

MICROPROCESSORS AND INTERFACING

PROGRAMMING AND HARDWARE

SECOND EDITION

DOUGLAS V. HALL



McGRAW-HILL INTERNATIONAL EDITIONS
Computer Science Series

8086 • 80286 • 80386 • 80486

CONTENTS

Preface xi

CHAPTER 1

Computer Number Systems, Codes, and Digital Devices	1
Computer Number Systems and Codes	1
Arithmetic Operations on Binary, Hex, and BCD Numbers	6
Basic Digital Devices	10

CHAPTER 2

Computers, Microcomputers, and Microprocessors—An Introduction	19
Types of Computers	19
How Computers and Microcomputers Are Used—An Example	20
Overview of Microcomputer Structure and Operation	23
Execution of a Three-Instruction Program	24
Microprocessor Evolution and Types	26
The 8086 Microprocessor Family—Overview	27
8086 Internal Architecture	28
Introduction to Programming the 8086	32

CHAPTER 3

8086 Family Assembly Language Programming—Introduction	37
Program Development Steps	37
Constructing the Machine Codes for 8086 Instructions	47
Writing Programs for Use with an Assembler	53
Assembly Language Program Development Tools	59

CHAPTER 4

Implementing Standard Program Structures in 8086 Assembly Language	65
Simple Sequence Programs	65
Jumps, Flags, and Conditional Jumps	71
If-Then, If-Then-Else, and Multiple If-Then-Else Programs	77
While-Do Programs	82
Repeat-Until Programs	84
Instruction Timing and Delay Loops	91

CHAPTER 5

Strings, Procedures, and Macros	95
The 8086 String Instructions	95
Writing and Using Procedures	99
Writing and Using Assembler Macros	127

CHAPTER 6

8086 Instruction Descriptions and Assembler Directives	131
Instruction Descriptions	131
Assembler Directives	158

CHAPTER 7

8086 System Connections, Timing, and Troubleshooting	163
A Basic 8086 Microcomputer System	163
Using a Logic Analyzer to Observe Microprocessor Bus Signals	168
An Example Minimum-Mode System, the SDK-86	173
Troubleshooting a Simple 8086-Based Microcomputer	201

CHAPTER 8

8086 Interrupts and Interrupt Applications	207
8086 Interrupts and Interrupt Responses	207
Hardware Interrupt Applications	216
8254 Software-Programmable Timer/Counter	221
8259A Priority Interrupt Controller	232
Software Interrupt Applications	240

CHAPTER 9

Digital Interfacing	245
Programmable Parallel Ports and Handshake Input/Output	245
Interfacing a Microprocessor to Keyboards	260
Interfacing to Alphanumeric Displays	267
Interfacing Microcomputer Ports to High-Power Devices	277
Optical Motor Shaft Encoders	283

CHAPTER 10

Analog Interfacing and Industrial Control	290
Review of Operational-Amplifier Characteristics and Circuits	290
Sensors and Transducers	295
D/A Converter Operation, Interfacing, and Applications	301
A/D Converter Specifications, Types, and Interfacing	304
A Microcomputer-Based Scale	307
A Microcomputer-Based Industrial Process-Control System	317
An 8086-Based Process-Control System	320
Developing the Prototype of a Microcomputer-Based Instrument	331
Robotics and Embedded Control	332
Digital Signal Processing and Digital Filters	336

CHAPTER 11

DMA, DRAMs, Cache Memories, Coprocessors, and EDA Tools	345
Introduction	346
The 8086 Maximum Mode	346
Direct Memory Access (DMA) Data Transfer	348
Interfacing and Refreshing Dynamic RAMs	353
A Coprocessor—The 8087 Math Coprocessor	365
Computer-Based Design and Development Tools	379

CHAPTER 12

C, a High-Level Language for System Programming	389
Introduction—A Simple C Program Example	389
Program Development Tools for C	391
Programming in C	395

CHAPTER 13

Microcomputer System Peripherals	435
System-Level Keyboard Interfacing	435
Microcomputer Displays	439
Computer Mice and Trackballs	462
Computer Vision	463
Magnetic-Disk Data-Storage Systems	465
Optical Disk Data Storage	478
Printer Mechanisms and Interfacing	479
Speech Synthesis and Recognition with a Computer	481
Digital Video Interactive	483

CHAPTER 14

Data Communications and Networks	487
Introduction to Asynchronous Serial Data Communication	487
Serial-Data Transmission Methods and Standards	493
Asynchronous Communication Software on the IBM PC	506
Synchronous Serial-Data Communication and Protocols	518
Local Area Networks	522
The GPIB, HPIB, IEEE488 Bus	529

CHAPTER 15

The 80286, 80386, and 80486 Microprocessors	534
Multiuser/Multitasking Operating System Concepts	535
The Intel 80286 Microprocessor	543
The Intel 80386 32-Bit Microprocessor	547
The Intel 80486 Microprocessor	568
New Directions	570

BIBLIOGRAPHY 577

APPENDIX A iAPX 86/10 16-BIT HMOS MICROPROCESSOR 579

APPENDIX B INSTRUCTIONS: 8086/8088, 186, 8087 592

INDEX 607

CHAPTER

Computers, Microcomputers, and Microprocessors—An Introduction

We live in a computer-oriented society, and we are constantly bombarded with a multitude of terms relating to computers. Before getting started with the main flow of the book, we will try to clarify some of these terms and to give an overview of computers and computer systems.

OBJECTIVES

At the conclusion of this chapter, you should be able to:

1. Define the terms *microcomputer*, *microprocessor*, *hardware*, *software*, *firmware*, *timesharing*, *multitasking*, *distributed processing*, and *multiprocessing*.
2. Describe how a microcomputer fetches and executes an instruction.
3. List the registers and other parts in the 8086/8088 execution unit and bus interface unit.
4. Describe the function of the 8086/8088 queue.
5. Demonstrate how the 8086/8088 calculates memory addresses.

TYPES OF COMPUTERS

Mainframes

Computers come in a wide variety of sizes and capabilities. The largest and most powerful are often called **mainframes**. Mainframe computers may fill an entire room. They are designed to work at very high speeds with large data words, typically 64 bits or greater, and they have massive amounts of memory. Computers of this type are used for military defense control, for business data processing (in an insurance company, for example), and for creating computer graphics displays for science fiction movies. Examples of this type of computer are the IBM 4381, the Honeywell DPS8, and the Cray Y-MP/832. The fastest and most powerful mainframes are called **supercomputers**. Figure 2-1a, p. 20, shows a photograph of a Cray Y-MP/832 supercom-

puter, which contains eight central processors and 32 million 64-bit words of memory.

Minicomputers

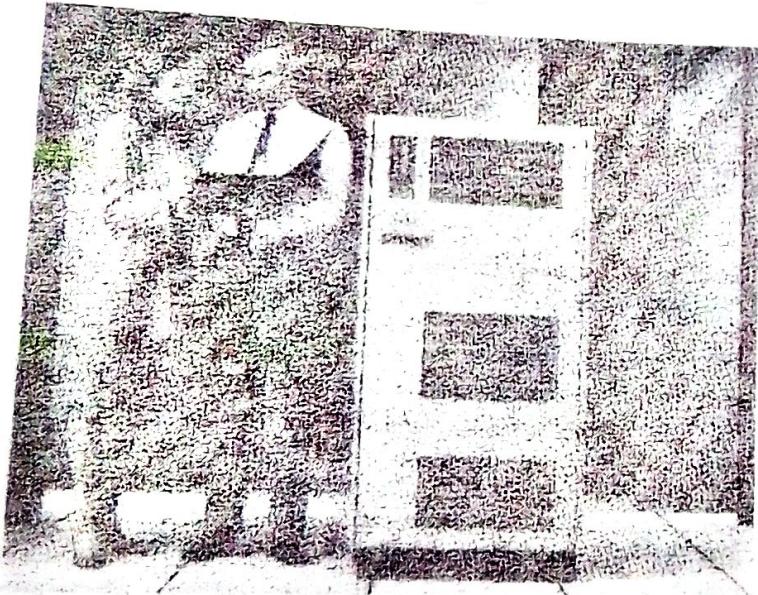
Scaled-down versions of mainframe computers are often called **minicomputers**. The main unit of a minicomputer usually fits in a single rack or box. A minicomputer runs more slowly, works directly with smaller data words (often 32-bit words), and does not have as much memory as a mainframe. Computers of this type are used for business data processing, industrial control (for an oil refinery, for example), and scientific research. Examples of this type of computer are the Digital Equipment Corporation VAX 6360 and the Data General MV/8000II. Figure 2-1b shows a photograph of a Digital Equipment Corporation's VAX 6360 minicomputer.

Microcomputers

As the name implies, **microcomputers** are small computers. They range from small controllers that work directly with 4-bit words and can address a few thousand bytes of memory to larger units that work directly with 32-bit words and can address billions of bytes of memory. Some of the more powerful microcomputers have all or most of the features of earlier minicomputers. Therefore, it has become very hard to draw a sharp line between these two types. One distinguishing feature of a microcomputer is that the CPU is usually a single integrated circuit called a **microprocessor**. Older books often used the terms **microprocessor** and **microcomputer** interchangeably, but actually the microprocessor is the CPU to which you add ROM, RAM, and ports to make a **microcomputer**. A later section in this chapter discusses the evolution of different types of microprocessors. Microcomputers are used in everything from smart sewing machines to computer-aided design systems. Examples of microcomputers are the Intel 8051 single-chip controller; the SDK-86, a single-board computer design kit; the IBM Personal Computer (PC); and the Apple Macintosh computer. The Intel 8051 microcontroller is contained in a single 40-pin chip. Figure 2-2a, p. 21, shows the SDK-86 board, and Figure 2-2b shows the Compaq 386/25 system.



(a)



(b)

FIGURE 2-1 (a) Photograph of Cray Y-MP/832 computer. (Courtesy Cray Research, Inc., and photographer, Paul Shambroom.) (b) Photograph of VAX 6360 minicomputer. (Courtesy Digital Equipment Corp.)

Computerizing an Electronics Factory—Problem

Now, suppose that we want to "computerize" an electronics company. By this we mean that we want to make computer use available to as many people in the company as possible as cheaply as possible. We want the engineers to have access to a computer which can help them design circuits. People in the drafting department should have access to a computer which can be used for computer-aided drafting. The accounting department should have access to a computer for doing all the financial bookkeeping. The warehouse should have access to a computer to help with inventory control. The manufacturing department should have access to a computer for controlling machines and testing finished products. The president, vice presidents, and supervisors should have access to a computer to help them with long-range planning. Secretaries should have access to a computer for word processing. Salespeople should have access to a computer to help them keep track of current pricing, product availability, and commissions. There are several ways to provide all the needed computer power. One solution is to simply give everyone an individual personal computer. The problem with this approach is that it makes it difficult for different people to access commonly needed data. In the next sections we show you two ways to provide computer power and common data to many users.

