (mask)} \\
0000\ 1010 & \text{A after masking}
\end{array}
$$

and then insert the new value:

$$
\begin{array}{ll}
0000\ 1010 & \text{A before} \\
\underline{1001\ 0000} & \text{B (insert)} \\
1001\ 1010 & \text{A after insertion}
\end{array}
$$

The mask operation is an AND microoperation and the insert operation is an OR microoperation.

The clear operation compares the words in A and B and produces an all 0's result if the two numbers are equal. This operation is achieved by an exclusive-OR microoperation as shown by the following example:

$$
\begin{array}{ll}
1010 & \text{A} \\
\underline{1010} & \text{B} \\
0000 & \text{A} \leftarrow \text{A} \oplus \text{B}
\end{array}
$$

when A and B are equal, the two corresponding bits are either both 0 or both 1. In either case the exclusive-OR operation produces a 0. The all-0's result is then checked to determine if the two numbers were equal.

## 2-6 SHIFT MICRO OPERATIONS

Shift micro operations are used for serial transfer of data. They are also used in conjunction with arithmetic, logic, and other data-processing operations. The contents of a register can be shifted to the left or the right. At the same time that the bits are shifted, the first flip-flop receives its binary information from the serial input. During a shift-left operation the serial input transfers a bit into the right most position. During a shift-right operation the serial input transfers a bit into the leftmost position. The information transferred through the serial input determines the type of shift. There are three types of shifts: logical, circular, and arithmetic.

A logical shift is one that transfers 0 through the serial input. We will adopt the symbols **shl** and **shr** for logical shift-left and shift-right microoperations. For example:

$$
\begin{array}{l}
\text{R1} \leftarrow \text{Sh1 R1} \\
\text{R2} \leftarrow \text{shr R2}
\end{array}
$$

are two microoperations that specify a 1-bit shift to the left of the content of register R1 and a 1-bit shift to the right of the content of register R2. The register symbol must be the same on both sides of the arrow. The bit transferred to the end position through the serial input is assumed to be 0 during a logical shift.

The circular shift (also known as a rotate operation) circulates the bits of the register around the two ends without loss of information. This is accomplished by connecting the serial output of the shift register to its serial input. We will use the symbols cil and cir for the circular shift left and right, respectively. The symbolic notation for the shift microoperation is shown in Table 2-7.

TABLE 2-7 Shift Microoperations

| Symbolic designation | Description |
|---|---|
| $R \leftarrow$ shl R | Shift-left register R |
| $R \leftarrow$ shr R | Shift-right register R |
| $R \leftarrow$ cil R | Circular shift-left register R |
| $R \leftarrow$ cir R | Circular shift-right register R |
| $R \leftarrow$ ashl R | Arithmetic shift-left R |
| $R \leftarrow$ ashr R | Arithmetic shift-right R |

An arithmetic shift is a microoperation that shifts a signed binary number to the left or right. An arithmetic shift-left multiplies a signed binary number by 2. An arithmetic shift-right divides the number by 2. Arithmetic shifts must leave the sign bit unchanged because the sign of the number remains the same when it is multiplied or
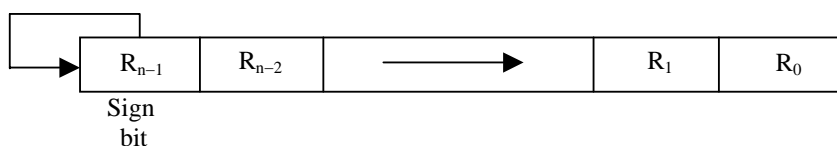


Figure 2-11 Arithmetic shift right.

divided by 2. The leftmost bit in a register holds the sign bit, and the remaining bits hold the number. The sign bit is 0 for positive and 1 for negative. Negative numbers are in 2's complement form. Figure 4-11 shows a typical register of n bits. Bit $R_{n-1}$ in the leftmost position holds the sign bit. $R_{n-2}$ is the most significant bit of the number and $R_0$ is the least significant bit. The arithmetic shift-right leaves the sign bit unchanged and shifts the number (including the sign bit) to the right. Thus $R_{n-1}$ remains the same, $R_{n-2}$ receives the bit from $R_{n-1}$, and so on tfor the other bits in the register. The bits in $R_0$ is lost.

The arithmetic shift-left inserts a 0 into $R_0$, and shifts all other bits to the left. The initial bit of $R_{n-1}$ is lost and replaced by the bit from $R_{n-2}$. A sign reversal occurs if the bit in $R_{n-1}$ changes in value after the shift. This happens if the multiplication by 2 causes an overflow. An overflow occurs after an arithmetic shift left if initially, before the shift, $R_{n-1}$ is not equal to $R_{n-2}$. An overflow flip-flop $V_s$ can be used to detect an arithmetic shift-left overflow.

$$V_s = R_{n-1} \oplus R_{n-2}$$

If $V_s = 0$, there is no overflow, bit if $V_s = 1$, there is an overflow and a sign reversal after the shift. $V_s$ must be transferred into the overflow flip-flop with the same clock pulse that shifts the register.

## SUMMARY

1. A digital system is an interconnection of digital hardware modules that accomplish a specific information-processing task.

2. Digital system design invariably uses a modular approach.

3. The symbolic notation used to describe the micro operation transfers among registers is called a register transfer language.

4. It is possible to specify the sequence of micro operations in a computer by explaining every operation in words, but this procedure usually involves a lengthy descriptive explanation.

5. Computer registers are designated by capital letters (sometimes followed by numerals) to denote the function of the register.

6. Information transfer from one register to another is designated in symbolic form by means of a replacement operator.

7. A typical digital computer has many registers, and paths must be provided to transfer in formation form one register to another. The number of wires will be excessive if separate lines are used between each register and all other registers in the system.

8. A bus system can be constructed with three-state gates instead of multiplexes. A three-state gate is a digital circuit that exhibits three states.

9. The transfer of information from a memory word to the outside environment is called a read operation. The transfer of new information to be stored into the memory is called a write operation. A memory word will be symbolized by the letter M.

10. A micro operation is an elementary operation performed with the data stored in registers.

11. To implement the add micro operation with hardware, we need the registers that hold the data and the digital component that performs the arithmetic addition.

12. The hardware implementation of logic micro operations requires that logic gates be inserted for each bit or pair of bits in the registers to perform the required logic function.

13. Logic micro operations are very useful for manipulating individual bits or a portion of a word stored in a register.

14. Shift micro operations are used for serial transfer of data. They are also used in conjunction with arithmetic, logic and other data-processing operations.

## SELF ASSESSMENT

1. What is typical digital computer? Explain it with the help of block diagram and its various parts.

2. Explain the significance of the shift micro operations.

3. What are program interrupts? Explain their significance?

4. Discuss the input and output operations of the computers system.

5. Show the block diagram of the hardware that implements the following register transfer statement:

$$yT_2; R2 \leftarrow R1, \quad R1 \leftarrow R2$$

6. Represent the following conditional control statement by two register transfer statements with control functions.

$$\text{If } (P = 1) \text{ then } (R1 \leftarrow R2) \text{ else if } (Q = 1) \text{ then } (R1 \leftarrow R3)$$

7. What has to be done to the bus system of Fig. 2-3 to be able to transfer information from any register to any other register? Specifically, show the connections that must be included to provide a path from the outputs of register C to the inputs of register A.

8. Draw a diagram of a bus system similar to the one shown in Fig. 2-3, but use three-state buffers and a decoder instead of the multiplexers.

9. A digital computer has a common bus system for 16 registers of 32 bits each. The bus is constructed with multiplexers.

   a. How many selection inputs are there in each multiplexer?

   b. What size of multiplexers are needed ?

   c. How many multiplexers are there in the bus?

10. The following transfer statements specify a memory. Explain the memory operation in each case.

    a. $R2 \leftarrow M[AR]$

    b. $M[AR] \leftarrow R3$

    c. $R5 \leftarrow M[R5]$

11. Draw the block diagram for the hardware that implements the following statements:

$$x + yz:\ AR \leftarrow AR + BR$$

where AR and BR are two n-bit registers and x, y, and z are control variables. Include the logic gates for the control function. (Remember that the symbol + designates an OR operation in a control or Boolean function but that it represents an arithmetic plus in a microoperation.)

12. Show the hardware that implements the following statement. Include the logic gates for the control function and a block diagram for the binary counter with a count enable input.

$$xyT_0 + T_1 + y'T_2:\ AR \leftarrow AR + 1$$

13. Consider the following register transfer statements for two 4-bit registers R1 and R2.

$$xT:\quad R1 \leftarrow R1 + R2$$

$$x'T:\ R1 \leftarrow R2$$

Every time that variable T = 1, either the content of R2 is added to the content of R1 if x = 1, or the content of R2 is transferred to R1 if x = 0. Draw a diagram showing the hardware implementation of the two statements. Use block diagrams for the two 4-bit registers, a 4-bit adder and a quadruple 2-to-1-line multiplexer that selects the inputs to R1. In the diagram, show how the control variables x and T select the inputs of the multiplexer and the load input of register R1.

# CHAPTER –III

# COMPUTER INSTRUCTIONS

**Author: Dr. Manoj Duhan**                                    **Vetter:  Dr. Pradeep Bhatia**

## 3.1 INTRODUCTION

The basic computer has three instruction code formats, as shown in Fig. 3-1. Each format has 16 bits. The operation code (opcode) part of the instruction contains three bits and the meaning of the remaining 13 bits depends on the operation code encountered. A memory-reference instruction uses 12 bits to specify an address and one bit to specify the addressing mode I. I is equal to 0 for direct address and to 1 for indirect address. The register-reference instructions are recognized by the operation code 111 with a 0 in the leftmost bit (bit 15) of the instruction. A register-reference instruction specifies an operation on or a test of the AC register. An operand from memory is not needed; therefore, the other 12 bits are used to specify the operation or test to be executed. Similarly, an input-output instruction does not need a reference  to memory and is recognized by the operation code 111 with a 1 in the leftmost bit of the instruction. The remaining 12 bits are used to specify the type of input-output operation or test performed.

The type of instruction is recognized by the computer control from the four bits in positions 12 through 15 of the instruction. If the three op code bits in positions 12 though 14 are not equal to 111, the instruction is a memory-reference type and the bit in position 15 is taken as the addressing mode I. If the 3-bit op code is equal to 111, control then inspects the bit in position 15. If this bit is 0, the instruction is a register-reference type. If the bit is 1, the instruction is an
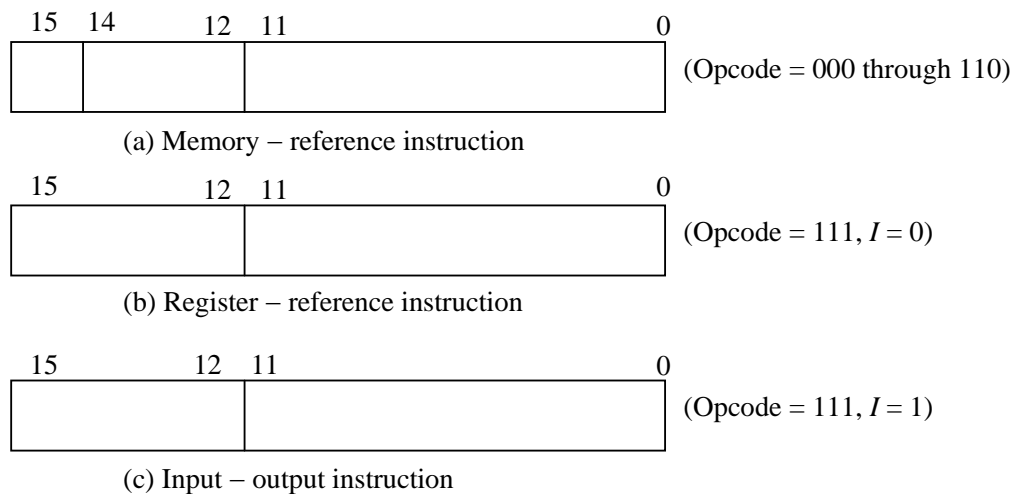
| 15 | 14 | 12 | 11 | | 0 | |
|----|----|----|----|----|----|----|

(Opcode = 000 through 110)

(a) Memory – reference instruction

| 15 | | 12 | 11 | | 0 | |
|----|----|----|----|----|----|----|

(Opcode = 111, $I = 0$)

(b) Register – reference instruction

| 15 | | 12 | 11 | | 0 | |
|----|----|----|----|----|----|----|

(Opcode = 111, $I = 1$)

(c) Input – output instruction

**Figure 3.1** Basic computer instruction formats.

input-output type. Note that the bit in position 15 of the instruction code is designated by the symbol I but is not used as a mode bit when the operation code is equal to 111.

Only three bits of the instruction are used for the operation code. It may seem that the computer is restricted to a maximum of eight distinct operations. However, since register-reference and input-output instructions use the remaining 12 bits as part of the operation code, the total number of instruction chosen for the basic computer is equal to 25.

The instructions for the computer are listed in Table 3-2. The symbol designation is a three-letter word and represents an abbreviation intended for programmers and users. The hexadecimal

**TABLE 3-1** Basic Computer Instruction

| Symbol | Hexadecimal code | | Description |
| | $I = 0$ | $I = 1$ | |
| --- | --- | --- | --- |
| AND | 0xxx | 8xxx | AND memory word to AC |
| ADD | 1xxx | 9xxx | Add memory word to AC |
| LDA | 2xxx | Axxx | Load memory word to AC |
| STA | 3xxx | Bxxx | Store content of AC in memory |
| BUN | 4xxx | Cxxx | Branch unconditionally |
| BSA | 5xxx | Dxxx | Branch and save return address |
| ISZ | 6xxx | Exxx | Increment and skp if zero |
| CLA | 7800 | | Clear AC |
| CLE | 7400 | | Clear E |
| CMA | 7200 | | Complement AC |
| CME | 7100 | | Complement E |
| CIR | 7080 | | Circulate right AC and E |
| CIL | 7040 | | Circulate left AC and E |
| INC | 7020 | | Increment AC |
| SPA | 7010 | | Skip next instruction if AC |
| SNA | 7008 | | positive |
| SZA | 7004 | | Skip next instruction if AC |
| SZE | 7002 | | negative |
| HLT | 7001 | | Skip next instruction if AC zero |
| | | | Skip next instruction if E is 0 |
| | | | Halt computer |
| INP | F800 | | Input character to AC |
| OUT | F400 | | Output character from AC |
| SKI | F200 | | Skip on input flag |
| SKO | F100 | | Skip on output flag |
| ION | F080 | | Interrupt on |
| IOF | F040 | | Interrupt off |

code is equal to the equivalent hexadecimal number of the binary code used for the instruction. By using the hexadecimal equivalent we reduced the 16 bits of an instruction code to four digits with each hexadecimal digit being equivalent to four bits. A memory-reference instruction has an address part of 12 bits. The address part is denoted by three x's and stand for the three hexadecimal digits corresponding to the 12-bit address. The last bit of the instruction is designated by the symbol I. When I = 0, the last four bits of an instruction have a hexadecimal digit equivalent from 0 to 6 since the last bit is 0. When I = 1, the hexadecimal digit equivalent of the last four bits of the instruction ranges from 8 to E since the last bit is I.

Register-reference instructions use 16 bits to specify an operation. The leftmost four bits are always 0111, which is equivalent to hexadecimal 7. The other three hexadecimal digits give the binary equivalent of the remaining 12 bits. The input-output instructions also use all 16 bits to specify an operation. The last four bits are always 1111, equivalent to hexadecimal F.

### 3.1.1   INSTRUCTION SET COMPLETENESS

Before investigating the operations performed by the instructions, let us discuss the type of instructions that must be included in a computer. A computer should have a set of instructions so that the user can construct machine language programs to evaluate any function that is known to be computable. The set of instructions are said to be complete if the computer includes a sufficient number of instructions in each of the following categories:

1.      Arithmetic, logical, and shift instructions

2.      Instructions for moving information to and from memory and processor registers

3.      Program control instructions together with instructions that check statues conditions

4.      Input and output instructions

Arithmetic, logical, and shift instructions provide computational capabilities for processing the type of data that the user may wish to employ. The bulk of the binary information in a digital computer is stored in memory, but all computations are done in processor registers. Therefore, the user must have the capability of moving information between these two units. Decision-making capabilities are an important aspect of digital computers. For example, two numbers can be compared, and if the first is greater than the second, it may be necessary to proceed differently than if the second is greater than the first. Program control instructions such as branch instructions are used to change the sequence in which the program is executed. Input and output instructions are needed for communication between the computer are the user. Programs and data must be transferred into memory and results of computations must be transferred back to the user.

The instructions listed in Table 3-1 constitute a minimum set that provides all the capabilities mentioned above. There is one arithmetic instruction, ADD, and two related instructions, complement AC (CMA) and increment AC (INC). With these three instructions we can add and subtract binary numbers when negative numbers are in signed-2's complement representation. The circulate instructions, CIR ad CIL, can be used for arithmetic shifts as well as any other type of shifts desired. Multiplication and division can be performed using addition, subtraction, and shifting. There are three logic operations: AND, complement AC (CMA), and clear AC (CLA). The AND and complement provide a NAND operation. It can be shown that with the NAND operation it is possible to implement all the other logic operations with two variables (listed in Table 3-1). Moving information from memory to AC is accomplished with the load AC (LDA) instruction. Storing information from AC into memory is done with the store AC (STA) instruction. The branch instructions BUN, BSA, and ISZ, together with the four skip

instructions, provide capabilities for program control and checking of status conditions. The input (INP) and output (OUT) instructions cause information to be transferred between the computer and external devices.

Although the set of instructions for the basic computer is complete, it is not efficient because frequently used operations are not performed rapidly. An efficient set of instructions will include such instructions as subtract, multiply, OR, and exclusive-OR. These operations must be programmed in the basic computer. The programs are presented in Chap. 6 together with other programming examples for the basic computer. By using a limited number of instructions it is possible to show the detailed logic design of the computer. A more complete set of instructions would have made the design too complex. In this way we can demonstrate the basic principles of computer organization and design without going into excessive complex details. In Chap. 8 we present a complete list of computer instructions that are included in most commercial computers.

The function of each instruction listed in Table 3-2. We can delay this discussion because we must first consider the control unit and understand its internal organization.

## 3.2    TIMING AND CONTROL

The timing for all registers in the basic computer is controlled by a master clock generator. The clock pulses are applied to all flip-flops and registers in the system, including the flip-flops and registers in the control unit. The clock pulses do not change the state of a register unless the register is enabled by a control signal. The control signals are generated in the control unit and provide control inputs for the multiplexers in the common bus, control inputs in processor registers, and microoperations or the accumulator.

There are two major types of control organization: hardwired control and micro programmed control. in the hardwired organization, the control logic is implemented with gates, flip-flops, decoders, and other digital circuits. It has the advantage that it can be optimized to produce a fast mode of operation. In the microprogrammed organization, the control information is stored in a control memory. The control memory is programmed to initiate the required sequence of cicrooperations. A hardwired control, as the name implies, requires changes in the wiring among the various components if the design has to be modified or changed. In the microprogrammed control, and required changes or modifications can be done by updating the microprogram in control memory. A hardwired control for the basic computer is presented in this section.

The block diagram of the control unit is shown in Fig. 3-2. It consists of two decoders, a sequence counter, and a number of control logic gates. An instruction read from memory is placed in the instruction register (IR): .The instruction register is shown again in Fig. 3.2 where it is divided into three parts: the I. bit, the operation code, and bits 0 through 11. The operation code in bits 12 through 14 are decoded with a $3 \times 8$ decoder. The eight outputs of the decoder are designated by the symbols $D_0$ through $D_7$. The subscripted decimal number is equivalent to the binary value of the corresponding operation code. Bits 15 of the instruction is transferred to a flip-flop designated by the symbol I. Bits 0 through 11 are applied to the control logic gates. The 4-bit sequence counter can count in binary from 0 through 15. The outputs of the counter are decoded into 16 timing signals $T_0$ through $T_{15}$. The internal logic of the control gates will be derived later when we consider the design of the computer in detail.

The sequence counter SC can be incremented or cleared synchronously. Most of the time, the counter is incremented to provide the sequence of timing signals out of the $4 \times 16$ decoder. Once in a while, the counter is cleared to 0, causing the next active timing signal to be $T_0$. As an example, consider the case where SC is incremented to provide timing signals $T_0$, $T_1$, $T_2$, $T_3$, and

$T_4$ in sequence. At time $T_4$, SC is cleared to 0 if decoder output $D_3$ is active. This is expressed symbolically by the statement

$$D_3 T_4: \quad SC \leftarrow 0$$

The timing diagram of Fig. 3-3 shows the time relationship of the control signals. The sequence counter SC responds to the positive transition of the clock. Initially, the CLR input of SC is active. The first positive transition of the
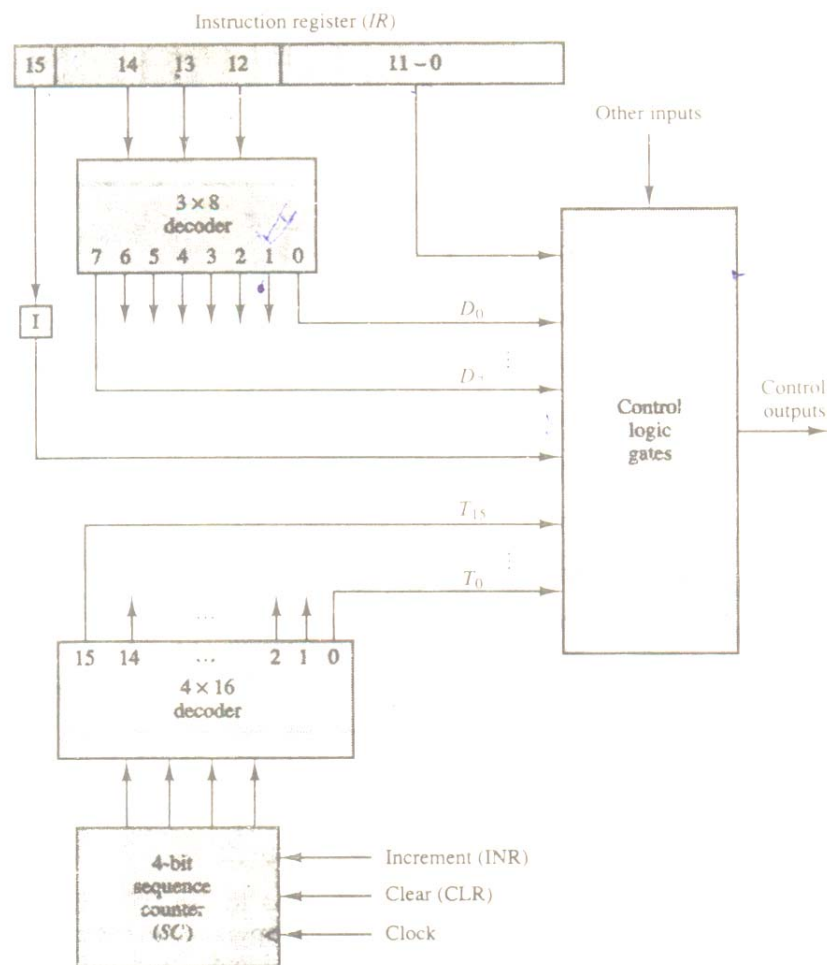


**Figure 3-2** Control unit basic computer

Clock clears SC to 0, which in turn activates the timing signal $T_0$ out of the decoder. $T_0$ is active during one clock cycle. The positive clock transition labeled $T_0$ in the diagram will trigger only those registers whose control inputs are connected to timing signal $T_0$. SC is incremented with every positive clock transition, unless its CLR input is active. This produces the sequence of timing signals $T_0$, $T_1$, $T_2$, $T_3$, $T_4$, and so on, as shown in the diagram (Note the relationship between the timing signal and its corresponding positive clock transition.) If SC is not cleared, the timing signals will continue with $T_5$, $T_6$, up to $T_{15}$ and back to $T_0$.
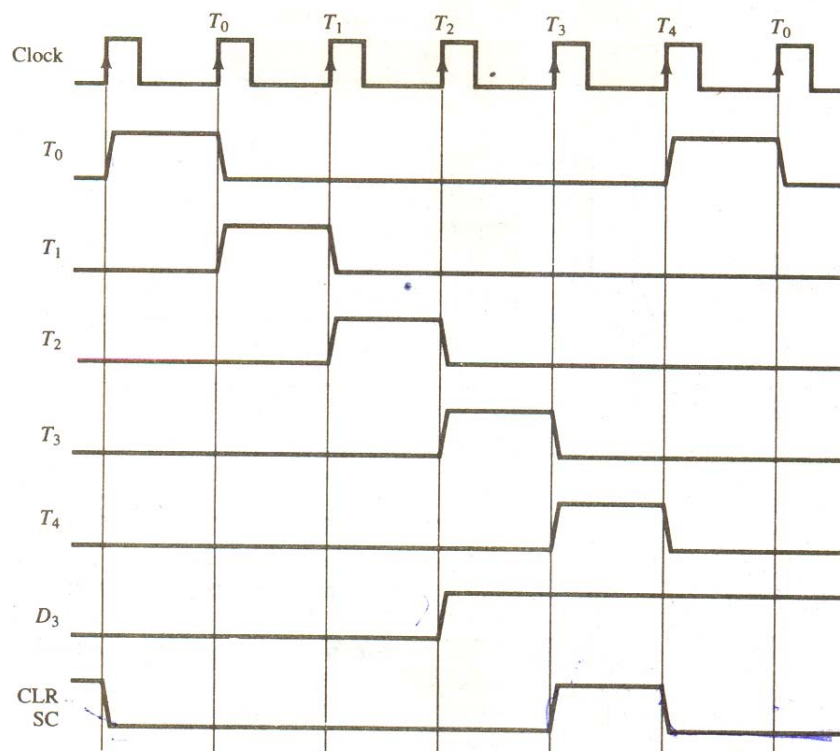
**Figure 3.3** Example of control timing signals

The last three waveforms in Fig 3.3 show how SC is cleared when $D_3T_4 = 1$. Output $D_3$ from the operation decoder becomes active at the end of timing signal $T_2$. When timing signal $T_4$ becomes active, the output of the AND gate that implements the control function $D_3D_4$ becomes active. This signal is applied to the CLR input of SC. On the next positive clock transition (the one marked $T_4$ in the diagram) the counter is cleared to 0. This causes the timing signal $T_0$ to become active instead of $T_5$ that would have been active if SC were incremented instead of cleared.

A memory read or write cycle will be initiated with the rising edge of a timing signal. It will be assumed that a memory cycle time is less than the clock cycle time. According to this assumption, a memory read or write cycle initiated by a timing signal will be completed by the time the next clock gives through its positive transition. The clock transition will then be used to load the memory word into a register. This timing relationship is not valid in many computers because the memory cycle time is usually longer than the processor clock cycle. In such a case it is necessary to provide wait cycles in the processor until the memory word is available. To facilitate the presentation, we will assume that a wait period is not necessary in the basic computer.

To fully comprehend the operation of the computer, it is crucial that one understands the timing relationship between the clock transition and the timing signals. For example, the register transfer statement

$$T_0: \quad AR \leftarrow PC$$

specifies a transfer of the content of PC into AR if timing signal $T_0$ is active during an entire clock cycle interval. During this time the content of PC is placed onto the bus (with $S_2$ $S_1$ $S_0$ = 010) and the LD (load) input of AR is enabled. The actual transfer does not occur until the end of the clock cycle when the clock goes through a positive transition. This same positive clock transition increments the sequence counter SC from 0000 to 0001. The next clock cycle has $T_1$ active and $T_0$ inactive.

## 3.3     INSTRUCTION CYCLE

A program residing in the memory unit of the computer consists of a sequence of instructions. The program is executed in the computer by going through a cycle for each instruction. Each instruction cycle in turn is subdivided into a sequence of sub cycles or phases. In the basic computer each instruction cycle consists or the following phases:

1. Fetch an instruction from memory.

2. Decode the instruction

3. Read the effective address from memory if the instruction has an indirect address.

4. Execute the instruction.

Upon the completion of step 4, the control goes back to step 1 to fetch, decode, and execute the next instruction. This process continues indefinitely unless a HALT instruction is encountered.

## 3.3.1   FETCH AND DECODE

Initially, the program counter PC is loaded with the address of the first instruction in the program. The sequence counter SC is cleared to 0, providing a decoded timing signal $T_0$. After each clock pulse, SC is incremented by one so that the timing signals go through a sequence $T_0$, $T_1$, $T_2$, and so on. The microoperations for the fetch and decode phases can be specified by the following register transfer statements.

$T_0$:   AR $\leftarrow$ PC

$T_1$:  IR  $\leftarrow$  M[AR], PC $\leftarrow$ PC + 1

$T_2$:  $D_0$,..., $D_1$ $\leftarrow$ Decode IR (12-14), AR $\leftarrow$ IR(0 –11), 1 $\leftarrow$ IR (15)

Since only AR is connected to the address inputs of memory, it is necessary to transfer the address from PC to AR during the clock transition associated with timing signal $T_0$. The instruction read from memory is then placed in the instruction register IR with the clock transition associated with timing signal $T_1$. At the same time, PC is incremented by one to prepare it for the address of the next instruction in the program. At time $T_2$, the operation code in IR is decoded, the indirect bit is transferred to flip-flop I, and the address part of the instruction is transferred to AR. Note that SC is incremented after each clock pulse to produce the sequence $T_0$, $T_1$, and $T_2$.
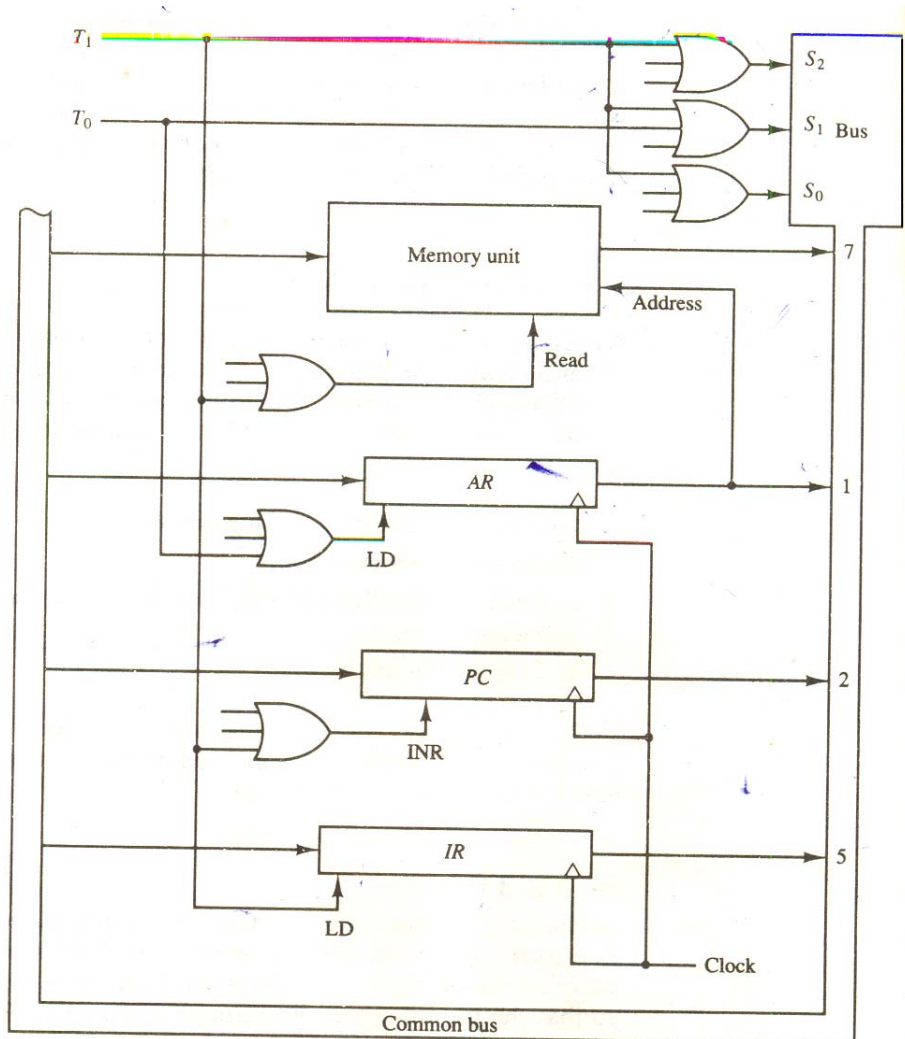
**Figure 3.4** Register transfers for the fetch phase.

Figure 3.4 shows how the first two register transfer statements are implemented in the bus system. To provide the data path for the transfer of PC to AR we must apply timing signal $T_0$ to achieve the following connection.:

1. Place the content of PC onto the bus by making the bus selection inputs $S_2S_1S_0$ equal to 010.

2. Transfer the content of the bus to AR y enabling the LD input of AR.

The next clock transition initiates the transfer form PC to AR since $T_0 = 1$. In order to implement the second statement

$$T_1: \quad IR \leftarrow M[AR], \quad PC \leftarrow PC + 1$$

it is necessary to use timing signal $T_1$ to provide the following connections in the bus system.

1. Enable the read input of memory.

2. Place the content of memory onto the bus by making $S_2S_1S_0 = 111$.

3. Transfer the content of the bus to IR by enabling the LD input of IR.

4. Increment PC by enabling the INR input of PC.

The next clock transition initiates the read and increment operations since $T_1 = 1$.

Figure 3.4 duplicates a portion of the bus system and shows how $T_0$ and $T_1$ are connected to the control inputs of the registers, the memory, and the bus selection inputs. Multiple input OR gates are included in the diagram because there are other control functions that will initiate similar operations.

### 3.3.2   DETERMINE THE TYPE OF INSTRUCTION

The timing signal that is active after the decoding is $T_3$. During time $T_3$, the control unit determines the type of instruction that was just read from memory. The flowchart of Fig. 3.5 presents an initial configuration for the instruction cycle and shows how the control determines the instruction type after the decoding.

Decoder output $D_7$ is equal to 1 if the operation code is equal to binary 111. We determine that if $D_1 = 1$, the instruction must be a register-reference or input-output type. If $D_7 = 0$, the operation code must be one of the other seven values 000 through 110, specifying memory-reference instruction. Control then inspects the value of the first bit of the instruction, which is now available in flip-flop I. If $D_7 = 0$ and $I = 1$, we have a memory-reference instruction with an indirect address. It is then necessary to read the effective address from memory. The micro operation for the indirect address condition can be symbolized by the register transfer statement

$$AR \leftarrow M [AR]$$

Initially, AR holds the address part of the instruction. This address is used during the memory read operation. The word at the address given by AR is read from memory and placed on the common bus. The LD input of AR is then enabled to receive the indirect address that resided in the 12 least significant bits of the memory word.

The three instruction types are subdivided into four separate paths. The selected operation is activated with the clock transition associated with timing signal $T_3$. This can be symbolized as follows:

| | | | |
|---|---|---|---|
| $D_1 = 0$ | $I = 1$ | $D'_7 IT_3$: | $AR \leftarrow M [AR]$ |
| 0 | 0 | $D'_7 I'T_3$; | Nothing |
| 1 | 0 | $D_7 I' T_3$: | Execute a register-reference instruction |
| 1 | 1 | $D_7 IT_3$: | Execute an input-output instruction |

When a memory-reference instruction with $I = 0$ is encountered, it is not necessary to do anything since the effective address is already in AR. However, the sequence counter SC must be incremented when $D'_7 T_3 = 1$, so that the execution of the memory-reference instruction can be continued with timing variable $T_4$. a register-reference or input-output instruction can be executed with the clock associated with timing signal $T_3$. After the instruction is executed, SC is cleared to 0 and control returns to the fetch phase with $T_0 = 1$.
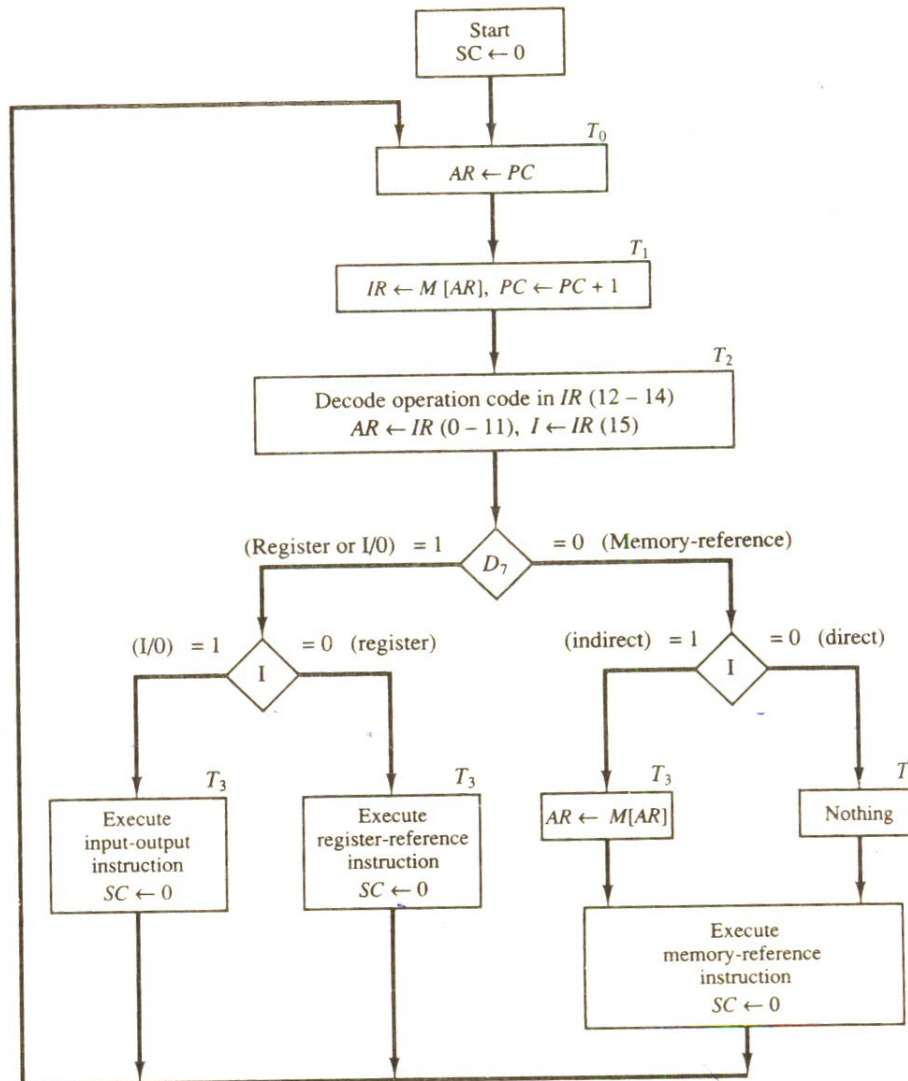
**Figure 3.5** Flowchart for instruction cycle (initial configuration).

Note that the sequence counter SC is either incremented or cleared to 0 with every positive clock transition. We will adopt the convention that if SC is incremented, we will not write the statement SC ← SC + 1, but it will be implied that the control goes to the next timing signal in sequence. When SC is to be cleared, we will include the statement Sc ← 0.

The register transfers needed for the execution of the register-reference instructions are presented in this section. The memory-reference instructions are explained in the next section.

### 3.3.3 REGISTER-REFERENCE INSTRUCTIONS

Register-reference instructions are recognized by the control when $D_7 = 1$ and I = 0. These instructions use bits 0 through 11 of the instruction code to specify one of 12 instruction. These 12 bits are available in IR (0-11). They were also transferred to AR during time $T_2$.

The control functions and microoperations for the register-reference instructions are listed in Table 3.2. These instructions are executed with the clock transition associated with timing variable $T_3$. Each control function needs the Boolean relation $D_7 I' T_3$, which we designate for convenience by the symbol r. The control function is distinguished by one of the bits in IR (0-11). By assigning the symbol $B_i$ to bit i of IR, all control functions can be simply denoted by $rB_i$. For example, the instruction CLA has the hexadecimal code 7800 (see Table 3.1), which gives the binary equivalent 0111 1000 0000 0000. The first bit is a zero and is equivalent to I'. The next three bits constitute the operation code and are recognized from decoder output $D_7$. Bit 11 in IR is 1 and is recognized from $B_{11}$. The control function that initiates the microoperation for this instruction is $D_7 I' T_3 B_{11} = rB_{11}$. The execution of a register-reference instruction is completed at time $T_3$. The sequence counter SC is cleared to 0 and the control goes back to fetch the next instruction with timing signal $T_0$.

The first seven register-reference instructions perform clear, complement, circular shift, and increment microoperations on the AC or E registers. The next four instructions cause a skip of the next instruction in sequence when a stated condition is satisfied. The skipping of the instruction is achieved by incrementing PC once again (in addition, it is being incremented during the fetch phase at time $T_1$). The condition control statements must be recognized as part o the control conditions. The AC is positive when the sign bit in $AC(15) = 0$; it is negative when $AC(15) = 1$. The content of AC is zero $(AC = 0)$ if all the flip-flops of the register are zero. The HLT instruction clears a start-stop flip-flops S and stops the sequence counter from counting. To restore the operation of the computer, the start-stop flip-flop must be set manually.

**TABLE 3.2** Execution of Register-Reference Instruciton

$D_1 I' T_3 = r$ (common to all register-reference instructions)
$\quad$ IR(i) = $B_i$ [bit in IR(0−11) that specifies the operation]

| | | | |
|---|---|---|---|
| | r: | $SC \leftarrow 0$ | Clear SC |
| CLA | $rB_{11}$: | $AC \leftarrow 0$ | Clear AC |
| CLE | $rB_{10}$: | $E \leftarrow 0$ | Clear E |
| CMA | $rB_9$: | $AC \leftarrow \overline{AC}$ | Complement |
| CME | $rB_8$: | $E \leftarrow \overline{E}$ | AC |
| CIR | $rB_7$: | $AC \leftarrow$ shr AC, $AC(15) \leftarrow$ | Complement E |
| CIL | $rB_6$: | $E, E \leftarrow AC(0)$ | Circulate right |
| INC | $rB_5$: | $AC \leftarrow$ shl AC, $AC(0) \leftarrow E$, | Circulate left |
| SPA | $rB_4$: | $E \leftarrow AC(15)$ | Increment AC |
| SNA | $rB_3$: | $AC \leftarrow AC + 1$ | Skip if positive |
| SZA | $rB_2$: | If $(AC(15) = 0)$ then $(PC$ | Skip if negative |
| SZE | $rB_1$ | $\leftarrow PC + 1)$ | Skip if AC zero |
| HLT | $rB_0$: | If $(AC(15) = 1)$ then $(PC$ | Skip if E zero |
| | | $\leftarrow PC + 1)$ | Halt computer |
| | | If $(AC = 0)$ then $PC \leftarrow PC$ | |

+ 1)

If (E = 0) the (PC ← PC +

1)

S ← 0 (S is a start-stop

flip-flop)

## 3.4    MEMORY-REFERENCE INSTRUCTIONS

In order to specify the microoperations needed for the execution of each instruction, it is necessary that the function that they are intended to perform be defined precisely.  Looking back to Table 3.1, where the instructions are listed, we find that some instructions have an ambiguous description. This is because the explanation of an instruction in words is usually lengthy, and not enough space is available in the table for such a lengthy explanation.  We will now show that the function of the memory-reference instructions can be defined precisely by means of register transfer notation.

Table 3.3 lists the seven memory-reference instructions.  The decoded output $D_i$ for i = 0, 1, 2, 3, 4, 5, and 6 from the operation decoder that belongs to each instruction is included in the table.  The effective address of the instruction is in the address register AR and was placed there during timing signal $T_2$ when I = 0, or during timing signal $T_3$ when I = 1.  The execution of the memory-reference instructions starts with timing signal $T_4$ the symbolic description of each instruction is specified in the table in terms of register transfer notation.  The actual execution of the instruction in the bus system will require a sequence of microoperations.  This is because data stored in memory cannot be processed directly.  The data must be read from memory to a register where they can be operated on with logic circuits.  We now explain the operation of each instruction and list the control functions and microoperations needed for their execution.  A flowchart that summarizes all the microoperations is presented at the end of this section.

**Table 3.3** Memory-Reference Instructions

| Symbol | Operation Decoder | Symbolic description |
|---|---|---|
| AND | $D_0$ | AC ← AC Λ M[AR] |
| ADD | $D_1$ | AC ← AC + M[AR], E ←$C_{out}$ |
| LDA | $D_2$ | AC ← M[AR] |
| STA | $D_3$ | M[AR] ← AC |
| BUN | $D_4$ | PC ← AR |
| BSA | $D_5$ | M[AR ← PC, PC ← AR + 1 |
| ISZ | $D_6$ | M[AR] ← M[AR] + 1, If M[AR] + 1 = 0 then PC ← PC + 1 |

**AND to AC**

This is an instruction that performs the AND logic operation on pairs of bits in AC and the memory word specified by the effective address.  The result of the operation is transferred to AC. The microoperations that execute this instruction are :

$$D_0T_4 : \quad DR \leftarrow M[AR]$$

$$D_0T_5: \quad AC \leftarrow AC \wedge DR, \quad SC \leftarrow 0$$

The control function for this instruction uses the operation decoder $D_0$ since this output of the decoder is active when the instruction has an AND operation whose binary code value is 000. Two timing signals are needed to execute the instruction. The clock transition associated with timing signal $T_4$ transfers the operand from memory into DR. The clock transition associated with the next timing signal $T_5$ transfers to AC the result of the AND logic operation between the contents of DR and AC. The same clock transition clears SC to 0, transferring control to timing signal $T_0$ to start a new instruction cycle.

### ADD to AC

The instruction adds the content of the memory word specified by the effective address to the value of AC. The sum is transferred into AC and the output carry $C_{out}$ is transferred to the E (extended accumulator) flip-flop. The microoperations needed to execute this instruction are

$$D_1T_4: \quad DR \leftarrow M[AR]$$

$$D_1T_5: \quad AC \leftarrow AC + DR, \quad E \leftarrow C_{out}, \quad SC \leftarrow 0$$

The same two timing signals, $T_4$ and $T_5$, are used again but with operation decoder $D_1$ instead of $D_0$, which was used for the AND instruction. After the instruction is fetched from memory and decoded, only one output of the operation decoder will be active, and that output determines the sequence of microoperations that the control follows during the execution of a memory-reference instruction.

### LDA : Load to AC

This instruction transfers the memory word specified by the effective address to AC. The microoperations needed to execute this instruction are

$$D_2T_4: \quad DR \leftarrow M[AR]$$

$$D_2T_5: \quad AC \leftarrow DR, SC \leftarrow 0$$

From the bus diagram we note that there is no direct path from the bus into AC. The adder and logic circuit receive information from DR which can be transferred into AC. Therefore, it is necessary to read the memory word into DR first and then transfer the content of DR into AC. The reason for not connecting the bus to the inputs of AC is the delay encountered in the adder and logic circuit. It is assumed that the time it takes to read from memory and transfer the word through the bus as well as the adder and logic circuit is more than the time of one clock cycle. By not connecting the bus to the inputs of AC we can maintain one clock cycle per microoperation.

### STA : Store AC

This instruction stores the content of AC into the memory word specified by the effective address. Since the output of AC is applied to the bus and the data input of memory is connected to the bus, we can execute this instruction with one microoperation:

$$D_3T_4: M[AR] \leftarrow AC, \quad SC \leftarrow 0$$

### BUN : Branch Unconditionally

This instruction transfers the program to the instruction specified by the effective address. Remember that PC holds the address of the instruction to be read from memory in the next instruction cycle. PC is incremented at time $T_1$ to prepare it for the address of the next instruction

in the program sequence. The BUN instruction allows the programmer to specify an instruction out of sequence and we say that the program branches (or jumps) unconditionally. The instruction is executed with one microoperation:

$$D_4T_4: \quad PC \leftarrow AR, \ SC \leftarrow 0$$

The effective address from AR is transferred through the common bus to PC. Resetting SC to 0 transfers control to $T_0$. The next instruction is then fetched and executed from the memory address given by the new value in PC.

### BSA : Branch and Save Return Address

This instruction is useful for branching to a portion of the program called a subroutine or procedure. When executed, the BSA instruction stores the address of the next instruction in sequence (which is available in PC) into a memory location specified by the effective address. The effective address plus one is then transferred to PC to serve as the address of the first instruction in the subordinate. This operation was specified in Table 5-4 with the following register transfer:

$$M[AR] \leftarrow PC, \quad PC \leftarrow AR + 1$$

A numerical example that demonstrates how this instruction is used with a subordinate is shown in Fig. 3.6. The BSA instruction is assumed to be in memory at address 20. The I bit is 0 and the address part of the instruction has the binary equivalent of 135. After the fetch and decode phases, PC contains 21, which is the address of the next instruction in the program (referred to as the return address). AR holds the effective address 135. This is shown in part (a) of the figure. The BSA instruction performs the following numerical operation:

$$M[135] \leftarrow 21, PC \leftarrow 135 + 1 = 136$$

The result of this operation is shown in part (b) of the figure. The return address 21 is stored in memory location 135 and control continues with the subordinate program starting from address 136. The return to the original program (at address 21) is accomplished by means of an indirect BUN instruction placed at the end of the subroutine. When this instruction is executed, control goes to the indirect phase to read the effective address at location 135, where it finds the previously saved address 21. When the BUN instruction is executed, the effective address 21 is transferred to PC. The next instruction cycle finds PC with the value 21, so control continues to execute the instruction at the return address.

The BSA instruction performs the function usually referred to as a subroutine call. The indirect BUN instruction at the end of the subroutine performs the function referred to as a subroutine return. In most commercial computers, the return address associated with a subroutine is stored in either a processor register or in a portion of memory called a stack.

It is not possible to perform the operation of the BSA instruction in one clock cycle when we use the bus system of the basic computer. To use the memory and the bus properly, the BSA instruction must be executed with a sequence of two microoperations:

$$D_5T_4: \quad M[AR] \leftarrow PC, AR \leftarrow AR + 1$$

$$D_5T_5: \quad PC \leftarrow AR, SC \leftarrow 0$$

Timing signal $T_4$ initiates a memory write operation, places the content of PC onto the bus, and enables the INR input of AR. The memory write operation is completed and AR is incremented by the time the next clock transition occurs. The bus is used at $T_5$ to transfer the content of AR to PC.

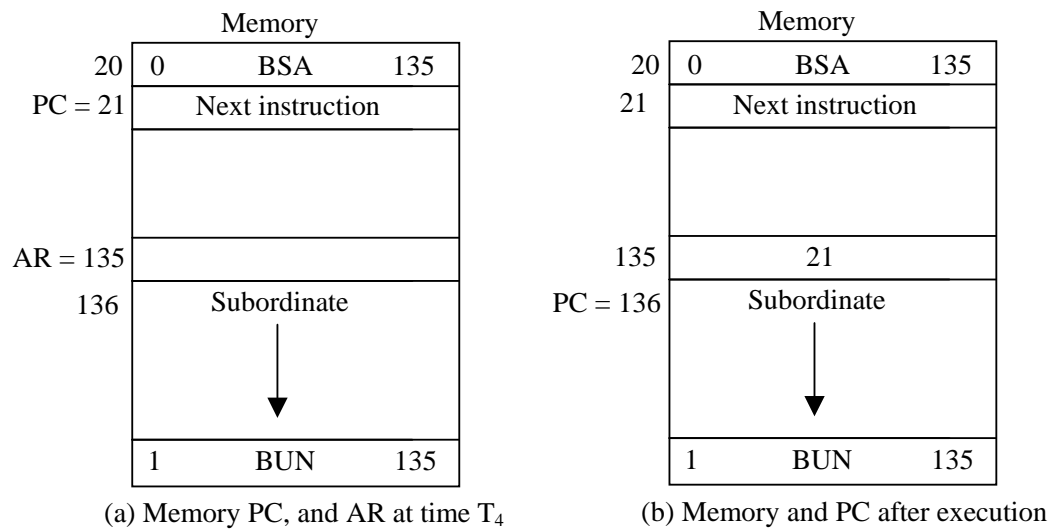(a) Memory PC, and AR at time $T_4$    (b) Memory and PC after execution

**Figure 3.6** Example of BSA instruction execution.

It is not possible to perform the operation of the BSA instruction in one clock cycle when we use the bus system of the basic computer. To use the memory and the bus properly, the BSA instruction must be executed with a sequence of two microoperations:

$$D_5T_4: \quad M[AR] \leftarrow PC, AR \leftarrow AR + 1$$

$$D_5T_5: \quad PC \leftarrow AR, SC \leftarrow 0$$

Timing signal $T_4$ initiates a memory write operation, places the content of PC onto the bus, and enables the INR input of AR. The memory write operation is completed and AR is incremented by the time the next clock transition occurs. The bus is used at $T_5$ to transfer the content of AR to PC.

**ISZ : Increment and Skip if Zero**

This instruction increments the word specified by the effective address, and if the incremented value is equal to 0, PC is incremented by 1. The programmer usually stores a negative number (in 2's complement) in the memory word. As this negative number is repeatedly incremented by one, it eventually reaches the value of zero. At that time PC is incremented by one in order to skip the next instruction in the program.

Since it is not possible to increment a word inside the memory, it is necessary to read the word into DR, increment DR, and store the word back into memory. This is done with the following sequence of microoperations:

$$D_6T_4: \quad DR \leftarrow M[AR]$$

$$D_6T_5: \quad DR \leftarrow DR + 1$$

$$D_6T_6: \quad M[AR] \leftarrow DR, \text{ if } (DR = 0) \text{ then } (PC \leftarrow PC + 1), SC \leftarrow 0$$

**3.4.1   CONTROL FLOWCHART**

A flowchart showing all microoperations for the execution of the seven memory-reference instructions is shown in Fig. 3.7. The control functions are indicated on top of each box. The microoperations that are performed during time $T_4$, $T_5$, or $T_6$ depend on the operation code value.

This is indicated in the flowchart by six different paths, one of which the control takes after the instruction is decoded. The sequence counter SC is cleared to 0 with the last timing signal in each case. This causes a transfer of control to timing signal $T_0$ to start the next instruction cycle.

Note that we need only seven timing signals to execute the longest instruction (ISZ). The computer can be designed with a 3-bit sequence counter. The reason for using a 4-bit counter for SC is to provide additional timing signals for other instructions that are presented in the problems section.
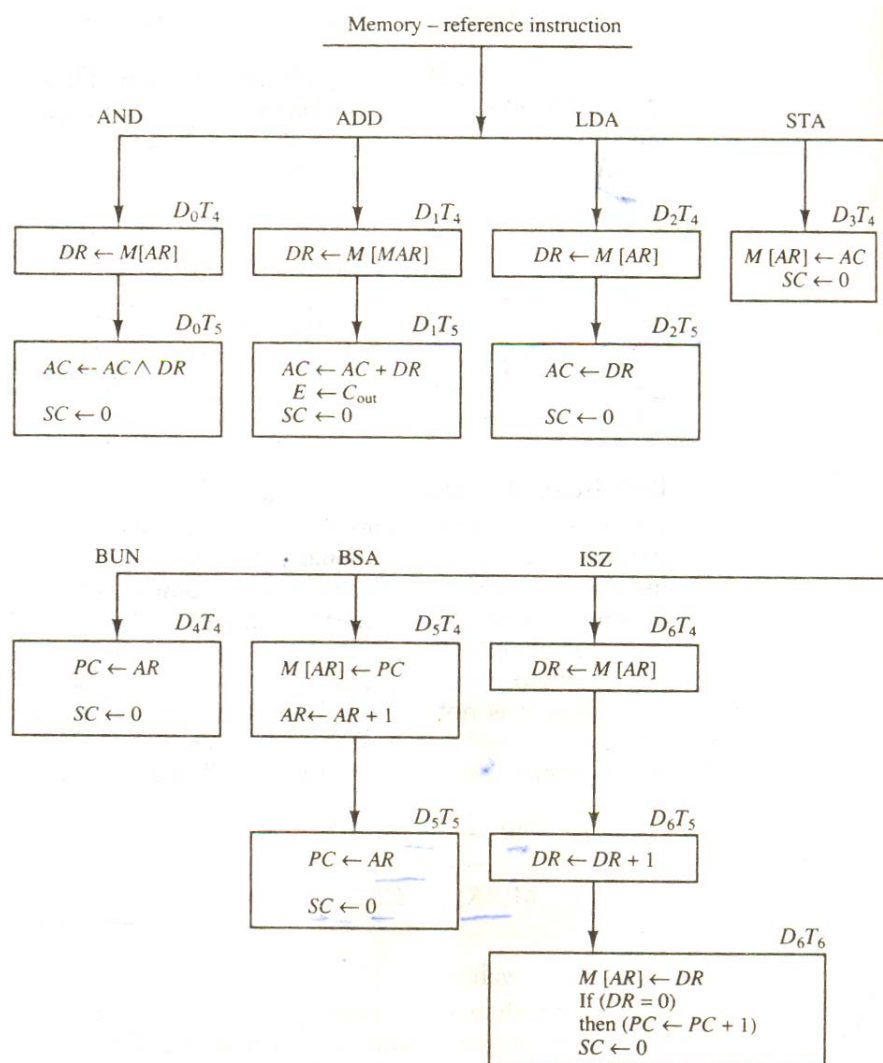


**Figure 3.7**     Flowchart for memory-reference instructions

## 3.5     INPUT-OUTPUT AND INTERRUPT

A computer can serve no useful purpose unless it communicates with the external environment. Instructions and data stored in memory must come from some input device. Computational results must be transmitted to the user through some output device. Commercial computers include many types of input and output devices. To demonstrate the most basic requirements for

input and output communication, we will use as an illustration a terminal unit with a keyboard and printer. Input-output organization is discussed further in Chap. 11.

### 3.5.1    INPUT-OUTPUT CONFIGURATION

The terminal sends and receives serial information. Each quantity of information has eight bits of an alphanumeric code. The serial information from the keyboard is shifted into the input register INPR. The serial information for the printer is stored in the output register OUTR. These two registers communicate with a communication interface serially and with the AC in parallel. The input-output configuration is shown in Fig. 3.8. The transmitter interface receives serial information from the keyboard and transmits it to INPR. The receiver interface receives information from OUTR and sends it to the printer serially.

The input register INPR consists of eight bits and holds an alphanumeric input information. The 1-bit input flag FGI is a control flip-flop. The flag bit is set to 1 when new
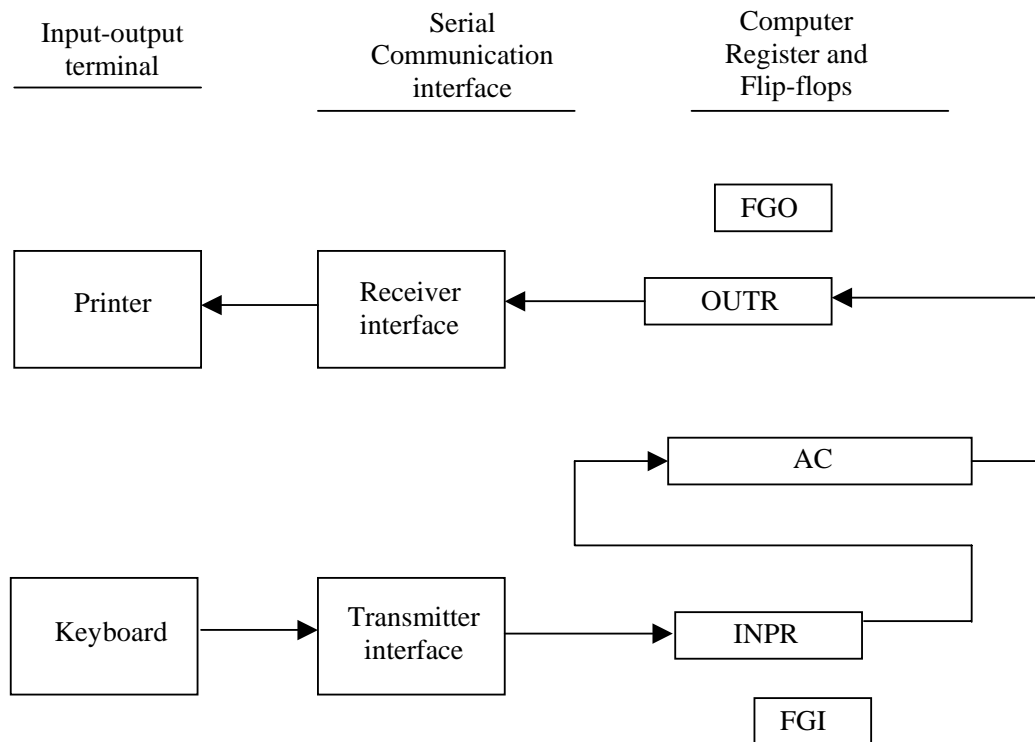


**Figure 3.8 Input Output Configuration**

Information is available in the input device and is cleared to 0 when the information is accepted by the computer. The flag is needed to synchronize the timing rate difference between the input device and the computer. The process of information transfer is as follows. Initially, the input flag FGI is cleared to 0. When a key is struck in the keyboard, an 8-bit alphanumeric code is shifted into INPR and the input flag FGI is set to 1. As long as the flag is set, the information in INPR cannot be changed by striking another key. The computer checks the flag bit; if it is 1, the information from INPR is transferred in parallel into AC and FGI is cleared to 0. Once the flag is cleared, new information can be shifted into INPR by striking another key.

The output register OUTR works similarly but the direction of information flow is reversed. Initially, the output flag FGO is set to 1. The computer checks the flag bit, if it is 1, the information from AC is transferred in parallel to OUTR and FGO is cleared to 0. The output device accepts the coded information, prints the corresponding character, and when the operation is completed, it sets FGO to the computer does not load a new character into OUTR when FGO is 0 because this condition indicates that the output device is in the process of printing the character.

### 3.5.2   INPUT OUTPUT INSTRUCTIONS

Input and output instructions are needed for transferring information to and from AC register, for checking the flag bits, and for controlling the interrupt facility. Input-output instructions have an operation code 1111 and are recognized by the control when $D_7 =$ and $I = 1$. The remaining bits of the instruction specify the particular operation. The control functions and microoperations for the input-output instructions are listed in Table 3.4. These instructions are executed with the clock transition associated with timing signal $T_3$. Each control function needs a Boolean relation $D\text{-}IT_3$, which we designate for convenience by the symbol p. The control function is distinguished by one of the bits in IR. By assigning the symbol $B_i$ to bit i of IR, all control functions can be

**Table 3.4** Input-Output Instructions

$D_7 IT_3 = p$ (common to all input-output instructions)

$IR(i) = B_i$ [bit in IR (6−11) that specifies the instruction]

|      |            |                                          |                     |
|------|------------|------------------------------------------|---------------------|
|      | p :        | $SC \leftarrow 0$                        | Clear SC            |
| INP  | $pB_{11}$: | $AC(0\text{−}7) \leftarrow INPR, FGI \leftarrow 0$ | Input character     |
| OUT  | $pB_{10}$: | $OUTR \leftarrow AC (0\text{−}7), FGO \leftarrow 0$ | Output character    |
| SKI  | $pB_9$:    | If (FGI = 1) then (PC $\leftarrow$ PC + 1) | Skip on input flag  |
| SKO  | $pB_8$:    | If (FGO = 1) then (PC $\leftarrow$ PC + 1) | Skip on output flag |
| ION  | $pB_7$:    | $IEN \leftarrow 1$                       | Interrupt enable on |
| IOF  | $pB_6$:    | $IEN \leftarrow 0$                       | Interrupt enable off|

denoted by $pB_1$ for i = 6 though 11. The sequence counter SC is cleared to 0 when $p = D_7 IT_3 = 1$.

The INP instruction transfers the input information from INPR into the eight low-order bits of AC and also clears the input flag to 0. The OUT instruction transfers the eight least significant bits of AC into the output register OUTR and clears the output flag to 0. The next two instructions in Table 3.4 check the status of the flags and cause a skip of the next instruction if the flag is 1. The instruction that is skipped will normally be a branch instruction to return and check the flag again. The branch instruction is not skipped if the flag is 0. If the flag is 1, the branch instruction is skipped and an input or output instruction is executed. The last two instructions set and clear an interrupt enable flip-flop IEN. The purpose of IEN is explained in conjunction with the interrupt operation.

### 3.5.3 PROGRAM INTERRUPT

The process of communication just described is referred to as programmed control transfer. The computer keeps checking the flag bit, and when it finds it set, it initiates an information transfer. The difference of information flow rate between the computer and that of the input-output device makes this type of transfer inefficient. To see why this is inefficient, consider a computer that can go through an instruction cycle in 1 $\mu$s. Assume that the input-output device can transfer information at a maximum rate of 10 characters per second. This is equivalent to one character every 100,000 $\mu$s. Two instructions are executed when the computer checks the flag bit and decides not to transfer the information. This means that at the maximum rate, the computer will check the flag 50,000 times between each transfer. The computer is wasting time while checking the flag instead of doing some other useful processing task.

An alternative to the programmed controlled procedure is to let the external device inform the computer when it is ready for the transfer. In the meantime the computer can be busy with other tasks. This type of transfer uses the interrupt facility. While the computer is running a program, it does not check the flags. However, when a flag is set, the computer is momentarily interrupted from proceeding with the current program and is informed of the fact that a flag has been set. The computer is wasting time while checking the flag instead of doing some other useful processing task.

An alternative to the programmed controlled procedure is to let the external device inform the computer when it is ready for the transfer. In the meantime the computer can be busy with other tasks. This type of transfer uses the interrupt facility. While the computer is running a program, it does not check the flags. However, when a flag is set, the computer is momentarily interrupted from proceeding with the current program and is informed of the fact that a flag has been set. The computer deviates momentarily from what it is doing to take care of the input or output transfer. It then returns to the current program to continue what it was doing before the interrupt.

The interrupt enable flip-flop IEN can be set and cleared with two instructions. When IEN is cleared to 0 (with the IOF instruction), the flags cannot interrupt the computer. When IEN is set to 1 (with the ION instruction), the computer can be interrupted. These two instructions provide the programmer with the capability of making a decision as to whether or not to use the interrupt facility.

The way that the interrupt is handled by the computer can be explained by means of the flowchart of Fig. 5-13. An interrupt flip-flop R is included in the computer. When R = 0, the computer goes through an instruction cycle. During the execute phase of the instruction cycle IEN is checked by the control. If it is 0, it indicates that the programmer does not want to use the interrupt, so control continues with the next instruction cycle. If IEN is 1, control checks the flag bits. If both flags are 0, it indicates that neither the input nor the output registers are ready for transfer of information. In this case, control continues with the next instruction cycle. If either flag is set to 1 while IEN = 1, flip-flop R is set to 1. At the end of the execute phase, control checks the value of R, and if it is equal to 1, it goes to an interrupt cycle instead of an instruction cycle.

The interrupt cycle is a hardware implementation of a branch and save return address operation. The return address available in PC is stored in a specific location where it can be found later when the program returns to the instruction at which it was interrupted. This location may be a processor register, a memory stack, or a specific memory location. Here we choose the memory location at address 0 as the place for storing the return address. Control then inserts

address 1 into PC and clears IEN and R so that no more interruptions can occur until the interrupt request from the flag has been serviced.
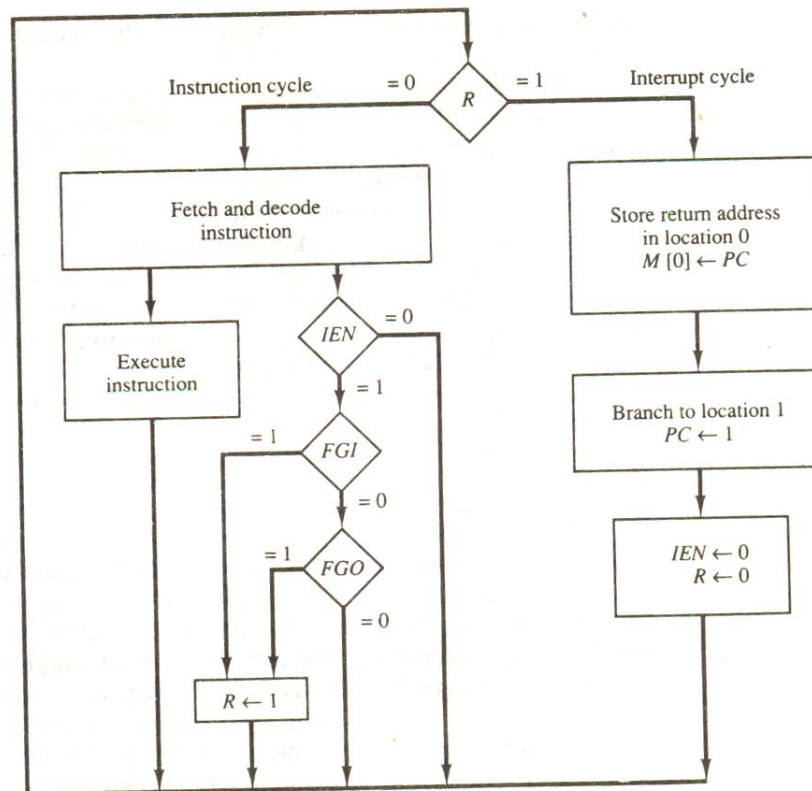


**Figure 3.9** Flowchart for interrupt cycle

An example that shows what happens during the interrupt cycle is shown in Fig. 3.10. Suppose that an interrupt occurs and R is set to 1 while the control is executing the instruction at address 255. At this time, the return address 256 is in PC. The programmer has previously placed an input-output service program in memory starting from address 1120 and a BUN 1120 instruction at address 1. This is shown in Fig. 3.10 (a).

When control reaches timing signal $T_0$ and finds that R = 1, it proceeds with the interrupt cycle. The content of PC (256) is stored in memory location 0, PC is set to 1, and R is cleared to 0. At the beginning of the next instruction cycle, the instruction that is read from memory is in address 1 since this is the content of PC. The branch instruction at address 1 causes the program to transfer to the input-output service program at address 1120. This program checks the flags, determines which flag is set, and then transfers the required input or output information. One this is done, the instruction ION is executed to set IEN to 1 (to enable further interrupts), and the program returns to the location where it was interrupted. This is shown in Fig. 3.10 (b).

The instruction that returns the computer to the original place in the main program is a branch indirect instruction with an address part of 0. This instruction is placed at the end of the I/O service program. After this instruction is read from memory during the fetch phase, control goes to the indirect phase (because I = 1) to read the effective address. The effective address is in

location 0 and is the return address that was stored there during the previous interrupt cycle. The execution of the indirect BUN instruction results in placing into PC the return address from location 0.
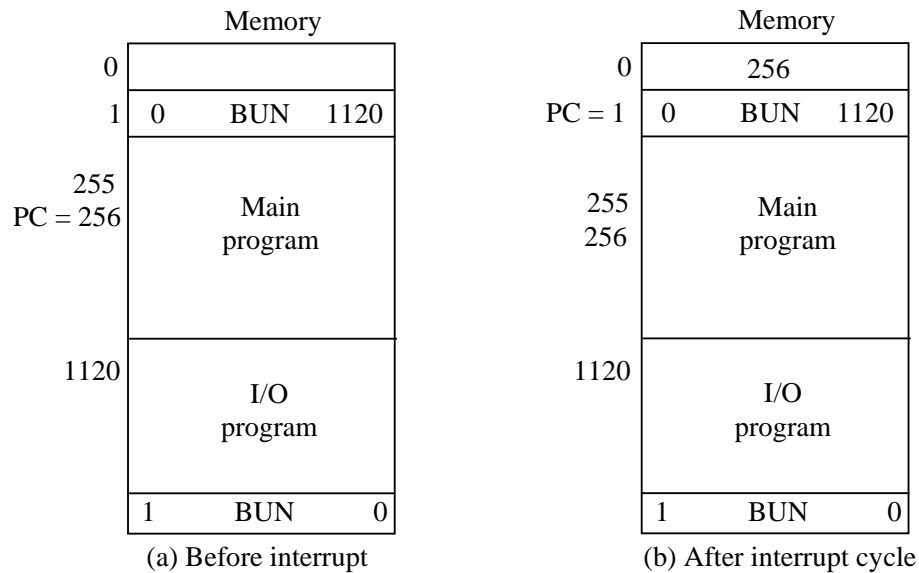


| | Memory | | | | | Memory | | |
| (a) Before interrupt | | | | | (b) After interrupt cycle | | | |

**Figure 3.10** Demonstration of the interrupt cycle

### 3.5.4   INTERRUPT CYCLE

We are now ready to list the register transfer statements for the interrupt cycle. The interrupt cycle is initiated after the last execute phase if the interrupt flip-flop R is equal to 1. This flip-flop is set to 1 if IEN = 1 and either FGI or FGO are equal to 1. This can happen with any clock transition except when timing signals $T_0$, $T_1$, or $T_2$ are active. The condition for setting flip-flop R to 1 can  be expressed with the following register transfer statement :

$$T'_0 T'_1 T'_2 \text{ (IEN) (FGI + FGO):} \quad R \leftarrow 1$$

The symbol + between FGI and FGO in the control function designates a logic OR operation. This is ANDed with IEN and $T'_0 T'_1 T'_2$ .

We now modify the fetch and decode phases of the instruction cycle. Instead of using only timing signals $T_0$, $T_1$, and $T_2$ (as shown in Fig. 3.5) we will AND the three timing signals with R′ so that the fetch and decode phases will be recognized from the three control functions $R'T_0$, $R'T_1$, and $R'T_2$. The reason for this is that after the instruction is executed and SC is cleared to 0, the control will go through a fetch phase only if R = 0. Otherwise, if R = I, the control will go through an interrupt cycle. The interrupt cycle stores the return address (available in PC) into memory location 0, branches to memory location 1, and clears IEN, R, and SC to 0. This can be done with the following sequence of micro operations.

$$RT_0: \quad AR \leftarrow 0, \qquad TR \leftarrow PC$$

$$RT_1: \quad M[AR] \leftarrow TR, PC \leftarrow 0$$

$$RT_2: \quad PC \leftarrow PC + 1, IEN \leftarrow 0, R \leftarrow 0, SC \leftarrow 0$$

During the first timing signal AR is cleared to 0, and the content of PC is transferred to the temporary register TR. With the second timing signal, the return address is stored in memory at location 0 and PC is cleared to 0. The third timing signal increments PC to 1, clears IEN and R, and control goes back to $T_0$ by clearing SC to 0. The beginning of the next instruction cycle has the condition $R'T_0$ and the content of PC is equal to 1. The control then goes through an instruction cycle that fetches and executes the BUN instruction in location 1.

## SUMMARY

1. The operation code (op code) part of the instruction contains three bits and the meaning of the remaining 13 bits depends on the operation code encountered.

2. The type of instruction is recognized by the computer control from the four bits in positions 12 through 15 of the instruction.

3. Before investigating the operations performed by the instructions, let us discuss the type of instructions that must be included in a computer.

4. Arithmetic, logical, and shift instructions provide computational capabilities for processing the type of data that the user may wish to employ.

5. A program residing in the memory unit of the computer consists of a sequence of instructions.

6. There are two major types of control organization: hardwired control and micro programmed control. in the hardwired organization, the control logic is implemented with gates, flip-flops, decoders, and other digital circuits.

7. The first seven register-reference instructions perform clear, complement, circular shift, and increment micro operations on the AC or E registers.

8. In order to specify the micro operations needed for the execution of each instruction, it is necessary that the function that they are intended to perform be defined precisely.

9. A computer can serve no useful purpose unless it communicates with the external environment. Instructions and data stored in memory must come from some input device.

10. The process of communication just described is referred to as programmed control transfer.

11. The process of communication just described is referred to as programmed control transfer. The computer keeps checking the flag bit, and when it finds it set, it initiates an information transfer.

## SELF ASSESSMENT

1. What is instruction set? Explain any five out of the available set of instructions.

2. What are program interrupts? What is their significance?

3. A computer uses a memory unit with 256K words of 32 bits each. A binary instruction code is stored in one word of memory. The instruction has four parts: an indirect bit, an operation code, a register code part to specify one of 64 registers, and an address part.

    a. How many bits are there in the operation code, the register code part, and the address part?

    b. Draw the instruction word format and indicate the number of bits in each part.

    c. How many bits are there in the data and address inputs of the memory?

4. What is the difference between a direct and an indirect address instruction? How many references to memory are needed for each type of instruction to bring an operand into a processor register?

5. Explain why each of the following microoperations cannot be executed during a single clock pulse in the system shown in Fig. 3-4. Specify a sequence of microoperations that will perform the operation.

    a. $IR \leftarrow M[PC]$

    b. $AC \leftarrow AC + TR$

    c. $DR \leftarrow DR + AC$ (AC does not change)

6. What are the two instructions needed in the basic computer in order to set the E flip-flop to 1?

7. Draw a timing diagram similar to Fig. 3-3 assuming that SC is cleared to 0 at time $T_3$ if control signal $C_7$ is active.

$$C_7 T_3: SC \leftarrow 0$$

$C_7$ is activated with the positive clock transition associated with $T_1$.

8. The content of AC in the basic computer is hexadecimal A937 and the initial value of E is 1. Determine the contents of AC, E, PC, AR, and IR in hexadecimal after the execution of the CLA instruction. Repeat 11 more times, starting from each one of the register-reference instructions. The initial value of PC is hexadecimal 021.

# CHAPTER-IV

# Micro Programmed Control

**Author: Dr. Manoj Duhan**                                    **Vetter: Dr. Pradeep Bhatia**

## 4.1 INTRODUCTION

The function of the control unit in a digital computer is to initiate sequences of microoperations. The number of different types of microoperations that are available in a given system is finite. The complexity of the digital system is derived from the number of sequences of microoperations that are performed. When the control signals are generated by hardware using conventional logic design techniques, the control unit is said to be hardwired. Microprogramming is a second alternative for designing the control unit of a digital computer. The principle of microprogramming is an elegant and systematic method for controlling the micro operation sequence in a digital computer.

The control function that specifies a microoperation is a binary variable. When it is in one binary state, the corresponding microoperation is executed. A control variable in the opposite binary state dies not change the state of the registers in the system. The active state of a control variable may be either the 1 state or the 0 state, depending on the application. In a bus-organized system, the control signals that specify micro operation are groups of bits that select the paths in multiplexers, decoders, and arithmetic logic units.