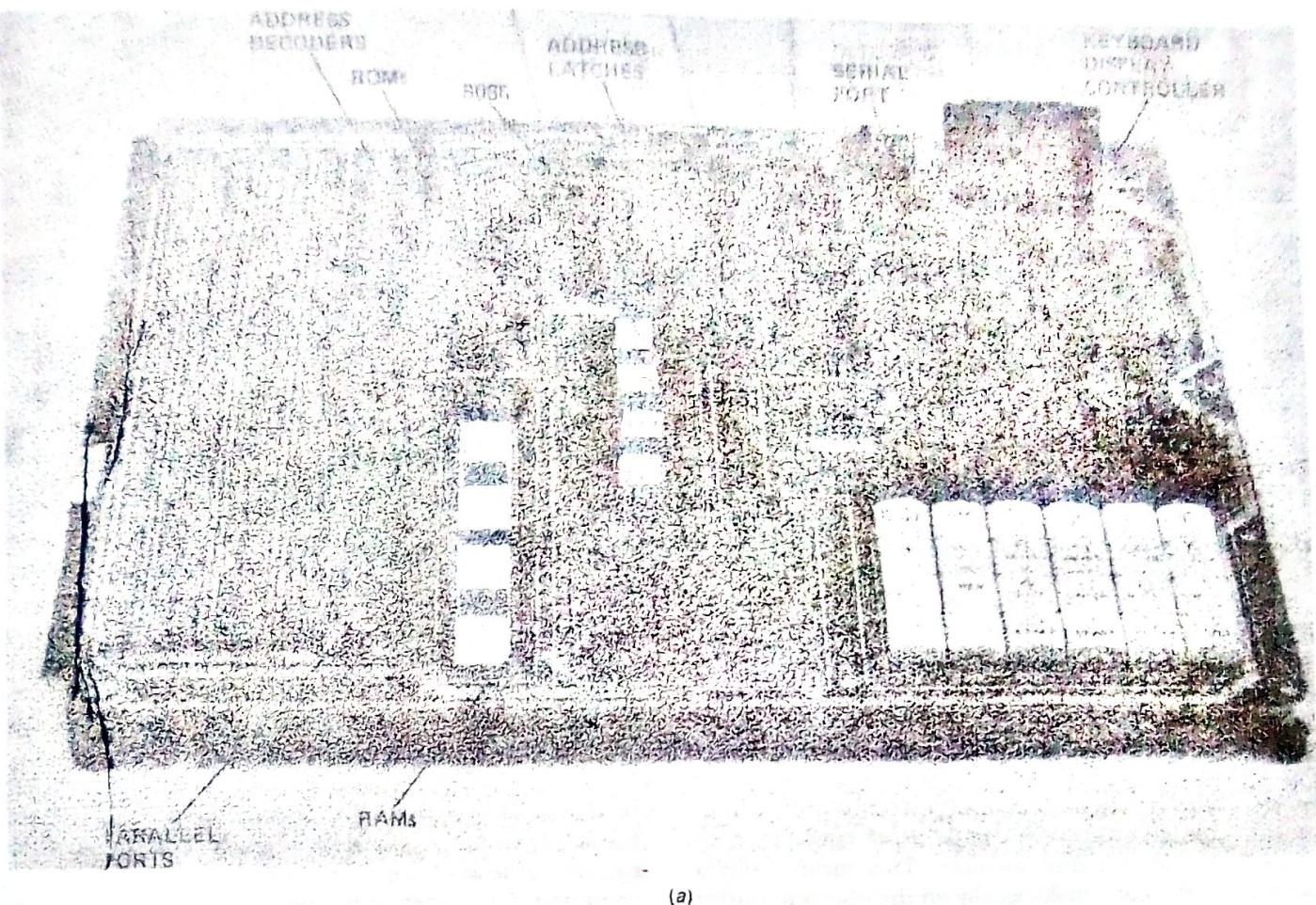
TIMESHARING AND MULTITASKING SYSTEMS

One common method of providing computer access is a *timesharing* system such as shown in Figure 2-3, p. 22. Several video terminals are connected to the computer through direct wires or through telephone lines. The terminal can be on the user's desk or even in the user's home. The rate at which a user usually enters data is very slow compared with the rate at which a computer can process the data. Therefore, the computer can serve many users by dividing its time among them in small increments. In other words, the computer works on user 1's program for perhaps 20 milliseconds (ms), then works on user 2's program for 20 ms, then works on user 3's program for 20 ms, and so on, until all the users have had a turn. In a few milliseconds the computer will get back to user 1 again and repeat the cycle. To each user it will appear as if he or she has exclusive use of the computer because the computer processes data as fast as the user enters it. A timesharing system such as this allows several users to interact with the computer at the same time. Each user can get information from or store information in the large memory attached to the computer. Each user can have an inexpensive printer attached to the terminal or can direct program or data output to a high-speed printer attached directly to the computer.

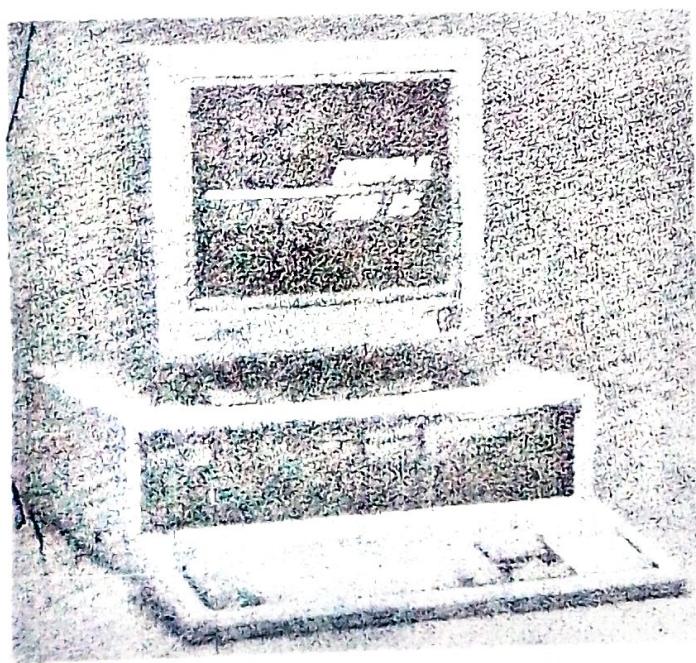
An airline ticket reservation computer might use a timesharing system such as this to allow users from all over the country to access flight information and make reservations. A time-multiplexed or time-sliced system such as this can also allow a computer to control many machines or processes in a factory. A computer is much faster than the machines or processes. Therefore, it can

HOW COMPUTERS AND MICROCOMPUTERS ARE USED—AN EXAMPLE

The following sections are intended to give you an overview of how computers are interfaced with users to do useful work. These sections should help you understand many of the features designed into current microprocessors and where this book is heading.



(a)



(b)

FIGURE 2-2 (a) Photograph of Intel SDK-86 board. (Intel Corp.) (b) Photograph of Compaq 386/25. (Compaq Corp.)

check and adjust many pressures, temperatures, motor speeds, etc., before it needs to get back and recheck the first one. A system such as this is often called a **multitasking system** because it appears to be doing many tasks at the same time.

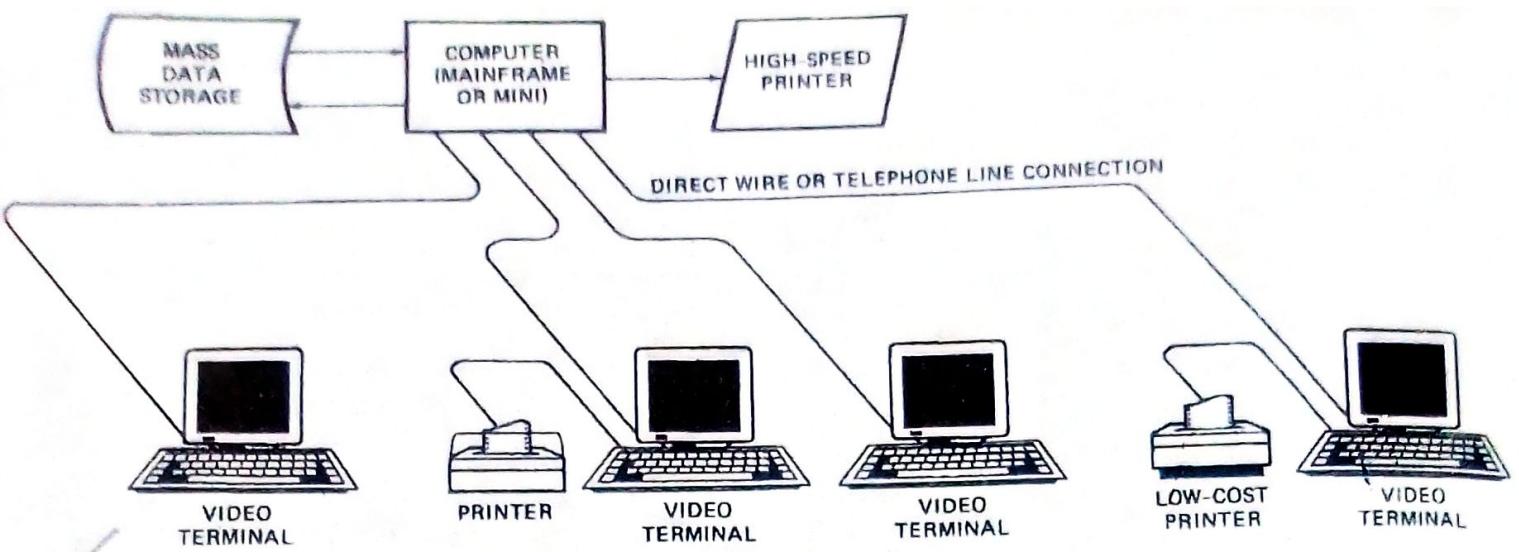
Now let's take another look at our problem of computerizing the electronics company. We could put a powerful computer in some central location and run wires from it to video display terminals on users' desks. Each user could then run the program needed to do a particular task. The accountant could run a ledger program, the secretary could run a word processing program, etc. Each user could access the computer's large data memory. Incidentally, a large collection of data stored in a computer's memory is often referred to as a *data base*. For a small company a system such as this might be adequate. However, there are at least two potential problems.

✓ The first potential problem is, "What happens if the computer is not working?" The answer to this question is that everything grinds to a halt. In a situation where people have become dependent on the computer, not much gets done until the computer is up and running again. The old saying about putting all your eggs in one basket comes to mind here.

✓ The second potential problem of the simple timesharing system is saturation. As the number of users increases, the time it takes the computer to do each user's task increases also. Eventually the computer's response time to each user becomes unreasonably long. People get very upset about the time they have to wait.