The control unit initiates a series of sequential steps of microoperations. During any given time, certain microoperations are to be initiated, while others remain idle. ;the control variables at any given time can be represented by a string of 1's and 0' a called a control word. As such, control words can be programmed to perform various operations on the components are stored in memory is called a micro programmed control unit. Each word in control memory contains within it a microinstruction. The microinstruction specifies one or more micro operations fro the system. A sequence of microinstructions constitutes a microprogram. Since alterations of the microprogram are not need once the control unit is in operation, the control memory can be a read-only memory (ROM). The content of the words in ROM are fixed and cannot be altered by simple programming since no writing capability is available in the ROM. ROM words are made permanent during he hardware production of the unit. The use of a microprogram involves placing all control variables in words of ROM for use by the control unit through successive read operations. The content of the word in ROM at a given address specifies a microinstruction.

A more advanced development known as dynamic microprogramming permits a microprogram to be loaded initially from an auxiliary memory such as a magnetic disk. Control units that use dynamic microprogramming employ a writable control memory. This type of memory can be used for writing (to change the microprogram) but is used mostly for reading. A memory that is part of a control unit is referred to as a control memory.

A computer that employs a micro programmed control unit will have two separate memories: a main memory and a control memory. The main memory is available to the user for storing the program. The contents of main memory may alter when the data are manipulated and every time that the program is changed. The user's program  in main memory consists of machine instructions and data. In contrast, the control memory holds a fixed microprogram that cannot be altered by the occasional user. The microprogram consists of microinstructions that specify various internal control signals for execution of register microoperations. Each machine instruction initials a series of microinstructions in control memory. These microinstructions

generate the rations to fetch the instruction from main memory; to evaluate the effective address, to execute the operation specified by the instruction, and to return control to the fetch phase in order to repeat the cycle for the next instruction.

The general configuration of a microprogrammed control unit is demonstrated in the block diagram of Fig. 4-1. The control memory is assumed to be a ROM, within which all control information is permanently stored. The
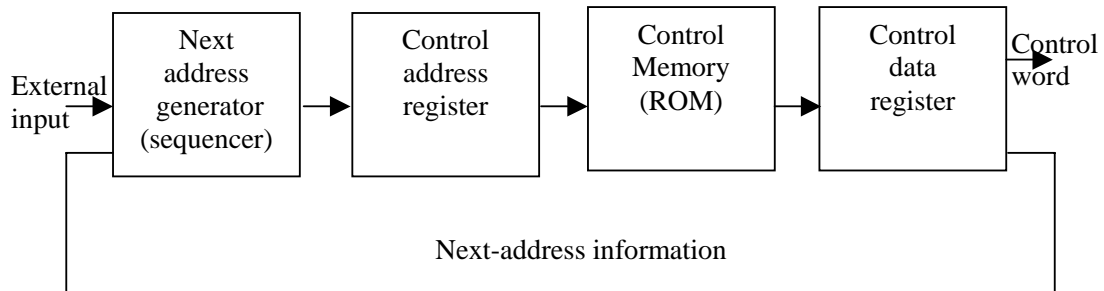


**Fig. 4-1 Microprogrammed control organization**

control memory address register specifies the address of the microinstruction, and the control data register holds the microinstruction read from memory. The microinstruction contains a control word that specifies one or more microoperations for the data processor. Once these operations are executed, the control must determine the nest address. The location of the next microinstruction may be the one next in sequence, or it may be located somewhere else in the control memory. For this reason it is necessary to use some bits of the present microinstruction to control the generation of the address of the next microinstruction. The next address may also be a function of external input conditions. While the microoperations are being executed, the next address is computer in the next address generator circuit and then transferred into the control address register to read the next microinstruction. Thus a microinstruction contains bits for initiating microoperations in the data processor part and bits that determine the address sequence for the control memory.

The next address generator is sometimes called a microprogram sequencer, as it determines the address sequence that is read from control memory. ;the address of the next microinstruction can be specified in several ways, depending on the sequencer inputs. Typical functions of a microprogram sequencer are incrementing the control address register by one, loading into the control address register an address from control memory, transferring an external address, or loading an initial address to start the control operations.

The control data register holds the present microinstruction while the next address is computed and read from memory; the data register is sometimes called a pipeline register. It allows the execution of the microoperations specified by the control word simultaneously with the generation of the next microinstruction. This configuration requires a two-phase clock., with one clock applied to the address register and the other to the data register.

The system can operate without the control data register by applying a single-phase clock to the address register. The control word and next-address information are taken directly from the control memory. ;it must be realized that a ROM operates as a combinational circuit, with the address value as the input and the corresponding word as the output. The content of the specified word in ROM remains in the output wires as long as its address value remains in the address register. No read signal is needed as in a random-access memory.Each clock pulse will execute the microoperations specified by the control word and also transfer a new address to the control address register. In the example that follows we assume a single-phase clock and therefore we do

not use a control data register. In this way the address register is the only component in the control system that receives clock pulses. The other two components: the sequencer and the control memory are combinational circuits and do not need a clock.

The main advantage of the micro programmed control is the fact that once the hardware configuration is established; there should be no need for further hardware or wiring changes. If we want to establish a different control sequence for the system, all we need to do is specify a different set of microinstruction for control memory. The hardware configuration should not be changed for different operations; the only thing that must be changed is the microprogram residing in control memory.

It should be mentioned that must computers based on the reduced instruction set computer (RISC) architecture concept, use hardwired control rather than a control memory with a microprogram.

## 4.2 ADDRESS SEQUENCING

Microinstructions are stored in control memory in groups, with each group specifying routine. Each computer instruction has its own microprogram routine in control memory to generate the microoperations that execute the instruction. The hardware that controls the address sequencing of the control memory must be capable of sequencing the microinstructions within a routine and be able to branch from one routine to another. To appreciate the address sequencing in a microprogram control unit, let us enumerate the steps that the control must undergo during the execution of a single computer instruction.

An initial address is loaded into the control address register when power is turned on in the computer. This address is usually the address of the first microinstruction that activates the instruction fetch routine. The fetch routine may be sequenced by incrementing the control address register through the rest of its microinstructions. At the end of the fetch routine, the instruction is in the instruction register of the computer.

The control memory next must go through the routine that determines the effective address of the operand. A machine instruction may have bits that specify various addressing modes, such as indirect address and index registers. The effective address computation routine in control memory can be reached through a branch microinstruction, which is conditioned on the status of the mode bits of the instruction. When the effective address computation routine is completed, the address of the operand is available in the memory address register.

The next step is to generate the microoperations that execute the instruction fetched from memory. The microoperation steps to be generated in processor register depend on the operation code part of the instruction. Each instruction has its own microprogram routine stored in a given location of control memory. The transformation from the instruction code bits to an address in control memory where the routine is located is referred to as a mapping process. A mapping procedure is a rule that transforms the instruction code into a control memory address. Once the required routine is reached, the microinstructions that execute the instruction may be sequenced by incrementing the control address register, but sometimes the sequence of microoperations will depend on values of certain status bits in processor registers. Micro programs that employ subroutines will require an external register for storing the return address. Return addresses cannot be stored in ROM  because the unit has no writing capability.

When the execution of the instruction is completed, control must return to the fetch routine. This is accomplished by executing an unconditional branch microinstruction to the first address of the fetch routine. In summary, the address sequencing capabilities required in control memory are:

1.       Incrementing of the control address register.

2.     Unconditional branch or conditional branch, depending on statues bit conditions.

3.     A mapping process from the bits of the instruction to an address for control memory.

4.     A facility for subroutine call and return.

Figure 4-2 shows a block diagram of a control memory and the associated hardware needed for selecting the next microinstruction. The microinstruction in control memory contains a set of bits to initiate micro oprerations in computer registers and other bits to specify the method by which the next address is obtained. The diagram shows four different paths from which the control address register (CAR) receive the address. The incremented increments the content of the control address register by one, to select the next microinstruction in sequence. Branching is achieved by specifying the branch address in one of the fields of the microinstruction. Conditional branching is obtained by using part of the microinstruction to select a specific statues bit in order to determine its condition. An external address is transferred into control memory via a mapping logic circuit. The return address for a subroutine is stored in a special register whose value is then used when the microprogram wishes to return form the subroutine.

### 4.2.1     CONDITIONAL BRANCHING

The branch logic provides decision-making capabilities in the control unit. The status conditions are special bits in the system that provide parameter information such as the carry-out of an adder, the sign bit of a number, the mode bits of an instruction, and input or output status conditions. Information in these bits can be tested and actions initiated based on their condition: whether their value is 1 or 0. The status bits, together with the field in the microinstruction that specifies a branch address, control the conditional branch decisions generated in the branch logic.

The branch logic hardware may be implemented in a variety of ways. The simplest way is to test the specified condition and branch to the indicated address if the condition is met; otherwise, the address register is incremented.

This can be implemented with a multiplexer. Suppose that there are eight status bit conditions in the system. Three bits in the microinstruction are used to specify any one of eight status bit conditions. These three bits provide the selection variables for the multiplexer. If the selected status bits in the 1 state, the output of the multiplexer is 1; otherwise, it is 0. A 1 output in the multiplexer generates a control signal to transfer the branch address from the microinstruction into the control address register. A 0 output in the multiplexer causes the address register to be incremented. In this configuration, the microprogram follows one of two possible paths, depending on the value of the selected status bit.

An unconditional branch microinstruction can be implemented by loading the branch address from control memory into the control address register. This can be accomplished by fixing the value of one status bit at the input of the multiplexer, so it is always equal to 1. A reference to this bit by the status bit select lines from control memory causes the branch address to loaded into the control address register unconditionally.
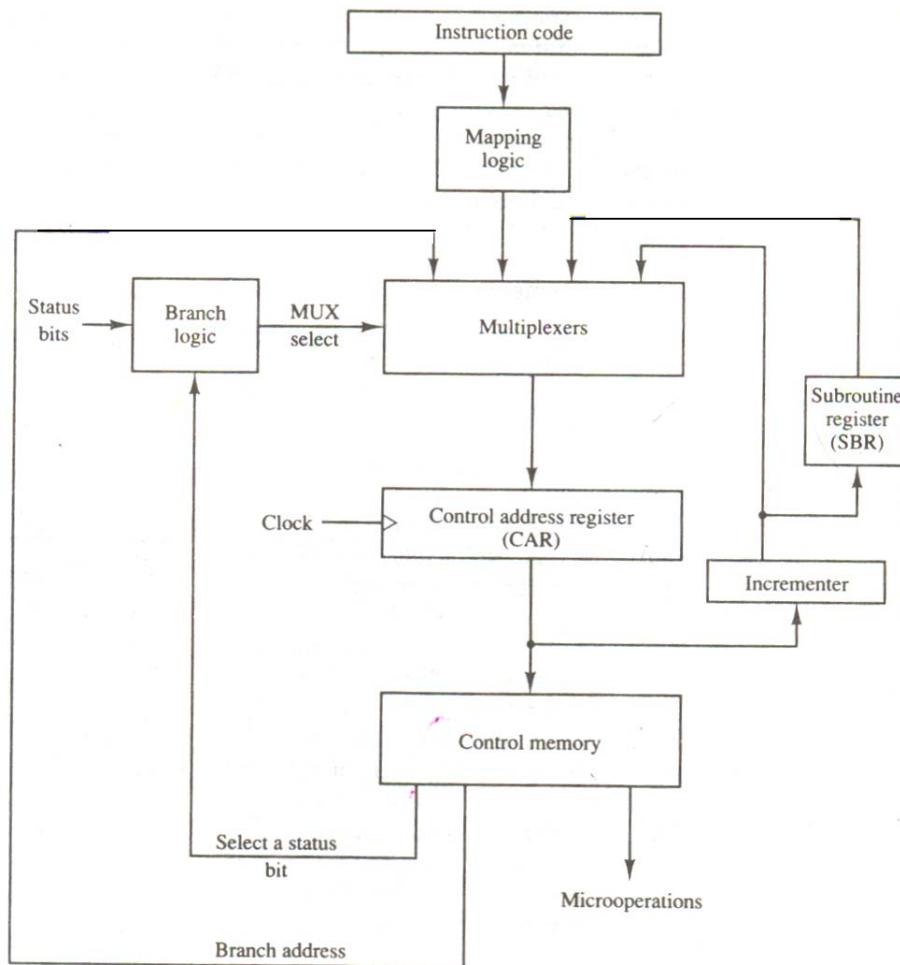
**Fig. 4.2** Selection of Address for control memory.

### 4.2.3 MAPPING OF INSTRUCTION

A special type of branch exists when a microinstruction specifies a branch to the first word in control memory where a microprogram routine for an instruction is located. The status bits for this type of branch are the bits in the operation code part of the instruction. For example, a computer with a simple instruction format as shown in Fig. 4-3 has an operation code of four bits which can specify up to 16 distinct instructions. Assume further that the control memory has 128 words, requiring an address of seven bits. For each operation code there exists a micro program routine in control memory that executes the instruction. One simple mapping process that converts the 4-bit operation code to a 7-bit address for control memory is shown in fig. 4-3. This mapping consists of placing a 0 in the most significant bit of the address, transferring the four operation code bits, and clearing the two lest significant bits of the control address register. This provides for each computer instruction a microprogram routine with a capacity of four microinstructions. If the routine needs more than four microinstructions, it can use addresses 1000000 through 1111111. If it uses fewer than four microinstructions, the unused memory cautions would be available for other routines.

One can extend this concept to a more general mapping rule by using a ROM to specify the mapping function. In this configuration, the bits of the instruction specify the address of a mapping ROM. The contents of the mapping ROM give the bits for the control address register. In this way the microprogram routine that executes the instruction can be placed in any desired location in control memory. The mapping concept provides flexibility for adding instruction for control memory as the need arises.
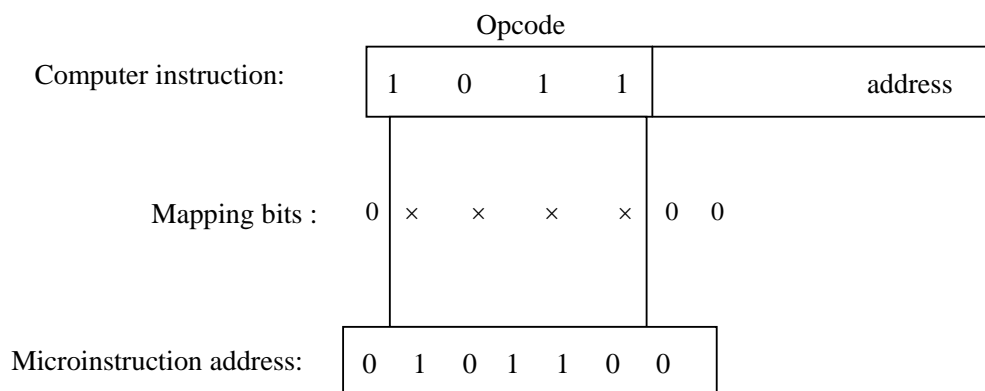
Opcode

| Computer instruction: | 1 | 0 | 1 | 1 | address |

| Mapping bits : | 0 | × | × | × | × | 0 | 0 |

| Microinstruction address: | 0 | 1 | 0 | 1 | 1 | 0 | 0 |

**Figure 4-3** Mapping from instruction code to microinstruction address.

The mapping function is sometimes implemented by means of an integrated circuit called programmable logic device or PLD. A PLD is similar to ROM in concept except that it uses AND and OR gates with internal electronic fuses. The interconnection between inputs, AND gates, OR gates, and outputs can be programmed as in ROM. A mapping function that can be expressed in terms of Boolean expressions can be implemented conveniently

## 4.2.4   SUBROUTINES

Subroutines are programs that are used by other routines to accomplish a particular task. A subroutine can be called from any point within the main body of the microprogram. Frequently, many micro programs contain identical sections of code. Micro instructions can be saved by employing subroutines that use common sections of microcode. For example, the sequence of microoperation needed to generate the effective address of the operand for an instruction is common to all memory reference instructions. This sequence could be a subroutine that is called from within many other routines to execute the effective address computation.

Micro programs that use subroutines must have a provision for storing the return address during a subroutine call and restoring the address during a subroutine return. This may be accomplished by placing the incremented output form the control address register into a subroutine register and branching to the beginning of the subroutine. The subroutine register can then become the source for transferring the address for the return to the main routine. The best way to structure a register file that stores addresses for subroutines is to organize the registers in a last-in, first-out (LIFO) stack.

## 4.3     MICROPROGRAM EXAMPLE

Once the configuration of a computer and its micro programmed control unit is established, the designer's task is to generate the microcode for the control memory. This code generation is called microprogramming and is a process similar to conventional machine language programming. To appreciate this process, we present here a simple digital computer and show how it is micro programmed. The computer used here is similar but not identical to the basic computer.

### 4.3.1 COMPUTER CONFIGURATION

The block diagram of the computer is shown in Fig 4-4. It consists of two memory units: a main memory for storing instructions and data, and a control memory for storing he microprogram. Four register are associated with the processor unit and two with the control unit. The processor registers are program counter PC, address register AR, data register DR, and accumulator register
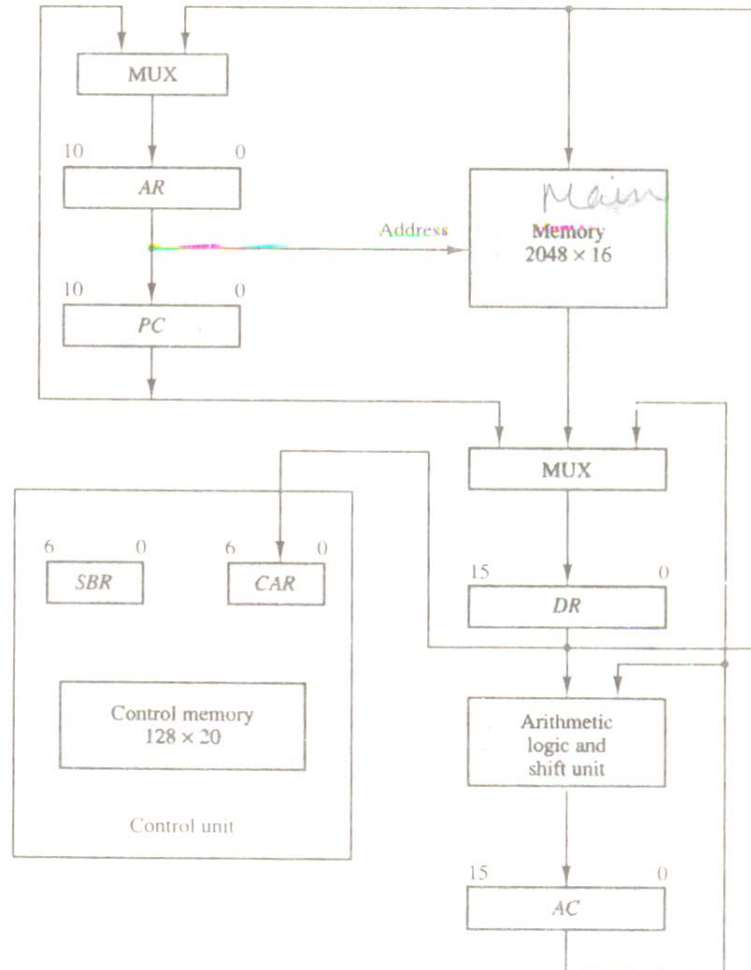


**Figure-4-4** Computer hardware configuration

DR, and accumulator register AC. The function of these registers is similar to the basic computer. The control unit has a control address register CAR and a subroutine register SBR. The control memory and its registers are organized as a micro programmed control unit, as shown in Fig. 4-2.

The transfer of information among the register in the processor is done through multiplexers rather than a common bus. DR can receive information from AC, PC, or memory. AR can receive information from PC or DR, PC can receive information only from AR. The arithmetic, logic, and shift unit performs micro opreations with data from AC and DR and places the result in AC. Note that memory receive its address from AR. Input data written to memory come from DR, and data read from memory can go only to DR.

The computer instruction format is depicted in Fig. 4-5(a). It consists of three fields: a 1-bit field for indirect addressing symbolized by I, a 4-bit operation code (opcode), and an 11-bit

address field. Figure 4-5(b) lists four of the 16 possible memory-reference instructions. The ADD instruction adds the content of the operand found in the effective address to the content of AC. The BRANCH instruction causes a branch to the effective address if the operand in AC is negative. The program proceeds with the next consecutive instruction if AC is not negative. The AC is negative if its sign bit (the bit in the leftmost position of the register) is a 1. The STORE instruction transfers the content of AC into the memory word specified by the effective address. The EXCHANGE instruction swaps the data between AC and the memory word specified by the effective address.

It will be shown subsequently that each computer instruction must be micro programmed. In order not to complicate the microprogramming example, only four instructions are considered here. It should be realized that 12 other instructions can be included and each instruction must be micro programmed by the procedure outlined below.

### 4.3.2 MICROINSTRUCTION FORMAT

The microinstruction format for the control memory is shown in Fig. 4-6. The 20 bits of the microinstruction are divided into four functional parts. The three fields F1, F2, and F3 specify microoperations for the computer. The CD field
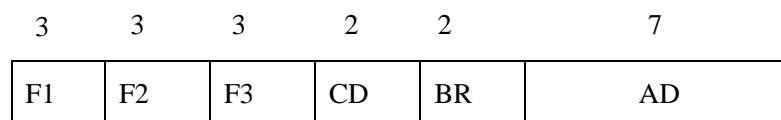
| 1 | Opcode | Address |
|---|--------|---------|

(a) Instruction format

| Symbol | Opcode | Description |
|--------|--------|-------------|
| ADD | 0000 | $AC \leftarrow AC + M[EA]$ |
| BRANCH | 0001 | If $(AC < 0)$ then $(PC \leftarrow EA)$ |
| STORE | 0010 | $M[EA] \leftarrow AC$ |
| EXCHANGE | 0011 | $AC \leftarrow M[EA], M[EA] \leftarrow AC$ |

EA is the effective address

(b) Four computer instructions

**Fig 4-5 Computer instructions**

| 3 | 3 | 3 | 2 | 2 | 7 |
|---|---|---|---|---|---|
| F1 | F2 | F3 | CD | BR | AD |

F1, F2 F3: Microoperation fields

CD: Condition for branching

BR : Branch field

AD : Address field

**Figure 4.6** Microinstruction code format (20 bits).

selects status bit conditions. The BR field specifies the type of branch to be used. The AD field contains a branch address. The address field is seven bits wide, since the control memory has 128 = $2^7$ words.

The microoperations are subdivided into three fields of three bits each. The three bits in each field are encoded to specify seven distinct microoperations as listed in Table 4-1. This gives a total of 21 microoperations. No more than three microoperations can be chosen for a microinstruction, one from each field. If fewer than three microoperations are used, one or more of the fields will use the binary code 000 for no operation. As an illustration, a microinstruction can specify two simultaneous micro operation from F2 and F3 and none from F1.

$$DR \leftarrow M [AR] \quad \text{with F2} = 100$$

$$\text{and} \quad PC \leftarrow PC + 1 \quad \text{with F3} = 101$$

The nine bits of the microoperation fields will than be 000 100 101. It is important to realize that two or more conflicting microoperations cannot be specified simultaneously. For example, a microoperation field 010 001 000 has no meaning because it specifies the operations to clear AC to 0 and subtract DR from AC at the same time.

Each microoperation in Table 4-1 is defined with a register transfer statement and is assigned a symbol for use in a symbolic microprogram. All transfer-type microoperations symbols use five letters. The first two letters designate the source register, the third letters is always a T, and the last two letters designate the destination register. For example, the microoperation that specifies the transfer AC $\leftarrow$ DR (F1 = 100) has the symbol DRTAC, which stands for a transfer from DR to AC.

The CD (condition) field consists of two bits which are encoded to specify four status bit conditions as listed in Table 7-1. The first condition is always a 1, so that a reference to CD = 00 (or the symbol U) will always find the condition to be true. When this condition is used in conjunction with the BR (branch) field, it provides an unconditional branch operation. The indirect bit

**Table 4-1** Symbols and Binary Code for Microinstruction Fields

| F1 | Microoperation | Symbol |
|-----|-----|-----|
| 000 | None | NOP |
| 001 | AC $\leftarrow$ AC + DR | ADD |
| 010 | AC $\leftarrow$ 0 | CLRAC |
| 011 | AC $\leftarrow$ AC + 1 | INCAC |
| 100 | AC $\leftarrow$ DR | DRTAC |
| 101 | AR $\leftarrow$ DR(0-10) | DRTAR |
| 110 | AR $\leftarrow$ PC | PCTAR |
| 111 | M [AR] $\leftarrow$ DR | WRITE |

| F2 | Microoperation | Symbol |
|-----|-----|-----|
| 000 | None | NOP |

| | | |
|---|---|---|
| 001 | $AC \leftarrow AC - DR$ | SUB |
| 010 | $AC \leftarrow AC \vee DR$ | OR |
| 011 | $AC \leftarrow AC \wedge DR$ | AND |
| 100 | $DR \leftarrow M [AR]$ | READ |
| 101 | $DR \leftarrow AC$ | ACTDR |
| 110 | $DR \leftarrow DR + 1$ | INCDR |
| 111 | $DR (0\text{-}10) \leftarrow PC$ | PCTDR |

| F3 | Microoperation | Symbol |
|---|---|---|
| 000 | None | NOP |
| 001 | $AC \leftarrow AC \oplus DR$ | XOR |
| 010 | $AC \leftarrow \overline{AC}$ | COM |
| 011 | $AC \leftarrow shl\ AC$ | SHL |
| 100 | $AC \leftarrow shr\ AC$ | SHR |
| 101 | $PC \leftarrow PC + 1$ | INCPC |
| 110 | $PC \leftarrow AR$ | ARTPC |
| 111 | Reserved | |

| CD | Condition | Symbol | Comments |
|---|---|---|---|
| 00 | Always = 1 | U | Unconditional branch |
| 01 | DR(15) | I | Indirect address bit |
| 10 | AC(15) | S | Sign bit of AC |
| 11 | AC = 0 | Z | Zero value in AC |

| BR | Symbol | Function |
|---|---|---|
| 00 | JMP | $CAR \leftarrow AD$ if condition = 1 |
| | | $CAR \leftarrow CAR + 1$ if condition = 0 |
| 01 | CALL | $CAR \leftarrow AD, SBR \leftarrow CAR + 1$ if condition = 1 |
| | | $CAR \leftarrow CAR + 1$ if condition = 0 |
| 10 | RET | $CAR \leftarrow SBR$ (Return from subroutine) |
| 11 | MAP | $CAR (2\text{-}5) \leftarrow DR (11\text{-}14), CAR (0,1,6) \leftarrow 0$ |

I is available from bit 15 of DR after an instruction is read from memory. The sign bit of AC provides the next status bit. The zero value, symbolized by Z, is a binary variable whose value is

equal to 1 if all the bits in AC are equal to zero. We will use the symbols U, I, S, and Z for the four status bits when we write micro programs in symbolic form.

The BR (branch) field consists of two bits. It is used, in conjunction with address field AD, to choose the address of the next microinstruction. As shown in Table 4-1, when BR = 00, the control performs a jump (JMP) operation (which is similar to a branch), and when BR = 01, it performs a call to subroutine (CALL) operation. The two operations are identical except that a call microinstruction stores the return address in the subroutine register SBR. The jump and call operations depend on the value of the CD field. If the status bit condition specified in the CD field is equal to 1, the next address in the AD field is transferred to the control address register CAR. Otherwise, CAR is incremented by 1.

The return from subroutine is accomplished with a BR field equal to 10. This causes the transfer of the return address from SBR to CAR. The mapping from the operation code bits of the instruction to an address for CAR is accomplished when the BR field is equal to 11. This mapping is as depicted in Fig. 4-3. the bits of the operation code are in DR after an instruction is read from memory. Note that the last two conditions in the BR field are independent of the values in the CD and AD fields.

### 4.3.3 SYMBOLIC MICROINSTRUCTIONS

The symbols defined in Table 4-1 can be used to specify microinstructions in symbolic form. A symbolic micro program can be translated into its binary equivalent by means of an assembler. A microprogram assembler is similar in concept to a conventional computer assembler. The simplest and most straightforward way to formulate an assembly language for a microprogram is to define symbols for each field of the microinstruction and to give users the capability for defining their own symbolic addresses.

Each line of the simply language micro program defines a symbolic microinstruction. Each symbolic microinstruction is divided into five fields: label, microoperation, CD, BR, and AD. The fields specify the following information.

1. The label field may be empty or it may specify a symbolic address. A label is terminated with a colon (:)

2. The microoperations field consists of one, two, or three symbols, separated by commas, from those defined in Table 4-1. There may be no more than one symbol from each F field. The NOP symbol is used when the microinstruction has no microoperations. This will be translated by the assembler to nine zeros.

3. The CD field has one of the letters U, I,S, or Z.

4. The BR field contains one of the four symbols defined in Table 4-1.

5. The AD field specifies a value for the address field of the microinstruction in one of three possible ways:

a. With a symbolic address, which must also appear as a label.

b. With the symbol NEXT to designate the next address in sequence.

c. When the BR field contains a RET or MAP symbol, the AD field is lift empty and is converted to seven zeros by the assembler.

We will use also the pseudo instruction ORG to define the origin, or first address, of a microprogram routine. Thus the symbol ORG 64 informs the assembler to place the next microinstruction in control memory at decimal address 64, which is equivalent to the binary address 1000000.

### 4.3.4 THE FETCH ROUTINE

The control memory has 128 words, and each word contains 20 bits. To microprogram the control memory, it is necessary to determine the bit values of each of the 128 words. The first 64 words (addresses 0 to 63) are to be occupied by the routines for the 16 instructions. The last 64 words may be used for any other purpose. A convenient starting location for the fetch routine is address 64. the microinstructions needed for the fetch routine are

$$AR \leftarrow PC$$

$$DR \leftarrow M [AR], PC \leftarrow PC + 1$$

$$AR \leftarrow DR (0\text{-}10), CAR (2\text{-}5) \leftarrow DR (11\text{-}14), CAR 90,1,6) \leftarrow 0$$

The address of the instruction is transferred from PC to AR and the instruction is then read from memory into DR. Since no instruction register is available, the instruction code remains in DR. The address part is transferred to AR and then control is transferred to one of 16 routines by mapping the operation code part of the instruction from DR into CAR.

The fetch routine needs three microinstructions, which are placed in control memory at addresses 64, 65, and 66. Using the assembly language conventions defined previously, we can write the symbolic microprogram for the fetch routine as follows:

|  | ORG 64 |  |  |  |
|---|---|---|---|---|
| FETCH: | PCTAR | U | JMP | NEXT |
|  | READ, INCPE | U | JMP | NEXT |
|  | DRTAR | U | MAP |  |

The translation of the symbolic microprogram to biary produces the following binary microprogram. The bit values are obtained from Table 4-1.

| Binary Address | F1 | F2 | F3 | CD | BR | AD |
|---|---|---|---|---|---|---|
| 1000000 | 110 | 000 | 000 | 00 | 00 | 1000001 |
| 1000001 | 000 | 100 | 101 | 00 | 00 | 1000010 |
| 1000010 | 101 | 101 | 000 | 00 | 11 | 0000000 |

The three microinstructions that constitute the fetch routine have been listed in three different representations. The register transfer representation shows the internal register transfer operations that each microinstruction implements. The symbolic representation is useful for writing microprograms in an assembly language format. The binary representation is the actual internal content that must be stored in control memory. It is customary to write microprograms in symbolic form and then use an assembler program to obtain a translation to binary.

### 4.3.5 SYMBOLIC MICROPROGRAM

The execution of the third (MAP) microinstruction in the fetch routine results in a branch to address 0××××0, where ×××× are the four bits of the operation code. For example, if the instruction is an ADD instruction whose operation code is 0000, the MAP microinstruction will

transfer to CAR the address 0000000, which is the start address for the ADD routine in control memory. The first address for the BRANCH and STORE routines are 0 0001 00 (decimal 4) and 0 0010 00 (decimal 8), respectively. The first address for the other 13 routines are at address values 12, 16, 20,…., 60. This gives four words in control memory for each routine.

In each routine we must provide microinstructions for evaluating the effective address and for executing the instruction. The indirect address mode is associated with all memory-reference instructions. A saving in the number of control memory words may be achieved if the microinstructions for the indirect address are stored as a subroutine. This subroutine, symbolized by INDRCT, is located right after the fetch routine, as shown in Table 7-2. the table also shows the symbolic microprogram for the fetch routine and the microinstruction in the ADD routine calls subroutine INDRCT, conditioned on status bit I. If I = 1 , a branch to INDRCT occurs and the return address (address 1 in this case) is stored in the subroutine register SBR. The INDRCT subroutine has two microinstructions:

|  | | | |
|---|---|---|---|
| INDRCT: | READ | U | JMP | NEXT |
| | KRTAR | U | RET | |

**Table 4.2** Symbolic Micro program

| Label | Micro operations | CD | BR | AD |
|---|---|---|---|---|
| | ORG 0 | | | |
| ADD | NOP | I | CALL | INDRCT |
| | READ | U | JMP | NEXT |
| | ADD | U | JMP | FETCH |
| | ORG 4 | | | |
| BRANCH: | NOP | S | JMP | OVER |
| | NOP | U | JMP | FETC |
| OVER: | NOP | I | CALL | INDRCT |
| | ARTPC | U | JMP | FETCH |
| | ORG 8 | | | |
| STORE: | NOP | I | CALL | INDRCT |
| | ACTDR | U | JMP | NEXT |
| | WRITE | U | JMP | FETCH |
| | ORG 12 | | | |
| EXCHANGE: | NOP | I | CALL | INDRCT |

|  | READ | U | JMP | NEXT |
|---|---|---|---|---|
|  | ACTDR, DRTAC | U | JMP | NEXT |
|  | WRITE | U | JMP | FETCH |

|  | ORG 64 | | | |
|---|---|---|---|---|
| FETCH: | PCTAR | U | JMP | NEXT |
|  | READ, INCPC | U | JMP | NEXT |
|  | DRTAR | U | MAP | |
| INDRCT: | READ | U | JMP | NEXT |
|  | DRTAR | U | RET | |

Remember that an indirect address considers the address part of the instruction as the address where the effective address is stored rather than the address of the operand. Therefore, the memory has to be accessed to get the effective address, which is then transferred to AR. The return from subroutine (RET) transfers the address from SBR to CAR, thus returning to the second microinstruction of the ADD routine.

The execution of the ADD instruction is carried out by the microinstruction at addresses 1 and 2. The first microinstruction reads the operand from memory into DR. the second microinstruction performs an add microoperation with the content of DR and AC and then jumps back to the beginning of the fetch routine.

The BRANCH instruction should cause a branch to the effective address if AC < 0. The AC will be less than zero if its sign is negative, which is detected from status bit S being a 1. The BRANCH routine in Table 4-2 starts by checking the value of S. If S is equal to 0, no branch occurs and the next microinstruction causes a jump back to the fetch routine without altering the content of PC. If S is equal to 1, the first JMP microinstruction transfers control to location OVER. The microinstruction at this location calls the INDRCT subroutine if I = 1. The effective address is then transferred from AR to PC and the microprogram jumps back to the fetch routine.

The STORE routine again uses the INDRCT subroutine if I = 1. The content of AC is transferred into DR. A memory write operation is initiated to store the content of DR in a location specified by the effective address in AR.

The EXCHANGE routine reads the operand from the effective address and places it in DR and AC are interchanged in the third microinstruction. This interchange is possible when the registers are of the edge-triggered type. The original content of AC that is now in now in DR is stored back in memory.

Note that Table 4-2 contains a partial list of the microprogram. Only four out of 16 possible computer instructions have been micro programmed. Also control memory words at locations 69 to 127 have not be4en used. Instructions such as multiply, divide, and others that

require a long sequence of microoperations will need more than four microinstructions for their execution. Control memory words 69 to 127 can be used for this purpose.

### 4.3.6 BINARY MICROPROGRAM

The symbolic microprogram is a convenient form for writing micro programs in a way that people can read and understand. But this is not the way that the microprogram is stored in memory. The symbolic microprogram must be translated to binary either by means of an assembler program or by the user if the microprogram is simple enough as in this example.

The equivalent binary form of the microprogram is listed in Table 4-3. The addresses for control memory are given in both decimal and binary. The binary content of each microinstruction is derived from the symbols and their equivalent binary values as defined in Table 4-1.

Note that address 3 has no equivalent in the symbolic microprogram since the ADD routine has only three microinstructions at addresses 0,1, and 2. The next routine starts at address 4. Even though address 3 is not used. Some binary value must be specified for each word in control memory. We could have specified all 0's in the word since this location will never be used. However, if some unforeseen error occurs, or if a noise signal sets CAR to the value of 3, it will be wise to jump to jump to address 64, which is the beginning of the fetch routine.

The binary micro program listed in Table 4-3 specifies the word content of the control memory. When a ROM is used for the control memory, the

| Micro Routine | Address | | Binary Microinstruction | | | | | |
|---|---|---|---|---|---|---|---|---|
| | Decimal | Binary | F1 | F2 | F3 | CD | BR | AD |
| ADD | 0 | 0000000 | 000 | 000 | 000 | 01 | 01 | 1000011 |
| | 1 | 0000001 | 000 | 100 | 000 | 00 | 00 | 0000010 |
| | 2 | 0000010 | 001 | 000 | 000 | 00 | 00 | 1000000 |
| | 3 | 0000011 | 000 | 000 | 000 | 00 | 00 | 1000000 |
| BRANCH | 4 | 0000100 | 000 | 000 | 000 | 10 | 00 | 0000110 |
| | 5 | 0000101 | 000 | 000 | 000 | 00 | 00 | 1000000 |
| | 6 | 0000110 | 000 | 000 | 000 | 01 | 01 | 1000011 |
| | 7 | 0000111 | 000 | 000 | 110 | 00 | 00 | 1000000 |
| STORE | 8 | 0001000 | 000 | 000 | 000 | 01 | 01 | 1000011 |
| | 9 | 0001001 | 000 | 101 | 000 | 00 | 00 | 0001010 |
| | 10 | 0001010 | 111 | 000 | 000 | 00 | 00 | 1000000 |
| | 11 | 0001011 | 000 | 000 | 000 | 00 | 00 | 1000000 |

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| EXCHANGE | 12 | 0001100 | 000 | 000 | 000 | 01 | 01 | 1000011 |
| | 13 | 0001101 | 001 | 000 | 000 | 00 | 00 | 0001110 |
| | 14 | 0001110 | 100 | 101 | 000 | 00 | 00 | 0001111 |
| | 15 | 0001111 | 111000 | 000 | 000 | 00 | 00 | 1000000 |
| FETCH | 64 | 1000000 | 110 | 000 | 000 | 00 | 00 | 1000001 |
| | 65 | 1000001 | 000 | 100 | 101 | 00 | 00 | 1000010 |
| | 66 | 1000010 | 101 | 000 | 000 | 00 | 11 | 0000000 |
| INDRCT | 67 | 1000011 | 000 | 100 | 000 | 00 | 00 | 1000100 |
| | 68 | 1000100 | 101 | 000 | 000 | 00 | 10 | 0000000 |

microprogram binary list provides the truth table for fabricating the unit. This fabrication is a hardware process and consists of creating a mask for the ROM so as to produce the 1's and 0's for each word. The bits of ROM are fixed once the internal links are fused during the hardware production. The ROM is made of IC packages that can be removed if necessary and replaced by other packages. To modify the instruction set of the computer, it is necessary to generate a new microprogram and mask a new ROM. The old one can be removed and the new one inserted in its place.