DISTRIBUTED PROCESSING OR MULTIPROCESSING

A partial solution for the two potential problems of a simple timesharing system is to use a distributed

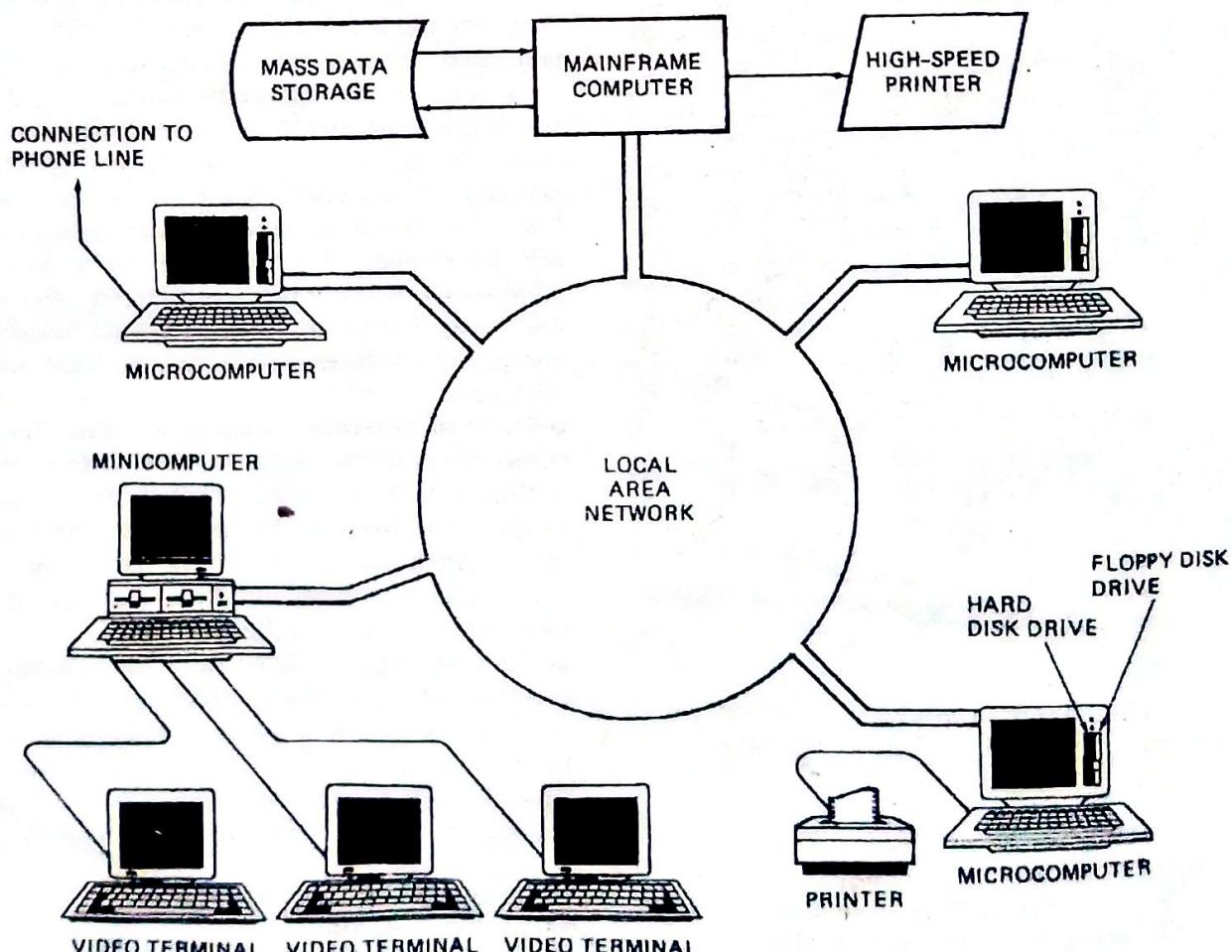


✓ FIGURE 2-3 Block diagram of a computer timesharing system.

processing system. Figure 2-4 shows a block diagram for such a system. The system has a powerful central computer with a large memory and a high-speed printer, as does the simple timesharing system described previously. However, in this system each user has a microcomputer instead of simply a video display terminal. In other words, each user station is an independently functioning microcomputer with a CPU, ROM, RAM, and probably magnetic or optical disk memory. This means that a person can do many tasks locally on the microcomputer

without having to use the large computer at all. Since the microcomputers are connected to the large computer through a network, however, a user can access the computing power, memory, or other resources of the large computer when needed.

✓ Distributing the processing to multiple computers or processors in a system has several advantages. First, if the large computer goes down, the local microcomputers can continue working until they need to access the large computer for something. Second, the burden on the



✓ FIGURE 2-4 Block diagram of a distributed processing computer system.

large computer is reduced greatly, because much of the computing is done by the local microcomputers. Finally, the distributed processing approach allows the system designer to use a local microcomputer that is best suited to the task it has to do.

COMPUTERIZED ELECTRONICS COMPANY OVERVIEW

Distributed processing seems to be the best way to go about computerizing our electronics factory. Engineers can have personal computers or engineering workstations on their desks. With these they can use available programs to design and test circuits. They can access the large computer if they need data from its memory. Through the telephone lines, the engineer with a personal computer can access data in the memory of other computers all over the world. The drafting people can have personal computers for simple work, or large computer-aided design systems for more complex work. Completed work can be stored in the memory of the large computer. The production department can have networked computers to keep track of product flow and to control the machines which actually mount components on circuit boards, etc. The accounting department can use personal computers with spreadsheet programs to work with financial data kept in the memory of the large computer. The warehouse supervisor can likewise use a personal computer with an inventory program to keep personal records and those in the large computer's memory updated. Corporate officers can have personal computers tied into the network. They then can interact with any of the other systems on the network. Salespeople can have portable personal computers that they can carry with them in the field. They can communicate with the main computer over the telephone lines using a modem. Secretaries doing word processing can use individual word processing units or personal computers. Users can also send messages to one another over the network. The specifics of a computer system such as this will obviously depend on the needs of the individual company for which the system is designed.

SUMMARY AND DIRECTION FROM HERE

The main concepts that you should take with you from this section are timesharing or multitasking and distributed processing or multiprocessing. As you work your way through the rest of this book, keep an overview

of the computerized electronics company in the back of your mind. The goal of this book is to teach you how the microcomputers and other parts of a system such as this work, how the parts are connected together, and how the system is programmed at different levels.

OVERVIEW OF MICROCOMPUTER STRUCTURE AND OPERATION

Figure 2-5 shows a block diagram for a simple microcomputer. The major parts are the central processing unit or CPU, memory, and the input and output circuitry or I/O. Connecting these parts are three sets of parallel lines called buses. The three buses are the address bus, the data bus, and the control bus. Let's take a brief look at each of these parts.

Memory

The memory section usually consists of a mixture of RAM and ROM. It may also have magnetic floppy disks, magnetic hard disks, or optical disks. Memory has two purposes. The first purpose is to store the binary codes for the sequences of instructions you want the computer to carry out. When you write a computer program, what you are really doing is writing a sequential list of instructions for the computer. The second purpose of the memory is to store the binary-coded data with which the computer is going to be working. This data might be the inventory records of a supermarket, for example.

Input/Output

The input/output or I/O section allows the computer to take in data from the outside world or send data to the outside world. Peripherals such as keyboards, video display terminals, printers, and modems are connected to the I/O section. These allow the user and the computer to communicate with each other. The actual physical devices used to interface the computer buses to external systems are often called ports. Ports in a computer function just as shipping ports do for a country. An input port allows data from a keyboard, an A/D converter, or some other source to be read into the computer under control of the CPU. An output port is used to send data from the computer to some peripheral, such as a video display terminal, a printer, or a D/A converter. Physically,

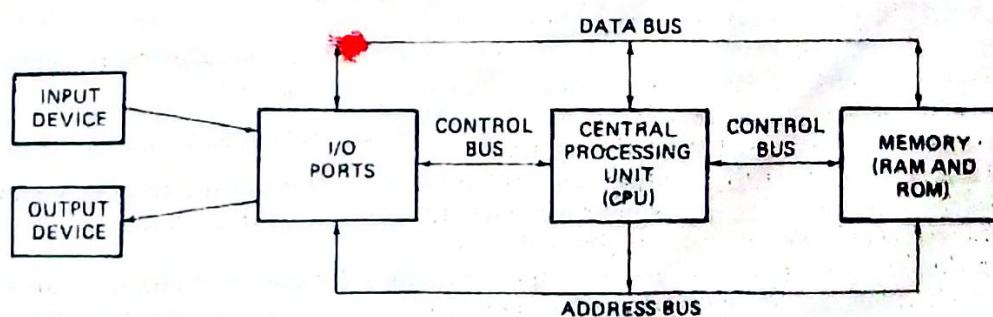


FIGURE 2-5 Block diagram of a simple microcomputer.

the simplest type of input or output port is just a set of parallel D flip-flops. If they are being used as an input port, the D inputs are connected to the external device, and the Q outputs are connected to the data bus which runs to the CPU. Data will then be transferred through the latches when they are enabled by a control signal from the CPU. In a system where they are being used as an output port, the D inputs of the latches are connected to the data bus, and the Q outputs are connected to some external device. Data sent out on the data bus by the CPU will be transferred to the external device when the latches are enabled by a control signal from the CPU.

Central Processing Unit

The central processing unit or CPU controls the operation of the computer. In a microcomputer the CPU is a microprocessor, as we discussed in an earlier section of the chapter. The CPU fetches binary-coded instructions from memory, decodes the instructions into a series of simple actions, and carries out these actions in a sequence of steps.

The CPU also contains an address counter or instruction pointer register, which holds the address of the next instruction or data item to be fetched from memory; general-purpose registers, which are used for temporary storage of binary data; and circuitry, which generates the control bus signals.

12

Address Bus

The address bus consists of 16, 20, 24, or 32 parallel signal lines. On these lines the CPU sends out the address of the memory location that is to be written to or read from. The number of memory locations that the CPU can address is determined by the number of address lines. If the CPU has N address lines, then it can directly address 2^N memory locations. For example, a CPU with 16 address lines can address 2^{16} or 65,536 memory locations, a CPU with 20 address lines can address 2^{20} or 1,048,576 locations, and a CPU with 24 address lines can address 2^{24} or 16,777,216 locations. When the CPU reads data from or writes data to a port, it sends the port address out on the address bus.

Data Bus

The data bus consists of 8, 16, or 32 parallel signal lines. As indicated by the double-ended arrows on the data bus line in Figure 2-5, the data bus lines are bidirectional. This means that the CPU can read data in from memory or from a port on these lines, or it can send data out to memory or to a port on these lines. Many devices in a system will have their outputs connected to the data bus, but only one device at a time will have its outputs enabled. Any device connected on the data bus must have three-state outputs so that its outputs can be disabled when it is not being used to put data on the bus.

Control Bus

The control bus consists of 4 to 10 parallel signal lines. The CPU sends out signals on the control bus to enable the outputs of addressed memory devices or port devices. Typical control bus signals are Memory Read, Memory Write, I/O Read, and I/O Write. To read a byte of data from a memory location, for example, the CPU sends out the memory address of the desired byte on the address bus and then sends out a Memory Read signal on the control bus. The Memory Read signal enables the addressed memory device to output a data word onto the data bus. The data word from memory travels along the data bus to the CPU.

Hardware, Software, and Firmware

When working around computers, you hear the terms hardware, software, and firmware almost constantly. *Hardware* is the name given to the physical devices and circuitry of the computer. *Software* refers to the programs written for the computer. *Firmware* is the term given to programs stored in ROMs or in other devices which permanently keep their stored information.

Summary of Important Points So Far

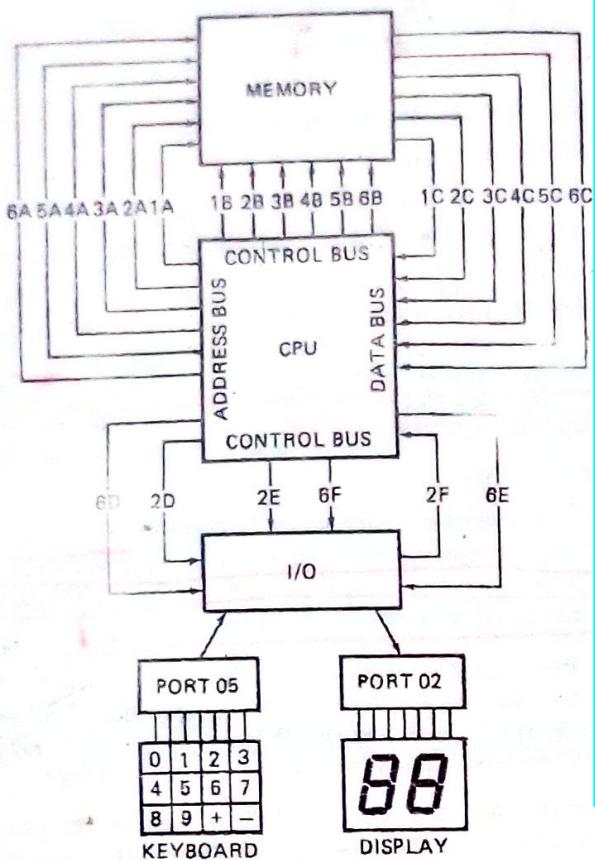
- A computer or microcomputer consists of memory, a CPU, and some input/output circuitry.
- These three parts are connected by the address bus, the data bus, and the control bus.
- The sequence of instructions or program for a computer is stored as binary numbers in successive memory locations.
- The CPU fetches an instruction from memory, decodes the instruction to determine what actions must be done for the instruction, and carries out these actions.

EXECUTION OF A THREE-INSTRUCTION PROGRAM

To give you a better idea of how the parts of a microcomputer function together, we will now describe the actions a simple microcomputer might go through to carry out (execute) a simple program. The three instructions of the program are

1. Input a value from a keyboard connected to the port at address 05H.
2. Add 7 to the value read in.
3. Output the result to a display connected to the port at address 02H.

Figure 2-6 shows in diagram form and sequential list form the actions that the computer will perform to execute these three instructions.



PROGRAM

1. INPUT A VALUE FROM PORT 05.
2. ADD 7 TO THIS VALUE.
3. OUTPUT THE RESULT TO PORT 02.

SEQUENCE

- 1A CPU SENDS OUT ADDRESS OF FIRST INSTRUCTION TO MEMORY.
- 1B CPU SENDS OUT MEMORY READ CONTROL SIGNAL TO ENABLE MEMORY.
- 1C INSTRUCTION BYTE SENT FROM MEMORY TO CPU ON DATA BUS.
- 2A ADDRESS NEXT MEMORY LOCATION TO GET REST OF INSTRUCTION.
- 2B SEND MEMORY READ CONTROL SIGNAL TO ENABLE MEMORY.
- 2C PORT ADDRESS BYTE SENT FROM MEMORY TO CPU ON DATA BUS.
- 2D CPU SENDS OUT PORT ADDRESS ON ADDRESS BUS.
- 2E CPU SENDS OUT INPUT READ CONTROL SIGNAL TO ENABLE PORT.
- 2F DATA FROM PORT SENT TO CPU ON DATA BUS.
- 3A CPU SENDS ADDRESS OF NEXT INSTRUCTION TO MEMORY.
- 3B CPU SENDS MEMORY READ CONTROL SIGNAL TO ENABLE MEMORY.
- 3C INSTRUCTION BYTE FROM MEMORY SENT TO CPU ON DATA BUS.
- 4A CPU SENDS NEXT ADDRESS TO MEMORY TO GET REST OF INSTRUCTION.
- 4B CPU SENDS MEMORY READ CONTROL SIGNAL TO ENABLE MEMORY.
- 4C NUMBER 07H SENT FROM MEMORY TO CPU ON DATA BUS.
- 5A CPU SENDS ADDRESS OF NEXT INSTRUCTION TO MEMORY.
- 5B CPU SENDS MEMORY READ CONTROL SIGNAL TO ENABLE MEMORY.
- 5C INSTRUCTION BYTE FROM MEMORY SENT TO CPU ON DATA BUS.
- 6A CPU SENDS OUT NEXT ADDRESS TO GET REST OF INSTRUCTION.
- 6B CPU SENDS OUT MEMORY READ CONTROL SIGNAL TO ENABLE MEMORY.
- 6C PORT ADDRESS BYTE SENT FROM MEMORY TO CPU ON DATA BUS.
- 6D CPU SENDS OUT PORT ADDRESS ON ADDRESS BUS.
- 6E CPU SENDS OUT DATA TO PORT ON DATA BUS.
- 6F CPU SENDS OUT OUTPUT WRITE SIGNAL TO ENABLE PORT.

(a) Machine Lang		
MEMORY ADDRESS	CONTENTS (BINARY)	CONTENTS (HEX)
00100H	11100100	E4
00101H	00000101	05
00102H	00000100	04
00103H	00000111	07
00104H	11100110	E6
00105H	00000010	02

(b)

FIGURE 2-6 (a) Execution of a three-step computer program. (b) Memory addresses and memory contents for a three-step program.

For this example, assume that the CPU fetches instructions and data from memory 1 byte at a time, as is done in the original IBM PC and its clones. Also assume that the binary codes for the instructions are in sequential memory locations starting at address 00100H. Figure 2-6b shows the actual binary codes that would be required in successive memory locations to execute this program on an IBM PC-type microcomputer.

The CPU needs an instruction before it can do anything, so its first action is to fetch an instruction byte from memory. To do this, the CPU sends out the address of the first instruction byte, in this case 00100H, to memory on the address bus. This action is represented by line 1A in Figure 2-6a. The CPU then sends out a Memory Read signal on the control bus (line 1B in the figure). The Memory Read signal enables the memory to output the addressed byte on the data bus. This action is represented by line 1C in the figure. The CPU reads in this first instruction byte (E4H) from the data bus and decodes it. By decode we mean that the CPU determines from the binary code read in what actions it is supposed to take. If the CPU is a microprocessor, it selects the sequence of microinstructions needed to

carry out the instruction read from memory. For the example instruction here, the CPU determines that the code read in represents an Input instruction. From decoding this instruction byte, the CPU also determines that it needs more information before it can carry out the instruction. The additional information the CPU needs is the address of the port that the data is to be input from. This port address part of the instruction is stored in the next memory location after the code for the input instruction.

To fetch this second byte of the instruction, the CPU sends out the next sequential address (00101H) to memory, as shown by line 2A in the figure. To enable the addressed memory device, the CPU also sends out another Memory Read signal on the control bus (line 2B). The memory then outputs the addressed byte on the data bus (line 2C). When the CPU has read in this second byte, 05H in this case, it has all the information it needs to execute the instruction.

To execute the input instruction, the CPU sends out the port address (05H) on the address bus (line 2D) and sends out an I/O Read signal on the control bus (line 2E). The I/O Read signal enables the addressed port

device to put a byte of data on the data bus (line 2F). The CPU reads in the byte of data and stores it in an internal register. This completes the fetching and execution of the first instruction.

Having completed the first instruction, the CPU must now fetch its next instruction from memory. To do this, it sends out the next sequential address (00102H) on the address bus (line 3A) and sends out a Memory Read signal on the control bus (line 3B). The Memory Read signal enables the memory device to put the addressed byte (04H) on the data bus (line 3C). The CPU reads in this instruction byte from the data bus and decodes it. From this instruction byte the CPU determines that it is supposed to add some number to the number stored in the internal register. The CPU also determines from decoding this instruction byte that it must go to memory again to get the next byte of the instruction, which contains the number that it is supposed to add. To get the required byte, the CPU will send out the next sequential address (00103H) on the address bus (line 4A) and another Memory Read signal on the control bus (line 4B). The memory will then output the contents of the addressed byte (the number 07H) on the data bus (line 4C). When the CPU receives this number, it will add it to the contents of the internal register. The result of the addition will be left in the internal register. This completes the fetching and executing of the second instruction.

The CPU must now fetch the third instruction. To do this, it sends out the next sequential address (00104H) on the address bus (line 5A) and sends out a Memory Read signal on the control bus (line 5B). The memory then outputs the addressed byte (E6H) on the data bus (line 5C). From decoding this byte, the CPU determines that it is now supposed to do an Output operation to a port. The CPU also determines from decoding this byte that it must go to memory again to get the address of the output port. To do this, it sends out the next sequential address (00105H) on the address bus (line 6A), sends out a Memory Read signal on the control bus (line 6B), and reads in the byte (02H) put on the data bus by the memory (line 6C). The CPU now has all the information that it needs to execute the Output instruction.

To output a data byte to a port, the CPU first sends out the address of the desired port on the address bus (line 6D). Next it outputs the data byte from the internal register on the data bus (line 6E). The CPU then sends out an I/O Write signal on the control bus (line 6F). This signal enables the addressed output port device so that the data from the data bus lines can pass through it to the LED displays. When the CPU removes the I/O Write signal to proceed with the next instruction, the data will remain latched on the output pins of the port device. The data will remain latched on the port until the power is turned off or until a new data word is output to the port. This is important because it means that the computer does not have to keep outputting a value over and over in order for it to remain on the output.

All the steps described above may seem like a great deal of work just to input a value from a keyboard, add 7 to it, and output the result to a display. Even a simple

microcomputer, however, can run through all these steps in a few microseconds.

Summary of Simple Microcomputer Bus Operation

1. A microcomputer fetches each program instruction in sequence, decodes the instruction, and executes it. (Memory read)
2. The CPU in a microcomputer fetches instructions or reads data from memory by sending out an address on the address bus and a Memory Read signal on the control bus. The memory outputs the addressed instruction or data word to the CPU on the data bus. (Memory write)
3. The CPU writes a data word to memory by sending out an address on the address bus, sending out the data word on the data bus, and sending a Memory Write signal to memory on the control bus.
4. To read data from a port, the CPU sends out the port address on the address bus and sends an I/O Read signal to the port device on the control bus. Data from the port comes into the CPU on the data bus.
5. To write data to a port, the CPU sends out the port address on the address bus, sends out the data to be written to the port on the data bus, and sends an I/O Write signal to the port device on the control bus.

MICROPROCESSOR EVOLUTION AND TYPES

As we told you in the preceding section, a microprocessor is used as the CPU in a microcomputer. There are now many different microprocessors available, so before we dig into the details of a specific device, we will give you a short microprocessor history lesson and an overview of the different types.

Microprocessor Evolution

A common way of categorizing microprocessors is by the number of bits that their ALU can work with at a time. In other words, a microprocessor with a 4-bit ALU will be referred to as a 4-bit microprocessor, regardless of the number of address lines or the number of data bus lines that it has. The first commercially available microprocessor was the Intel 4004, produced in 1971. It contained 2300 PMOS transistors. The 4004 was a 4-bit device intended to be used with some other devices in making a calculator. Some logic designers, however, saw that this device could be used to replace PC boards full of combinational and sequential logic devices. Also, the ability to change the function of a system by just changing the programming, rather than redesigning the hardware, is very appealing. It was these factors that pushed the evolution of microprocessors.

In 1972 Intel came out with the 8008, which was capable of working with 8-bit words. The 8008, however,

required 20 or more additional devices to form a functional CPU. In 1974 Intel announced the 8080, which had a much larger instruction set than the 8008 and required only two additional devices to form a functional CPU. Also, the 8080 used NMOS transistors, so it operated much faster than the 8008. The 8080 is referred to as a second-generation microprocessor.

Soon after Intel produced the 8080, Motorola came out with the MC6800, another 8-bit general-purpose CPU. The 6800 had the advantage that it required only a +5-V supply rather than the -5-V, +5-V, and +12-V supplies required by the 8080. For several years the 8080 and the 6800 were the top-selling 8-bit microprocessors. Some of their competitors were the MOS Technology 6502, used as the CPU in the Apple II microcomputer, and the Zilog Z80, used as the CPU in the Radio Shack TRS-80 microcomputer.

As designers found more and more applications for microprocessors, they pressured microprocessor manufacturers to develop devices with architectures and features optimized for doing certain types of tasks. In response to the expressed needs, microprocessors have evolved in three major directions during the last 15 years.

Dedicated or Embedded Controllers

One direction has been dedicated or embedded controllers. These devices are used to control "smart" machines, such as microwave ovens, clothes washers, sewing machines, auto ignition systems, and metal lathes. Texas Instruments has produced millions of their TMS-1000 family of 4-bit microprocessors for this type of application. In 1976 Intel introduced the 8048, which contains an 8-bit CPU, RAM, ROM, and some I/O ports all in one 40-pin package. Other manufacturers have followed with similar products. These devices are often referred to as microcontrollers. Some currently available devices in this category—the Intel 8051 and the Motorola MC6801, for example—contain programmable counters and a serial port (UART) as well as a CPU, ROM, RAM, and parallel I/O ports. A more recently introduced single-chip microcontroller, the Intel 8096, contains a 16-bit CPU, ROM, RAM, a UART, ports, timers, and a 10-bit analog-to-digital converter.