If a writable control memory is employed, the ROM is replaced by a RAM. The advantage of employing RAM for the control memory is that the microprogram can be altered simply by writing a new pattern of 1's and 0's without resorting to hardware procedures. A writable control memory possesses the flexibility of choosing the instruction set a computer dynamically by changing the microprogram under processor control. However, most microprogrammed systems use a ROM for the control memory because it is cheaper and faster than a RAM and also to prevent the occasional user from changing the architecture of the system.

## 4.4 DESIGN OF CONTROL UNIT

The bits of the microinstruction are usually divided into fields, with each field defining a distinct, separate function. The various fields encountered in instruction formats provide control bits to initiate microoperations in the system, special bits to specify the way that the next address is to be evaluated, and an address field for branching. The number of control bits that initiate micro operation can be reduced by grouping mutually exclusive variables into fields and encoding the k bits in each field to provide 2 microoperations. Each field requires a decoder to produce the corresponding control signals. This method reduces the size of the microinstruction bits but requires additional hardware external to the control memory. It also increases the delay time of the control signals because they must propagate through the decoding circuits.

The encoding of control bits was demonstrated in the programming example of the preceding section. The nine bits of the microoperation field are divided into three subfields of three bits each. The control memory output of each subfield must be decoded to provide the distinct microoperations. The outputs of the decoders are connected to the appropriate inputs in the processor unit.

Figure 4-7 shows the three decoders and some of the connections that must be made from their outputs. Each of the three fields of the microinstruction presently available in the output of control memory are decoded with a 3 × 8 decoder to provide eight outputs. Each of these outputs must be connected to the proper circuit to initiate the corresponding microoperation as specified in Table 7-1. For example, when FI = 101 (binary 5), the next clock pulse transition transfers the content of DR (0-10) to AR (symbolized by DRTAR in Table 7-1). Similarly, when F1 = 101 (binary 6) there is a transfer from PC to AR (symbolized by PCTAR). As shown in Fig. 4-7, outputs 5 and 6 of decoder F1 are connected to the load input of AR so that when either one of these outputs is active, information from the multiplexers is transferred to AR. The multiplexers select the information from DR when output 5 is active and from PC when output 5 in inactive. The transfer into AR occurs with a clock pulse transition only when output 5 or output 6 of the decoder are active. The other outputs of the decoders that initiate transfers between registers must be connected in a similar fashion.

The arithmetic logic shift unit can be designed, instead of using gates to generate the control signals marked by the symbols AND, ADD, and DR in Fig 4.7, these inputs will now come from the outputs of the decoders associated with the symbols AND, ADD, and DRTAC, respectively



**Figure 4-7** Decoding of microoperation fields

as shown in Fig. 4-7. the other output of the decoders that are associated with an AC operation must also be connected to the arithmetic logic shift unit in a similar fashion .

### 4.4.1    MICROPROGRAM SEQUENCER

The basic components of a micro programmed control unit are the control unit  are the control memory and the circuits that select the next address. The address selection part is called a

microprogram sequencer. A  microprogram sequencer can be constructed with digital functions to suit a particular application. However, just as there are large ROM units available in integrated circuit packages, so are general purpose sequencers suited for the construction of microprogram control units. To guarantee a wide range of acceptability, an integrated circuit sequencer must provide an internal organization that can be adapted to a wide range of applications.

The purpose of a microprogram sequencer is to present is to present an address to the control memory so that a microinstruction may be read and executed. The next-address logic of the sequencer determines the specific address source to be loaded into the control address register. The choice of the address source is guided by the next-address information bits that the sequencer receives from the present microinstruction. Commercial sequencers include within the unit an internal register stack used for temporary storage of addresses during microprogram looping and subroutine calls. ;some sequencers provide an output register which can function as the address register for the control memory.

To illustrate the internal structure of a typical microprogram sequencer we will show a particular unit that is suitable for use in the microprogram computer example developed in the preceding section. The block diagram of the microprogram sequencer is shown in Fig. 4-8. The control memory is included in the diagram to show the interaction between the sequencer and the memory attached to it. There are two multiplexers in the circuit. The first multiplexer selects an address from one of four sources and routes it into a control address register CAR. The second multiplexer tests the value of a selected status bit and the result of the test is applied to an input logic circuit. The output from CAR is incremented and applied to one of the multiplexer inputs and to the subroutine register SBR. The other three inputs to multiplexer number 1 come from the address field of the present microinstruction, from the output of SBR, and from an external source that maps the instruction. Although the diagram shows a single subroutine register, a typical sequencer will have a register stack about four to eight levels deep. In this way, a number of subroutines can be active at the same time. A push and pop operation, in conjunction with a stack pointer, stores and retrieves the return address during the call and return microinstructions.

The CD (condition) field of the microinstruction selecting one of the status bits in the second multiplexer. If the bit selected is equal to 1, the T (test)variable is equal to 1; otherwise, it is equal to 0. The T value together with the two bits from the BR (branch) field go to an input logic circuit. The input logic in a particular sequencer will determine the type of operations that are available in the unit. Typical sequencer operations are: increment, branch or jump, call and return from subroutine, load an external address, push or pop the stack, and other address sequencing operations. With three inputs, the sequencer  can provide up to eight address sequencing operations. Some commercial sequencers have three or four inputs in addition to the T input and thus provide a wider range of operations.

The input logic circuit in Fig. 4-8 has three inputs, $I_0$, $I_1$, and T, and three outputs, $S_0$, $S_1$ and L. variables $S_0$ and $S_1$ select one of the source addresses for CAR. Variable L enables the load input in SBR. The binary values of the two selection variables determine the path in the multiplexer. For example, with $S_1 S_0 = 10$, multiplexer input number 2 is selected and establishes a transfer
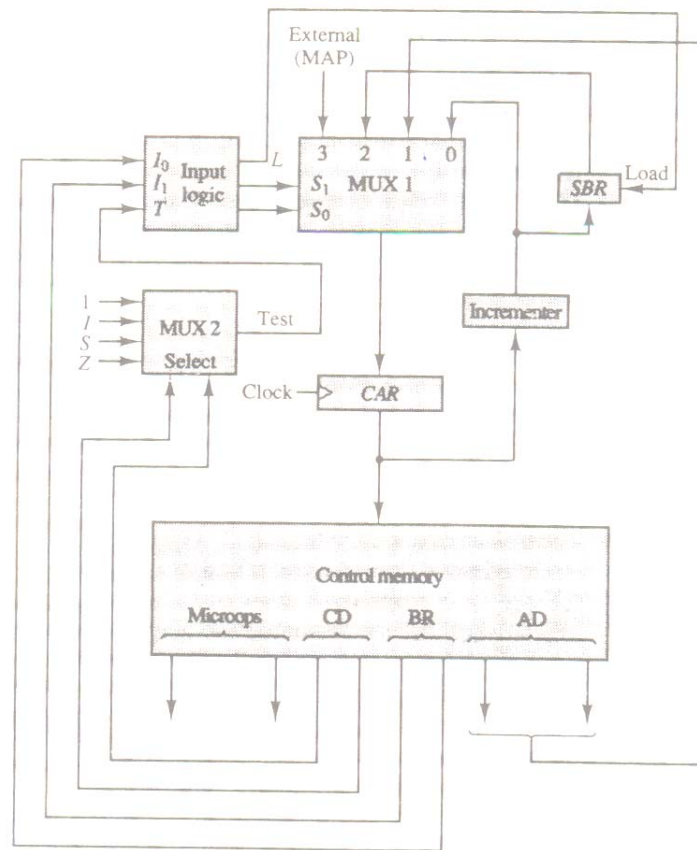
**Figure 4-8** Microprogram sequencer for a control memory. Page 234.

path from SBR to CAR. Note that each of the four inputs as well as the output of MUX 1 contains a 7-bit address.

The truth table for the input logic circuit is shown in Table 4-4. inputs $I_1$ and $I_0$ are identical to the bit values in the BR field. The function listed in each entry was defined in Table 4-1. The bit values for $S_1$ and $S_0$ are determined from the stated function and the path in the multiplexer that establishes the required transfer. The subroutine register is loaded with the incremented value of CAR during a call microinstruction (BR = 01) provided that the status bit condition is satisfied (T = 1). The truth table can be used to obtain the simplified Boolean functions for the input logic:

$$S_1 = I_1$$
$$S_0 = I_1 I_0 + I'_1 T$$
$$L = I'_1 I_0 T$$

**Table 4-4**       Input Logic Truth Table for Microprogram Sequencer

| BR | Field | $I_1$ | $I_0$ | T | $S_1$ | $S_0$ | Load SBR L |
|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 |
| 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 1 | 1 | 0 | 1 | 1 |
| 1 | 0 | 1 | 0 | × | 1 | 0 | 0 |
| 1 | 1 | 1 | 1 | × | 1 | 1 | 0 |

The circuit cab be constructed with three AND gates, an OR gate, and an inverter.

Note that the incremented circuit in the sequencer of Fig. 4-8 is not a counter constructed with flip-flops but rather a combinational circuit constructed with gates. A combinational circuit incrementer can be designed by cascading a series of half-adder circuits. The output carry from one state must be applied to the input of the next stage. One input in the first least significant stage must be equal to 1 to provide the increment-by-one operation.

## SUMMARY

1. The function of the control unit in a digital computer is to initiate sequences of micro operations. The number of different types of micro operations that are available in a given system is finite.

2. The control function that specifies a micro operation is a binary variable. When it is in one binary state, the corresponding micro operation is executed.

3. A computer that employs a micro programmed control unit will have two separate memories: a main memory and a control memory.

4. Microinstructions are stored in control memory in groups, with each group specifying routine. Each computer instruction has its own micro program routine in control memory to generate the micro operations that execute the instruction.

5. The branch logic provides decision-making capabilities in the control unit.

6. A special type of branch exists when a microinstruction specifies a branch to the first word in control memory where a micro program routine for an instruction is located.

7. Subroutines are programs that are used by other routines to accomplish a particular task.

8. Once the configuration of a computer and its micro programmed control unit is established, the designer's task is to generate the microcode for the control memory.

9. Each line of the simply language micro program defines a symbolic microinstruction. Each symbolic microinstruction is divided into five fields: label, micro operation, CD, BR, and AD. The fields specify the following information.

10. The purpose of a micro program sequencer is to present is to present an address to the control memory so that a microinstruction may be read and executed.

## SELF ASSESSMENT

1. What is the difference between a microprocessor and a microprogram? Is it possible to design a microprocessor without a microprogram? Are all micro programmed computers also micro processors?

2. Explain the difference between hardwired control and micro programmed control. Is it possible to have a hardwired control associated with a control memory?

3.   Define the following: (a) micro operation; (b) micro instruction; (c) micro program; (d) microcode.

4.   The micro programmed control organization showing in Fig. 4.1 has the following propagation delay times. 20 ns to generate the next address, 10 ns to transfer the address into the control address register, 20 ns to access the control memory ROM, 5 ns to transfer the microinstruction into the control data register, and 20 ns to perform the required microoperations specified by the control word. What is the maximum clock frequency that the control can use? What would the clock frequency be if the control data register is not used?

5.   The system shown in Fig. 4.2 uses a control memory of 1024 words of 32 bits each. The microinstruction has three fields as shown in the diagram. The micro operations field has 16 bits.

   a.   How many bits are there in the branch address field and the select field?

   b.   If there are 16 status bits in the system, how many bits of the branch logic are used to select a status bit?

   c.   How many bits are left to select an input for the multiplexers?

6.   The control memory in Fig. 4.2 has 4096 words of 24 bits each.

   a.   How many bits are there in the control address register?

   b.   How many bits are there in each of the four inputs shown going into the multiplexers?

   c.   What are the number of inputs in each multiplier and how many multiplexers are needed?

7.   Using the mapping procedure described in Fig. 4.3, give the first microinstruction address for the following operation code: (a) 0010; (b) 1011; (c) 1111.

8.   Formulate a mapping procedure that provides eight consecutive microinstructions for each routine. The operation code has six bits and the control memory has 2048 words.

9.   Explain how the mapping from an instruction code to a microinstruction address can be done by means of a read-only memory. What is the advantage of this method compared to the one in Fig. 4-3?

10.  Why do we need the two multiplexers in the computer hardware configuration showing in Fig. 4.4? Is there another way that information from multiple sources can be transferred to a common destination?

# CHAPTER-V

# CENTRAL PROCESSING UNIT

**Author: Dr. Manoj Duhan**                    **Vetter: Dr. Pradeep Bhatia**

## 5.1 INTRODUCTION

The part of the computer that performs the bulk of data-processing operations is called the central processing unit and is referred to as the CPU. The CPU is made up of three major parts, as shown in Fig.5-1. The register set stores intermediate data used during the execution of the instructions. The arithmetic logic unit (ALU) performs the required microoperations for executing the instructions. The control unit supervises the transfer of information among the registers and instructs the ALU as to which operation to perform.

The CPU performs a variety of functions dictated by the type of instructions that are incorporated in the computer. Computer architecture is sometimes defined as the computer structure and behavior as seen by the programmer that uses machine language instructions. This includes the instruction formats, addressing modes, the instruction set, and the general organization of the CPU registers.

One boundary where the computer designer ad the computer programmer see the same machine is the part of the CPU associated with the instruction set. From the designer's point of view, the computer instruction set provides the specifications for the design of the CPU. The design of a CPU is



**Figure 5-1** Major components of CPU.

a task that in large part involves choosing the hardware for implementing the machine instructions. The user who programs the computer in machine or assembly language must be aware of the register set, the memory structure, instruction performs.

Design examples of simple CPUs are carried out in previous chapters. This chapter describes the organization and architecture of the CPU with an emphasis on the user's view of the computer. We briefly describe how the registers communicate with the ALU through buses and explain the operation of the memory stack. We then present the type of instruction formats available, the addressing modes used to retrieve data from memory, and typical instructions commonly incorporated in computers. The last section presents the concept of reduced instruction set computer (RISC).
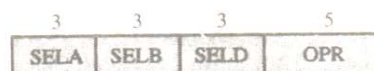
## 5.2 GENERAL REGISTER ORGANIZATION

In the programming examples, we have shown that memory locations are needed for storing pointers, counters, return addresses, temporary results, and partial products during multiplication. Having to refer to memory locations for such applications is time consuming because memory access is the most time-consuming operation in a computer. It is more

convenient and more efficient to store these intermediate values in processor registers. When a large number of registers are included in the CPU, it is most efficient to connect them



(a) Block diagram

(b) Control word

**Figure 5-2** Register set with common ALU.

through a common bus system. The registers communicate with each other not only for direct data transfers, but also while performing various microoperations. Hence it is necessary to provide a common unit that can perform all the arithmetic, logic, and shift microoperations in the processor.

A bus organization for seven CPU registers is shown in Fig. 5-2. The output of each register is connected to two multiplexers (MUX) to form the two buses A and B. The selection lines in each multiplexer select one register or the input data for the particular bus. The A and B

buses form the inputs to a common arithmetic logic unit (ALU). The operation selected in the ALU determines the arithmetic or logic microoperation that is to be performed. The result of the microoperation is available for output data and also goes into the inputs of all the registers. The register that receives the information from the output bus is selected by a decoder. The decoder activates one of the register load inputs, thus providing a transfer path between the data in the output bus and the inputs of the selected destination register.

The control unit that operates the CPU bus system directs the information flow through the registers and ALU by selecting the various components in the system. For example, to perform the operation.

$$R1 \leftarrow R2 + R3$$

the control must provide bi9nary selection variables to the following selector inputs:

1. MUX A selector (SELA): to place the content of R2 into bus A.

2. MUX B selector (SELB): to place the content of R3 into bus B.

3. ALU operation selector (OPR): to provide the arithmetic addition A + B.

4. Decoder destination selector (SELD): to transfer the content of the output bus into R1.

The four control selection variables are generated in the control unit and must be available at the beginning of a clock cycle. The data from the two source registers propagate through the gates in the multiplexers and the ALU, to the output bus, and into the inputs of the destination register, all during the clock cycle interval. Then, when the next clock transition occurs, the binary information from the output bus is transferred into R1. To achieve a fast response time, the ALU is constructed with high-speed circuits. The buses are implemented with multiplexers or three-state gates.

## 5.2.1 CONTROL WORD

There are 14 binary selection inputs in the unit, and their combined value specifies a control word. The 14-bit control word is defined in Fig. 5-2. It consists of four fields. Three fields contain three bits each, and one field has five bits. The three bits of SELA select a source register for the A input of the ALU. The three bits of SELB select a register for the B input of the ALU. The three bits of SELD select a destination register using the decoder and its seven load output. The five bits of OPR select one of the operations in the ALU. The 14-bit control word when applied to the selection inputs specify a particular microoperation.

The encoding of the register selections is specified in Table 5-1. The 3-bit

**Table 5.1** Encoding of Register Selection Fields

| Binary Code | SELA | SELB | SELD |
|---|---|---|---|
| 000 | Input | Input | None |
| 001 | R1 | R1 | R1 |
| 010 | R2 | R2 | R2 |
| 011 | R3 | R3 | R3 |
| 100 | R4 | R4 | R4 |
| 101 | R5 | R5 | R5 |

| | | | |
|---|---|---|---|
| 110 | R6 | R6 | R6 |
| 111 | R7 | R7 | R7 |

binary code listed in the first column of the table specifies the binary code for each of the three fields. The register selected by fields SELA, SELB, and SELD is the one whose decimal number is equivalent to the binary number in the code. When SELA or SELB is 000, the corresponding multiplexer selects the external input data. When SELD = 000, no destination register is selected but the contents of the output bus are available in the external output.

The ALU provides arithmetic and logic operations. In addition, the CPU must provide shift operations. The shifter may be placed in the input of the ALU to provide a preshift capability, or at the output of the ALU to provide post shifting capability. In some cases, the shift operations are included with the ALU. The encoding of the ALU operations for the CPU is specified in Table 5.2. The OPR field has five bits and each operation is designated with a symbolic name.

**Table 5.2** Encoding of ALU Operations

| OPR Select | Operation | Symbol |
|---|---|---|
| 00000 | Transfer A | TSFA |
| 00001 | Increment A | INCA |
| 00010 | Add A + B | ADD |
| 00101 | Subtract A − B | SUB |
| 00110 | Decrement A | DECA |
| 01000 | AND A and B | AND |
| 01010 | OR A and B | OR |
| 01100 | XOR A and B | XOR |
| 01110 | Complement A | COMA |
| 10000 | Shift right A | SHRA |
| 11000 | Shift left A | SHLA |

## 5.2.2 EXAMPLE OF MICRO OPERATIONS

A control word of 14 bits is needed to specify a microoperation in the CPU. The control word for a given microoperation can be derived from the selection variables. For example, the subtract microoperation given by the statement.

$$R1 \leftarrow R2 - R3$$

specifies R2 for the A input of the ALU, R3 for the B input of the ALU, R1 for the destination register, and an ALU operation to subtract A − B. Thus the control word is specified by the four fields and the corresponding binary value for each field is obtained from the encoding listed in Tables 5-1 and 5-2. The binary control word for the subtract microoperation is 010  011  001  00101 and is obtained as follows :

| Field : | SELA | SELB | SELD | OPR |
|---|---|---|---|---|

| Symbol : | R2 | R3 | R1 | SUB |
|---|---|---|---|---|
| Control word: | 010 | 011 | 001 | 00101 |

The control word for this micro operation and a few others are listed in Table5-3.

The increment and transfer microoperations do not use the B input of the ALU. For these cases, the B field is marked with a dash. We assign 000 to any unused field when formulating the binary control word, although any other binary number may be used. To place the content of a register into the output terminals we place the content of the register into the A input of the ALU, but none of the registers are selected to accept the data. The ALU operation TSFA places the data from the register, through the ALU, into the output terminals. The direct transfer from input to output is accomplished with a control word of all 0's (making the B field 000). A register can be

**Table 5-3** Examples of Microoperations for the CPU

| Microoperation | Symbolic Designation | | | | Control Word | | | |
|---|---|---|---|---|---|---|---|---|
| | SELA | SELB | SELD | OPR | | | | |
| R1←R2–R3 | R2 | R3 | R1 | SUB | 010 | 011 | 001 | 00101 |
| R4←R4 ∨ R5 | R4 | R5 | R4 | OR | 100 | 101 | 100 | 01010 |
| R6←R6 + 1 | R6 | — | R6 | INCA | 110 | 000 | 110 | 00001 |
| R7 ← R1 | R1 | — | R7 | TSFA | 001 | 000 | 111 | 00000 |
| Output ←R2 | R2 | — | None | TSFA | 010 | 000 | 000 | 00000 |
| Output← Input | Input | — | None | TSFA | 000 | 000 | 000 | 00000 |
| R4← sh1 R4 | R4 | — | R4 | SHLA | 100 | 000 | 100 | 11000 |
| R5 ← 0 | R5 | R5 | R5 | XOR | 101 | 101 | 101 | 01100 |

cleared to 0 with an exclusive-OR operation. This is because $x \oplus x = 0$.

It is apparent from these examples that from these ex apples that many other microoperations can be generated in the CPU. The most efficient way to generate control words with a large number of bits is to store them in a memory unit. A memory unit that stores control words is referred to as a control memory. Y reading consecutive control words from memory, it is possible to initiate the desired sequence of microoperation for the CPU. This type of control is referred to as microprogrammed control. A microprogrammed control unit is shown in Fig. 7-8. The binary control word for the CPU will come from the outputs of the control memory marked "micro-ops."

## 5-3 STACK ORGANIZATION

A useful feature that is included in the CPU of most computers is a stack or last-in, first-out (LIFO) list. A stack is a storage device that stores information in such a manner that the item stored last is the first item retrieved. The operation of a stack can be compared to a stack of trays. The last tray placed on top of the stack is the first to be taken off.

The stack in digital computers is essentially a memory unit with an address register that can count only (after an initial value is loaded into it). The register that holds the address for the

stack is called a stack pointer (SP) because its value always points at the top item in the stack. Contrary to a stack of trays where the tray itself may be taken out or inserted, the physical registers of a stack are always available for reading or writing. It is the content of the word that is inserted or deleted.

The two operations of a stack are the insertion and deletion of items. The operation of insertion is called push (or push-down) because it can be through of as the result of pushing a new item on top. The operation of deletion is called pop (or pop-up) because it can be thought of as the result of removing one item so hat the stack pops up. However, nothing is pushed or popped in a computer stack. these operation are simulated by incrementing or decrementing the stack pointer register.

### 5.3.1 REGISTER STACK

A stack can be placed in portion of a large memory or it can be organized as a collection of a finite number of memory words or registers. Figure 5-3 shows the organization of a 64-word register stack. The stack pointer register SP contains a binary number whose value is equal to the address of the word that is currently on top of the stack. Three items are placed in the stack: A, B, and C, in that order. Item C is on top of the stack so that the content of SP is now 3. To remove the top item, the stack is popped by reading the memory word

Figure 5-3 Block diagram of a 64-word stack

at address 3 and decrementing the content of SP. Item B is now on top the stack since SP holds address 2. To insert a new item, the stack is pushed by incrementing SP and writing a word in the next-higher location in the stack. Note that item C has been read out but not physically removed. This does not matter because when the stack is pushed, a new item is written in its place.

In a 64-word stack, the stack pointer contains 6 bits because $2^6 = 64$. Since SP has only six bits, it cannot exceed a number greater than 63 (111111 in binary). When 63 is incremented by 1, the result is 0 since $111111 + 1 = 1000000$ in binary, but SP can accommodate only the six

least significant bits. Similarly, when 000000 is decremented by 1, the result is 111111. The one-bits register FULL is set to 1 when the stack is full, and the one-bit register EMTY is set to 1 when the stack is empty of items. DR is the data register that holds the binary data to be written into or read out of the stack. .

Initially, SP is cleared to 0, EMTY is set to 1, and FULL is cleared to 0, so that SP points to the word at address 0 and the stack is marked empty and not full. If the stack is not full (if FULL = 0), a new item is inserted with a push operation. The push operation is implemented with the following sequence of microoperations:

| | |
|---|---|
| $SP \leftarrow SP + 1$ | Increment stack pointer |
| $M[SP] \leftarrow DR$ | Write item on top of the stack |
| If $(SP = 0)$ then $(FULL \leftarrow 1)$ | Check if stack is full |
| $EMTY \leftarrow 0$ | Mark the stack not empty |

The stack pointer is incremented so that it points to the address of the next-higher word. A memory write operation inserts the word form DR into the top of the stack. Note that SP holds the address of the top of the stack and that M[SP] denotes the memory word specified by the address presently available in SP. The first item stored in the stack is at address 1. The last item is stored at address 0. If SP reaches 0, the stack is full of items so FULL is set to 1. This condition is reached if the top item prior to the last push was in location 63 and, after incrementing SP, the last item is stored in location 0.Once an item is stored in location 0, there are no more empty registers in the stack. If an item is written in the stack, obviously the stack cannot be empty, so EMTY is cleared to 0.

A new item is deleted from the stack if the stack is not empty (if EMTY = 0). The pop operation consists of the following sequence of microoperations:

| | |
|---|---|
| $DR \leftarrow M[SP]$ | Read item from the top of stack |
| $SP \leftarrow SP - 1$ | Decrement stack pointer |
| If $(SP = 0)$ then $(EMTY \leftarrow 1)$ | Check if stack is empty |
| $FULL \leftarrow 0$ | Mark the stack not full |

The top item is read from the stack into DR. the stack pointer is then decremented. If its value reaches zero, the stack is empty, so EMTY is set to 1. This condition is reached if the item read was in location 1. Once this item is read out, SP is decremented and reaches the value 0, which is the initial value of SP. Note that if a pop operation reads the item from location 0 and then SP is decremented, SP changes to 111111, which is equivalent to decimal 63. In this configuration, the word in address 0 receives the last item in the stack. Note also that an erroneous operation will result if the stack is pushed when FULL = 1 or popped when EMTY = 1.

## 5.3.2   MEMORY STACK

A stack can exist as a stand-alone unit as in Fig. 5-3 or can be implemented in a random-access memory attached to a CPU. The implementation of a stack in the CPU is done by assigning a portion of memory to a stack operation and using a processor register as a stack pointer. Figure 8-4 shows a portion of computer memory partitioned into three segments: program, data, and stack. The program counter PC points at the address of the next instruction in the program. The address register AR points at an array of data. The stack pointer.
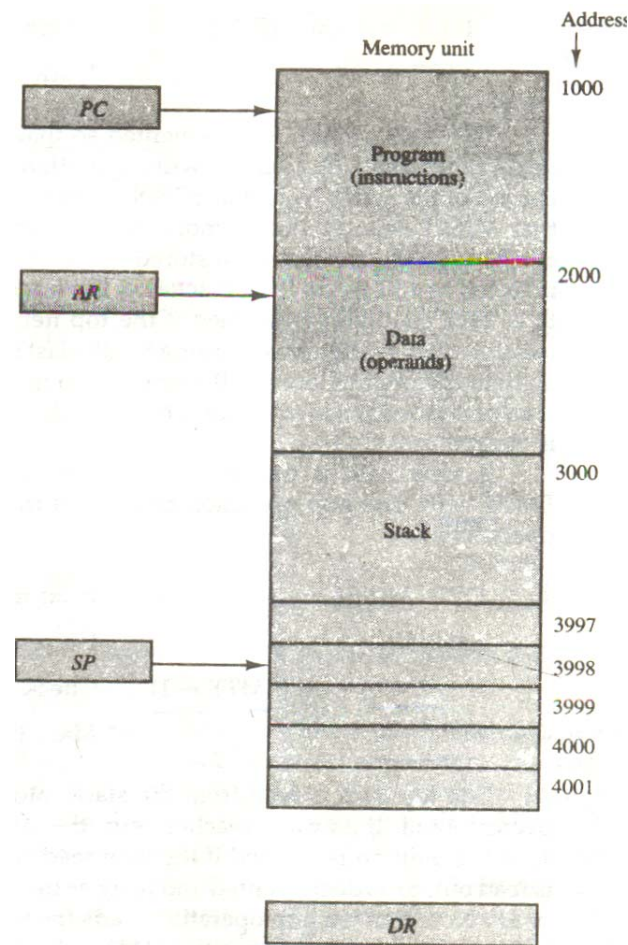
**Figure 5-4** Computer memory with program, data, and stack segments.

SP points at the top of the stack. The three registers are connected to a common address bus, and either one can provide an address for memory. PC is used during the fetch phase to read an instruction. AR is used during the execute phase to read an operand. SP is used to push or pop items into or from the stack

As shown in Fig. 5-4, the initial value of SP is 4001 and the stack grows with decreasing addresses. Thus the first item stored in the stack is at address 4000, the second item is stored at address 3999, and the last address that can be used for the stack is 3000. No provisions are available for stack limit checks.

We assume that the items in the stack communicate with a data register DR. A new item is inserted with the push operation as follow:

$$SP \leftarrow SP - 1$$

$$M[SP] \leftarrow DR$$

The stack pointer is decremented so that it points at the address of the next word. A memory write operation inserts the word form DR into the top of the stack. A new item is deleted with a pop operation as follows:

$$DR \leftarrow M[SP]$$

$$SP \leftarrow SP + 1$$

The top is read form the stack into DR. The stack pointer is then incremented to point at the next item in the stack.

Most computers do not provide hardware to check for stack overflow (full stack) or underflow (empty stack). The stack limits can be checked by using two processor registers: one to hold the upper limit (3000 in this case), and the other to hold the lower limit (4001 in this case). After a push operation, SP is compared with the upper-limit register and after a pop operation, SP is compared with the lower-limit register.

The two microoperations needed for either the push or pop are (1) an access to memory through SP, and (2) updating SP. Which of the two microoperations is dine first and whether SP is updated by incrementing or decrementing depends on the organization of the stack. In Fig. 5-4 the stack grows by decreasing the memory address. The stack may be constructed to grow by increasing the memory address as in Fig. 5-3. In such a case, SP is incremented for the push operation and decremented for the pop operation. A stack may be constructed so that SP points at the next empty location above the top of the stack. In this case the sequence of microoperations must be interchanged.

A stack pointer is loaded with an initial value. This initial value must be the bottom address of an assigned stack in memory. Henceforth, SP is automatically decremented or incremented with every push or pop operation. The advantage of a memory stack is that the CPU can refer to it without having to specify an address, since the address is always available and automatically updated in the stack pointer.

### 5.3.3   REVERSE POLISH NOTATION

A stack organization is very effective for evaluating arithmetic expressions. The common mathematical method of writing arithmetic expression imposes difficulties when evaluated by a computer. The common arithmetic expressions are written in infix notation, with each operation written between the operands. Consider the simple arithmetic expression

$$A * B + C * D$$

The star (denoting multiplication) is placed between two operands A and B or C and D. The plus is between the two products. To evaluate this arithmetic expression it is necessary to compute the product A * B, store this product while computing C * D, and then sum the two products. From this example we see that to evaluate arithmetic expressions in infix notation it is necessary to scan back and forth along the expression to determine the next operation to be performed.

The polish mathematician Lukasiewicz showed that arithmetic expressions can be represented in prefix notation. This representation often referred to as Polish notation, places the operator before the operands. The postifix notation, referred to as reverse polish notation (RPN), places the operator after the operands. The following examples demonstrate the three representations:

|  |  |
|---|---|
| A + B | Infix notation |
| + AB | Prefix or Polish notation |
| AB + | Postfix or reverse Polish notation |

The reverse Polish notation is in a form suitable for stack manipulation. The expression.

$$A * B + C * D$$

is written is reverse Polish notation as

$$AB * CD * +$$

and is evaluated as follows: scan the expression from left to right. When an operator is reached, perform the operation with the two operands found on the left side of the operator. Remove the two operands and the operator and replace them by the number obtained form the result of the operation. Continue to scan the expression and repeat the procedure for every operator encountered until there are no more operators.

For the expression above we find the operator $*$ after A and B. We perform the operation A $*$ B and replace A, B, and $*$ by the product to obtain

$$(A * B)\, CD * +$$

where (A $*$ B) is a single quantity obtained from the product. The next operator is a $*$ and its previous two operands are C and D, so we perform C $*$ D and obtain an expression with two operands and one operator:

$$(A * B)\,(C * D) +$$

the next operator is + and the two operands to be added are the two products, so we add the two quantities to obtain the result.

The conversion from infix notation to reverse Polish notation must take into consideration the operational hierarchy adopted for infix notation. This hierarchy dictates that we first perform all arithmetic inside inner parentheses, then inside outer parentheses, and do multiplication and division operations before addition and subtraction operations. Consider the expression

$$(A + B) * [C * (D + E) + F]$$

to evaluate the expression we must first perform the arithmetic inside the parentheses (A + B) and (D + E). Next we must calculate the expression inside the square brackets. The multiplication of C $*$ (D + E) must be done prior to the addition of F since multiplication has precedence over addition. The last operation is the multiplication of the two terms between the parentheses and brackets. The expression can be converted to reverse Polish notation, without the use of parentheses, by taking into consideration the operation hierarchy. The converted expression is

$$AB + DE + C * F + *$$

Proceeding from left to right, we first add A and B, then add D and E. At this point we are left with

$$(A + B)\,(D + E)\, C * E + *$$

where (A + B) and (D + E) are each a single number obtained from the sum. The two operands for the next $*$ are C and (D + E). These two numbers are multiplied and the product added to F. The final $*$ causes the multiplication of the two terms.

### 5.3.4 EVALUATION OF ARITHMETIC EXPRESSIONS

Reverse Polish notation, combined with a stack arrangement of registers, is the most efficient way known for evaluating arithmetic expressions. This procedure is employed in some electronic calculators and also in some computers. The stack is particularly useful for handling long,

complex problems involving chain calculations. It is based on the fact that any arithmetic expression can be expressed in parentheses-free Polish notation.

The procedure consists of first converting the arithmetic expression into its equivalent reverse Polish notation. The operands are pushed into the stack in the order in which they appear. The initiation of an operation depends on whether we have a calculator or a computer. In a calculator, the operators are entered through the keyboard. In a computer, they must be initiated by instructions that contain an operation field (no address field is required). The following microoperations are executed with the stack when an operation is entered in a calculator or issued by the control in a computer: (1) the two topmost operands in the stack are used for the operation, and (2) the stack is popped and the result of the operation replaces the lower operand. By pushing the operands into the stack continuously and performing the operations as defined above, the expression is evaluated in the proper order and the final result remains on top of the stack.

The following numerical example may clarify this procedure. Consider the arithmetic expression

$$(3 * 4) + (5 * 6)$$

In reverse Polish notation, it is expressed as

$$34 * 56 * +$$

Now consider the stack operating shown in Fig. 8-5. Each box represents one represents one stack operation and the arrow always points to the top of the stack. Scanning the expression from left to right, we encounter two operands. First the number 3 is pushed into the stack, then the number 4. The next symbol is the multiplication operator $*$. This causes a multiplication of the two topmost items in the stack. The stack is then popped and the product is placed on top of the stack, replacing the two original operands. Next we encounter the two operands 5 and 6, so they are pushed into the stack. The stack operation that results from the next $*$ replaces these two numbers by their product. The last operation causes an arithmetic addition of the two topmost numbers in the stack to produce the final result of 42.

Scientific calculators that employ an internal stack require that the user convert the arithmetic expressions into reverse Polish notation. Computers that use a stack-organized CPU provide a system program to perform the conversion for the user. Most compilers, irrespective
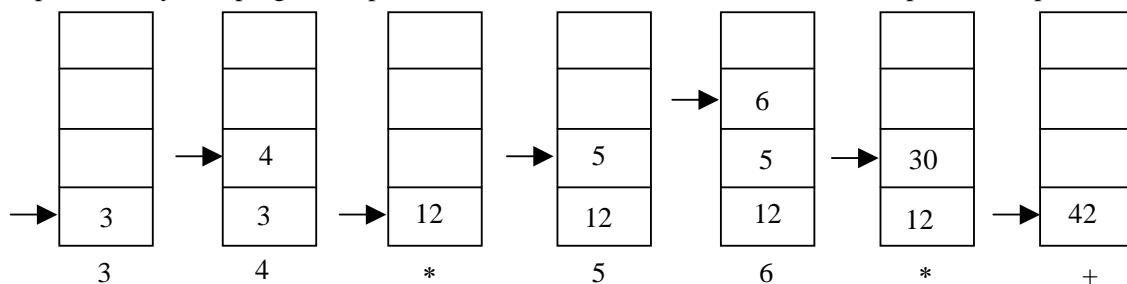


**Figure 5.5** Stack operations to evaluate (3 $*$ 4 + 5 $*$ 6)

of their CPU organization, convert all arithmetic expressions into Polish notation anyway because this is the most efficient method for translating arithmetic expressions into machine language instructions. So in essence, a stack-organized CPU may be more efficient in some applications than a CPU without a stack.

## 5-4 INSTRUCTION FORMATS

The physical and logical structure of computers is normally described in reference manuals provided with the system. Such manuals explain the internal construction of the CPU, including the processor registers available and their logical capabilities. They list all hardware-implemented instructions, specify their binary code format, and provide a precise definition of each instruction. A computer will usually have a variety of instruction code formats. It is the function of the control unit within the CPU to interpret each instruction code and provide the necessary control functions needed to process the instruction.

The format of an instruction is usually depicted in a rectangular box symbolizing the bits of the instruction as they appear in memory words or in a control register. The bits of the instruction are divided into groups called fields. The most common fields found in instruction formats are:

1.      An operation code field that specifies the operation to be performed.

2.      An address field that designates a memory address or a processor register.

3.      A mode field that specifies the way the operand or the effective address is determined.

Other special fields are sometimes employed under certain circumstances, as for example a field that gives the number of shifts in a shift-type instruction.

The operation code field of an instruction is a group of bits that define various processor operations, such as add, subtract, complement, and shift. The bits that define the mode field of an instruction code specify a variety of alternatives for choosing the operands from the given address. The various addressing modes that have been formulated for digital computers are presented in Sec. 5.5. In this section we are concerned with the address field of an instruction format and consider the effect of including multiple address fields is an instruction.

Operations specified by computer instructions are executed on some data stored in memory or processor registers, Operands residing in processor registers are specified with a register address. A register address is a binary number of k bits that defines one of $2^k$ registers in the CPU. Thus a CPU with 16 processor registers R0 through R15 will have a register address field of four bits. The binary number 0101, for example, will designate register R5.

Computers may have instructions of several different lengths containing varying number of addresses. The number of address fields in the instruction format of a computer depends on the internal organization of its registers. Most computers fall into one of three types of CPU organizations:

1.      Single accumulator organization.

2.      General register organization.

3.      Stack organization.

An example of an accumulator-type organization is the basic computer presented in Chap. 5. All operations are performed with an implied accumulator register. The instruction format in this type of computer uses one address field. For example, the instruction that specifies an arithmetic addition is defined by an assembly language instruction as

ADD    X

where X is the address of the operand.  The ADD instruction in this case results in the operation AC ← AC + M[X].  AC is the accumulator register and M[X] symbolizes the memory word located at address X.

An example of a general register type of organization was presented in Fig. 7.1.  The instruction format in this type of computer needs three register address fields.  Thus the instruction for an arithmetic addition may be written in an assembly language as

ADD   R1,    R2,    R3

to denote the operation R1 ← R2 + R3.  The number of address fields in the instruction can be reduced from three to two if the destination register is the same as one of the source registers. Thus the instruction

ADD   R1,    R2

would denote the operation R1 ← R1 + R2.  Only register addresses for R1 and R2 need be specified in this instruction.

Computers with multiple processor registers use the move instruction with a mnemonic MOV to symbolize a transfer instruction. Thus the instruction

MOV   R1,    R2

denotes the transfer R1 ← R2 (or R2 ← R1, depending on the particular computer). Thus transfer-type instructions need tow address fields to specify the source and the destination.

General register-type computers employ two or three address fields in their instruction format. Each address field may specify a processor register or a memory word. An instruction symbolized by

ADD    R1,   X

would specify the operation R1 ← R + M [X]. It has two address fields, one for register R1 and the other for the memory address X.

The stack-organized CPU was presented in Fig. 8-4. Computers with stack organization would have PUSH and POP instructions which require an address field. Thus the instruction

PUSH    X

will push the word at address X to the top of the stack. The stack pointer is updated automatically. Operation-type instructions do not need an address field in stack-organized computers. This is because the operation is performed on the two items that are on top of the stack. The instruction

ADD

in a stack computer consists of an operation code only with no address field. This operation has the effect of popping the two top numbers from the stack, adding the numbers, and pushing the sum into the stack. There is no need to specify operands with an address field since all operands are implied to be in the stack.

Most computers fall into one of the three types of organizations that have just been described. Some computers combine features from more than one organization structure. For example, the Intel 808- microprocessor has seven CPU registers, one of which is an accumulator register. As a consequence, the processor has some  of the characteristics of a general register type and some of the characteristics of a accumulator type. All arithmetic and logic instruction, as well as the load and store instructions, use the accumulator register, so these instructions have

only one address field. On the other hand, instructions that transfer data among the seven processor registers have a format that contains two register address fields. Moreover, the Intel 8080 processor has a stack pointer and instructions to push and pop from a memory stack. The processor, however, does not have the zero-address-type instructions which are characteristic of a stack-organized CPU.

To illustrate the influence of the number of addresses on computer programs, we will evaluate the arithmetic statement

$$X = (A + B) * (C + D)$$

using zero, one, two, or three address instruction. We will use the symbols ADD, SUB, MUL, and DIV for the four arithmetic operations; MOV for the transfer-type operation; and LOAD and STORE for transfers to and from memory and AC register. We will assume that the operands are in memory addresses A, B, C, and D, and the result must be stored in memory at address X.

## 5.4.1  THREE-ADDRESS INSTRUCTIONS

Computers with three-address instruction formats can use each address field to specify either a processor register or a memory operand. The program in assembly language that evaluates X = (A + B) * (C + D) is shown below, together with comments that explain the register transfer operation of each instruction.

| | | |
|---|---|---|
| ADD | R1, A, B | R1 ← M [A] + M [B] |
| ADD | R2, C, D | R2 ← M [C] + M [D] |
| MUL | X, R1, R2 | M [X] ← R1 * R2 |

It is assumed that the computer has two processor registers, R1 and R2. The symbol M [A] denotes the operand at memory address symbolized by A.

The advantage of the three-address format is that it results in short programs when evaluating arithmetic expressions. The disadvantage is that the binary-coded instructions require too many bits to specify three addresses. An example of a commercial computer that uses three-address instructions is the Cyber 170. The instruction formats in the Cyber computer are restricted to either three register address fields or two register address fields and one memory address field.

## 5.4.2  TWO-ADDRESS INSTRUCTIONS

Two address instructions are the most common in commercial computers. Here again each address field can specify either a processor register or a memory word. The program to evaluate X = (A + B) * (C + D) is as follows:

| | | |
|---|---|---|
| MOV | R1, A | R1 ← M [A] |
| ADD | R1, B | R1 ← R1 + M [B] |
| MOV | R2, C | R2 ← M [C] |
| ADD | R2, D | R2 ← R2 + M [D] |
| MUL | R1, R2 | R1 ← R1 * R2 |
| MOV | X, R1 | M [X] ← R1 |

The MOV instruction moves or transfers the operands to and from memory and processor registers. The first symbol listed in an instruction is assumed to be both a source and the destination where the result of the operation is transferred.

### 5.4.3   ONE-ADDRESS INSTRUCTIONS

One-address instructions use an implied accumulator (AC) register for all data manipulation. For multiplication and division there is a need for a second register. However, here we will neglect the second and assume that the AC contains the result of tall operations. The program to evaluate X = (A + B) ∗ (C + D) is

| LOAD | A | AC ← M [A] |
|------|---|-----------|
| ADD | B | AC ← A [C] + M [B] |
| STORE | T | M [T] ← AC |
| LOAD | C | AC ← M [C] |
| ADD | D | AC ← AC + M [D] |
| MUL | T | AC ← AC ∗ M [T] |
| STORE | X | M [X] ← AC |

All operation are done between the AC register and a memory operand. T is the address of a temporary memory location required for storing the intermediate result.

### 5.4.5   ZERO-ADDRESS INSTRUCTIONS

A stack-organized computer does not use an address field for the instructions ADD and MUL. The PUSH and POP instructions, however, need an address field to specify the operand that communicates with the stack. The following program shows how X = (A + B) ∗ (C + D) will be written for a stack organized computer. (TOS stands for top of stack)

| PUSH | A | TOS ← A |
|------|---|---------|
| PUSH | B | TOS ← B |
| ADD | | TOS ← (A + B) |
| PUSH | C | TOS ← C |
| PUSH | D | TOS ← D |
| ADD | | TOS ← (C + D) |
| MUL | | TOS ← (C + D) ∗ (A + B) |
| POP | X | M [X] ← TOS |

To evaluate arithmetic expressions in a stack computer, it is necessary to convert the expression into reverse Polish notation. The name "zero-address" is given to this type of computer because of the absence of an address field in the computational instructions,

### 5.4.6   RISC INSTRUCTIONS

The reduced instruction set computer (RISC) architecture have several advantages. The instruction set of a typical of a typical RISC processor is restricted to the use of load and store instructions when communicating between memory and CPU. All other instruction are executed within the registers of the CPU without referring to memory. A program for a RISC-type CPU

consists of LOAD and STORE instructions that have one memory and one register address, and computational-type instructions that have three addresses with all three specifying processor registers. The following is a program to evaluate X = (A + B) ∗ (C + D).

| | | |
|------|-----------|---------------------|
| LOAD | R1, A | R1 ← M [A] |
| LOAD | R2, B | R2 ← M [B] |
| LOAD | R3, C | R3 ← M [C] |
| LOAD | R4, D | R4 ← M [D] |
| ADD | R1, R1, R2 | R1 ← R1 + R2 |
| ADD | R3, R3, R4 | R3 ← R3 + R4 |
| MUL | R1, R1, R3 | R1 ← R1 ∗ R3 |
| STORE | X, R1 | M [X] ← R1 |

The load instruction transfer the operands form memory to CPU registers. The add and multiply operations are executed with data in the registers without accessing memory. The result of the computations is then stored in memory with a store instruction.

## 5-5  ADDRESSING MODES

The operation field of an instruction specifies the operation to be performed. This operation must be executed on some data stored in computer registers or memory words. The way the operands are chosen during program execution in dependent on the addressing mode of the instruction. The addressing mode of the instruction. The addressing mode specifies a rule for interpreting or modifying the address field of the instruction before the operand is actually referenced. Computers use addressing mode techniques for the purpose of accommodating one or both of the following provisions:

1.      To gave programming versatility to the user by providing such facilities as pointers to memory, counters for loop control, indexing of data, and program relocation.

2.      To reduce the number of bits in the addressing field of the instruction.

3.       The availability of the addressing modes gives the experienced assembly language programmer flexibility for writing programs that are more efficient with respect to the number of instructions and execution time.

To understand the various addressing modes to be presented in this section, it is imperative that we understand the basic operation cycle of the computer.  The control unit of a computer is designed to go through an instruction cycle that is divided into three major phases:

1.      Fetch the instruction from memory

2.      Decode the instruction.

3.      Execute the instruction.

There is one register in the computer called the program counter of PC that keeps track of the instructions in the program stored in memory.  PC holds the address of the instruction to be executed next and is incremented each time an instruction is fetched from memory.  The decoding done in step 2 determines the operation to be performed, the addressing mode of the instruction

and the location of the operands. The computer then executes the instruction and returns to step 1 to fetch the next instruction in sequence.

In some computers the addressing mode of the instruction is specified with a distinct binary code, just like the operation code is specified. Other computers use a single binary code that designates both the operation and the mode of the instruction. Instructions may be defined with a variety of addressing modes, and sometimes, two or more addressing modes are combined in one instruction.

An example of an instruction format with a distinct addressing mode field is shown in Fig. 5-6. The operation code specified the operation to be performed. The mode field is sued to locate the operands needed for the operation. There may or may not be an address field in the instruction. If there is an address field, it may designate a memory address or a processor register. Moreover, as discussed in the preceding section, the instruction may have more than one address field, and each address field may be associated with its own particular addressing mode.

Although most addressing modes modify the address field of the instruction, there are two modes that need no address field at all. These are the implied and immediate modes.

**5.5.1    Implied Mode:** In this mode the operands are specified implicitly in the definition of the instruction.    For example, the instruction "complement accumulator" is an implied-mode instruction because the operand in the accumulator register is implied in the definition of the instruction. In fact, all register reference instructions that sue an accumulator are implied-mode instructions.

| Op code | Mode | Address |
|---------|------|---------|

**Figure  5.6** Instruction format with mode field

Zero-address instructions in a stack-organized computer are implied-mode instructions since the operands are implied to be on top of the stack.

**5.5.2    Immediate Mode :** In this mode the operand is specified in the instruction itself. In other words, an immediate-mode instruction has an operand field rather than an address field. The operand field contains the actual operand to be used in conjunction with the operation specified in the instruction. Immediate-mode instructions are useful for initializing registers to a constant value.

It was mentioned previously that the address field of an instruction may specify either a memory word or a processor register. When the address field specifies a processor register, the instruction is said to be in the register mode.

**5.5.3    Register Mode :** In this mode the operands are in registers that reside within the CPU. The particular register is selected from a register field in the instruction. A k-bit field can specify any one of $2^k$ registers.

**5.5.4    Register Indirect Mode :** In this mode the instruction specifies a register in the CPU whose contents give the address of the operand in memory. In other words, the selected register contains the address of the operand rather than the operand itself. Before using a register indirect mode instruction, the programmer must ensure that the memory address fo the operand is placed in the processor register with a previous instruction. A reference to the register is then equivalent to specifying a memory address. The advantage of a register indirect mode instruction is that the address field of the instruction sues fewer bits to select a register than would have been required to specify a memory address directly.

**5.5.5    Autoincrement or Autodecrement Mode:**  This is similar to the register indirect mode except that the register is incremented or decremented after (or before) its value is used to access memory.  When the address stored in the register refers to a table of data in memory, it is necessary to increment or decrement the register after every access to the table.  This can be achieved by using the increment or decrement instruction.  However, because it is such a common requirement, some computers incorporate a special mode that automatically increments or decrements the content of the register after data access.

The address field of an instruction is used by the control unit in the CPU to obtain the operand from memory. Sometimes the value given in the address field is the address of the operand, but sometimes it is just an address from which the address of the operand is calculated. To differentiate among the various addressing modes it is necessary to distinguish between the address part of the instruction and the effective address used by the control when executing the instruction.  The effective address is defined to be the memory address obtained from the computation dictated by the given addressing mode.  The effective address is the address of the operand in a computational-type instruction.  It is the address where control branches in response to a branch-type instruction.  We have already defined two addressing modes in previous chapter.

**5.5.6    Direct Address Mode:** In this mode the effective address is equal to the address part of the instruction.  The operand resides in memory and its address is given directly by the address field of the instruction.  In a branch-type instruction the address field specifies the actual branch address.

**5.5.7    Indirect Address Mode:** In this mode the address field of the instruction gives the address where the effective address is stored in memory.  Control fetches the instruction from memory and uses its address part to access memory again to read the effective address.

**5.5.8    Relative Address Mode:** In this mode the content of the program counter is added to the address part of the instruction in order to obtain the effective address.  The address part of the instruction is usually a signed number (in 2's complement representation) which can be either positive or negative.  When this number is added to the content of the program counter, the result produces an effective address whose position in memory is relative to the address of the next instruction.  To clarify with an example, assume that the program counter contains the number 825 and the address part of the instruction contains the number 24.  The instruction at location 825 is read from memory during the fetch phase and the program counter is then incremented by one to 826 + 24 = 850.  This is 24 memory locations forward from the address of the next instruction.  Relative addressing is often used with branch-type instructions when the branch address is in the area surrounding the instruction word itself.  It results in a shorter address field in the instruction format since the relative address can be specified with a smaller number of bits compared to the number of bits required to designate the entire memory address.

**5.5.9    Indexed Addressing Mode:** In this mode the content of an index register is added to the address part of the instruction to obtain the effective address. The index register is a special CPU register that contains an index value.  The address field of the instruction defines the beginning address of a data array in memory. Each operand in the array is stored in memory relative to the beginning address. The distance between the beginning address and the address of the operand is the index value stores in the index register.  Any operand in the array can be accessed with the same instruction provided that the index register contains the correct index value.  The index register can be incremented to facilitate access to consecutive operands.  Note that if an index-type instruction does not include an address field in its format, the instruction converts to the register indirect mode of operation.

Some computers dedicate one CPU register to function solely as an index register. This register is involved implicitly when the index-mode instruction is used. In computers with many processor registers, any one of the CPU registers can contain the index number. In such a case the register must be specified explicitly in a register field within the instruction format.

**5.5.10 Base Register Addressing Mode:** In this mode the content of a base register is added to the address part of the instruction to obtain the effective address. This is similar to the indexed addressing mode except that the register is now called a base register instead of an index register. The difference between the two modes is in the way they are used rather than in the way that they are computed. An index register is assumed to hold an index number that is relative to the address part of the instruction. A base register is assumed to hold a base address and the address field of the instruction gives a displacement relative to this base address. The base register addressing mode is used in computers to facilitate the relocation of programs in memory. When programs and data are moved from one segment of memory to another, as required in multiprogramming systems, the address values of the base register requires updating to reflect the beginning of a new memory segment.

## SUMMARY

1. The part of the computer that performs the bulk of data-processing operations is called t he central processing unit and is referred to as the CPU.

2. The CPU performs a variety of functions dictated by the type of instructions that are incorporated in the computer.

3. The registers communicate with each other not only for direct data transfers, but also while performing various micro operations. Hence it is necessary to provide a common unit that can perform all the arithmetic, logic, and shift micro operations in the processor.

4. The control unit that operates the CPU bus system directs the information flow through the registers and ALU by selecting the various components in the system.

5. A stack can be placed in portion of a large memory or it can be organized as a collection of a finite number of memory words or registers.

6. Reverse Polish notation, combined with a stack arrangement of registers, is the most efficient way known for evaluating arithmetic expressions.

7. A stack organization is very effective for evaluating arithmetic expressions. The common mathematical method of writing arithmetic expression imposes difficulties when evaluated by a computer.

8. The physical and logical structure of computers is normally described in reference manuals provided with the system.

9. The operation code field of an instruction is a group of bits that define various processor operations, such as add, subtract, complement, and shift.

10. The instruction set of a typical of a typical RISC processor is restricted to the use of load and store instructions when communicating between memory and CPU.

11. The operation field of an instruction specifies the operation to be performed. This operation must be executed on some data stored in computer registers or memory words. The way the operands are chosen during program execution in dependent on the addressing mode of the instruction.

**SELF ASSESSMENT**

1. What is the difference between a microprocessor and a microprogram? Is it possible to design a microprocessor without a microprogram? Are all micro programmed computers also microprocessors?

2. Explain the difference between hardwired control and micro programmed control. Is it possible to have a hardwired control associated with a control memory?

3. Define the following : (a) microoperation; (b) microinstruction; (c) microprogram; (d) microcode.

4. What are different instruction formats we are using?

5. What is addressing mode? How many addressing modes are there? Use particular cases to explain the concept of the addressing modes.

6. The micro programmed control organization showing in Fig. 5.1 has the following propagation delay times. 40 ns to generate the next address, 10 ns to transfer the address into the control address register, 40 ns to access the control memory ROM, 10 ns to transfer the microinstruction into the control data register, and 40 ns to perform the required micro operations specified by the control word. What is the maximum clock frequency that the control can use? What would the clock frequency be if the control data register is not used?

7. The system shown in Fig. 5.2 uses a control memory of 1024 words of 32 bits each. The microinstruction has three fields as shown in the diagram. The micro operations field has 16 bits.

    a. How many bits are there in the branch address field and the select field?

    b. If there are 16 status bits in the system, how many bits of the branch logic are used to select a status bit?

    c. How many bits are left to select an input for the multiplexers?

8. A bus-organized CPU similar to Fig. 5-2 has 16 registers with 32 bits in each, ALU, and a destination decoder.

    a. How many multiplexers are there in the A bus, and what is the size of each multiplexer?

    b. How many selection inputs are needed for MUX A and MUX B?

    c. How many inputs and outputs are there in the decoder?

    d. How many inputs and outputs are there in the ALU for data, including input and output carries?

    e. Formulate a control word for the system assuming that the ALU has 35 operations.

9. The bus system of Fig. 5-2 has the following propagation delay times: 30 ns for the signals to propagate through the multiplexers, 80 ns to perform the ADD operation in the ALU, 20 ns delay in he destination decoder, and 10 ns to clock the data into the destination register. What is the minimum cycle time that can be used for the clock?

10. Specify the control word that must be applied to the processor of Fig. 5-2 to implement the following micro operations.

    a. R1 ← R2 + R3

    b. R4 ← R4

    c. R5 ← R5 − 1

    d. R6 ← sh1 R1

    e. R7 ← input

11. Determine the microoperations that will be executed in the processor of Fig. 5-2 when the following 14-bit control words are applied.

    1. 00101001100101

    2. 00000000000000

    3. 01001001001100

    4. 11110001110000

# CHAPTER –VI

# INPUT-OUTPUT ORGANIZATION

**Author: Dr. Manoj Duhan**                          **Vetter: Mr. Ravinder Rathee**

## 6.1 Peripheral Devices

The input-output subsystem of a computer, referred to as I/O, provides an efficient mode of communication between the central system and the outside environment. Programs and data must be entered into computer memory for processing and results obtained from computations must be recorded or displayed for the user. A computer serves no useful purpose without the ability to receive information from an outside source and to transmit results in a meaningful form.

The most familiar means of entering information into a computer is through a typewriter-like keyboard that allows a person to enter alphanumeric information directly. Every time a key is depressed, the terminal sends a binary coded character to the computer. The fastest possible speed for entering information this way depends on the person's typing speed. On the other hand, the central processing unit is an extremely fast device capable of performing operation this way depends on the person's typing speed. On the other hand, the central processing unit is an extremely fast device capable of performing operations at very high speed. When input information is transferred to the processor via a slow keyboard, the processor will be idle most of the time while waiting for the information to arrive. To use a computer efficiently, a large amount of programs and data must be prepared in advance and transmitted into a storage medium such as magnetic tapes or disks. The information in the disk is then transferred into computer memory at a rapid rate. Results of programs are also transferred into a high-speed storage, such as disks, from which they can be transferred later into a printer to provide a printed output of results.

Devices that are under the direct control of the computer are said to be connected on-line. These devices are designed to read information into or out of the memory unit upon command from the CPU and are considered to be part of the total computer system. Input or output devices attached to the computer are also called peripherals. Among the most common peripherals are keyboards, display units, and printers. Peripherals that provide auxiliary storage for the system are magnetic disks and tapes. Peripherals are electro-mechanical and electromagnetic devices of some complexity. Only a very brief discussion of their function will be given here without going into detail of their internal construction.

Video monitors are the most commonly used peripherals. They consist of a keyboard as the input device and a display unit as the output device. There are different types of video monitors, but the most popular use a cathode ray tube (CRT). The CRT contains an electronic gun that sends an electronic beam to a phosphorescent screen in front of the tube. The beam can be deflected horizontally and vertically. To produce a pattern on the screen, a grid inside the CRT receives a variable voltage that causes the beam to hit the screen and make it glow at selected spots. Horizontal and vertical signals deflect the beam and make it sweep across the tube, causing the visual pattern to appear on the screen. A characteristic feature of display devices is a cursor that marks the position in the screen where the next character will be inserted. The cursor can be moved to any position in the screen, to a single character, the beginning of a word, or to any line. Edit dyes add or delete information based on the cursor position. The display terminal can operate in a single-character mode where all character entered on the screen through the keyboard are transmitted to the computer simultaneously. In the block mode, the edited text is first stored in a logical memory inside the terminal. The text is transferred to the computer as a block of data.

Printers provide a permanent record on paper of computer output data or text. There are three basic types of character printers: daisywheel, dot matrix, and laser printers. The daisywheel printer contains a wheel with the characters placed along the circumference. To print a character, the wheel rotates to the proper position and an energized magnet then presses the letter against the ribbon. The dot matrix printer contains a set of dots along the printing mechanism. For example, a $5 \times 7$ dot matrix printer that prints 80 characters per line has seven horizontal lines, each consisting of $5 \times 7$ dot matrix printer that prints 80 characters per line has seven horizontal lines, each consisting of $5 \times 80 = 400$ dots. Each dot can be printed or not, depending on the specific characters that are printed on the line. The laser printer uses a rotating photographic drum that is used to imprint the character images. The pattern is then transferred onto paper in the same manner as a copying machine.

Magnetic tapes are used mostly for storing files of data: for example, a company's payroll record. Access is sequential and consists of records that can be accessed one after other as the tape moves along a stationary read-write mechanism. It is one of the cheapest and slowest methods for storage and has the advantage that tapes can be recovered when not in use. Magnetic disks have high-speed rotational surfaces coated with magnetic material. Access is achieved by moving a read-write mechanism to a track in the magnetized surface. Disks are used mostly for bulk storage of programs and data. Tapes and disks are discussed further in Sec. 12-1 in conjunction with their role as auxiliary memory.

Other input and output devices encountered in computer systems are digital incremental plotters, optical and magnetic character readers, analog-to-digital converters, and various data acquisition equipment. Not all input comes from people, and not all output is intended for people. Computers are used to control various processes in real time, such as machine tooling, assembly line procedures, and chemical and industrial processes. For such applications, a method must be provided for sensing status condition in the process and sending control signals to the process being controlled.

The input-output organization of a computer is a function of the size of the computer and the devices connected to it. The difference between a small and a large system is mostly dependent on the amount of hardware the computer has available for communicating with peripheral units and the number of peripherals connected to the system. Since each peripheral behaves differently from any other, it would be prohibitive to dwell on the detailed interconnections needed between the computer and each peripheral. Certain techniques common to most peripherals are presented in this chapter.

### 6.1.1 ASCII ALPHANUMERIC CHARACTERS

Input and output devices that communicate with people and the computer are usually involved in the transfer of alphanumeric information to and from the device and the computer is ASCII (American Standard Code for Information Interchange). It uses seven bits to code 128 characters as shown in Table 6-1. The seven bits of the code are designated by $b_1$ through $b_7$ being the most significant bit. The letter A, for example, is represented in ASCII as 1000001 (column 100, row 0001). The ASCII code contains 94 characters that can be printed and 34 nonprinting characters used for various control functions. The printing characters consist of the 26 uppercase letters A through Z, the 26 lowercase letters, the 10 numerals 0 through 9, and 32 special printable characters such as %, *, and $.

The 34 control characters are designated in the ASCII table with abbreviated names. They are listed again below the table with their functional names. The control characters are used for routing data and arranging the printed text into a prescribed format. There are three types of control characters: format effectors, information separators, and communication control

characters. Format effectors are characters that control the layout of printing. They include the familiar typewriter controls, such as backspace (BS), horizontal tabulation (HT), and carriage return (CR). Information separators are used to separate the data into divisions like paragraphs and pages. They include characters such as record separator (RS) and file separator (FS). The communication control characters are useful during the transmission of text between remote terminals. Examples of communication control characters are STX (start of text) and ETX (end of text), which are used to frame a text message when transmitted through a communication medium.

ASCII is a 7 bit code, but most computer manipulate an 8-bit quantity as a single unit called a byte. Therefore, ASCII characters most often are stored one per byte.  Therefore,  ASCII characters most often are stored one per byte.  The extra bit is sometimes used for other purposes, depending on the application.  For example, some printers recognize 8-bit ASCII characters with the most significant bit set to 0.  Additional 128 8-bit characters with the most significant bit set to 1 are used for other symbols, such as the Greek alphabet or italic type font.  When used in data communication, the eighth bit may be employed to indicate the parity of the binary-coded character.

## 6.2 INPUT-OUTPUT INTERFACE

Input-output interface provides a method for transferring information between internal storage and external I/O devices. Peripherals connected to a computer need special communication links for interfacing them with the central processing unit. The purpose of the communication link is to resolve the differences that exist between the central computer and each peripheral. The major differences are:

1.    Peripherals are electromechanical and electromagnetic devices and their manner of operation is different from the operation of the CPU and memory, which are electronic devices. Therefore, a conversion of signal values may be required.

2.    The data transfer rate of peripherals is usually slower than the transfer rate of the CPU, and consequently, a synchronization mechanism may be need.

3.    Data codes and formats in peripherals differ form the word format in the CPU and memory.

4.    The operating modes of peripherals are different from each other and each must be controlled so as not to disturb the operation of other peripherals connected to the CPU.

To resolve these differences, computer systems include special hardware components between the CPU and peripherals to supervise and synchronize all input and output transfers. These components are called interface units because they interface between the processor bus and the peripheral device. In addition, each device may have its own controller that supervises the operations of the particular mechanism in the peripheral.

## 6.2.1    I/O BUS AND INTERFACE MODULES

A typical communication link between the processor and several peripherals is shown in Fig. 6-1 The I/O bus consists of data lines, address lines, and control lines. The magnetic disk, printer, and terminal are employed in practically any general-purpose computer. The magnetic tape is used in some computers for backup storage. Each peripheral device has associated with it an interface unit. Each interface decodes the address and control received from the I/O bus, interprets them for the peripheral, and provides signals for the peripheral controller. It also synchronizes the data flow and supervises the transfer between peripheral and processor. Each peripheral has its own controller that operates the particular electromechanical device. For example, the printer

controller controls the paper motion, the print timing, and the selection of printing characters. A controller may be housed separately or may be physically integrated with the peripheral.

The I/O bus from the processor is attached to all peripheral interfaces. To communicate with a particular device, the processor places a device address on the address lines. Each interface attached to the I/O bus contains an address decoder that monitors the address lines. When the interface detects its own address, it activates the path between the bus lines and the device that it controls. All peripherals whose address does not correspond to the address in the bus are disabled their interface.

At the same time that the address is made available in the address lines, the processor provides a function code in the control lines. The interface



Figure 6-1 Connection of I/O bus to input devices.

selected responds to the function code and proceeds to execute it. The function code is referred to as an I/O command and is in essence an instruction that is executed in the interface and its attached peripheral unit. The interpretation of the command depends on the peripheral that the processor is addressing. There are four types of commands that an interface may receive. They are classified as control, status, status, data output, and data input.

A control command is issued to activate the peripheral and to inform it what to do. For example, a magnetic tape unit may be instructed to backspace the tape by one record, to rewind the tape, or to start the tape moving in the forward direction. The particular control command issued depends on the peripheral, and each peripheral receives its own distinguished sequence of control commands, depending on its mode of operation.

A status command is used to test various status conditions in the interface and the peripheral. For example, the computer may wish to check the status of the peripheral before a transfer is initiated. During the transfer, one or more errors may occur which are detected by the interface. These errors are designated by setting bits in a status register that the processor can read at certain intervals.

A data output command causes the interface to respond by transferring data from the bus into one of its registers. Consider an example with a tape unit. The computer starts the tape moving by issuing a control command. The processor then monitors the status of the tape by means of a status command. When the tape is in the correct position, the processor issues a data output command. The interface responds to the address and command and transfers the information from the data lines in the bus to its buffer register. The interface then communicates with the tape controller and sends the data to be stored on tape.

The data input command is the opposite of the data output. In this case the interface receives an item of data from the peripheral and places it in its buffer register. The processor checks if data are available by means of a status command and then issues a data input command. The interface places the data on the data lines, where they are accepted by the processor.

## 6.2.2    I/O VERSUS MEMORY BUS

In addition to communicating with I/O, the processor must communicate with the memory unit. Like the I/O bus, the memory bus contains data, address, and read/write control lines. There are three ways that computer buses can be used to communicate with memory and I/O:

1.    Use two separate buses, one for memory and the other for I/O.

2.    Use one common bus for both memory and I/O but have separate control lines for each.

3.    Use one common bus for memory and I/O with common control lines.

In the first method, the computer has independent sets of data, address, and control buses, one for accessing memory and the other for I/O. This is done in computers that provide a separate I/O processor (IOP) in addition to the central processing unit (CPU). The memory communicates with both the CPU and the IOP through a memory bus. The IOP communicates also with the input and output devices through a separate I/O bus with its own address, data and control lines. The purpose of the IOP is to provide an independent pathway for the transfer of information between external devices and internal memory.

## 6.2.3    ISOLATED VERSUS MEMORY-MAPPED I/O

Many computers use one common bus to transfer information between memory or I/O and the CPU. The distinction between a memory transfer and I/O transfer is made through separate read and write lines. The CPU specifies whether the address on the address lines is for a memory word or for an interface register by enabling one of two possible read or write lines. The I/O read and I/O write control lines are enabled during an I/O transfer. The memory read and memory write control lines are enabled during a memory transfer. This configuration isolates all I/O interface addresses from the addresses assigned to memory and is referred to as the isolated I/O method for assigning addresses in a common bus.

In the isolated I/O configuration, the CPU has distinct input and output instructions, and each of these instructions is associated with the address of an interface register. When the CPU fetches and decodes the operation code of an input or output instruction, it places the address associated with the instruction into the common address lines. At the same time, it enables the I/O read (for input) or I/O write (for output) control line. This informs the external components that are attached to the common bus that the address in the address lines is for an interface register and not for a memory word.  On the other hand, when the CPU is fetching an instruction or an operand from memory, it places the memory address on the address lines and enables the memory read or memory write  control line. This informs the external components that the address is for a memory word and not for an I/O interface.

The  isolated I/O method isolates memory word and not for an I/O addresses so that memory address values are not affected by interface address assignment since each has its own address space. The other alternative is to use the same address space for both memory and I/O. This is the case in computers that employ only one set of read and write signals and do not distinguish between memory and I/O addresses. This configuration is referred to as memory-mapped I/O. The computer treats an interface register as being part of the memory system. The assigned addresses for interface registers cannot be used for memory words, which reduces the memory address range available.

In a memory-mapped I/O organization there are no specific input or output instructions. The CPU can manipulate I/O data residing in interface registers with the same instructions that are used to manipulate memory words. Each interface is organized as a set of registers that respond to read and write requests in the normal address space. Typically, a segment of the total address space is reserved for interface registers, but in general, they can be located at any address as long as there is not also a memory word that responds to the same address.

Computers with memory-mapped I/O can use memory-type instructions to access I/O data. It allows the computer to use the same instructions for either input-output transfers or for memory transfers. The advantage is that the load and store instructions used for reading and writing from memory can be used to input and output data from I/O registers. In a typical computer, there are more memory-reference instructions than I/O instructions. With memory-mapped I/O all instructions that refer to memory are also available for I/O.

## 6.2.4    EXAMPLE OF I/O INTERFACE

An example of an I/O interface unit is shown in block diagram form in Fig. 6-2. It consists of two data registers called ports, a control register, a status register, bus buffers, and timing and control circuits. The interface communicates with the CPU through the data bus. The chip select and register select inputs determine the address assigned to the interface. The I/O read and write are two control lines that specify an input or output, respectively. The four registers communicate directly with the I/O device attached to the interface.

The I/O data to and from the device can be transferred into either port A or Port B. The interface may operate with an output device or with an input device, or with a device that requires both input and output. If the interface is connected to a printer, it will only output data, and if it services a character reader, it will only input data. A magnetic disk unit transfers data in both directions but not at the same time, so the interface can use bidirectional lines. A command is passed to the I/O device by sending a word to the appropriate interface register. In a system like this, the function code in the I/O bus is not needed because control is sent to the control register, status information is received from the status register, and data are transferred to and from ports A and B registers. Thus the transfer of data, control, and status information is always via the common data bus. The distinction between data, control, or status information is determined from the particular register with which the CPU communicates.

The control register receives control information from the CPU. By loading appropriate bits into the control register, the interface and the I/O device attached to it can be placed in a variety of operating modes. For example, port A may be defined as an input port and port B as an output port. A magnetic tape unit may be instructed to rewind the tape or to start the tape moving in the forward direction. The bits in the status register are used for status conditions and for recording errors that may occur during the data transfer. For example, a status bit may indicate that port A has received a new data item from the I/O device. Another bit in the status register may indicate that a parity error has occurred during the transfer.

The interface registers communicate with the CPU through the bidirectional data bus. The address bus selects the interface unit through the chip select and the two register select inputs. A circuit must be provided externally (usually, a decoder) to detect the address assigned to the interface registers. This circuit enables the chip select (CS) input when the interface is selected by the address bus. The two register select inputs RS1 and RS0 are usually connected to the two least significant lines of the lines address bus. These two inputs select one of the four

| CS | RS1 | RS0 | Register selected |
|----|-----|-----|-------------------|
| 0 | × | × | None: data bus in high-impedance |
| 1 | 0 | 0 | Port A register |
| 1 | 0 | 1 | Port B register |
| 1 | 1 | 0 | Control register |
| 1 | 1 | 1 | Status register |

**Figure 6-2** Example of I/O interface unit.

registers in the interface as specified in the table accompanying the diagram. The content of the selected register is transfer into the CPU via the data bus when the I/O read signal is enabled. The CPU transfers binary information into the selected register via the data bus when the I/O write input is enabled.

## 6.3  ASYNCHRONOUS DATA TRANSFER

The internal operations in a digital system are synchronized by means of clock pulses supplied by a common pulse generator. Clock pulses are applied to all registers within a unit and all data transfers among internal registers occur simultaneously during the occurrence of a clock pulse. Two units, such as a CPU and an I/O interface, are designed independently of each other. If the registers in the interface share a common clock with the CPU registers, the transfer between the two units is said to be synchronous. In most cases, the internal timing in each unit is independent from the other in that each uses its own private clock for internal registers. In that case, the two units are said to be asynchronous to each other. This approach is widely used in most computer systems.

Asynchronous data transfer between two independent units requires that control signals be transmitted between the communicating units to indicate the time at which data is being transmitted. One way of achieving this is by means of a strobe pulse supplied by one of the units to indicate to the other unit when the transfer has to occur. Another method commonly used is to accompany each data item being transferred with a control signal that indicates the presence of

data in the bus. The unit receiving the data item responds with another control signal to acknowledge receipt of the data. This type of agreement between two independent units is referred to as handshaking.

The strobe pulse method and the handshaking method of asynchronous data transfer are not restricted to I/O transfers. In fact, they are used extensively on numerous occasions requiring the transfer of data between two independent units. In the general case we consider the transmitting unit as the source and the receiving unit as the destination. For example, the CPU is the source unit during an output or a write transfer and it is the destination unit during an input or a read transfer. It is customary to specify the asynchronous transfer between two independent units by means of a timing diagram that shows the timing relationship that must exist between the control signals and the data in buses. The sequence of control during an asynchronous transfer depends on whether the transfer is initiated by the source or by the destination unit.

## 6.3.1   STROBE CONTROL

The strobe control method of asynchronous data transfer employs a single control line to time each transfer. The strobe may be activated by either the source or the destination unit. Figure 6-3(a) shows a source-initiated transfer.



(a) Block diagram

(b) Timing diagram

Figure 6-3 Source-initiated strobe for data transfer.

The data bus carries the binary information from source unit to the destination unit. Typically, the bus has multiple lines to transfer an entire byte or word. The strobe is a single line that informs the destination unit when a valid data word is available in the bus.

As shown in the timing diagram of Fig. 6-3(b), the source unit first places the data on the data bus. After a brief delay to ensure that the data settle to a steady value, the source activates the strobe pulse. The information on the data bus and the strobe signal remain in the active state for a sufficient time period to allow the destination unit to receive the data. Often, the destination unit uses the falling edge of the strobe pulse to transfer the contents of the data bus into one of its internal registers. The source removes the data from the bus a brief period after it disables its strobe pulse. Actually, the source does not have to change the information in the data bus. The fact that the strobe signal is disabled indicates that the data bus does not contain valued data. New valid data will be available only after the strobe is enabled again.

Figure 6-4 shows a data transfer initiated by the destination unit. In this case the destination unit activates the strobe pulse, informing the source to provide the data. The source unit responds by placing the requested binary information on the data bus. The data must be valid and remain in the bus long enough for the destination unit to accept it. The falling edge of the

strobe pulse can be used again to trigger a destination register. The destination unit then disables the strobe. The source removes the data from the bus after a predeter-mined time interval.

In many computers the strobe pulse is actually controlled by the clock pulses in the CPU. The CPU is always in control of the buses and informs  the external units how to transfer data. For example, the strobe of Fig. 6-3 could be a memory-write control signal from the CPU to a memory unit. The source, being the CPU, places a word on the data bus and informs the memory units



(a) Block diagram

(b) Timing diagram

Figure 6-4  Destination-initiated strobe for data transfer.

which is the destination, that this is a write operation. Similarly, the strobe of fig. 6-4 could be a memory-read control signal from the CPU to a memory unit. The destination, the CPU, initiates the read operation to inform the memory, which is the source, to place a selected word into the data bus.

The transfer of data between the CPU and an interface unit is similar to the strobe transfer just described. Data transfer between an interface and an I/O device is commonly controlled by a set of handshaking lines.

## 6.3.2   HANDSHAKING

The disadvantage of the strobe method is that the source unit that initiates the transfer has no way of knowing whether the destination unit has actually received the data item that was placed in the bus. Similarly, a destination unit that initiates the transfer has no way of knowing whether the source unit has actually placed the data on the bus,. The handshake method solves this problem by introducing a second control signal that provides a reply to the unit that initiates the transfer. The basic principle of the two-write handshaking method of data transfer is as follows. One control line is in the same direction as the data flow in the bus from the source to the destination. It is used by the source unit to inform the destination unit whether there are valued data in the bus. The other control line is in the other direction from the destination to the source. It is used by the destination unit to inform the source whether it can accept data. The sequence of control during the transfer depends on the unit that initiates the transfer.