Bit-Slice Processors

A second direction of microprocessor evolution has been bit-slice processors. For some applications, general-purpose CPUs such as the 8080 and 6800 are not fast enough or do not have suitable instruction sets. For these applications, several manufacturers produce devices which can be used to build a custom CPU. An example is the Advanced Micro Devices 2900 family of devices. This family includes 4-bit ALUs, multiplexers, sequencers, and other parts needed for custom-building a CPU. The term slice comes from the fact that these parts can be connected in parallel to work with 8-bit words, 16-bit words, or 32-bit words. In other words, a designer can add as many slices as needed for a particu-

lar application. The designer not only custom-designs the hardware of the CPU, but also custom-makes the instruction set for it using microcode.

General-Purpose CPUs

The third major direction of microprocessor evolution has been toward general-purpose CPUs which give a microcomputer most or all of the computing power of earlier minicomputers. After Motorola came out with the MC6800, Intel produced the 8085, an upgrade of the 8080 that required only a +5-V supply. Motorola then produced the MC6809, which has a few 16-bit instructions, but is still basically an 8-bit processor. In 1978 Intel came out with the 8086, which is a full 16-bit processor. Some 16-bit microprocessors, such as the National PACE and the Texas Instruments 9900 family of devices, had been available previously, but the market apparently wasn't ready. Soon after Intel came out with the 8086, Motorola came out with the 16-bit MC68000, and the 16-bit race was off and running. The 8086 and the 68000 work directly with 16-bit words instead of with 8-bit words; they can address a million or more bytes of memory instead of the 64 Kbytes addressable by the 8-bit processors, and they execute instructions much faster than the 8-bit processors. Also, these 16-bit processors have single instructions for functions such as multiply and divide, which required a lengthy sequence of instructions on the 8-bit processors.

The evolution along this last path has continued on to 32-bit processors that work with gigabytes (10^9 bytes) or terabytes (10^{12} bytes) of memory. Examples of these devices are the Intel 80386, the Motorola MC68020, and the National 32032.

Since we could not possibly describe in this book the operation and programming of even a few of the available processors, we confine our discussions primarily to one group of related microprocessors. The family we have chosen is the Intel 8086, 8088, 80186, 80188, 80286, 80386, 80486 family. Members of this family are very widely used in personal computers, business computer systems, and industrial control systems. Our experience has shown that learning the programming and operation of one family of microcomputers very thoroughly is much more useful than looking at many processors superficially. If you learn one processor family well, you will most likely find it quite easy to learn another when you have to.

THE 8086 MICROPROCESSOR FAMILY—OVERVIEW

The Intel 8086 is a 16-bit microprocessor that is intended to be used as the CPU in a microcomputer. The term 16-bit means that its arithmetic logic unit, its internal registers, and most of its instructions are designed to work with 16-bit binary words. The 8086 has a 16-bit data bus, so it can read data from or write data to memory and ports either 16 bits or 8 bits at a time. The 8086 has a 20-bit address bus, so it can address any one of 2^{20} , or 1,048,576, memory locations.

Each of the 1,048,576 memory addresses of the 8086 represents a byte-wide location. Sixteen-bit words will be stored in two consecutive memory locations. If the first byte of a word is at an even address, the 8086 can read the entire word in one operation. If the first byte of the word is at an odd address, the 8086 will read the first byte with one bus operation and the second byte with another bus operation. Later we will discuss this in detail. The main point here is that if the first byte of a 16-bit word is at an even address, the 8086 can read the entire word in one operation.

The Intel 8088 has the same arithmetic logic unit, the same registers, and the same instruction set as the 8086. The 8088 also has a 20-bit address bus, so it can address any one of 1,048,576 bytes in memory. The 8088, however, has an 8-bit data bus, so it can only read data from or write data to memory and ports 8 bits at a time. The 8086, remember, can read or write either 8 or 16 bits at a time. To read a 16-bit word from two successive memory locations, the 8088 will always have to do two read operations. Since the 8086 and the 8088 are almost identical, any reference we make to the 8086 in the rest of the book will also pertain to the 8088 unless we specifically indicate otherwise. This is done to make reading easier. The Intel 8088, incidentally, is used as the CPU in the original IBM Personal Computer, the IBM PC/XT, and several compatible personal computers.

The Intel 80186 is an improved version of the 8086, and the 80188 is an improved version of the 8088. In addition to a 16-bit CPU, the 80186 and 80188 each have programmable peripheral devices integrated in the same package. In a later chapter we will discuss these integrated peripherals. The instruction set of the 80186 and 80188 is a superset of the instruction set of the 8086. The term *superset* means that all the 8086 and 8088 instructions will execute properly on an 80186 or an 80188, but the 80186 and the 80188 have a few additional instructions. In other words, a program written for an 8086 or an 8088 is upward-compatible to an 80186 or an 80188, but a program written for an 80186 or an 80188 may not execute correctly on an 8086 or an 8088. In the instruction set descriptions in Chapter 6, we specifically indicate which instructions work only with the 80186 or 80188.

The Intel 80286 is a 16-bit, advanced version of the 8086 which was specifically designed for use as the CPU in a multilayer or multitasking microcomputer. When operating in its *real address mode*, the 80286 functions mostly as a fast 8086. Most programs written for an 8086 can be run on an 80286 operating in its real address mode. When operating in its *virtual address mode*, an 80286 has features which make it easy to keep users' programs separate from one another and to protect the system program from destruction by users' programs. In Chapter 15 we discuss the operation and use of the 80286. The 80286 is the CPU used in the IBM PC/AT personal computer.

The Intel 80386 is a 32-bit microprocessor which can directly address up to 4 gigabytes of memory. The 80386 contains more sophisticated features than the 80286 for use in multilayer and multitasking microcomputer

systems. In Chapter 15 we discuss the features of the 80386 and the 80486, which is an evolutionary step up from the 80386.

8086 INTERNAL ARCHITECTURE

Before we can talk about how to write programs for the 8086, we need to discuss its specific internal features, such as its ALU, flags, registers, instruction byte queue, and segment registers.

As shown by the block diagram in Figure 2-7, the 8086 CPU is divided into two independent functional parts, the *bus interface unit* or BIU, and the *execution unit* or EU. Dividing the work between these two units speeds up processing.

The BIU sends out addresses, fetches instructions from memory, reads data from ports and memory, and writes data to ports and memory. In other words, the BIU handles all transfers of data and addresses on the buses for the execution unit.

The execution unit of the 8086 tells the BIU where to fetch instructions or data from, decodes instructions, and executes instructions. Let's take a look at some of the parts of the execution unit.

The Execution Unit

CONTROL CIRCUITRY, INSTRUCTION DECODER, AND ALU

As shown in Figure 2-7, the EU contains control circuitry which directs internal operations. A decoder in the EU translates instructions fetched from memory into a series of actions which the EU carries out. The EU has a 16-bit *arithmetic logic unit* which can add, subtract, AND, OR, XOR, increment, decrement, complement, or shift binary numbers.

FLAG REGISTER

A *flag* is a flip-flop which indicates some condition produced by the execution of an instruction or controls certain operations of the EU. A 16-bit *flag register* in the EU contains nine active flags. Figure 2-8 shows the location of the nine flags in the flag register. Six of the nine flags are used to indicate some condition produced by an instruction. For example, a flip-flop called the *carry flag* will be set to a 1 if the addition of two 16-bit binary numbers produces a carry out of the most significant bit position. If no carry out of the MSB is produced by the addition, then the carry flag will be a 0. The EU thus effectively runs up a "flag" to tell you that a carry was produced.

The six conditional flags in this group are the *carry flag* (CF), the *parity flag* (PF), the *auxiliary carry flag* (AF), the *zero flag* (ZF), the *sign flag* (SF), and the *overflow flag* (OF). The names of these flags should give you hints as to what conditions affect them. Certain 8086 instructions check these flags to determine which of two alternative actions should be done in executing the instruction.

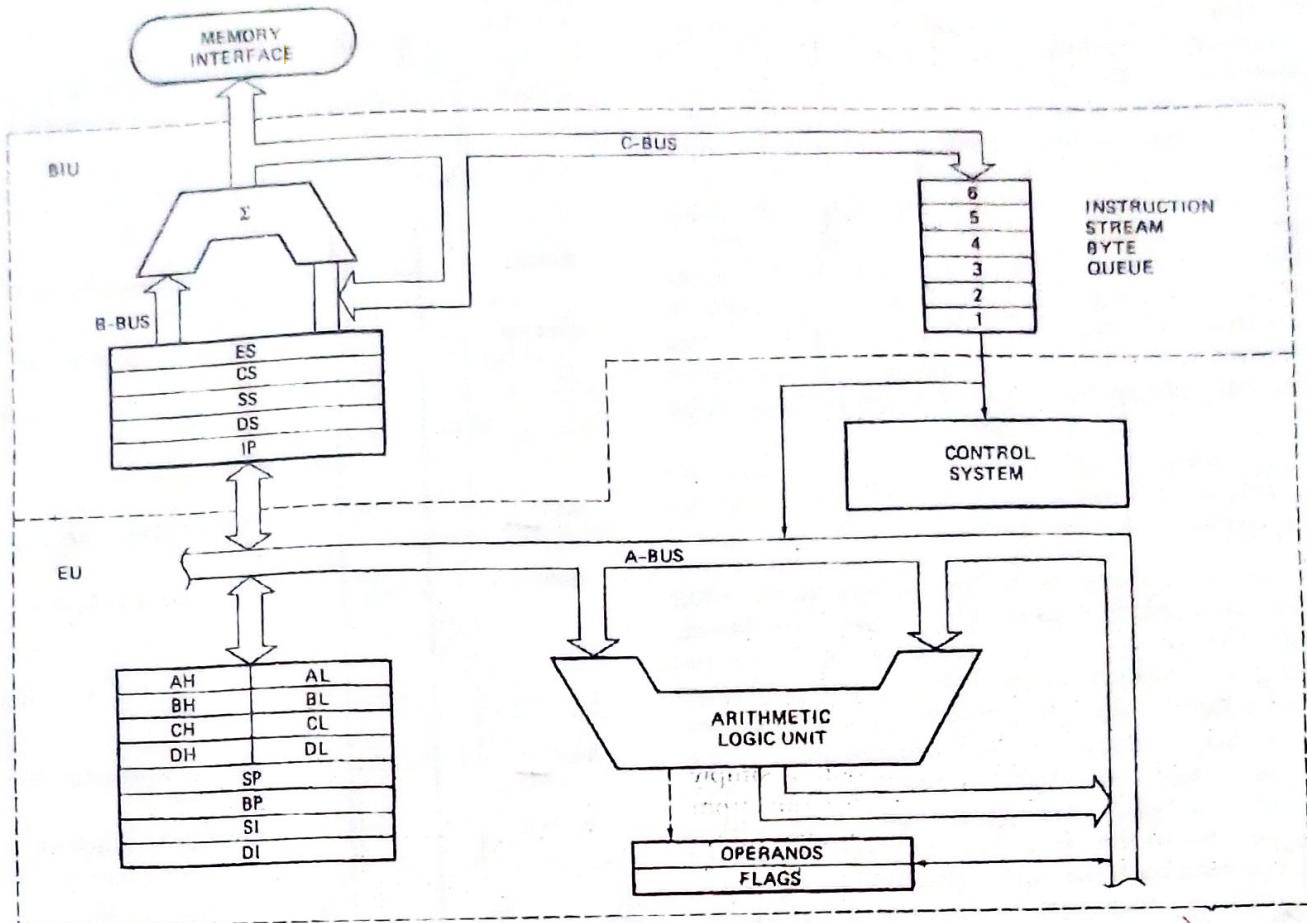


FIGURE 2-7 8086 internal block diagram. (Intel Corp.)