Figure 6-5 shows the data transfer procedure when initiated by the source. The two handshaking lines are data valid, which is generated by the source unit, and data accepted, generated by the destination unit. The timing diagram shows the exchange of signals between the two units. The sequence of events listed in part (c) shows the four possible states that the system
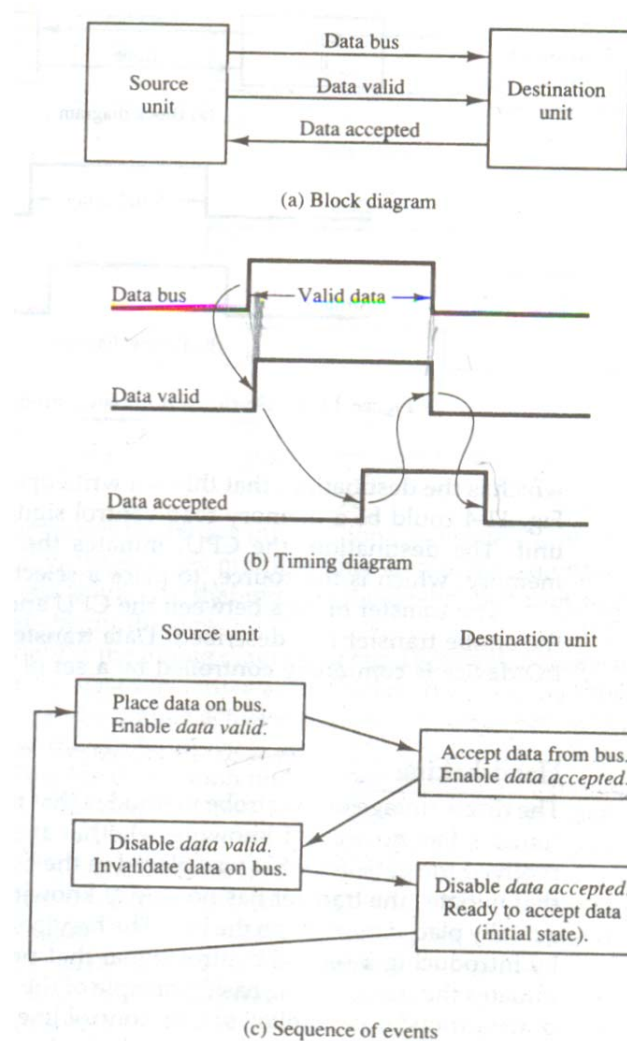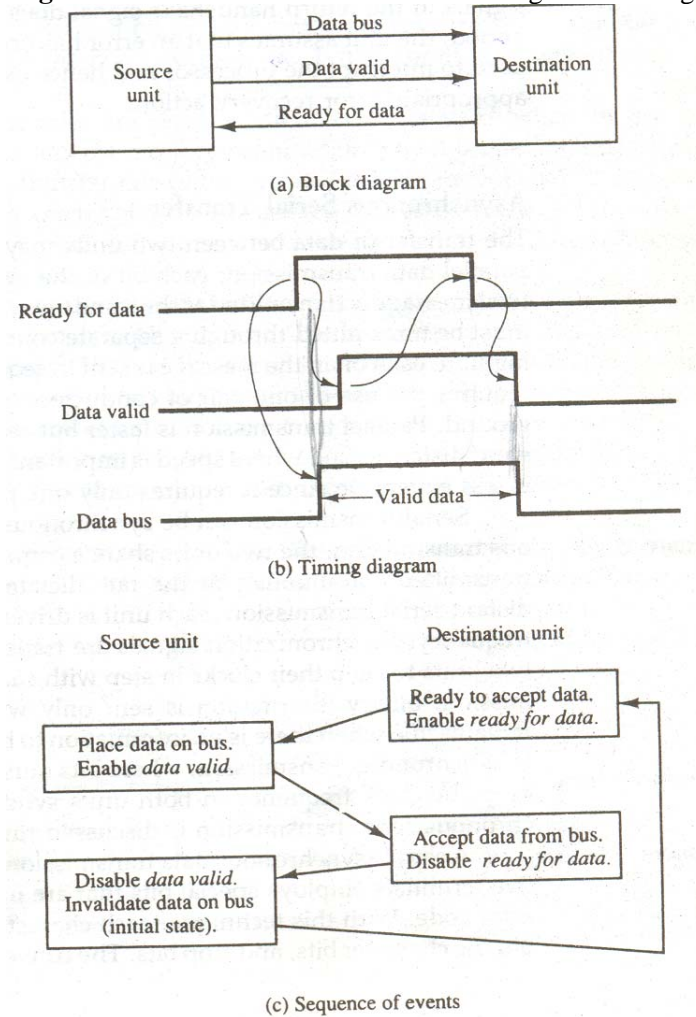
(a) Block diagram

(b) Timing diagram

(c) Sequence of events

**Figure 6-5** Source-initiated transfer using handshaking.

can be at any given time. The source unit initiates the transfer by placing the data on the bus and enabling its data valid signal. The data accepted signal is activated by the destination unit after it accepts the data from the bus. The source unit then disables its data valid signal, which invalidates the data on the bus. The destination unit then disables its data accepted signal and the system gose into its initial state. The source dies not send the next data item until after the destination unit shows its readiness to accept new data by disabling its data accepted signal. This scheme allows arbitrary delays from one state to the next and permits each unit to respond at its own data transfer rate. The rate of transfer is determined by the slowest unit.

The destination-initiated transfer using handshaking lines is shown in Fig. 6-6. Note that the name of the signal generated by the destination unit has been changed to ready fro data to reflect its new meaning. The source unit in this case dies not place data on the bus until after it receives the ready for data signal from the destination unit. From there on, the handshaking procedure follows the same pattern as in the source-initiated case. Note that the sequence of

**Figure 6-6** Destination-initiated transfer using handshaking.



(a) Block diagram

(b) Timing diagram

(c) Sequence of events

events in both cases would be identical if we consider the ready for data signal as the complement of data accepted. In fact, the only difference between the source-initiated and the destination-initiated transfer is in their choice of initial state.

The handshaking scheme provides a high degree of flexibility and reality because the successful completion of a data transfer relies on active participation by both units. If one unit is faulty, the data transfer will not be completed. Such an error can be detected by means of a timeout mechanism, which produces an alarm if the data transfer is not completed within a predetermined time. The timeout is implemented by means of an internal clock that starts counting time when the unit enables one of its handshaking control signals. If the return handshake signal does not respond within a given time period, the unit assumes that an error has

occurred. The timeout signal can be used to interrupt the processor and hence execute a service routine that takes appropriates error recovery action.

### 6.3.3 ASYNCHRONOUS SERIAL TRANSFER

The transfer of data between two units may be done in parallel or serial. In parallel data transmission, each bit of the message has its own path and the total message is transmitted at the same time. This means that an n-bit message must be transmitted through in seqatate conductor paths. In serial data transmission, each bit in the message is sent in sequence one at a time. This method requires the use of one pair of conductors or one conductor and a common ground. Paralled transmission is faster but requires many wires. It is used for short distances and where speed is important. Serial transmission is slower but is less expensive since it requires only one pair of conductors.

Serial transmission can be can be synchronous or asynchronous. In synchronous transmission, the two units share a common clock frequency and bits are transmitted continuously at the rate dictated by the clock pulses. In long-distant serial transmission, each unit is driven by a separate clock of the same frequency. Synchronization signals are transmitted periodically between the two units to keep their clocks in step with each other. In asynchronous transmission, binary information is sent only when it is available and the line remains idle when there is no information to be transmitted. This is in contrast to synchronous transmission, where bits must be transmitted continuously to deep the clock frequency in both units synchronized with each other.

Serial asynchronous data transmission technique used in many interactive terminals employs special bits that are inserted at both ends of the character code. With this technique, each character consists of three parts: a start bit, the character bits, and stop bits. The convention is that the transmitter rests at the 1-state when no characters are transmitted. The first bit, called the start bit, is always a 0 and is used to indicate the beginning of a character. The last bit called the stop bit is always a 1. An example of this format is shown in Fig. 6-7.

A transmitted character can be detected by the receiver from knowledge of the transmission rules:

1.      When a character is not being sent, the line is kept in the 1-state.

2.      The initiation of a character transmission is detected from the start bit, which is always 0.

3.      The character bits always follow the start bit.

4.      After the last bit of the character is transmitted, a stop bit is detected when the line returns to the 1-state for at least one bit time.

Using these rules, the receiver can detect the start bit when the line gives from 1 to 0. A clock in the receiver examines the line at proper bit times. The receiver knows the transfer rate of the bits and the number of character bits to accept. After the character bits are transmitted, one or two stop bits are sent. The stop bits are always in the 1-state and frame the end of the character to signify the idle or wait state.

At the end of the character the line is held at the 1-state for a period of at least one or two bit times so that both the transmitter and receiver can resynchronize. The length of time that the line stays in this state depends on the amount of time required for the equipment to resynchronize. Some older electromechanical terminals use two stop bits, but newer terminals use one stop bit. The line remains in the 1-state until another character is transmitted. The stop time ensures that a new character will not follow for one or two bit times.

As illustration, consider the serial transmission of a terminal whose transfer rate is 10 characters per second. Each transmitted character consists of a start bit, eight information bits, of
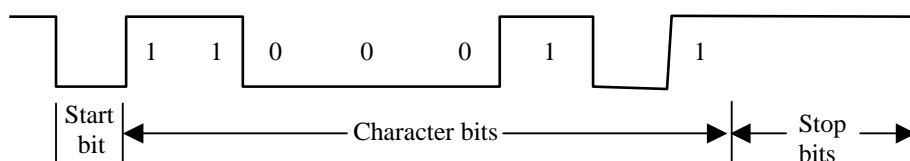
]Figure 6-7  Asynchronous serial transmission

 a start bit, eight information bits, and two stop bits, for a total of 11 bits. Ten characters per second means that each character takes   0. 1s for transfer. Since there are 11 bits to be transmitted, it follows that the bit time is 9.09 ms. The baud rate is defined as the rate at which serial information is transmitted and is equivalent to the data transfer in bits per second. Ten characters per second with an 11-bit format has a transfer rate of 110 baud.

The terminal has a keyboard and a printer. Every time a key is depressed, the terminal sends 11 bits serially along a wire. To print a character. To print a character in the printer, an 11-bit message must be received along another wire. The terminal interface consists of a transmitter and a receiver. The transmitter accepts an 8-bit character room the computer and proceeds to send a serial 11-bit message into the printer line. The receiver accepts a serial 11-bit message from the keyboard line and forwards the 8-bit character code into the computer. Integrated circuits are available which are specifically designed to provide the interface between computer and similar interactive terminals. Such a circuit is called an asynchronous communication interface or a universal asynchronous receiver-transmitter (UART).

## 6.4 MODES OF TRANSFER

 Binary information received from an external device is usually stored in memory for later processing. Information transferred from the central computer into an external device originates in the memory unit. The CPU merely executes the I/O instructions and may accept the data temporarily, but the ultimate source or destination is the memory   unit. Data transfer between the central computer and I/O devices may be handled in a variety of modes. Some modes use the CPU as an intermediate path; other transfer the data directly to and from the memory unit. Data transfer to and from peripherals may be handled in one of three possible modes:

1.      Programmed I/O

2.      Interrupt-initiated I/O

3.      Direct memory access (DMA)

Programmed I/O operations are the result of I/O instructions written in the computer program. Each data item transfer is initiated by an instruction in the program. Usually, the transfer is to and from a CPU register and peripheral. Other instructions are needed to transfer the data to and from CPU and memory. Transferring data under program control requires constant monitoring of the peripheral by the CPU. Once a data transfer is initiated, the CPU is required to monitor the interface to see when a transfer can again be made. It is up to the programmed instructions executed in the CPU to keep close tabs on everything that is taking place in the interface unit and the I/O device.

In the programmed I/O method, the CPU stays in a program loop until the I/O unit indicates that it is ready for data transfer. This is a time-consuming process since it keeps the processor busy needlessly. It can be avoided by using an interrupt facility and special commands to inform the interface to issue an interrupt request signal when the data are available from the

device. In the meantime the CU can proceed to execute another program. The interface meanwhile keeps monitoring the device. When the interface determines that the device is ready for data transfer, it generates an interrupt request to the computer. Upon detecting the external interrupt signal, the CPU momentarily stops the task it is processing, branches to a service program to process the I/O transfer, and then returns to the task it was originally performing

Transfer of data under programmed I/O is between CPU and peripheral. In direct memory access (DMA), the interface transfers data into and out of the memory unit through the memory bus. The CPU initiates the transfer by supplying the interface with the starting address and the number of words needed to be transferred and then proceeds to execute other tasks. When the transfer is made, the DMA requests memory cycles through the memory bus. When the request is granted by the memory controller, the DMA transfers the data directly into memory. The CPU merely delays its memory access operation to allow the direct memory I/O transfer. Since peripheral speed is usually slower than processor speed, I/O-memory transfers are infrequent compared to processor access to memory.

Many computers combine the interface logic with the requirements for direct memory access into one unit and call it an I/O processor (IOP). The IOP can handle many peripherals through a DMPA and interrupt facility. In such a system, the computer is divided into three separate modules: the memory unit, the CPU, and the IOP.

## 6.4.1   EXAMPLE OF PROGRAMMED I/O

In the programmed I/O method, the I/O device dies not have direct access to memory. A transfer from an I/O device to memory requires the execution of several instructions by the CPU, including an input instruction to transfer the data from the device to the CPU, and a store instruction to transfer the data from the CPU to memory. Other instructions may be needed to verify that the data are available from the device and to count the numbers of words transferred.

An example of data transfer from an I/O device through an interface into the CPU is shown in Fig. 6.8. The device transfers bytes of data one at a time as they are available. When a byte of data is available, the device places it in the I/O bus and enables its data valid line. The interface accepts the byte into its data register and enables the data accepted line. The interface sets a it in the status register that we will refer to as an F  or "flag" bit. The device can now disable the data valid line, but it will not transfer another byte until the data accepted line is disabled by the interface. This is according to the handshaking procedure established in Fig. 6-5.

A program is written for the computer to check the flag in the status register to determine if a byte has been placed in the data register by the I/O device. This is done by reading the status register into a CPU register and checking the value of the flag bit. If the flag is equal to 1, the CPU reads the data from the data register. The flag bit is then cleared to 0 by either the CPU or the interface, depending on how the interface circuits are designed. Once the flag is cleared, the interface disables the data accepted line and the device can then transfer the next data  byte.

A flowchart of the program that must be written for the CPU is shown in Fig. 6-9. It is assumed that the device is sending a sequence of bytes that must be stored in memory. The transfer of each byte requires three instructions:

1.   Read the status register.

2.   Check the status of the flag bit and branch to step 1 if not set or to step 3 if set.

3.   Read the data register.

Each byte is read into a CPU register and then transferred to memory with a store instruction. A common I/O programming task is to transfer a block of words form an I/O device and store them in a memory buffer. A program that stores input characters in a memory buffer using the

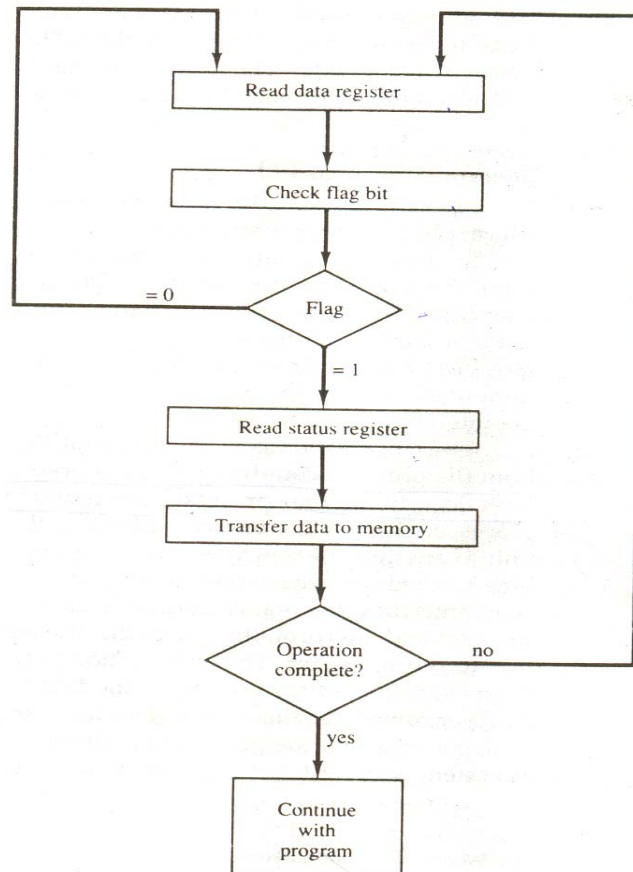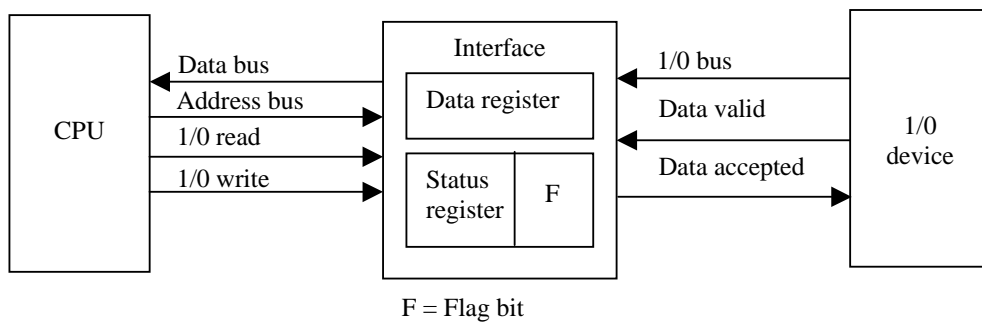**Figure 6-8** Data transfer form I/O device to CPU



F = Flag bit



**Figure 6-9** Flowchart for CPU program to input data.

instructions mentioned in the earlier chapter.

The programmed I/O method is particularly useful in small low-speed computers or in systems that are dedicated to monitor a device continuously. The difference in information

transfer rate between the CPU and the I/O device makes this type of transfer inefficient. To see why this is inefficient, consider a typical computer that can execute the two instructions that read the status register and check the flag in 1 πs. Assume that the input device transfers its data at an average rate of 100 bytes per second. This is equivalent to one byte every 10,000 πs. This means that the CPU will check the flag 10,000 times between each transfer. The CPU is wasting time while checking the flag instead of doing some other useful processing task.

### 6.4.2   INTERRUPT-INITIATED I/O

An alternative to the CPU constantly monitoring the flag is to let the interface inform the computer when it is ready to transfer data. This mode of transfer uses the interrupt facility. While the CPU is running a program, it does not check the flag. However, when the flag is set, the computer is momentarily interrupted from proceeding with the current program and is informed of the fact that the flag has been set. The CPU deviates from what it is doing to take care of the input or output transfer. After the transfer is completed, the computer returns to the previous program to continue what it was doing before the interrupt.

The CPU responds to the interrupt signal by storing the return address from the program counter into a memory stack and then control branches to a service routine that processes the required I/O transfer. The way that the processor chooses the branch address of the service routine varies from tone unit to another. In principle, there are two methods for accomplishing this. One is called vectored interrupt and the other, no vectored   interrupt. In a non vectored interrupt, the branch address is assigned to a fixed location in memory. In a vectored interrupt, the source that interrupts supplies the branch information to the computer. This information is called the interrupt vector. In some computers the interrupt vector is the first address of the I/O service routine. In other computers the interrupt vector is an address that points to a location in memory where the beginning address of the I/O service routine is stored.

### 6.4.3   SOFTWARE CONSIDERATIONS

The previous discussion was concerned with the basic hardware needed to interface I/O devices to a computer system. A computer must also have software routines for controlling peripherals and for transfer of data between the processor and peripherals. I/O routines must issue control commands to activate the peripheral and to check the device status to determine when it is ready for data transfer. Once ready, information is transferred item by item until all the data are transferred. In some cases, a control command is then given to execute a device function such as stop tape or print characters. Error checking and other useful steps often accompany the transfers. In interrupt-controlled transfers, the I/O software must  issue commands to the peripheral to interrupt when ready and to service the interrupt when it occurs. In DMA transfer, the I/O software must initiate the DMA channel to start its operation.

Software control of input-output equipment is a complex undertaking. For this reason I/O routines for standard peripherals are provided by the manufacturer as part of the computer system. They are usually included within the operating system. Most operating systems are supplied with a variety of I/O programs to support the particular line of peripherals offered for the computer. I/O routines are usually available as operating system procedures and the user refers to the established routines to specify the type of transfer required without going into detailed machine language programs.

### 6.5   PRIORITY INTERRUPT

Data transfer between the CPU and an I/O device is initiated by the CPU. However, the CPU cannot start the transfer unless the device is ready to communicate with the CPU. The readiness of the device can be determined from an interrupt signal. The CPU responds to the interrupt

request by storing the return address from PC into a memory stack and then the program branches to a service routine that processes the required transfer. Some processors also push the current PSW for the service routine. We neglect the PSW here in order not to complicate the discussion of I/O interrupts.

In a typical application a number of I/O devices are attached to the computer, with each device being able to originate an interrupt request. The first task of the interrupt system is to identify the source of the interrupt. There is also the possibility that several sources will request service simultaneously. In this case the system must also decide which device to service first.

A priority interrupts is a system that establishes a priority over the various sources to determine which condition is to be serviced first when two or more request arrive simultaneously. The system may also determine which conditions are permitted to interrupt the computer while another interrupt is being serviced. Higher-priority interrupt levels are assigned to request which, if delayed of interrupted, could have serious consequences. Devices with high-speed transfers such as keyboards receive low priority. When two devices interrupt the computer at the same time, the computer services the devices interrupt the computer at the same time, the computer services the device, with the higher priority first.

Establishing the priority of simultaneous interrupts can be done by software or hardware. A polling procedure is used to identify the highest-priority source by software means. In this method there is one common branch address for all interrupts. The program that takes care of interrupts begins at the branch address and polls the interrupt sources in sequence. The order in which they are tested determines the priority of each interrupt. The highest-priority source is tested first, and if its interrupt signal is on, control branches to a service routine for this source. Otherwise, the next-lower-priority source is tested, and so on. Thus the initial service routine for all interrupt consists of a program that tests the interrupt sources in sequence and branches to one of many possible service routines. The particular service routine reached belongs to the highest-priority device among all devices that interrupted the computer. The disadvantage of the soft ware method is that if there are many interrupts, the time required to poll them can exceed the time available to service the I/O device. In this situation a hardware priority-interrupt unit can be used to speed up the operation.

A hardware priority-interrupt unit functions as an overall manager in an interrupt system environment. It accepts interrupt requests from many sources, determines which of the incoming requests has the highest priority, and issues an interrupt request to the computer based on this determination. To speed up the operation, each interrupt source has its own interrupt vector to access its own service routine directly. Thus no polling is required because all the decisions are established by the hardware priority-interrupt unit. The hardware priority function can be established by either a serial or a parallel connection of interrupt lines. The serial connection is also known as the daisy chaining method.

### 6.5.1   DAISY-CHAINING PRIORITY

The daisy-chaining method of establishing priority consists of a serial connection of all devices that request an interrupt. The device with the highest priority is placed in the first position, followed by lower-priority devices up to the device with the lowest priority, which is placed last in the chain. This method of connection between three devices and the CPU is shown in Fig. 6-10 The interrupt request line is common to all devices and forms a wired logic connection. If any device has its interrupt signal in the low-level state, the interrupt line goes to the low-level state and enables the interrupt input in the CPU. When no interrupts are pending, the interrupt line stays in the high-level state and no interrupts are recognized by the CPU. This is equivalent to a negative logic OR operation. The CPU responds to an interrupt request by enabling the interrupt

acknowledge line. This signal is received by device 1 at its PI (priority in) input. The acknowledge signal passes on to the next device through the PO (priority out) output only if device 1 is not requesting an interrupt. If device 1 has a pending interrupt, it blocks the acknowledge signal from the next device by placing a 0 in the PO output. It then proceeds to insert its own interrupt vector address (VAD) into the data bus for the CPU to use during the interrupt cycle.

A device with a 0 in its PI input generates a 0 in its PO output to inform the next-lower-priority device that the acknowledge signal has been blocked. A device that is requesting an interrupt and has a 1 in its PI input will intercept the acknowledge signal by placing a 0 in its PO output. If the device does not have pending interrupts, it transmits the acknowledge signal to the next device



**Figure 6-10** Daisy-chain priority interrupt

by placing a 1 in its PO output. Thus the device with PI = 1 and PO = 0 is the one with the highest priority that is requesting an interrupt, and this device places its VAD on the data bus. The daisy chain arrangement gives the highest priority to the device that receives the interrupt acknowledge signal from the CPU. The farther the device is from the first position, the lower is its priority.

Figure 6.11 shows the internal logic that must be included with in each device when connected in the daisy-chaining scheme. The device sets its RF flip-flop when it wants to interrupt the CPU. The output of the RF flip-flop goes through an open-collector inverter, a circuit that provides the wired logic for the common interrupt line. If PI = 0, both PO and the enable line to VAD are equal to 0, irrespective of the value of RF. If PI = 1 and RF = 0, then PO = 1 and the vector address is disabled. This condition passes the acknowledge signal to the next device through PO. The device is active when PI = 1 and RF = 1. This condition places a 0 in PO and enables the vector address for the data bus. It is assumed that each device has its own distinct vector address. The RF flip-flop is reset after a sufficient delay to ensure that the CPU has received the vector address.

## 6.5.2  PARALLEL PRIORITY INTERRUPT

The parallel priority interrupt method uses a register whose bits are set separately by the interrupt signal from each device. Priority is established according to the position of the bits in the register. In addition to the interrupt register the circuit may include a mask register whose

purpose is to control the status of each interrupt request. The mask register can be programmed to disable
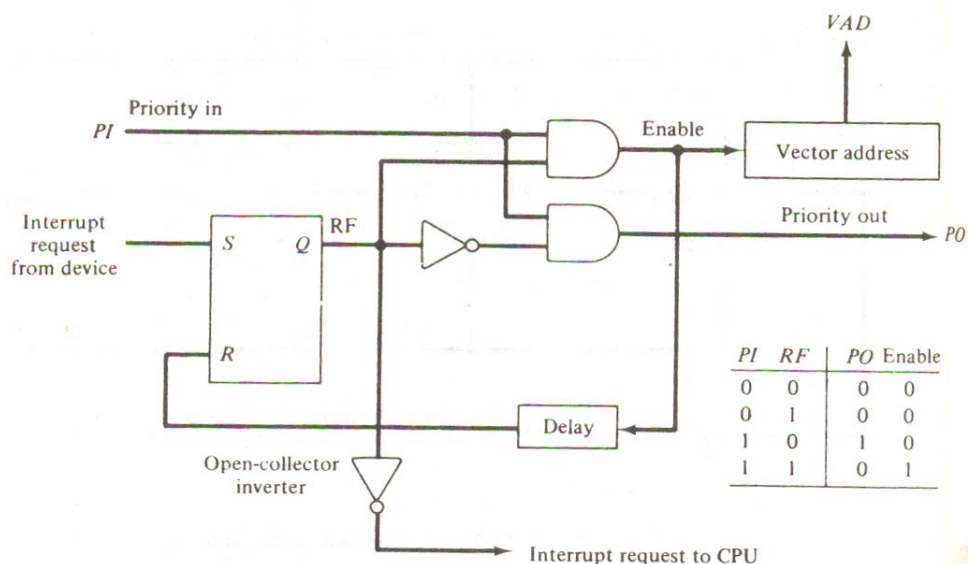


**Figure 6-11** One state of the daisy-chain priority arrangement.

lower-priority interrupts while a higher-priority device is being serviced. It can also provide a facility that allows a high-priority device to interrupt the CPU while a lower-priority device is being serviced.

The priority logic for a system of four interrupt sources is shown in Fig. 6.12. It consists of an interrupt register whose individual bits are set by  external conditions and cleared by program instructions. The magnetic disk, being a high-speed device, is given the highest priority. The printer has the next priority, followed by a character reader and a keyboard. The mask register has the same number of bits as the interrupt register. By means of program instructions, it is possible to set or reset any bit in the mask register. Each interrupt bit and its corresponding mask bit are applied to an AND gate to produce the four inputs to a priority encoder. In this way an interrupt is recognized only if its corresponding mask bit is set to 1 by the program. The priority encoder generates two bits of the vector address, which is transferred to the CPU.

Another output from the encoder sets an interrupt status flip-flop IST when an interrupt that is not masked occurs. The interrupt enable flip-flop IEN can be set or cleared by the program to provide an overall control over the interrupt system. The outputs of IST ANDed with IEN provide a common interrupt signal for the CPU. The interrupt acknowledge INTACK signal from the CPU enables the bus buffers in the output register and a vector address VAD is placed into the data bus. We will now explain the priority encoder circuit and then discuss the interaction between the priority interrupt controller and the CPU.
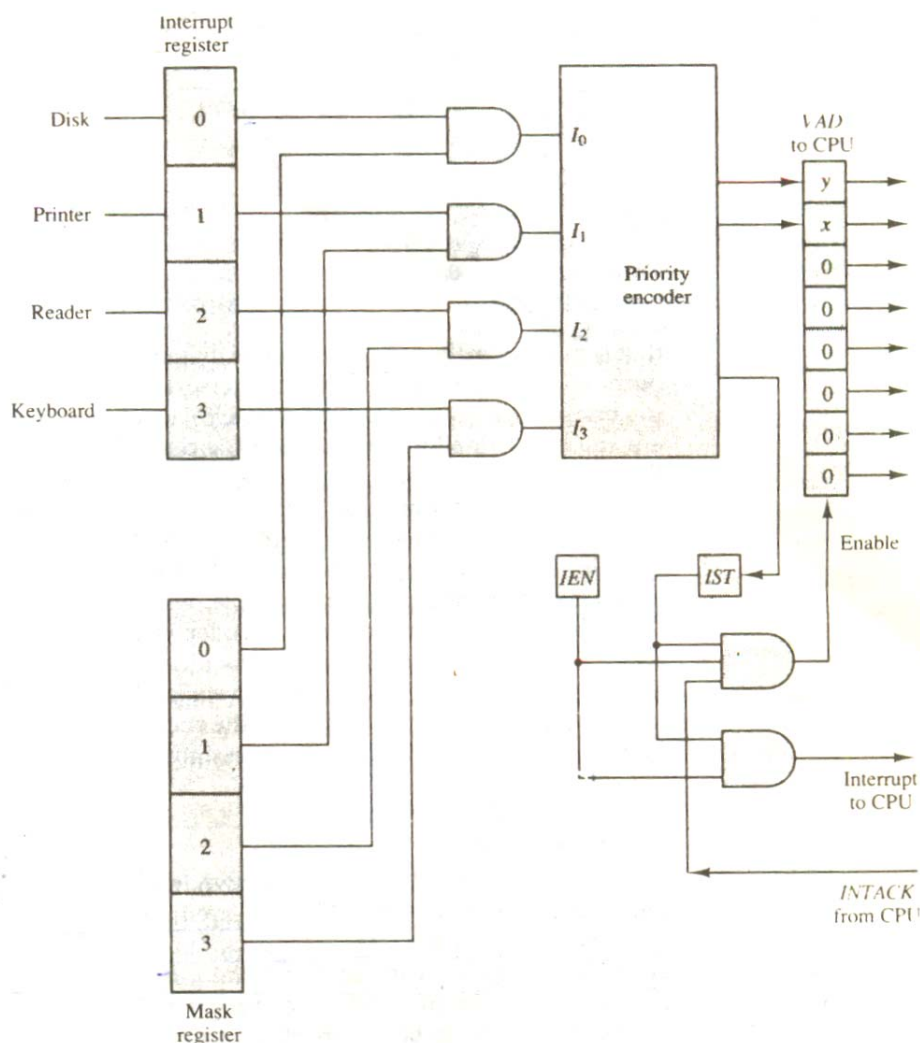
**Figure 6.12** Priority interrupt hardware.

## 6.6 DIRECT MEMORY ACCESS (DMA)

The transfer of data between a fast storage device such as magnetic disk and memory is often limited by the speed of the CPU. Removing the CPU from the path and letting the peripheral device manage the memory buses directly would improve the speed of transfer. This transfer technique is called direct memory access (DMA). During DMA transfer, the CPU is idle and has no control of the memory buses. A DMA controller takes over the buses to manage the transfer directly between the I/O device and memory.

The CPU may be placed in an idle state in a variety of ways. One common method extensively used in microprocessors is to disable the buses through special control signals. Figure 6.13 shows two control signals in the CPU that facilitate the DMA transfer. The bus request (BR) input is used by the DMA controller to request the CPU to relinquish control of the buses. When this input is active, the CPU terminates the execution of the current instruction and places the address bus, the data bus, and the read and write lines into a high-impedance state behaves like an open circuit, which means that the output is disconnected and dies not have a logic significance. The CPU activates the Bus grant (BG) output to inform the external DMA that the buses are in

the high-impedance state. The DMA that originated the bus request can now take control of the buses to conduct memory transfers without processor intervention. When the DMA terminates the transfer, it disables the bus request line. The CPU disables the bus grant, takes control of the buses, and returns to its normal operation.

When the DMA takes control of the bus system, it communicates directly with the memory. The transfer can be made in several ways. In DMA burst transfer, a block sequence consisting of a number of memory words is transferred in a continuous burst while the DMA controller is master of the memory buses. This mode of transfer is needed for fast devices such as magnetic disks, where data transmission cannot be stopped or slowed down until an entire block is transferred. An alternative technique called cycle stealing allows the DMA controller to transfer one data word at a time after which it must return control of the buses to the CPU. The CPU merely delays its operation for one memory cycle to allow the direct memory I/O transfer to "steal" one memory cycle.

### 6.6.1   DMA CONTROLLER

The DMA controller needs the usual circuits of an interface to communicate with the CPU and I/O device. In addition, it needs an address register, a word count register, and a set of address lines. The address register and address lines

**Figure 6.13** CPU bus signals for DMA transfer.



are used for direct communication with the memory. The word count register specifies the number of words that must be transferred. The data transfer may be done directly between the device and memory under control of the DMA.

Figure 6.14 shows the block diagram of a typical DMA controller. The unit communicates with the CPU via the data bus and control lines. The registers in the DMA are selected by the CPU through the address bus by enabling the  DS (DMA select) and RS (register select) inputs. The RD (read) and WR (write) inputs are bidirectional. When the BG (bus grant) input is 0, the CPU can communicate with the DMA registers through the data bus to read from or write to the DMA registers. When BG = 1, the CPU has relinquished the buses and the DMA can communicate directly with the memory by specifying an address in the address bus and activating the RD or WR control. ;the DMA communicates with the external peripheral through the request and acknowledge lines by using a prescribed handshaking procedure.

The DMA controller has three registers: an address register, a word count register, and a control register. The address register contains an address to specify the desired location in memory. The address bits go through bus buffers into the address bus. The address register is incremented after each word that is transferred to memory. The word count register is incremented after each word that is transferred to memory. The word count register holds the number of words to be transferred. This register is decremented by one after each word transfer and internally tested for zero. The control register specifies the mode of transfer. All registers in the DMA appear to the CPU as I/O interface registers. Thus the CPU can read from or write into the DMA registers under program control via the data bus.

Figure 6.14  Block diagram of DMA controller.



The DMA is first initialized by the CPU.  After that, the DMA starts and continues to transfer data between memory and peripheral unit until an entire block is transferred. The initialization process is essentially a program consisting of I/O instructions that include the address for selecting particular DMA registers. The CPU initializes the DMA by sending the following information through the data bus:

1.          The starting address of the memory block where data are available (for read) or where data are to be stored (for write)

2.          The word count, which is the number of words in the memory block

3.          Control to specify the mode of transfer such as read or write

4.          A control to start the DMA transfer

The starting address is stored in the address register. The word count is stored in the word count register, and the control information in the control register. Once the DMA is initialized, the CPU stops communicating with the DMA unless it receives an interrupt signal or if it wants to check how many words have been transferred.

### 6.6.2    DMA Transfer

The position of the DMA controller among the other components in a computer system is illustrated in Fig. 6.15. The CPU communicates with the DMA through the address and data buses as with any interface unit. The DMA has its own address, which activates the DS and RS lines. The CPU initializes the DMA through the data bus. Once the DMA receives the start control command, it can start the transfer between the peripheral device and the memory.

When the peripheral device sends a DMA request, the DMA controller activates the BR line, informing the CPU to relinquish the buses. The CPU responds with its BG line, informing the DMA that its buses are disabled. The DMA then puts the current value of its address register into the address bus, initiates the RD or WR signal, and sends a DMA acknowledge to the peripheral device. Note that the RD and WR lines in the DMA controller are bidirectional. The

direction of transfer depends on the status of the BG line. When BG line. When BG = 0, the RD and WR are input lines allowing the CPU to communicate with the internal DMA registers. When BG = 1, the RD and WR and output lines from the DMA controller to the random-access memory to specify the read or write operation for the data.

When the peripheral device receives a DMA acknowledge, it puts a word in the data us (for write) or receives a word from the data bus (for read). Thus the DMA controls the read or write operations and supplies the address for the memory. The peripheral unit can then communicate with memory through the data bus for direct transfer between the two units while the CPU is momentarily disabled.



**Figure 6.15** DMA transfer in a computer system.

For each word that is transferred, the DMA increments its address register and decrements its word count register. If the word count does not reach zero, the DMA checks the request line coming from the peripheral. For a high-speed device, the line will be active as soon as the previous transfer is completed. A second transfer is then initiated, and the process continues until the entire block is transferred. If the peripheral speed is slower, the DMA request line may come somewhat later. In this case the DMA disables the bus request line so that the CPU can continue

to execute its program. When the peripheral requests a transfer, the DMA requests the buses again.

It the word count register reaches zero, the DMA stops any further transfer and removes its bus request. It also informs the CPU of the termination by means of an interrupt. When the CPU responds to the interrupt, it reads the content of the word count register. The zero value of this register indicates that all the words were transferred successfully. The CPU can read this register at any time to check the number of words already transferred.

A DMA controller may have more than on channel. In this case, each channel has a request and acknowledges pair of control signals which are connected to separate peripheral devices. Each channel also has its own address register and word count register within the DMA controller. A priority among the channels may be established so that channels with high priority are serviced before channels with lower priority.

DMA transfer is very useful in many applications. It is used for fast transfer of information between magnetic disks and memory. It is also useful for updating the display in an interactive terminal. Typically, an image of the screen display of the terminal is kept in memory which can be updated under program control. The contents of the memory can be transferred to the screen periodically by means of DMA transfer.

## SUMMARY

1.  The input-output subsystem of a computer, referred to as I/O, provides an efficient mode of communication between the central system and the outside environment.

2.  Devices that are under the direct control of the computer are said to be connected on-line.

3.  Input and output devices that communicate with people and the computer are usually involved in the transfer of alphanumeric information to and from the device and the computer is ASCII

4.  Input-output interface provides a method for transferring information between internal storage and external I/O devices. Peripherals connected to a computer need special communication links for interfacing them with the central processing unit.

5.  The internal operations in a digital system are synchronized by means of clock pulses supplied by a common pulse generator.

6.  The transfer of data between two units may be done in parallel or serial. In parallel data transmission, each bit of the message has its own path and the total message is transmitted at the same time.

7.  Binary information received from an external device is usually stored in memory for later processing. Information transferred from the central computer into an external device originates in the memory unit.

8.  An alternative to the CPU constantly monitoring the flag is to let the interface inform the computer when it is ready to transfer data.

9.  Data transfer between the CPU and an I/O device is initiated by the CPU. However, the CPU cannot start the transfer unless the device is ready to communicate with the CPU. The readiness of the device can be determined from an interrupt signal.

10. The DMA controller needs the usual circuits of an interface to communicate with the CPU and I/O device. In addition, it needs an address register, a word count register, and a set of address lines.

**11.** The transfer of data between a fast storage device such as magnetic disk and memory is often limited by the speed of the CPU. Removing the CPU from the path and letting the peripheral device manage the memory buses directly would improve the speed of transfer.

**12.** A DMA controller may have more than on channel. In this case, each channel has a request and acknowledges pair of control signals which are connected to separate peripheral devices.

## SELF ASSESSMENT

1. Enlist and explain that how many types of I/O devices can be connected to the Computer**?**

2. What is I/O bus interface module? Explain any one in detail?

3. Compare and contrast I/O and memory bus.

4. What is asynchronous data transfer? What are the modes with the help of which we can transfer the data to the said destination?

5. What are the different modes of transfer?

6. What is a priority interrupt? What are the different types of the priorities in DMA?

7. What is a DMA ? Draw its block and IC diagram and also explain its working.

8. How DMA transfer takes place? Comment.

# CHAPTER-VII

# MEMORY ORGANIZATION

**Author: Dr. Manoj Duhan.**                                    **Vetter: Mr. Ravinder Rathee.**

## 7-1   MEMORY HIERARCHY

The memory unit is an essential component in any digital computer since it is needed for storing programs and data. A very small computer with a limited application may be able to fulfill its intended task without the need of additional storage capacity. Most general-purpose computers would run more efficiently if they were equipped with additional storage beyond the capacity of the main memory. There is just not enough space in one memory unit to accommodate all the programs used in a typical computer. Moreover, most computer users accumulate and continue to accumulate large amounts of data-processing software. Not all accumulated information is needed by the processor at the same time. Therefore, it is more economical to use low-cost storage devices to serve as a backup for storing the information that is not currently used by the CPU. The memory unit that communicates directly with the CPU is called the main memory. Devices that provide backup storage are called auxiliary memory. The most common auxiliary memory devices used in computer systems are magnetic disks and tapes. They are used for storing system programs, large data files, and other backup information. Only  programs and data currently needed by the processor reside in main memory. All other information is stored in auxiliary memory and transferred to main memory when needed.

   The total memory capacity of a computer can be visualized as being a hierarchy of components. The memory hierarchy system consists of tall storage devices employed in a computer system from the slow but high-capacity auxiliary memory to a relatively faster main memory, to an even smaller and faster cache memory accessible to the high-speed processing logic. Figure 7-1 illustrates the components in a typical memory hierarchy. At the bottom of the hierarchy are the relatively slow magnetic tapes used to store removable files. Next are the magnetic disks used as backup storage. The main memory occupies a central position by being able to communicate directly with the CPU and with auxiliary memory devices through an I/O processor. When programs not residing in main memory are needed by the CPU, they are brought in from auxiliary memory. Programs not currently needed in main memory are transferred into auxiliary memory to provide space for currently used programs and data.

   A special very-high speed memory called a cache is sometimes used to increase the speed of processing by making current programs and data available to the CPU at a rapid rate. The cache memory is employed in computer systems to compensate for the speed differential between main memory access time and processor logic. CPU logic is usually faster than main memory access time, with the result that processing speed is limited primarily by the speed of main memory. A technique used to compensate for the mismatch in operating speeds is to employ in extremely fast, small cache between the CPU and main memory whose access time is close to processor logic clock cycle time. The cache is used for storing segments of programs currently being executed in the CPU and temporary data frequently needed in the present calculations  by
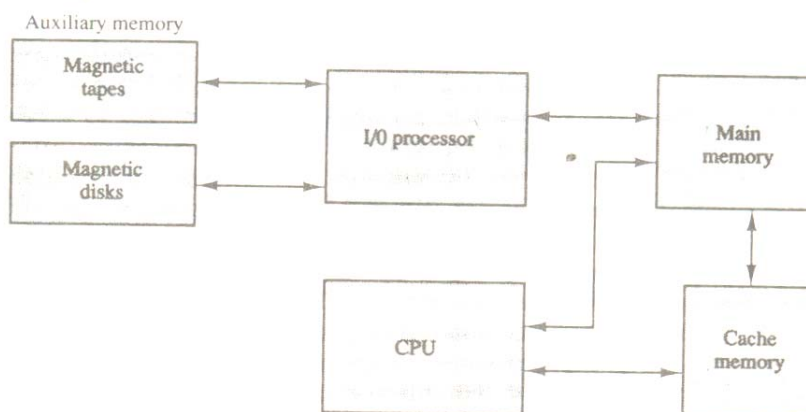
**Figure 7-1** Memory hierarchy in a computer system.

making programs and data available at a rapid rate, it is possible to increase the performance rate of the computer.

While the I/O processor manages data transfers between auxiliary memory and main memory, the cache organization is concerned with the transfer of information between main memory and CPU. Thus each is involved with a different level in the memory hierarchy system. The reason for having two or three levels of memory hierarchy is economics. As the storage capacity of the memory increases, the cost per bit for storing binary information decreases and the access time of the memory becomes longer. The auxiliary memory has a large storage capacity, is relatively inexpensive, but has low access speed compared to main memory. The cache memory is very small, relatively expensive, and has very high access speed. Thus as the memory access speed increases, so does its relative cost. The overall goal of using a memory hierarchy is to obtain the highest-possible average access speed while minimizing the total cost of the entire memory system.

Auxiliary and cache memories are used for different purposes. The cache holds those parts of the program and data that are most heavily used, while the auxiliary memory holds those parts that are not presently used by the CPU. Moreover, the CPU has direct access to both cache and main memory but not to auxiliary memory. The transfer from auxiliary to main memory is usually done by means of direct memory access of large blocks of data. The typical access time ratio between cache and main memory is about 1 to 7. For example, a typical cache memory may have an access time of 100ns, while main memory access time may be 700ns. Auxiliary memory average access time is usually 1000 times that of main memory. Block size in auxiliary memory typically ranges from256 to 2048 words, while cache block size is typically from 1 to 16 words.

Many operating systems are designed to enable the CPU to process a number of independent programs concurrently. This concept, called multiprogramming, refers to the existence of two or more programs indifferent parts of the memory hierarchy at the same time. In this way it is possible to keep all parts of the computer busy by working with several programs in sequence. For example, suppose that a program is being executed in the CPU and an I/O transfer is required. The CPU initiates the I/O processor to start executing the transfer. This leaves the CPU free to execute another program. In a multiprogramming system, when one program is waiting for input or output transfer, there is another program ready to utilize the CPU.

With multiprogramming the need arises for running partial programs, for varying the amount of main memory in use by a given program, and for moving programs around the memory hierarchy. Computer programs are sometimes too long to be accommodated in the total space available in main memory. Moreover, a computer system uses many programs and all the programs cannot reside in main memory at all times. A program with its data normally resides in auxiliary memory. When the program or a segment of the program is to be executed, it is transferred to main memory to be executed by the CPU. Thus one may think of auxiliary memory as containing the totality of information stored in a computer system. It is the task of the operating system to maintain in main memory a portion of this information that is currently active. The part of the computer system that supervises the flow of information between auxiliary memory and main memory is called the memory management system.

## 7.2    MAIN MEMORY

The main memory is the central storage unit in a computer system. It is a relatively large and fast memory used to store programs and data during the computer operation. The principal technology used for the main memory is based on semiconductor integrated circuits. Integrated circuit RAM chips are available in two possible operating modes, static and dynamic. The static RAM consists essentially of internal flip-flops that store the binary information. The stored information remains valid as long as power is applied to unit. The dynamic RAM stores the binary information in the form of electric charges that are applied to capacitors. The capacitors are provided inside the chip by MOS transistors. The stored charge on the capacitors tend to discharge with time and the capacitors must be periodically recharged by refreshing the dynamic memory. Refreshing is done by cycling through the words every few milliseconds to restore the decaying charge. The dynamic RAM offers reduced power consumption and larger storage capacity in a single memory chip. The static RAM is easier to use and has shorted read and write cycles.

Most of the main memory in a general-purpose computer is made up of RAM integrated circuit chips, but a portion of the memory may be constructed with ROM chips. Originally, RAM was used to refer to a random-access memory, but now it is used to designate a read/write memory to distinguish it from a read-only memory, although ROM is also random access. RAM is used for storing the bulk of the programs and data that are subject to change. ROM is used for storing programs that are permanently resident in the computer and for tables of constants that do not change in value one the production of the computer is completed.

Among other things, the ROM portion of main memory is needed for storing an initial program called a bootstrap loader. The bootstrap loader is a program whose function is to start the computer software operating when power is turned on. Since RAM is volatile, its contents are destroyed when power is turned off. The contents of ROM remain unchanged after power is turned off and on again. The startup of a computer consists of turning the power on and starting the execution of an initial program. Thus when power is turned on, the hardware of the computer sets the program counter to the first address of the bootstrap loader. The bootstrap program loads a portion of the operating system from disk to main memory and control is then transferred to the operating system, which prepares the computer fro general use.
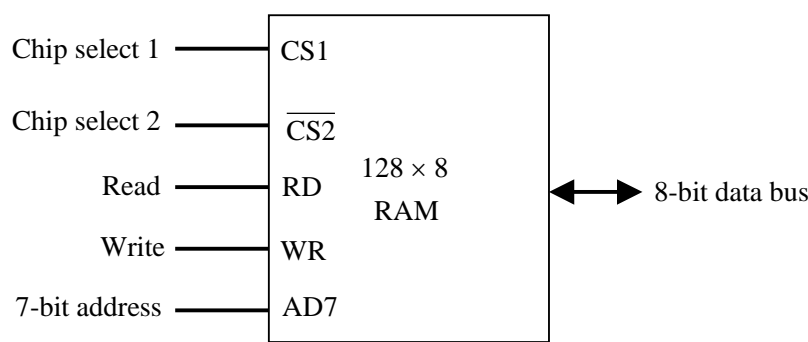
RAM and ROM chips are available in a variety of sizes. If the memory needed for the computer is larger than the capacity of one chip, it is necessary to combine a number of chips to form the required memory size. To demonstrate the chip interconnection, we will show an example of a $1024 \times 8$ memory constructed with $128 \times 8$ RAM chips and $512 \times 8$ ROM chips.

### 7.2.1   RAM AND ROM CHIPS

A RAM chip is better suited for communication with the CPU if it has one or more control inputs that select the chip only when needed. Another common feature is a bidirectional data bus that allows the transfer of data either from memory to CPU during a read operation or from CPU to memory during a write operation. A bidirectional bus can be constructed with three-state buffers. A three-state buffer output can be placed in one of three possible states: a signal equivalent to logic 1, a signal equivalent to logic 0, or a high-impedance state. The logic 1 and 0 are normal digital signals. The high-impedance state behaves like an open circuit, which means that the output does not carry a signal and has no logic significance.

The block diagram of a RAM chip is shown in Fig. 7-2. The capacity of the memory is 128 words of eight bits (one byte) per word. This requires a 7-bit

**Figure 7.2**  Typical RAM chip.



(a) Block diagram

| CSI | $\overline{CS2}$ | RD | WR | Memory function | State of data bus |
|-----|------|-----|-----|-----------------|-------------------|
| 0 | 0 | × | × | Inhibit | High-impedance |
| 0 | 1 | × | × | Inhibit | High-impedance |
| 1 | 0 | 0 | 0 | Inhibit | High-impedance |
| 1 | 0 | 0 | 1 | Write | Input data to RAM |
| 1 | 0 | 1 | × | Read | Output data from RAM |
| 1 | 1 | × | × | Inhibit | High-impedance |

(b) Function table

address and an 8-bit bidirectional data bus. The read and write inputs specify the memory operation and the two chips select (CS) control inputs are for enabling the chip only when it is selected by the microprocessor. The availability of more than one control input to select the chip facilitates the decoding of the address lines when multiple chips are used in the microcomputer. The read and write inputs are sometimes combined into one line labeled R/W. When the chip is selected, the two binary states in this line specify the two operations or read or write.

The function table listed in Fig. 7-2(b) specifies the operation of the RAM chip. The unit is in operation only when CSI = 1 and $\overline{CS2}$ = 0. The bar on top of the second select variable indicates that this input in enabled when it is equal to 0. If the chip select inputs are not enabled,

or if they are enabled but the read but the read or write inputs are not enabled, the memory is inhibited and its data bus is in a high-impedance state. When SC1 = 1 and $\overline{CS2}$ = 0, the memory can be placed in a write or read mode. When the WR input is enabled, the memory stores a byte from the data bus into a location specified by the address input lines. When the RD input is enabled, the content of the selected byte is placed into the data bus. The RD and WR signals control the memory operation as well as the bus buffers associated with the bidirectional data bus.

A ROM chip is organized externally in a similar manner. However, since a ROM can only read, the data bus can only be in an output mode. The block diagram of a ROM chip is shown in Fig. 7-3. For the same-size chip, it is possible to have more bits of ROM occupy less space than in RAM. For this reason, the diagram specifies a 512-byte ROM, while the RAM has only 128 bytes.

The nine address lines in the ROM chip specify any one of the 512 bytes stored in it. The two chip select inputs must be CS1 = 1 and $\overline{CS2}$ = 0 for the unit to operate. Otherwise, the data bus is in a high-impedance state. There is no need for a read or write control because the unit can only read. Thus when the chip is enabled by the two select inputs, the byte selected by the address lines appears on the data bus.

## 7.2.2 MEMORY ADDRESS MAP

The designer of a computer system must calculate the amount of memory required for the particular application and assign it to either RAM or ROM. The interconnection between memory and processor is then established form knowledge of the size of memory needed and the type of RAM and ROM chips available. The addressing of memory can be established by means of a table that specifies the memory address assigned to each chip. The table, called a memory address map, is a pictorial representation of assigned address space for each chip in the system.

To demonstrate with a particular example, assume that a computer system needs 512 bytes of RAM and 512 bytes of ROM. The RAM and ROM chips
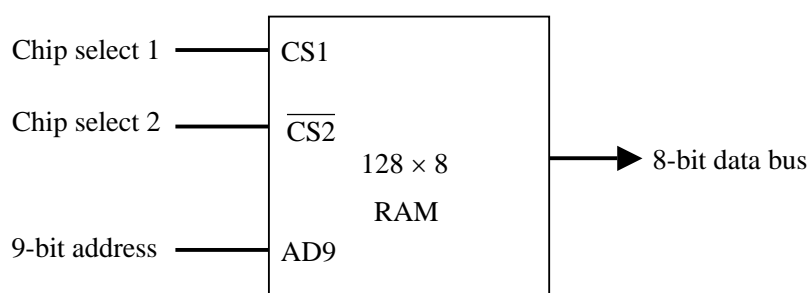


**Figure 7-3** Typical ROM chip.

to be used are specified in Figs. 7-2 and 7-3. The memory address map for this configuration is shown in Table 7-1. The component column specifies whether a RAM or a ROM chip is used. The hexadecimal address column assigns a range of hexadecimal equivalent addresses for each chip. The address bus lines are listed in the third column. Although there are 16 lines in the address bus, the table shows only 10 lines because the other 6 are not used in this example and are assumed to be zero. The small x's under the address bus lines designate those lines that must be connected to the address inputs in each chip. The RAM chips have 128 bytes and need seven address lines. The ROM chip has 512 bytes and needs 9 address lines. The x's are always assigned to the low-order bus lines: lines 1 through 7 for the RAM and lines 1 through 9 for the

ROM. It is now necessary to distinguish between four RAM chips by assigning to each a different address. For this particular example we choose bus lines 8 and 9 to represent four distinct binary combinations. Note that any other pair of unused bus lines can be chosen for this purpose. The table clearly shows that the nine low-order bus lines constitute a memory space fro RAM equal to $2^9 = 512$ bytes. The distinction between a RAM and ROM address is done with another bus line. Here we choose line 10 for this purpose. When line 10 is 0, the CPU selects a RAM, and when this line is equal to 1, it selects the ROM.

The equivalent hexadecimal address for each chip is obtained form the information under the address bus assignment. The address bus lines are subdivided into groups of four bits each so

TABLE 7-1 Memory Address Map for Microprocomputer

| Component | Hexadecimal address | Address bus | | |
|---|---|---|---|---|
| | | 10  9 | 8  7  6  5 | 4  3  2  1 |
| RAM 1 | 0000—007F | 0  0 | 0  x  x  x | x  x  x  x |
| RAM 2 | 0080—00FF | 0  0 | 1  x  x  x | x  x  x  x |
| RAM 3 | 0100—017F | 0  1 | 0  x  x  x | x  x  x  x |
| RAM 4 | 0180—01FF | 0  1 | 1  x  x  x | x  x  x  x |
| ROM | 0200—03FF | 1  x | x  x  x  x | x  x  x  x |

that each group can be represented with a hexadecimal digit. The first hexadecimal digit represents lines 13 to 16 and is always 0. The next hexadecimal digit represents lines 9 to 12, but lines 11 and 12 are always 0. The range of hexadecimal addresses for each component is determined from the x's associated with it. These x's represent a binary number that can range from an all-0's to an all-1's value.

### 7.2.3  MEMORY CONNECTION TO CPU

RAM and ROM chips are connected to a CPU through the data and address buses. The low-order lines in the address bus select the byte within the chips and other lines in the address bus select a particular chip through its chip select inputs. The connection of memory chips to the CPU is shown in Fig. 7-4. This configuration gives a memory capacity of 512 bytes of RAM and 512 bytes of ROM. It implements the memory map of Table 7-1. Each RAM receives the seven low-order bits of the address bus to select one of 128 possible bytes. The particular RAM chip selected is determined from lines 8 and 9 in the address bus. This is done through a $2 \times 4$ decoder whose outputs go to the SCI input in each RAM chip. Thus, when address lines 8 and 9 are equal to 00, the first RAM chip is selected. When 01, the second RAM chip is selected, and so on. The RD and WR outputs from the microprocessor are applied to the inputs of each RAM chip.

The selection between RAM and ROM is achieved through bus line 10. The RAMs are selected when the bit in this line is 0, and the ROM when the bit is 1. The other chip select input in the ROM is connected to the RD control line for the ROM chip to be enabled only during a read operation. Address bus lines 1 to 9 are applied to the input address of ROM without going through the decoder. This assigns addresses 0 to 511 to RAM and 512 to 1023 to ROM. The data bus of the ROM has only an output capability, whereas the data bus connected to the RAMs can transfer information in both directions.

The example just shown gives an indication of the interconnection complexity that can exist between memory chips and the CPU. The more chips that are connected, the more external decoders are required for selection among the chips. The designer must establish a memory map that assigns addresses to the various chips from which the required connections are determined.

## 7.3 AUXILIARY MEMORY

The most common auxiliary memory devices used in computer systems are magnetic disks and tapes. Other components used, but not as frequently, are magnetic drums, magnetic bubble memory, and optical disks. To understand fully the physical mechanism of auxiliary memory devices one must have a knowledge of magnetics, electronics, and electromechanical systems. Although the physical properties of these storage devices can be quite complex, their logical properties can be characterized and compared by a few parameters. The important characteristics of any device are its access mode, access time, transfer rate, capacity, and cost.

The average time required to reach a storage location in memory and obtain its contents is called the access time. In electromechanical devices with moving parts such as disks and tapes, the access time consists of a seek time required to position the read-write head to a location and a transfer time required to transfer data to or from the device. Because the seek time is usually much longer than the transfer time, auxiliary storage is organized in records or blocks. A record is a specified number of characters or words. Reading or writing is always done on entire records. The transfer rate is the number of characters or words that the device can transfer per second, after it has been positioned at the beginning of the record.

Magnetic drums and disks are quite similar in operation. Both consist of high-speed rotating surfaces coated with a magnetic recording medium. The rotating surface of the drum is a cylinder and that of the disk, a round flat plate. The recording surface rotates at uniform speed and is not stared or stopped during access operations. Bits are recorded as magnetic spots on the surface as it passes a stationary mechanism called a write head. Stored bits are detected by a change in magnetic field produced by a recorded spot on the surface as it passes through a read head. The amount of surface available for recording in a disk is greater than in a drum of equal physical size. Therefore, more information can be stored on a disk than on a drum of comparable size. For this reason, disks have replaced drums in more recent computers

### 7.3.1 MAGNETIC DISKS

A magnetic disk is a circular plate constructed of metal or plastic coated with magnetized material. Often both sides of the disk are used and several disks may be stacked on one spindle with read/write heads available on each surface. All disks rotate together at high speed and are not stopped or started from access purposes. Bits are stored in the magnetized surface in spots along concentric circles called tracks. The tracks are commonly divided into sections called sectors. In most systems, the minimum quantity of information which can be transferred is a sector. The sub division of tone disk surface into tracks and sectors.

Some units use a single read/write head from each disk surface. In this type of unit, the track address bits are used by a mechanical assembly to move the head into the specified track position before reading or writing. In other disk systems, separate read/write heads are provided for each track in each surface. The address can then select a particular track electronically through a decoder circuit. This type of unit is more expensive and is found only in very large computer systems.
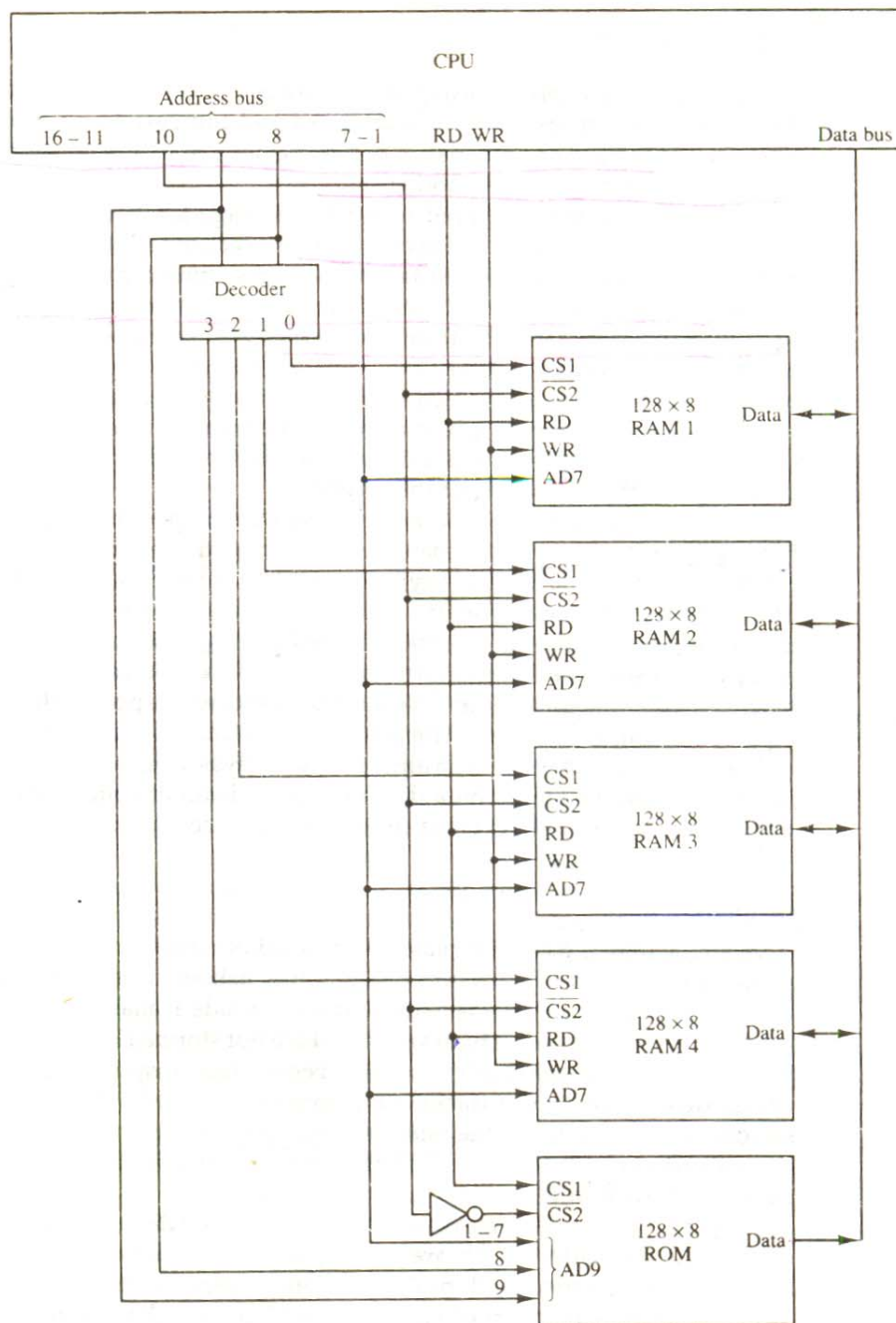
**Figure 7-4** Memory connection to the CPU.

Permanent timing tracks are used in disks to synchronize the bits and recognize the sectors. A disk system is addressed by address bits that specify the disk number, the disk surface, the sector number and the track within the sector. After the read/write heads are positioned in the specified track, the system has to wait until the rotating disk reaches the specified sector under the read/write head. Information transfer is very fast once the beginning of a sector has been reached.

Disks may have multiple heads and simultaneous transfer of bits from several tracks at the same time.

A track in a given sector near the circumference is longer than a track near the center of the disk. If bits are recorded with equal density, some tracks will contain more recorded bits than others. To make all the records in a sector of equal length, some disks use a variable recording density with higher density on tracks near the center than on tracks near the circumference. This equalizes the number of bits on all tracks of a given sector.

Disks that are permanently attached to the unit assembly and cannot be removed by the occasional user are called hard disks. A disk drive with removable disks is called a floppy disk. The disks used with a floppy disk drive are small removable disks made of plastic coated with magnetic recording material. There are two sizes commonly used, with diameters of 5.25 and 3.5 inches. The 3.5-inch disks are smaller and can store more data than can the 5.25-inch disks. Floppy disks are extensively used in personal computers as a medium for distributing software to computer users.

## 7.3.2   MAGNETIC TAPE

A magnetic tape transport consists of the electrical, mechanical, and electronic components to provide the parts and control mechanism for a magnetic-tape unit. The tape itself is a strip of plastic coated with a magnetic recording medium. Bits are recorded as magnetic spots on the tape along several tracks. Usually, seven or nine bits are recorded simultaneously to form a character together with a parity bit. Read/write heads are mounted one in each track so that data can be recorded and read as a sequence of characters.

Magnetic tape units can be stopped, started to move forward or in reverse, or can be rewound. However, they cannot be started or stopped fast enough between individual characters. For this reason, information is recorded in blocks referred to as records. Gaps of unrecorded tape are inserted between records where the tape can be stopped. The tape starts moving while in a gap and attains its constant speed by the time it reaches the next record. Each record on tape has an identification bit pattern at the beginning and end. By reading the bit pattern at the beginning, the tape control identifies the record number. By reading the bit pattern at the end of the record, the control recognizes the beginning of a gap. A tape unit is addressed by specifying the record number of characters in the record. Records may be of fixed or variable length.

## 7.4   ASSOCIATIVE MEMORY

Many data-processing applications require the search of items in a table stored in memory. An assembler program searches the symbol address table in order to extract the symbol's binary equivalent. An account number may be searched in a file to determine the holder's name and account status. The established way to search a table is to store all items where they can be addressed in sequence. The search procedure is a strategy for choosing a sequence of addresses, reading the content of memory at each address, and comparing the information read with the item being searched until a match occurs. The number of accesses to memory depends on the location of the item and the efficiency of the search algorithm. Many search algorithms have been developed to minimize the number of accesses while searching for an item in a random or sequential access memory.

The time required to find an item stored in memory can be reduced considerably if stored data can be identified for access by the content of the data itself rather than by an address. A memory unit accessed by content is called an associative memory or content addressable memory (CAM). This type of memory is accessed simultaneously and in parallel on the basis of data content rather than by specific address or location. When a word is written in an associative

memory, no address is given,. The memory is capable of finding an empty unused location to store the word. When a word is to be read from an associative memory, the content of the word, or part of the word, is specified. The memory locaters all words which match the specified content and marks them for reading.
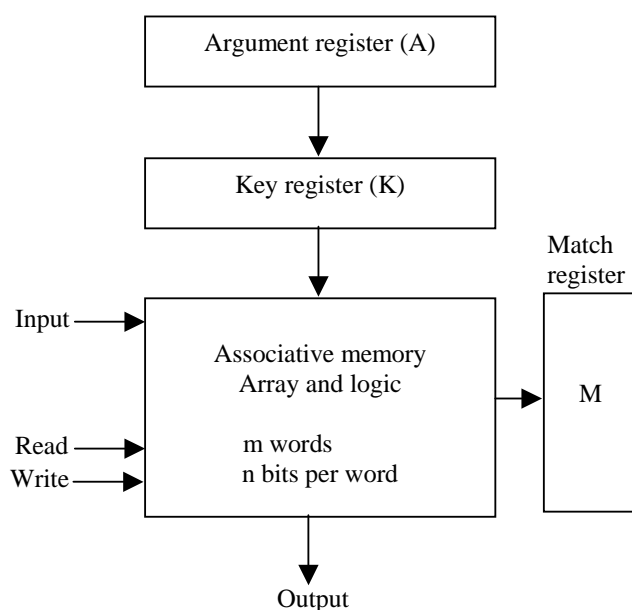
Because of its organization, the associative memory is uniquely suited to do parallel searches by data association. Moreover, searches can be done on an entire word or on a specific field within a word. An associative memory is more expensive then a random access memory because each cell must have storage capability as well as logic circuits for matching its content with an external argument. For this reason, associative memories are used in applications where the search time is very critical and must be very short.

### 7.4.1 HARDWARE ORGANIZATION

The block diagram of an associative memory is shown in Fig. 7-6. It consists of a memory array and logic from words with n bits per word. The argument register A and key register K each have n bits, one for each bit of a word. The match register M has m bits, one for each memory word. Each word in memory is compared in parallel with the content of the argument register. The words that match the bits of the argument register set a corresponding bit in the match register. After the matching process, those bits in the match register that have been set indicate the fact that their corresponding words have been matched. Reading is accomplished by a sequential access to memory for those words whose corresponding bits in the match register have been set.

The key register provides a mask for choosing a particular field or key in the argument word. The entire argument is compared with each memory word if the key register contains all 1's. Otherwise, only those bits in the argument that have 1's in their corresponding position of the key register are compared. Thus the key provides a mask or identifying piece of information which specifies how the reference to memory is made.

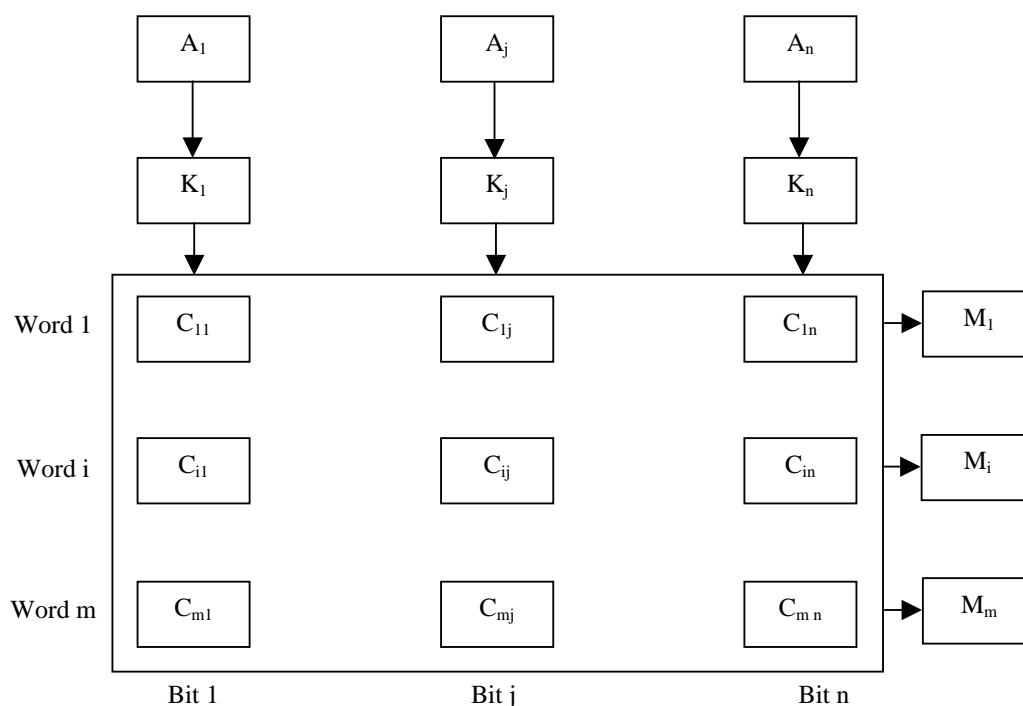**Figure 7-6** Block diagram of associative memory

To illustrate with a numerical example, suppose that the argument register A and the key register K have the bit configuration shown below. Only the three leftmost bits of A are compared with memory words because K has 1's in these positions.

| | | |
|---|---|---|
| A | 101 | 111100 |
| K | 111 | 000000 |
| Word 1 | 100 | 111100 no match |
| Word 2 | 101 | 000001 match |

Word 2 matches the unmasked argument field because the three leftmost bits of the argument and the word are equal.

The relation between the memory array and external registers in an associative memory is shown in Fig. 7-7. The cells in the array are marked by the letter C with two subscripts. The first subscript gives the word number and the second specifies the bit position in the word. Thus cell $C_{ij}$ is the cell for bit j in word i. A bit $A_j$ in the argument register is compared with all the bits in column j of the array provided that $K_j = 1$. This is done for all columns j = 1, 2,...,n. If a match occurs between all the unmasked bits of the argument and the bits in word i, the corresponding bit $M_i$ in the match register is set to 1. If one or more unmasked bits of the argument and the word do not match, $M_i$ is cleared to 0.

**Figure 7-7** Associative memory of m word, n cells per word



The internal organization of a typical cell $C_{ij}$ is shown in Fig. 7-8. It consists of a flip-flop storage element $F_{ij}$ and the circuits for reading, writing, and matching the cell. The input bit is transferred into the storage cell during a write operation. The bit stored is read out during a read operation. The match logic compares the content of the storage cell with the corresponding unmasked bit of the argument and provides an output for the decision logic that sets the bit in $M_i$.

## 7.4.2 MATCH LOGIC

The match logic for each word can be derived from the comparison algorithm for two binary numbers. First, we neglect the key bits and compare the argument in A with the bits stored in the cells of the words. Word i is equal to the argument in A if $A_j = F_{ij}$ for $j = 1, 2,\ldots, n$. Two bits are equal if they are both 1 or both 0. The equality of two bits can be expressed logically by the Boolean function
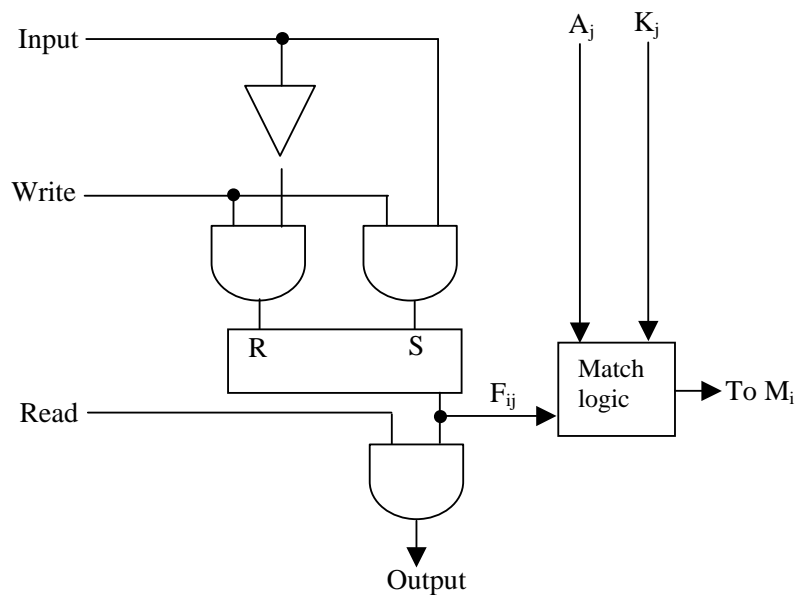
$$x_j = A_j F_{ij} + A_j^{'} F_{ij}^{'}$$

where $x_j = 1$ if the pair of bits in position j are equal; otherwise, $x_j = 0$.

For a word i to be equal to the argument in A we must have all $x_j$ variables equal to 1. This is the condition for setting the corresponding match bit $M_i$ to 1. The Boolean function for this condition is

$$M_i = x_1 x_2 x_3 \ldots x_n$$

and constitutes the AND operation of all pairs of matched bits in a word.

**Figure 7-8** One cell of associative memory.



We now include the key bit $K_j$ in the comparison logic. The requirement is that if $K_j = 0$, the corresponding bits of $A_j$ and $F_{ij}$ need no comparison. Only when $K_j = 1$ must they be compared. This requirement is achieved by ORing each term with $K_j^{'}$ , thus:

$$x_j + K'_j = \begin{cases} x_j & \text{if } K_j = 1 \\ 1 & \text{if } K_j = 0 \end{cases}$$

When $K_j = 1$, we have $K_j^{'} = 0$ and $x_j + 0 = x_j$. When $K_j = 0$, then $K_j^{'} = 1$ $x_j + 1 = 1$. A term $(x_j + K_j^{'})$ will be in the 1 state if its pair of bitsis not compared. This is necessary because each term is ANDed with all other terms so that an output of 1 will have no effect. The comparison of the bits has an effect only when $K_j = 1$.

The match logic for word i in an associative memory can now be expressed by the following Boolean function:

$$M_i = (x_1 + K_j^{'}) (x_2 + K_j^{'}) (x_3 + K_j^{'}) \dots (x_n + K_j^{'})$$

Each term in the expression will be equal to 1 if its corresponding $K_j^{'} = 0$. if $K_j = 1$, the term will be either 0 or 1 depending on the value of $x_j$. A match will occur and $M_i$ will be equal to 1 if all terms are equal to 1.

If we substitute the original definition of $x_j$. the Boolean function above can be expressed as follows:

$$M_i = \prod_{j=1}^{n} (A_j F_{ij} + A_j^{'} F_j^{'} + K_j^{'})$$

Where $\prod$ is a product symbol designating the AND operation of all n terms. We need m such functions, one for each word i = 1, 2, 3, …., m.

The circuit for catching one word is shown in Fig. 7-9. Each cell requires two AND gates and one OR gate. The inverters for $A_j$ and $K_j$ are needed once for each column and are used for all bits in the column. The output of all OR gates in the cells of the same word go to the input of a common AND gate to generate the match signal for $M_i$. $M_i$ will be logic 1 if a catch occurs and 0 if no match occurs. Note that if the key register contains all 0's, output $M_i$ will be a 1 irrespective of the value of A or the word. This occurrence must be avoided during normal operation.