The three remaining flags in the flag register are used to control certain operations of the processor. These flags are different from the six conditional flags described above in the way they are set or reset. The six conditional flags are set or reset by the EU on the basis of the results of some arithmetic or logic operation. The control flags, are deliberately set or reset with specific instructions you put in your program. The three control flags are the trap flag (TF), which is used for single stepping through a program; the interrupt flag (IF), which is used to allow or prohibit the interruption of a program; and the direction flag (DF), which is used with string instructions.

Later we will discuss in detail the operation and use of the nine flags.

GENERAL-PURPOSE REGISTERS

Observe in Figure 2-7 that the EU has eight general-purpose registers, labeled AH, AL, BH, BL, CH, CL, DH, and DL. These registers can be used individually for temporary storage of 8-bit data. The AL register is also called the accumulator. It has some features that the other general-purpose registers do not have.

Certain pairs of these general-purpose registers can be used together to store 16-bit data words. The acceptable

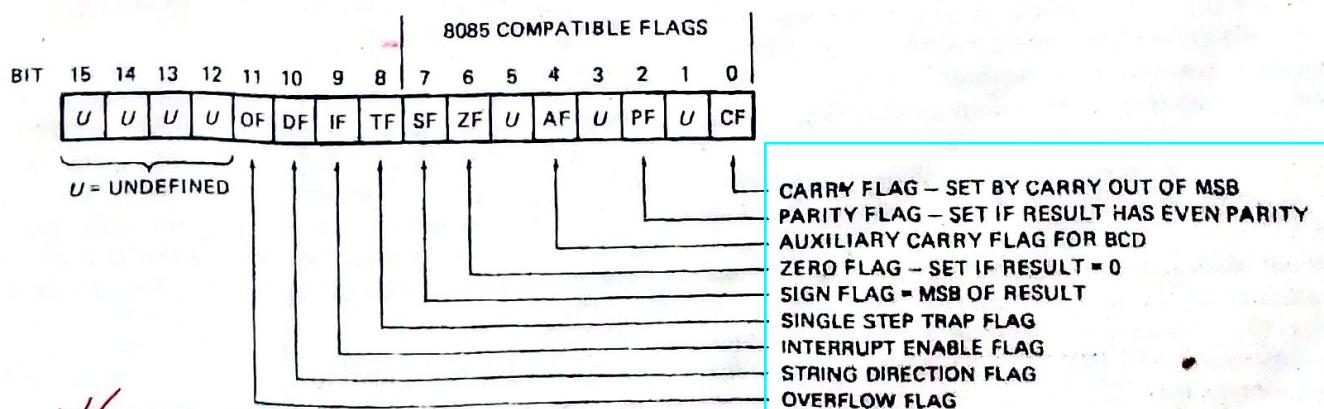


FIGURE 2-8 8086 flag register format. (Intel Corp.)

register pairs are AH and AL, BH and BL, CH and CL, and DH and DL. The AH–AL pair is referred to as the AX register, the BH–BL pair is referred to as the BX register, the CH–CL pair is referred to as the CX register, and the DH–DL pair is referred to as the DX register.

The 8086 general-purpose register set is very similar to those of the earlier-generation 8080 and 8085 microprocessors. It was designed this way so that the many programs written for the 8080 and 8085 could easily be translated to run on the 8086 or the 8088. The advantage of using internal registers for the temporary storage of data is that, since the data is already in the EU, it can be accessed much more quickly than it could be accessed in external memory. Now let's look at the features of the BIU.

The BIU

THE QUEUE

While the EU is decoding an instruction or executing an instruction which does not require use of the buses,

the BIU fetches up to six instruction bytes for the following instructions. The BIU stores these prefetched bytes in a first-in-first-out register set called a queue.

When the EU is ready for its next instruction, it simply reads the instruction byte(s) for the instruction from the queue in the BIU. This is much faster than sending out an address to the system memory and waiting for memory to send back the next instruction byte or bytes. The process is analogous to the way a bricklayer's assistant fetches bricks ahead of time and keeps a queue of bricks lined up so that the bricklayer can just reach out and grab a brick when necessary. Except in the cases of JMP and CALL instructions, where the queue must be dumped and then reloaded starting from a new address, this prefetch-and-queue scheme greatly speeds up processing. Fetching the next instruction while the current instruction executes is called *pipelining*.

SEGMENT REGISTERS

The 8086 BIU sends out 20-bit addresses, so it can address any of 2^{20} or 1,048,576 bytes in memory. However, at any given time the 8086 works with only four 65,536-byte (64-Kbyte) segments within this 1,048,576-byte (1-Mbyte) range. Four segment registers in the BIU are used to hold the upper 16 bits of the starting addresses of four memory segments that the 8086 is working with at a particular time. The four segment registers are the code segment (CS) register, the stack segment (SS) register, the extra segment (ES) register, and the data segment (DS) register.

Figure 2-9 shows how these four segments might be positioned in memory at a given time. The four segments can be separated as shown, or, for small programs which do not need all 64 Kbytes in each segment, they can overlap.

To repeat, then, a segment register is used to hold the upper 16 bits of the starting address for each of the segments. The code segment register, for example, holds the upper 16 bits of the starting address for the segment from which the BIU is currently fetching instruction code bytes. The BIU always inserts zeros for the lowest

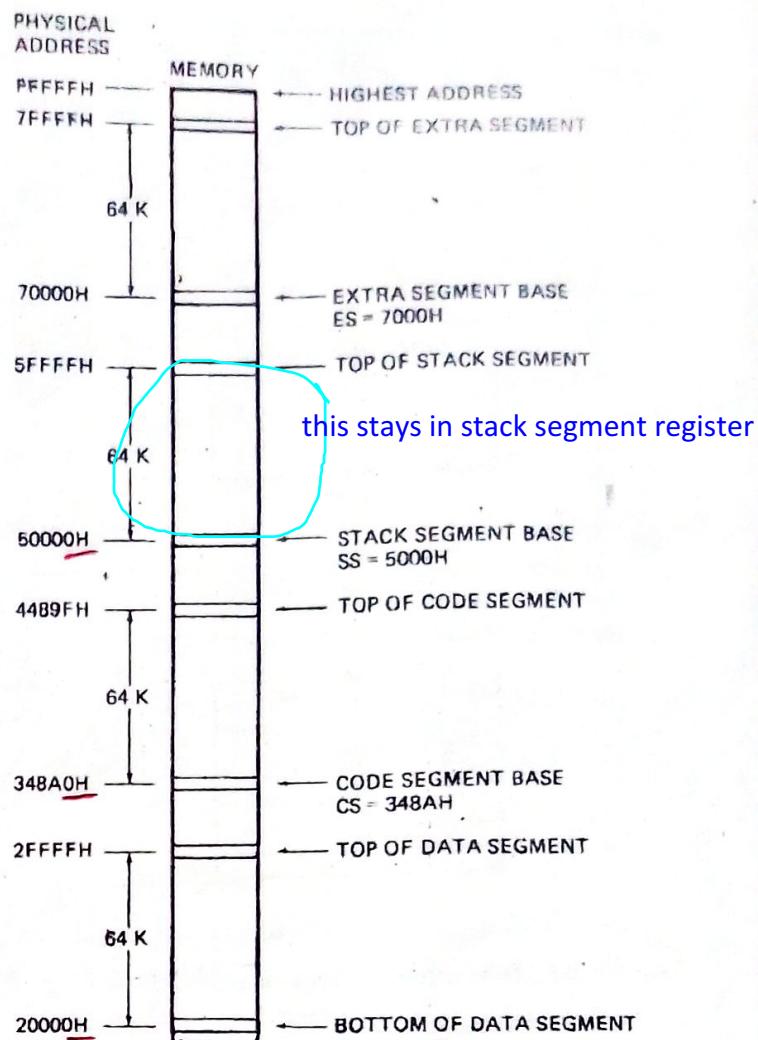


FIGURE 2-9 One way four 64-Kbyte segments might be positioned within the 1-Mbyte address space of an 8086.

4 bits (nibble) of the 20-bit starting address for a segment. If the code segment register contains 348AH, for example, then the code segment will start at address 348AOH. In other words, a 64-Kbyte segment can be located anywhere within the 1-Mbyte address space, but the segment will always start at an address with zeros in the lowest 4 bits. This constraint was put on the location of segments so that it is only necessary to store and manipulate 16-bit numbers when working with the starting address of a segment. The part of a segment starting address stored in a segment register is often called the *segment base*.

A *stack* is a section of memory set aside to store addresses and data while a *subprogram* executes. The stack segment register is used to hold the upper 16 bits of the starting address for the program stack. We will discuss the use and operation of a stack in detail later.

The extra segment register and the data segment register are used to hold the upper 16 bits of the starting addresses of two memory segments that are used for data.

INSTRUCTION POINTER

The next feature to look at in the BIU is the instruction pointer (IP) register. As discussed previously, the code

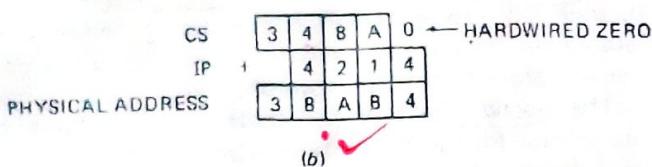
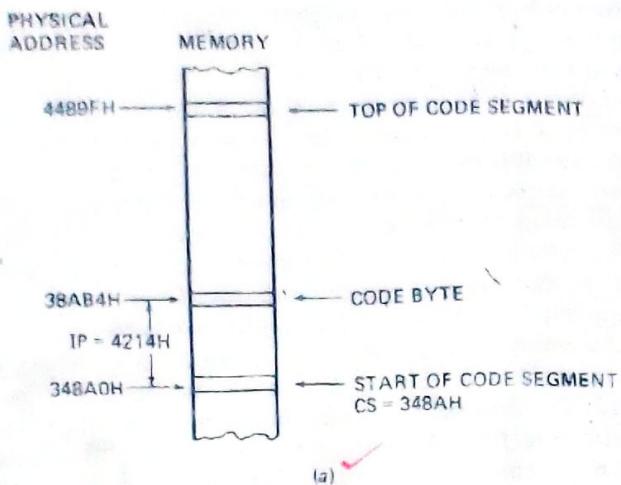


FIGURE 2-10 Addition of IP to CS to produce the physical address of the code byte. (a) Diagram. (b) Computation.

segment register holds the upper 16 bits of the starting address of the segment from which the BIU is currently fetching instruction code bytes. The instruction pointer register holds the 16-bit address, or offset, of the next code byte within this code segment. The value contained in the IP is referred to as an offset because this value must be offset from (added to) the segment base address in CS to produce the required 20-bit physical address sent out by the BIU. Figure 2-10a shows in diagram form how this works. The CS register points to the base, or start of the current code segment. The IP contains the distance or offset from this base address to the next instruction byte to be fetched. Figure 2-10b shows how the 16-bit offset in IP is added to the 16-bit segment base address in CS to produce the 20-bit physical address. Notice that the two 16-bit numbers are not added directly in line, because the CS register contains only the upper 16 bits of the base address for the code segment. As we said before, the BIU automatically inserts zeros for the lowest 4 bits of the segment base address.