### 7.4.3    READ OPERATION

If more than one word in memory matches the unmasked argument field, all the matched words will have 1's in the corresponding bit position of the catch register. It is then necessary to scan the bits of the match register on eat a time. The matched words are read in sequence by applying a read signal to each word line whose corresponding $M_i$ bit is a 1.
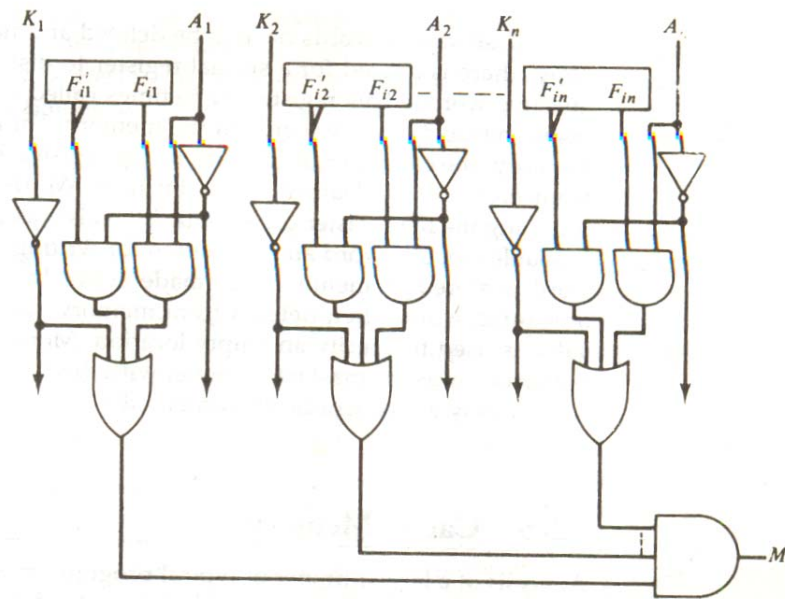


**Figure 7-9** Match logic for one word of associative memory

In most applications, the associative memory stores a table with no two identical items under a given key. In this case, only one word may match the unmasked argument field. By connecting output $M_i$ directly to the read line in the same word position (instead of the M register), the content of the matched word will be presented automatically at the output lines and no special read command signal is needed. Furthermore, if we exclude words having a zero content, an all-zero output will indicate that no match occurred and that the searched item is not available in memory.

### 7.4.4 WRITE OPERATION

An associative memory must have a write capability for storing the information to be searched. Writing in an associative memory can take different forms, depending on the application. If the entire memory is loaded with new information at once prior to a search operation then the writing can be done by addressing each location in sequence. This will make the device a random-access memory for writing and a content addressable memory for reading. The advantage here is that the address for input can be decoded as in a random-access memory. Thus instead of having m address lines, one for each word in memory, the number of address lines can be reduced by the decoder to d lines, where $m = 2^d$.

If unwanted words have to be deleted and new words inserted one at a time, there is a need for a special register to distinguish between active and inactive words. This register, sometimes called a tag register, would have as many bits as there are words in the memory. For every active word stored in memory, the corresponding bit in the tag register is set to 1. A word is deleted from memory by clearing its tag bit to 0. Words are stored in memory by scanning the tag register until the first 0 bit is encountered. This gives the first available inactive word and a position for writing a new word. After the new word is stored in memory it is made active by setting its tag bit to 1. An unwanted word when deleted from memory can be cleared to all 0's if this value is used to specify an empty location. Moreover, the words that have a tag bit of 0 must be masked (together with the $K_j$ bits) with the argument word so that only active words are compared.

### 7-5 CACHE MEMORY

Analysis of a large number of typical programs has shown that the references, to memory at any given interval of time tend to be confined within a few localized areas in memory. The phenomenon is known as the property of locality of reference. The reason for this property may be understood considering that a typical computer program flows in a straight-line fashion with program loops and subroutine calls encountered frequently. When a program loop is executed, the CPU repeatedly refers to the set of instructions in memory that constitute the loop. Every time a given subroutine is called, its set of instructions are fetched from memory. Thus loops and subroutines tend to localize the references to memory for fetching instructions. To a lesser degree, memory references to data also tend to be localized. Table-lookup procedures repeatedly refer to that portion in memory where the table is stored. Iterative procedures refer to common memory locations and array of numbers are confined within a local portion of memory. The result of all these observations is the locality of reference property, which states that over a short interval of time, the addresses generated by a typical program refer to a few localized areas of memory repeatedly, while the remainder of memory is accessed relatively frequently.

If the active portions of the program and data are placed in a fast small memory, the average memory access time can be reduced, thus reducing the total execution time of the program. Such a fast small memory is referred to as a cache memory. It is placed between the

CPU and main memory as illustrated in Fig. 7-10. The cache memory access time is less than the access time of main memory by a factor of 5 to 10. The cache is the fastest component in the memory hierarchy and approaches the speed of CPU components.

The fundamental idea of cache organization is that by keeping the most frequently accessed instructions and data in the fast cache memory, the average memory access time will approach the access time of the cache. Although the cache is only a small fraction of the size of main memory, a large fraction of memory requests will be found in the fast cache memory because of the locality of reference property of programs.

The basic operation of the cache is as follows. When the CPU needs to access memory, the cache is examined. If the word is found in the cache, it is read from the fast memory. If the word addressed by the CPU is not found in the cache, the main memory is accessed to read the word. A block of words containing the one just accessed is then transferred from main memory to cache memory. The block size may vary from one word (the one just accessed) to about 16 words adjacent to the one just accessed. In this manner, some data are transferred to cache so that future references to memory find the required words in the fast cache memory.

The performance of cache memory is frequently measured in terms of a quantity called hit ratio. When the CPU refers to memory and finds the word in cache, it is said to produce a hit. If the word is not found in cache, it is in main memory and it counts as a miss. The ratio of the number of hits divided by the total CPU references to memory (hits plus misses) is the hit ratio. The hit ratio is best measured experimentally by running representative programs in the computer and measuring the number of hits and misses during a given interval of time. Hit ratios of 0.9 and higher have been reported. This high ratio verifies the validity of the locality of reference property.

The average memory access time of a computer system can be improved considerably by use of a cache. If the hit ratio is high enough so that most of the time the CPU accesses the cache instead of main memory, the average access time is closer to the access time of the fast cache memory. For example, a computer with cache access time of 100 ns, a main memory access time of 1000 ns, and a hit ratio of 0.9 produces an average access time of 200 ns. This is a considerable improvement over a similar computer without a cache memory, whose access time is 1000 ns.

The basic characteristic of cache memory is its fast access time. Therefore, very little or no time must be wasted when searching for words in the cache. The transformation of data from main memory to cache memory is referred to as a mapping process. Three types of mapping procedures are of practical interest when considering the organization of cache memory:

1.     Associative mapping

2.     Direct mapping

3.     Set-associative mapping

To helping the discussion of these three mapping procedures we will use a specific example of a memory organization as shown in Fig. 7-10. The main memory can store 32K words of 12 bits each. The cache is capable of storing 512 of these words at any given time. For every word stored in cache, there is a duplicate copy in main memory.

The CPU communicates with both memories. It first sends a 15-bit address to cache. If there is a hit, the CPU accepts the 12-bit data from cache. If there is a miss, the CPU reads the word from main memory and the word is then transferred to cache.
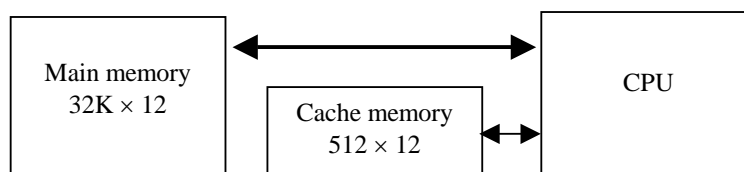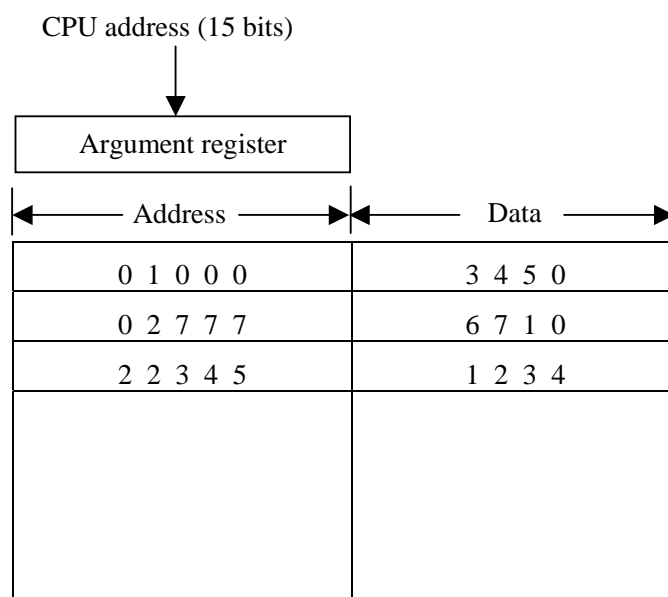
**Figure 7-10** Example of cache memory

## 7.5.1 ASSOCIATIVE MAPPING

The fasters and most flexible cache organization uses an associative memory. This organization is illustrated in Fig. 7-11. The associative memory stores both the address and content (data) of the memory word. This permits any location in cache to store any word from main memory. The diagram shows three words presently stored in the cache. The address value of 15 bits is shown as a five-digit octal number and its corresponding 12-bit word is shown as a four-digit octal number. A CPU address of 15 bits is placed in the argument register and the associative memory is searched for a matching address. If the address is found, the corresponding 12-bit data is read

**Figure 7-11** Associative mapping cache (all numbers in octal)

CPU address (15 bits)



and sent to the CPU. If no match occurs, the main memory is accessed for the word. The address-data pair is then transferred to the associative cache memory. If the cache is full, an address−data pair must be displaced to make room for a pair that is needed and not presently in the cache. The decision as to what pair is replaced is determined from the replacement algorithm that the designer chooses for the cache. A simple procedure is to replace cells of the cache in round-robin order whenever a new word is requested from main memory. This constitutes a first-in first-out (FIFO) replacement policy.

## 7.5.2 DIRECT MAPPING

Associative memories are expensive compared to random-access memories because of the added logic associated with each cell. The possibility of using a random-access memory for the cache is investigated in Fig. 7-12. The CPU address of 15 bits is divided into two fields. The

nine least significant bits constitute the index field and the remaining six bits form the tag and the index bits. The number of bits in the index field is equal to the number of address bits required to access the cache memory.

In the general case, there are $2^k$ words in cache memory and $2^n$ words in main memory. The n-bit memory address is divided into two fields: k bits for the index field and n − k bits for the tag field. The direct mapping cache organization uses the n-bit address to access the main memory and the k-bit index to access the cache. The internal organization of the words in the cache memory is as shown in Fig. 7–13(b). Each word in cache consists of the data word and its associated tag. When a new word is first brought into the cache, the tag bits are stored alongside the data bits. When the CPU generates a memory request, the index field is used for the address to access the cache.

**Figure 7-12** Addressing relationships between main and cache memories.
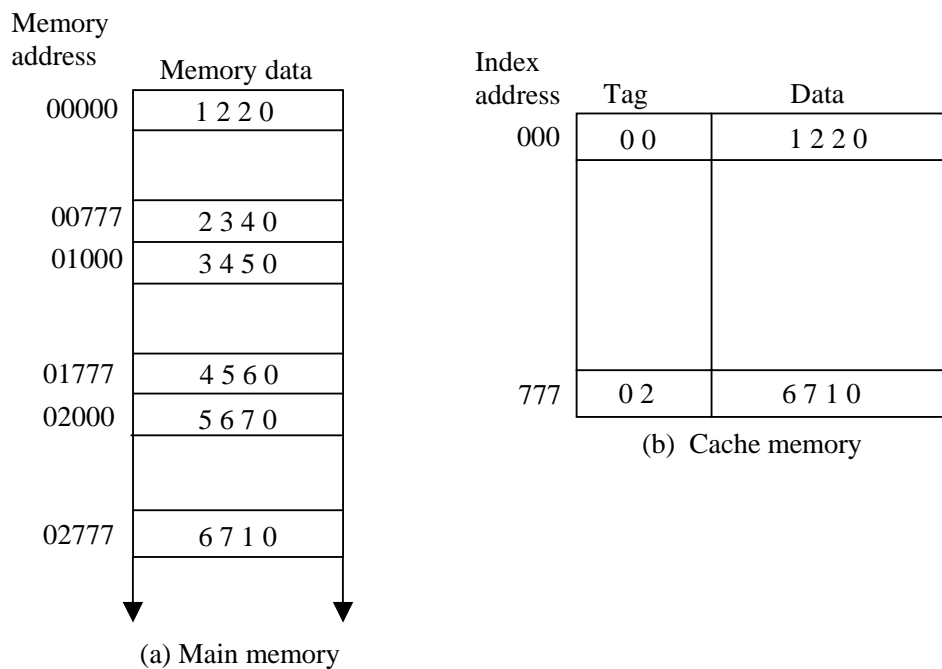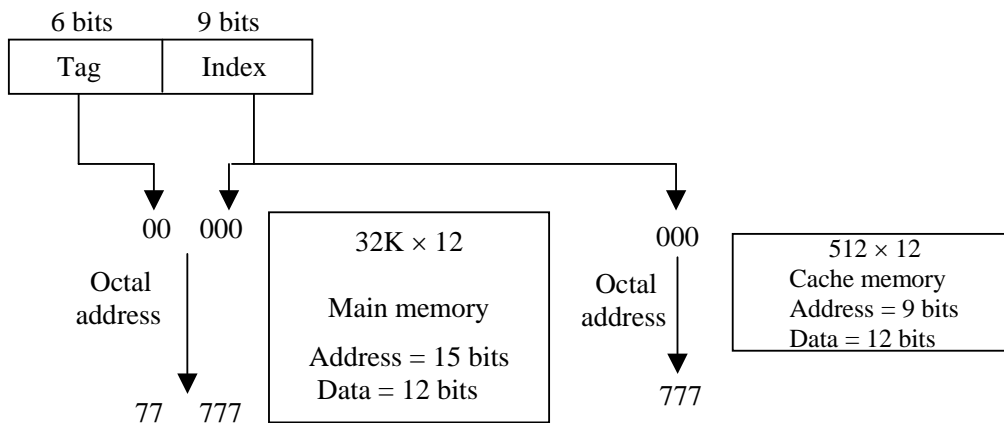




(a) Main memory

(b) Cache memory

**Figure 7-13** Direct mapping cache organization

The tag field of the CPU address is compared with the tag in the word read from the cache.  If the two tags match, there is a hit and the desired data word is in cache. If the two tags match, there is a hit and the desired data word is in cache.  If there is no match, there is a miss and the required word is read from main memory.  It is then stored in the cache together with the new tag, replacing the previous value.  The disadvantage of direct mapping is that the hit ratio can droop considerably if two or more words whose addresses have the same index but different tags are accessed repeatedly.  However, this possibility is minimized by the fact that such words are relatively far apart in the address range (multiples of 512 locations in this example).

To see how the direct-mapping organization operates, consider the numerical example shown in Fig. 7-13.  The word at address zero is presently stored in the cache (index = 000, tag = 00, data = 1220).  Suppose that the CPU now wants to access the word at address 02000.  The index address is 000, so it is sued to access the cache.  The two tags are then compared.  The cache tag is 00 but the address tag is 02, which does not produce a match.  Therefore, the main memory is accessed and the data word 5670 is transferred to the CPU.  The cache word at index address 000 is then replaced with a tag of 02 and data of 5670.

The direct-mapping example just described uses a block size of one word.  The same organization but using a block size of 8 words is shown in Fig. 7-14.  The index
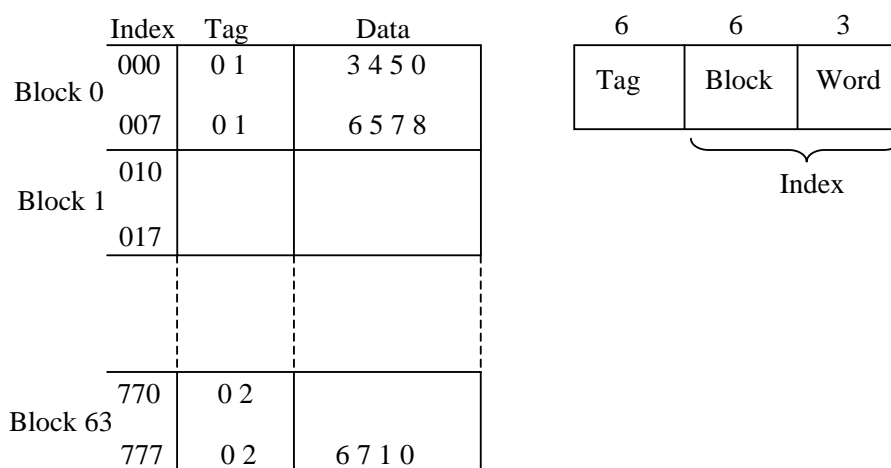


Fig 7.14

field is now divided into two parts: the block field and the word field. In a 512-word cache there are 64 block of 8 words each, since $64 \times 8 = 512$. The block number is specified with a 6-bit field and the word within the block is specified with a 3-bit field. The tag field stored within the cache is common to all eight words of the same block. Every time a miss occurs, an entire block of eight words must be transferred form main memory to cache memory. Although this takes extra time, the hit ratio will most likely improve with a larger block size because of the sequential nature of computer programs.

### 7.5.3   SET-ASSOCIATIVE MAPPING

It was mentioned previously that the disadvantage of direct mapping is that two words with the same index in their address but with different tag values cannot reside in cache memory at the same time. A third type of cache organization, called set-associative mapping, is an improvement over the direct-mapping organization in that each word of cache can store two or more words of

memory under the same index address. Each data word is stored together with its tag and the number of tag-data items in one word of cache is said to form a set. An example of a set-associative cache organization for a set size of two is shown in Fig. 7-15. Each index address refers to two data words and their associated tags. Each tag requires six bits and each data word has 12 bits, so the word length is 2(6 + 12) = 36 bits. An index address of nine bits can accommodate 512 words. Thus the size of cache memory is 512 × 36. It can accommodate 1024 words of main memory since each word of cache contains two data words. In general, a set-associative cache of set size k will accommodate k words of main memory in each word of cache.

| Index | Tag | Data | Tag | Data |
|---|---|---|---|---|
| 000 | 0 1 | 3 4 5 0 | 0 2 | 5 6 7 0 |
| | | | | |
| 777 | | 6 7 1 0 | 0 0 | 2 3 4 0 |

Figure 7-15 Two-way set-associative mapping cache.

The octal numbers listed in Fig. 7-15 are with reference to the main memory content illustrated in Fig. 7-13(a). The words stored at addresses 01000 and 02000 of main memory are stored in cache memory at index address 000. Similarly, the words at addresses 02777 and 00777 are stored in cache at index address 777. When the CPU generates a memory request, the index value of the address is used to access the cache. The tag field of the CPU address is then compared with both tags in the cache to determine if a catch occurs. The comparison logic is done by an associative search of the tags in the set similar to an associative memory search: thus the name "set-associative". The hit ratio will improve as the set size increases because more words with the same index but different tage can reside in cache. However, an increase in the set size increases the number of bit s in words of cache and requires more complex comparison logic.

When a miss occurs in a set-associative cache and the set is full, it is necessary to replace one of the tag-data items with a new value. The most common replacement algorithms used are: random replacement, first-in, first out (FIFO), and least recently used (LRU). With the random replacement policy the control chooses one tag-data item for replacement at random. The FIFO procedure selects for replacement the item that has been in the set the longest. The LRU algorithm selects for replacement the item that has been least recently used by the CPU. Both FIFO and LRU can be implemented by adding a few extra bits in each word of cache.

## 7.5.4  WRITING INTO CACHE

An important aspect of cache organization is concerned with memory write requests. When the CPU finds a word in cache during read operation, the main memory is not involved in the transfer. However, if the operation is a write, there are two ways that the system can proceed.

The simplest and most commonly used procedure is to up data main memory with every memory write operation, with cache memory being updated in parallel if it contains the word at the specified address. This is called the write-through method. This method has the advantage that main memory always contains the same data as the cache,. This characteristic is important in systems with direct memory access transfers. It ensures that the data residing in main memory are valid at tall times so that an I/O device communicating through DMA would receive the most recent updated data.

The second procedure is called the write-back method. In this method only the cache location is updated during a write operation. The location is then marked by a flag so that later when the words are removed form the cache it is copied into main memory. The reason for the write-back method is that during the time a word resides in the cache, it may be updated several times; however, as long as the word remains in the cache, it does not matter whether the copy in main memory is out of date, since requests from the word are filled from the cache. It is only when the word is displaced from the cache that an accurate copy need be rewritten into main memory. Analytical results indicate that the number of memory writes in a typical program ranges between 10 and 30 percent of the total references to memory.

## 7.5.5   CACHE INITIALIZATION

One more aspect of cache organization that must be taken into consideration is the problem of initialization. The cache is initialized when power is applied to the computer or when the main memory is loaded with a complete set of programs from auxiliary memory. After initialization the cache is considered to be empty, built in effect it contains some non-valid data. It is customary to include with each word in cache a valid bit to indicate whether or not the word contains valid data.

The cache is initialized by clearing all the valid bits to 0. The valid bit of a particular cache word is set to 1 the first time this word is loaded from main memory and stays set unless the cache has to be initialized again. The introduction of the valid bit means that a word in cache is not replaced by another word unless the valid bit is set to 1 and a mismatch of tags occurs. If the valid bit happens to be 0, the new word automatically replaces the invalid data. Thus the initialization condition has the effect of forcing misses from the cache until it fills with valid data.

## 7.6   VIRTUAL MEMORY

In a memory hierarchy system, programs and data are brought into main memory as they are needed by the CPU. Virtual memory is a concept used in some large computer systems that permit the user to construct programs as though a large memory space were available, equal to the totality of auxiliary memory. Each address that is referenced by the CPU goes through an address mapping from the so-called virtual address to a physical address in main memory. Virtual memory is used to give programmers the illusion that they have a very large memory at their disposal, even though the computer actually has a relatively small main memory. A virtual memory system provides a mechanism for translating program-generated addresses into correct main memory locations. This is done dynamically, while programs are being executed in the CPU. The translation or mapping is handled automatically by the hardware by means of a mapping table.

### 7.6.1 ADDRESS SPACE AND MEMORY SPACE

An address used by a programmer will be called a virtual address, and the set of such addresses the address space. An address in main memory is called a location or physical address. The set of such locations is called the memory space. Thus the address space is the set of addresses generated by programs as they reference instructions and data; the memory space consists of the actual main memory locations directly addressable for processing. In most computers the address and memory spaces are identical. The address space is allowed to be larger than the memory space in computers with virtual memory.

As an illustration, consider a computer with a main-memory capacity of 32K words (K = 1024). Fifteen bits are needed to specify a physical address in memory since $32K = 2^{15}$. Suppose that the computer has available auxiliary memory for storing $2^{20} = 1024K$ words. Thus auxiliary memory has a capacity for storing information equivalent to the capacity of 32 main memories. Denoting the address space by N and the memory space by M, we then have for this example N = 1024K and M = 32K.

In a multiprogram computer system, programs and data are transferred to and from auxiliary memory and main memory based on demands imposed by the CPU. Suppose that program 1 is currently being executed in the CPU. Program 1 and a portion of its associated data re moved from auxiliary memory into main memory as shown in Fig. 7-16. Portions of programs and data need not be in contiguous locations in memory since information is being moved in and out, and empty spaces may be available in scattered locations in memory.

In a virtual memory system, programmers are told that they have the total address space at their disposal. Moreover, the address field of the instruction code has a sufficient number of bits to specify all virtual addresses. In our example, the address field of an instruction code will consist of 20 bits but physical memory addresses must be specified with only 15 bits. Thus CPU will reference instructions and data with a 20-bit address, but the information at this address must be taken from physical memory because access to auxiliary storage for individual words will be prohibitively long. (Remember
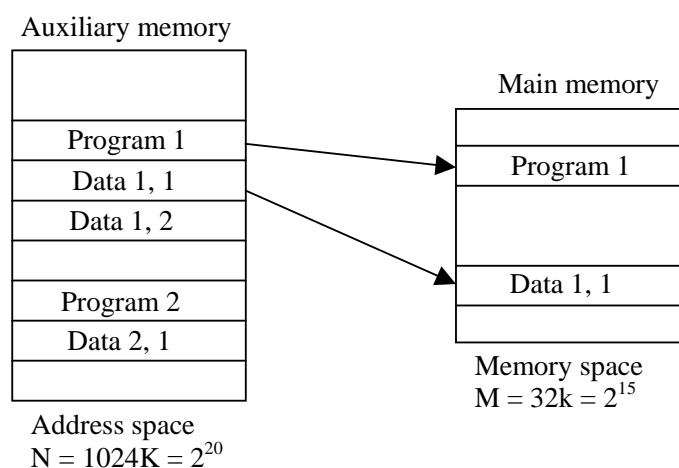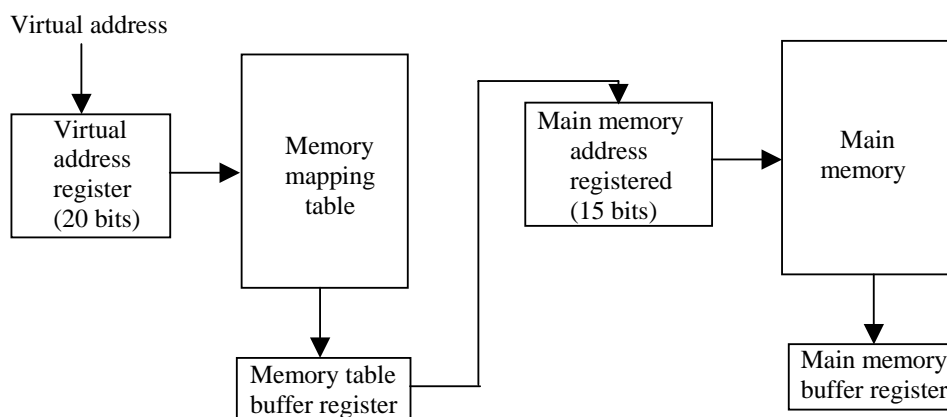


**Figure 7-16** Relation between address and memory space in a virtual memory system.

that for efficient transfers, auxiliary storage moves an entire record to the main memory). A table is then needed, as shown in Fig. 7-17, to map a virtual address of 20 bits to a physical address of

15 bits. The mapping is a dynamic operation, which means that every address is translated immediately as a word is referenced by CPU.

The mapping table may be stored in a separate memory as shown in Fig. 7-17 or in main memory. In the first case, an additional memory unit is required as well as one extra memory access time. In the second case, the table

Figure 7-17 Memory table for mapping a virtual address.



takes space from main memory and two accesses to memory are required with the program running at half speed. A third alternative is to use an associative memory as explained below.

## 7.6.2    ADDRESS MAPPING USING PAGES

The table implementation of the address mapping is simplified if the information in the address space and the memory space are each divided into groups of fixed size. The physical memory is broken down into groups of equal size called blocks, which may range from 64 to 4096 words each. The term page refers to groups of address space of the same size. For example, if a page or block consists of 1K words, then, using the previous example, address space is divided into 1024 pages and main memory is divided into 32 blocks. Although both a page and a block are split into groups of 1K words, a page refers to the organization of address space, while a block refers to the organization of memory space. The programs are also considered to be split into pages. Portions of programs are moved from auxiliary memory to main memory in records equal to the size of a page. The term "page frame" is sometimes used to denote a block.

Consider a computer with an address space of 8K and a memory space of 4K. If we split each into groups of 1K words we obtain eight pages and four blocks as shown in Fig. 7-18. At any given time, up to four pages of address space may reside in main memory in any one of the four blocks.

The mapping from address space to memory space is facilitated if each virtual address is considered to be represented by two numbers: a page number address and a line within the page. In a computer with $2^p$ words per page, p bits are used to specify a line address and the remaining high-order bits of the virtual address specify the page number. In the example of Fig. 7-18, a virtual address has 13 bits. Since each page consists of $2^{10} = 1024$ words, the high-order three bits of a virtual address will specify one of the eight pages and the low-order 10 bits give the line address within the page. Note that the line address in address space and memory space is the same; the only mapping required is from a page number to a block number.

The organization of the memory mapping table in a paged system is shown in Fig. 7-19. The memory-page table consists of eight words, one for each page. The address in the page table denotes the page number and the content of the word gives the block number where that page is stored in main memory. The table shows that pages 1, 2, 5 and 6 are now available in main memory in blocks 3, 0, 1, and 2, respectively. A presence bit in each location indicates whether the page has been transferred from auxiliary memory into main memory. A 0 in the presence bit indicates that this page is not available in main memory. The CPU references a word in memory with a virtual address of 13 bits. The three high-order bits of the virtual address specify a page number and also an address for the memory-page table. The content of the word in the memory
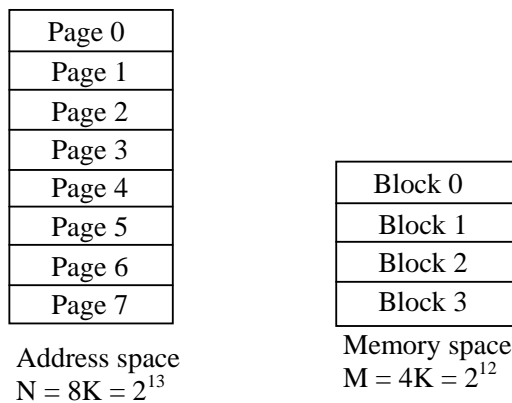


**Figure 7-18** Address space and memory space split into groups of 1K words.

page table at the page number address is read out into the memory table buffer register. If the presence bit is a 1, the block number thus read is transferred to the two high-order bits of the main memory address register. The line number from the virtual address is transferred into the 10 low-order bits of the memory address register. A read signal to main memory transfers the content of
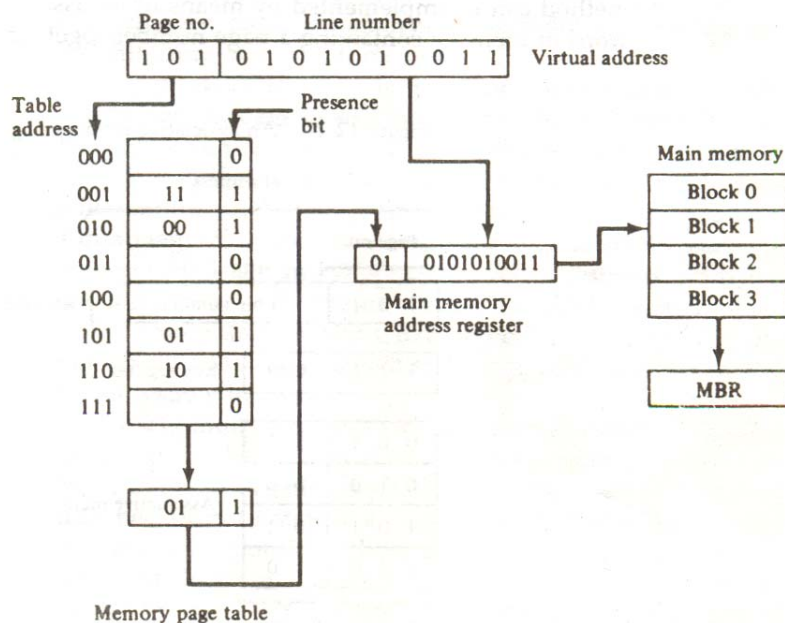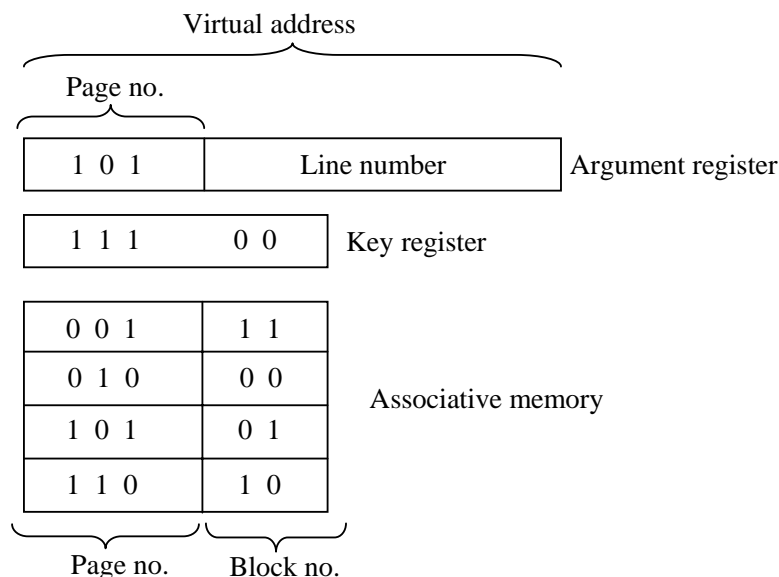


Figure 7-19 Memory table in a paged system.

the word to the main memory buffer register ready to be used by the CPU. If the presence bit in the word read from the page table is 0, it signifies that the content of the word referenced by the virtual address does not reside in main memory. A call to the operating system is then generated to fetch the required page from auxiliary memory and place it into main memory before resuming computation.

### 7.6.3   ASSOCIATIVE MEMORY PAGE TABLE

A random-access  memory page table is inefficient with respect to storage utilization. In the example of Fig. 7-19 we observe that eight words of memory are needed, one for each page, but at least four words will always be marked empty because main memory cannot accommodate more than four blocks. In general, system with n pages and m blocks would require a memory-page table of n locations of which up to m blocks will be marked with block numbers and all others will be empty. As a second numerical example, consider an address space of 1024K words and memory space of 32K words. If each page or block contains 1K words, the number of pages is 1024 and the number of blocks 32. The capacity of the memory-page table must be 1024 words and only 32 locations may have a presence bit equal to 1. At any given time, at least 992 locations will be empty and not in use.

A more efficient way to organize the page table would be to construct it with a number of words equal to the number of blocks in main memory. In this way the size of the memory is reduced and each location is fully utilized. This method can be implemented by means of an associative memory with each word in memory containing a page number together with its corresponding

**Figure 7-20**  An associative memory page table.



block number. The page field in each word is compared with the page number in the virtual address. If a match occurs, the word is read from memory and its corresponding block number is extracted.

Consider again the case of eight pages and four blocks as in the example of Fig. 7-19. We replace the random access memory-page table with an associative memory of four words as shown in Fig. 7-20. Each entry in the associative memory array consists of two fields. The first

three bits specify a field fro storing the page number. The last two bits constitute a field for storing the block number. The virtual address is placed in the argument register. The page number bits in the argument are compared with all page numbers in the page field of the associative memory. If the page number is found, the 5-bit word is read out from memory. The corresponding block number, being in the same word, is transferred to the main memory address register. If no match occurs, a call to the operating system is generated to bring the required page from auxiliary memory.

### 7.6.4 PAGE REPLACEMENT

A virtual memory system is a combination of hardware and software techniques. The memory management software system handles all the software operations for the efficient utilization of memory space. It must decide (1) which page in main memory ought to be removed to make room for a new page, (2) when a new page is to be transferred from auxiliary memory to main memory, and (3) where the page is to be placed in main memory. The hardware mapping mechanism and the memory management software together constitute the architecture of a virtual memory.

When a program start execution, one or more pages are transferred into main memory and the page table is set to indicate their position. The program is executed from main memory until it attempts to reference a page that is still in auxiliary memory. This condition is called page fault. When page fault occurs, the execution of the present program is suspended until the required page is brought into main memory. Since loading a page from auxiliary memory to main memory is basically an I/O operation, the operating system assigns this task to the I/O processor. In the meantime, controls transferred to the next program in memory that is waiting to be processed in the CPU. Later, when the memory block has been assigned and the transfer completed, the original program can resume its operation.

When a page fault occurs in a virtual memory system, it signifies that the page referenced by the CPU is not in main memory. A new page is then transferred from auxiliary memory to main memory. If main memory is full, it would be necessary to remove a page from a memory block to make room for the new page. The policy for choosing pages to remove is determined from the replacement algorithm that is used. The goal of a replacement policy is to try to remove the page least likely to be referenced in the immediate future.

Two of the most common replacement algorithms used are the first-in first-out (FIFO) and the least recently used (LRU). The FIFO algorithm selects for replacement the page the has been in memory the longest time. Each time a page is loaded into memory, its identification number is pushed into a FIFO stack. FIFO will be full whenever memory has no more empty blocks. When a new page must be loaded, the page least recently brought in is removed. The page to be removed is easily determined because its identification number is at the top of the FIFO stack. The FIFO replacement policy has the advantage of being easy to implement. It has the disadvantage that under certain circum-stances pages  are removed and loaded form memory too frequently.

The LRU policy is more difficult to implement but has been more attractive on the assumption that the least recently used page is a better candidate for removal than the least recently loaded pages in FIFO. The LRU algorithm can be implemented by associating a counter with every page that is in main memory. When a page is referenced, its associated counter is set to zero. At fixed intervals of time, the counters associated with all pages presently in memory are incremented by 1. The least recently used page is the page with the highest count. The counters are often called aging registers,  as their count indicates their age, that is, how long ago their associated pages have been referenced.

# SUMMARY

1. The memory unit is an essential component in any digital computer since it is needed for storing programs and data.

2. The total memory capacity of a computer can be visualized as being a hierarchy of components.

3. A special very-high speed memory called a cache is sometimes used to increase the speed of processing by making current programs and data available to the CPU at a rapid rate.

4. Auxiliary and cache memories are used for different purposes. The cache holds those parts of the program and data that are most heavily used, while the auxiliary memory holds those parts that are not presently used by the CPU.

5. The main memory is the central storage unit in a computer system. It is a relatively large and fast memory used to store programs and data during the computer operation.

6. A RAM chip is better suited for communication with the CPU if it has one or more control inputs that select the chip only when needed.

7. Most of the main memory in a general-purpose computer is made up of RAM integrated circuit chips, but a portion of the memory may be constructed with ROM chips.

8. RAM and ROM chips are connected to a CPU through the data and address buses. The low-order lines in the address bus select the byte within the chips and other lines in the address bus select a particular chip through its chip select inputs.

9. The designer of a computer system must calculate the amount of memory required for the particular application and assign it to either RAM or ROM.

10. The most common auxiliary memory devices used in computer systems are magnetic disks and tapes.

11. A magnetic disk is a circular plate constructed of metal or plastic coated with magnetized material.

12. Many data-processing applications require the search of items in a table stored in memory. An assembler program searches the symbol address table in order to extract the symbol's binary equivalent.

13. The match logic for each word can be derived from the comparison algorithm for two binary numbers.

14. If more than one word in memory matches the unmasked argument field, all the matched words will have 1's in the corresponding bit position of the catch register.

15. An associative memory must have a write capability for storing the information to be searched.

16. Analysis of a large number of typical programs has shown that the references, to memory at any given interval of time tend to be confined within a few localized areas in memory.

17. The fasters and most flexible cache organization uses an associative memory.

**18.** It was mentioned previously that the disadvantage of direct mapping is that two words with the same index in their address but with different tag values cannot reside in cache memory at the same time.

**19.** A virtual memory system is a combination of hardware and software techniques. The memory management software system handles all the software operations for the efficient utilization of memory space.

**20.** Two of the most common replacement algorithms used are the first-in first-out (FIFO) and the least recently used (LRU).

## SELF ASSESSMENT

1. Explain the memory hierarchy in the computer systems.

2. Explain different types of memory available with us.

3. What is memory address mapping? Explain the concept with the help of ROM chips.

4. How we can connect a memory with any CPU? Explain.

5. What is the concept of Associative memory? Explain.

6. Explain the significance on Match Logic.

7. What is cache memory? How it is fast as compared to conventional memory?

8. Explain the concept of virtual memory.

9. What is the concept of page replacement?