If the CS register, for example, contains 348AH, you know that the starting address for the code segment is 348A0H. When the BIU adds the offset of 4214H in the IP to this segment base address, the result is a 20-bit physical address of 38AB4H.

An alternative way of representing a 20-bit physical address is the segment base offset form. For the address of a code byte, the format for this alternative form will be CS:IP. As an example of this, the address constructed in the preceding paragraph, 38AB4H, can also be represented as 348A:4214.

To summarize, then, the CS register contains the upper 16 bits of the starting address of the code segment in the 1-Mbyte address range of the 8086. The instruction pointer register contains a 16-bit offset which

tells where in that 64-Kbyte code segment the next instruction byte is to be fetched from. The actual physical address sent to memory is produced by adding the offset contained in the IP register to the segment base represented by the upper 16 bits in the CS register.

Any time the 8086 accesses memory, the BIU produces the required 20-bit physical address by adding an offset to a segment base value represented by the contents of one of the segment registers. As another example of this, let's look at how the 8086 uses the contents of the stack segment register and the contents of the stack pointer register to produce a physical address.

STACK SEGMENT REGISTER AND STACK POINTER REGISTER

A stack, remember, is a section of memory set aside to store addresses and data while a subprogram is executing. The 8086 allows you to set aside an entire 64-Kbyte segment as a stack. The upper 16 bits of the starting address for this segment are kept in the stack segment register. The stack pointer (SP) register in the execution unit holds the 16-bit offset from the start of the segment to the memory location where a word was most recently stored on the stack. The memory location where a word was most recently stored is called the top of stack. Figure 2-11a shows this in diagram form.

The physical address for a stack read or a stack write is produced by adding the contents of the stack pointer register to the segment base address represented by the upper 16 bits of the base address in SS. Figure 2-11b shows an example. The 5000H in SS represents a segment base address of 50000H. When the FFE0H in the SP is added to this, the resultant physical address for the top of the stack will be 5FFEOH. The physical address can be represented either as a single number, 5FFEOH, or in SS:SP form as 5000:FFEOH.

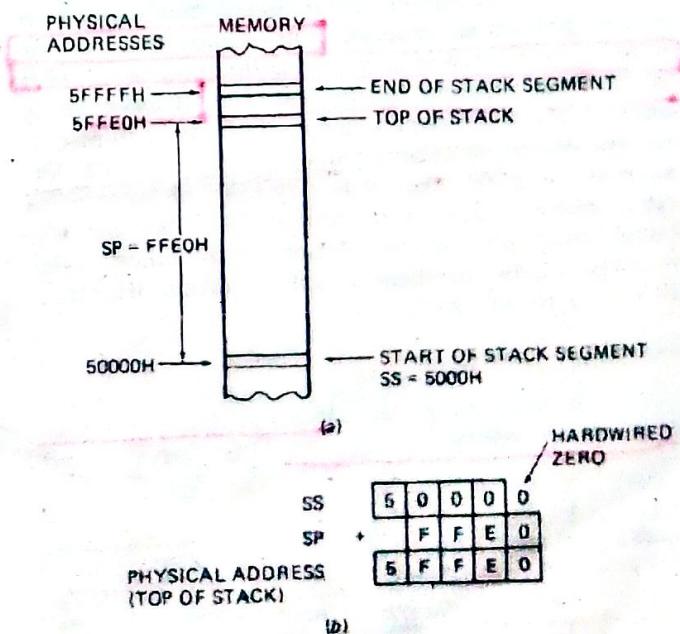


FIGURE 2-11 Addition of SS and SP to produce the physical address of the top of the stack. (a) Diagram. (b) Computation.

The operation and use of the stack will be discussed in detail later as need arises.

POINTER AND INDEX REGISTERS IN THE EXECUTION UNIT

In addition to the stack pointer register (SP), the EU contains a 16-bit base pointer (BP) register. It also contains a 16-bit source index (SI) register and a 16-bit destination index (DI) register. These three registers can be used for temporary storage of data just as the general-purpose registers described above. However, their main use is to hold the 16-bit offset of a data word in one of the segments. SI, for example, can be used to hold the offset of a data word in the data segment. The physical address of the data in memory will be generated in this case by adding the contents of SI to the segment base address represented by the 16-bit number in the DS register. After we give you an overview of the different levels of languages used to program a microcomputer, we will show you some examples of how we tell the 8086 to read data from or write data to a desired memory location.

INTRODUCTION TO PROGRAMMING THE 8086

Programming Languages

Now that you have an overview of the 8086 CPU, it is time to start thinking about how it is programmed. To run a program, a microcomputer must have the program stored in binary form in successive memory locations, as shown in Figure 2-12. There are three language levels that can be used to write a program for a microcomputer.

MACHINE LANGUAGE

You can write programs as simply a sequence of the binary codes for the instructions you want the microcomputer to execute. The three-instruction program in Figure 2-6b is an example. This binary form of the program is referred to as machine language because it is the form required by the machine. However, it is difficult, if not impossible, for a programmer to memorize the thousands of binary instruction codes for a CPU such as the 8086. Also, it is very easy for an error to occur when working with long series of 1's and 0's. Using hexadecimal representation for the binary codes might help some, but there are still thousands of instruction codes to cope with.

ASSEMBLY LANGUAGE

To make programming easier, many programmers write programs in assembly language. They then translate

the assembly language program to machine language so that it can be loaded into memory and run. Assembly language uses two-, three-, or four-letter mnemonics to represent each instruction type. A mnemonic is just a device to help you remember something. The letters in an assembly language mnemonic are usually initials or a shortened form of the English word(s) for the operation performed by the instruction. For example, the mnemonic for subtract is SUB, the mnemonic for Exclusive OR is XOR, and the mnemonic for the instruction to copy data from one location to another is MOV.

Assembly language statements are usually written in a standard form that has four fields, as shown in Figure 2-12. The first field in an assembly language statement is the label field. A label is a symbol or group of symbols used to represent an address which is not specifically known at the time the statement is written. Labels are usually followed by a colon. Labels are not required in a statement, they are just inserted where they are needed. We will show later many uses of labels.

The opcode field of the instruction contains the mnemonic for the instruction to be performed. Instruction mnemonics are sometimes called operation codes, or opcodes. The ADD mnemonic in the example statement in Figure 2-12 indicates that we want the instruction to do an addition.

The operand field of the statement contains the data, the memory address, the port address, or the name of the register on which the instruction is to be performed. Operand is just another name for the data item(s) acted on by an instruction. In the example instruction in Figure 2-12, there are two operands. AL and 07H, specified in the operand field. AL represents the AL register, and 07H represents the number 07H. This assembly language statement thus says, "Add the number 07H to the contents of the AL register." By Intel convention, the result of the addition will be put in the register or the memory location specified before the comma in the operand field. For the example statement in Figure 2-12, then, the result will be left in the AL register. As another example, the assembly language statement ADD BH, AL, when converted to machine language and run, will add the contents of the AL register to the contents of the BH register. The results will be left in the BH register.

The final field in an assembly language statement such as that in Figure 2-12 is the comment field, which starts with a semicolon. Comments do not become part of the machine language program, but they are very important. You write comments in a program to remind you of the function that an instruction or group of instructions performs in the program.

To summarize why assembly language is easier to use than machine language, let's look a little more closely at the assembly language ADD statement. The general format of the 8086 ADD instruction is

ADD destination, source

The source can be a number written in the instruction, the contents of a specified register, or the contents of a memory location. The destination can be a specified register or a specified memory location. However, the

LABEL FIELD	OP CODE FIELD	OPERAND FIELD	COMMENT FIELD
NEXT:	ADD	AL, 07H	; ADD CORRECTION FACTOR

FIGURE 2-12 Assembly language program statement format.

source and the destination in an instruction cannot both be memory locations.

A later section on 8086 addressing modes will show all the ways in which the source of an operand and the destination of the result can be specified. The point here is that the single mnemonic ADD, together with a specified source and a specified destination, can represent a great many 8086 instructions in an easily understandable form.

The question that may occur to you at this point is, "If I write a program in assembly language, how do I get it translated into machine language which can be loaded into the microcomputer and executed?" There are two answers to this question. The first method of doing the translation is to work out the binary code for each instruction a bit at a time using the templates given in the manufacturer's data books. We will show you how to do this in the next chapter, but it is a tedious and error-prone task. The second method of doing the translation is with an assembler. An assembler is a program which can be run on a personal computer or microcomputer development system. It reads the file of assembly language instructions you write and generates the correct binary code for each. For developing all but the simplest assembly language programs, an assembler and other program development tools are essential. We will introduce you to these program development tools in the next chapter and describe their use throughout the rest of this book.

HIGH-LEVEL LANGUAGES

Another way of writing a program for a microcomputer is with a *high-level language*, such as BASIC, Pascal, or C. These languages use program statements which are even more English-like than those of assembly language. Each high-level statement may represent many machine code instructions. An *interpreter program* or a *compiler program* is used to translate higher-level language statements to machine codes which can be loaded into memory and executed. Programs can usually be written faster in high-level languages than in assembly language because the high-level language works with bigger building blocks. However, programs written in a high-level language and interpreted or compiled almost always execute more slowly and require more memory than the same programs written in assembly language. Programs that involve a lot of hardware control, such as robots and factory control systems, or programs that must run as quickly as possible are usually best written in assembly language. Complex data processing programs that manipulate massive amounts of data, such as insurance company records, are usually best written in a high-level language. The decision concerning which language to use has recently been made more difficult by the fact that current assemblers allow the use of many high-level language features, and the fact that some current high-level languages provide assembly language features.

OUR CHOICE

For most of this book we work very closely with hardware, so assembly language is the best choice. In later chap-

ters, however, we do show you how to write programs which contain modules written in assembly language and modules written in the high-level language C. In the next chapter we introduce you to assembly language programming techniques. Before we go on to that, however, we will use a few simple 8086 instructions to show you more about accessing data in registers and memory locations.

How the 8086 Accesses Immediate and Register Data

In a previous discussion of the 8086 BIU, we described how the 8086 accesses code bytes using the contents of the CS and IP registers. We also described how the 8086 accesses the stack using the contents of the SS and SP registers. Before we can teach you assembly language programming techniques, we need to discuss some of the different ways in which an 8086 can access the data that it operates on. The different ways in which a processor can access data are referred to as its *addressing modes*. In assembly language statements, the addressing mode is indicated in the instruction. We will use the 8086 MOV instruction to illustrate some of the 8086 addressing modes.

The MOV instruction has the format

MOV destination, source

When executed, this instruction copies a word or a byte from the specified source location to the specified destination location. The source can be a number written directly in the instruction, a specified register, or a memory location specified in 1 of 24 different ways. The destination can be a specified register or a memory location specified in any 1 of 24 different ways. The source and the destination cannot both be memory locations in an instruction.

IMMEDIATE ADDRESSING MODE

Suppose that in a program you need to put the number 437BH in the CX register. The MOV CX, 437BH instruction can be used to do this. When it executes, this instruction will put the *immediate* hexadecimal number 437BH in the 16-bit CX register. This is referred to as *immediate addressing mode* because the number to be loaded into the CX register will be put in the two memory locations immediately following the code for the MOV instruction. This is similar to the way the port address was put in memory immediately after the code for the input instruction in the three-instruction program in Figure 2-6b.

A similar instruction, MOV CL, 48H, could be used to load the 8-bit immediate number 48H into the 8-bit CL register. You can also write instructions to load an 8-bit immediate number into an 8-bit memory location or to load a 16-bit number into two consecutive memory locations, but we are not yet ready to show you how to specify these.

REGISTER ADDRESSING MODE

Register addressing mode means that a register is the source of an operand for an instruction. The instruction

`MOV CX, AX`, for example, copies the contents of the 16-bit AX register into the 16-bit CX register. Remember that the destination location is specified in the instruction before the comma, and the source is specified after the comma. Also note that the contents of AX are just copied to CX, not actually moved. In other words, the previous contents of CX are written over, but the contents of AX are not changed. For example, if CX contains 2A84H and AX contains 4971H before the `MOV CX, AX` instruction executes, then after the instruction executes, CX will contain 4971H and AX will still contain 4971H. You can `MOV` any 16-bit register to any 16-bit register, or you can `MOV` any 8-bit register to any 8-bit register. However, you cannot use an instruction such as `MOV CX, AL` because this is an attempt to copy a byte-type operand (AL) into a word-type destination (CX). The byte in AL would fit in CX, but the 8086 would not know which half of CX to put it in. If you try to write an instruction like this and you are using a good assembler, the assembler will tell you that the instruction contains a type error. To copy the byte from AL to the high byte of CX, you can use the instruction `MOV CH, AL`. To copy the byte from AL to the low byte of CX, you can use the instruction `MOV CL, AL`.

Accessing Data in Memory

OVERVIEW OF MEMORY ADDRESSING MODES

The addressing modes described in the following sections are used to specify the location of an operand in memory. To access data in memory, the 8086 must also produce a 20-bit physical address. It does this by adding a 16-bit value called the effective address to a segment base address represented by the 16-bit number in one of the four segment registers. The effective address (EA) represents the displacement or offset of the desired operand from the segment base. In most cases, any of the segment bases can be specified, but the data segment is the one most often used. Figure 2-13a shows in graphic form how the EA is added to the data segment base to point to an operand in memory. Figure 2-13b shows how the 20-bit physical address is generated by the BIU. The starting address for the data segment in Figure 2-13b is 20000H, so the data segment register will contain 2000H. The BIU adds the effective address, 437AH, to the data segment base address of 20000H to produce the physical address sent out to memory. The 20-bit physical address sent out to memory by the BIU will then be 2437AH. The physical address can be represented either as a single number 2437AH or in the segment base:offset form as 2000:437AH.

The execution unit calculates the effective address for an operand using information you specify in the instruction. You can tell the EU to use a number in the instruction as the effective address, to use the contents of a specified register as the effective address, or to compute the effective address by adding a number in the instruction to the contents of one or two specified registers. The following section describes one way you can tell the execution unit to calculate an effective address. In later chapters we show other ways of specifying the effective address. Later we also show how the

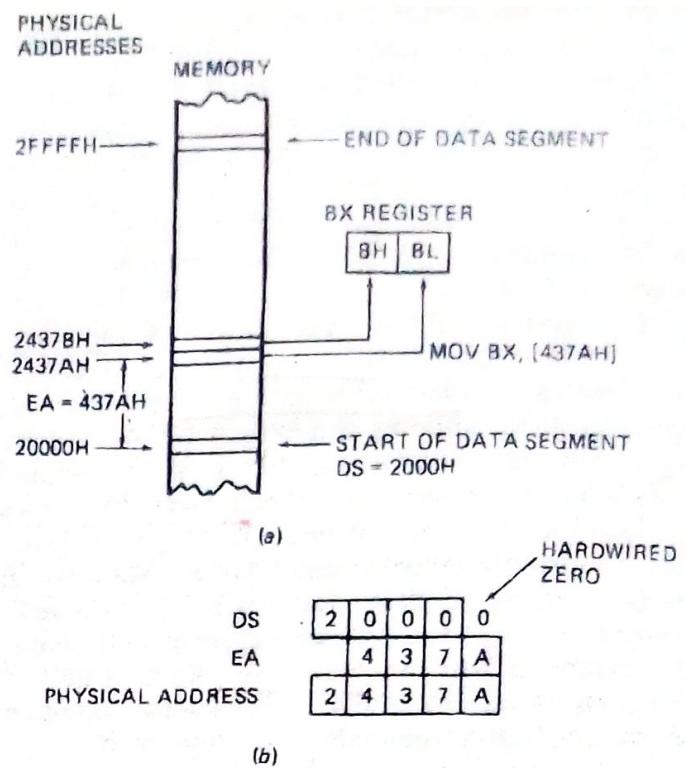


FIGURE 2-13 Addition of data segment register and effective address to produce the physical address of the data byte. (a) Diagram. (b) Computation.

addressing modes this provides are used to solve some common programming problems.

DIRECT ADDRESSING MODE

For the simplest memory addressing mode, the effective address is just a 16-bit number written directly in the instruction. The instruction `MOV BL, [437AH]` is an example. The square brackets around the 437AH are shorthand for "the contents of the memory location(s) at a displacement from the segment base of." When executed, this instruction will copy "the contents of the memory location at a displacement from the data segment base of" 437AH into the BL register, as shown by the rightmost arrow in Figure 2-13a. The BIU calculates the 20-bit physical memory address by adding the effective address 437AH to the data segment base, as shown in Figure 2-13b. This addressing mode is called direct because the displacement of the operand from the segment base is specified directly in the instruction. The displacement in the instruction will be added to the data segment base in DS unless you tell the BIU to add it to some other segment base. Later we will show you how to do this.

Another example of the direct addressing mode is the instruction `MOV BX, [437AH]`. When executed, this instruction copies a 16-bit word from memory into the BX register. Since each memory address of the 8086 represents a byte of storage, the word must come from two memory locations. The byte at a displacement of 437AH from the data segment base will be copied into BL, as shown by the right arrow in Figure 2-13a. The contents of the next higher address, displacement 437BH, will be copied into the BH register, as shown by the left arrow in Figure 2-13a. From the instruction

coding, the 8086 will automatically determine the number of bytes that it must access in memory.

An important point here is that an 8086 always stores the low byte of a word in the lower of the two addresses and stores the high byte of a word in the higher address. To stick this in your mind, remember:

Low byte-low address, high byte-high address

The previous two examples showed how the direct addressing mode can be used to specify the source of an operand. Direct addressing can also be used to specify the destination of an operand in memory. The instruction `MOV [437AH], BX`, for example, will copy the contents of the `BX` register to two memory locations in the data segment. The contents of `BL` will be copied to the memory location at a displacement of `437AH`. The contents of `BH` will be copied to the memory location at a displacement of `437BH`. This operation is represented by simply reversing the direction of the arrows in Figure 2-13a.

NOTE: When you are hand-coding programs using direct addressing of the form shown above, make sure to put in the square brackets to remind you how to code the instruction. If you leave the brackets out of an instruction such as `MOV BX, [437AH]`, you will code it as if it were the instruction `MOV BX, 437AH`. This second instruction will load the immediate number `437AH` into `BX`, rather than loading a word from memory at a displacement of `437AH` into `BX`. Also note that if you are writing an instruction using direct addressing such as this for an assembler, you must write the instruction in the form `MOV BL, DS:BYTE PTR [437AH]` to give the assembler all the information it needs. As we will show you in the next chapter, when you are using an assembler, you usually use a name to represent the direct address rather than the actual numerical value.

A FEW WORDS ABOUT SEGMENTATION

At this point you may be wondering why Intel designed the 8086 family devices to access memory using the segment:offset approach rather than accessing memory directly with 20-bit addresses. The segment:offset scheme requires only a 16-bit number to represent the base address for a segment, and only a 16-bit offset to access any location in a segment. This means that the 8086 has to manipulate and store only 16-bit quantities instead of 20-bit quantities. This makes for an easier interface with 8- and 16-bit-wide memory boards and with the 16-bit registers in the 8086.

The second reason for segmentation has to do with the type of microcomputer in which an 8086-family CPU is likely to be used. A previous section of this chapter described briefly the operation of a timesharing microcomputer system. In a timesharing system, several users share a CPU. The CPU works on one user's program for perhaps 20 ms, then works on the next user's program for 20 ms. After working 20 ms for each of the other users, the CPU comes back to the first user's program

again. Each time the CPU switches from one user's program to the next, it must access a new section of code and new sections of data. Segmentation makes this switching quite easy. Each user's program can be assigned a separate set of logical segments for its code and data. The user's program will contain offsets or displacements from these segment bases. To change from one user's program to a second user's program, all that the CPU has to do is to reload the four segment registers with the segment base addresses assigned to the second user's program. In other words, segmentation makes it easy to keep users' programs and data separate from one another, and segmentation makes it easy to switch from one user's program to another user's program. In Chapter 15 we tell you much more about the use of segmentation in multiuser systems.

CHECKLIST OF IMPORTANT TERMS AND CONCEPTS IN THIS CHAPTER

If you do not remember any of the terms or concepts in the following list, use the index to find them in the chapter.

Microcomputer, microprocessor

Hardware, software, firmware

Timesharing computer system

Multitasking computer system

Distributed processing system

Multiprocessing

CPU

Memory, RAM, ROM

I/O ports

Address, data, and control buses

Control bus signals

ALU

Segmentation

Bus interface unit (BIU)

Instruction byte queue, pipelining,
ES, CS, SS, DS registers, IP register

Execution unit (EU)

AX, BX, CX, DX registers, flag register,
ALU, SP, BP, SI, DI registers

Machine language, assembly language, high-level language

Mnemonic, opcode, operand, label, comment

Assembler, compiler

Immediate address mode, register address mode, direct address mode

Effective address

REVIEW QUESTIONS AND PROBLEMS

1. Describe the main advantages of a distributed processing computer system over a simple time-sharing system.
2. Describe the sequence of signals that occurs on the address bus, the control bus, and the data bus when a simple microcomputer fetches an instruction.
3. What determines whether a microprocessor is considered an 8-bit, a 16-bit, or a 32-bit device?
4.
 - a. How many address lines does an 8086 have?
 - b. How many memory addresses does this number of address lines allow the 8086 to access directly?
 - c. At any given time, the 8086 works with four segments in this address space. How many bytes are contained in each segment?
5. What is the main difference between the 8086 and the 8088? (?)
6.
 - a. Describe the function of the 8086 queue.
 - b. How does the queue speed up processing?
7.
 - a. If the code segment for an 8086 program starts at address 70400H, what number will be in the CS register?
 - b. Assuming this same code segment base, what physical address will a code byte be fetched from if the instruction pointer contains 539CH?
8. What physical address is represented by:
 - a. 4370:561EH
 - b. 7A32:0028H
9. What is the advantage of using a CPU register for temporary data storage over using a memory location? (?)
10. If the stack segment register contains 3000H and the stack pointer register contains 8434H, what is the physical address of the top of the stack?
11.
 - a. What is the advantage of using assembly language instead of writing a program directly in machine language?
 - b. Describe the operation an 8086 will perform when it executes ADD AX, BX.
12. What types of programs are usually written in assembly language?
13. Describe the operation that an 8086 will perform when it executes each of the following instructions:
 - a. MOV BX, 03FFH
 - b. MOV AL, 0DBH
 - c. MOV DH, CL
 - d. MOV BX, AX
14. Write the 8086 assembly language statement which will perform the following operations:
 - a. Load the number 7986H into the BP register.
 - b. Copy the BP register contents to the SP register.
 - c. Copy the contents of the AX register to the DS register.
 - d. Load the number F3H into the AL register.
15. If the 8086 execution unit calculates an effective address of 14A3H and DS contains 7000H, what physical address will the BIU produce?
16. If the data segment register (DS) contains 4000H, what physical address will the instruction MOVAL, [234BH] read?
17. If the 8086 data segment register contains 7000H, write the instruction that will copy the contents of DL to address 74B2CH.
18. Describe the difference between the instructions MOV AX, 2437H and MOV AX, [2437H].