



COMPUTER SYSTEM ARCHITECTURE

Lecture Slides

Dr. Nisha Chaurasia

WHAT IS COMPUTER SYSTEM ARCHITECTURE

- Computer architecture is concerned with the structure and behavior of the various functional modules of the computer and how they interact to provide the processing needs of the user.
- Computer organization is concerned with the way the hardware components are connected together to form a computer system.
- Computer design is concerned with the development of the hardware for the computer taking into consideration a given set of specifications.

CHAPTER-1

DIGITAL LOGIC CIRCUITS

- Introduces the fundamental knowledge needed for the design of digital systems constructed with individual gates and flip-flops.
- It covers Boolean algebra, combinational circuits, and sequential circuits, providing necessary background for understanding the digital circuits.

- The digital computer is a digital system that performs various computational tasks.
- The word digital implies that the information in the computer is represented by variables that take a limited number of discrete values.
- For e.g., the decimal digits 0, 1, 2, ..., 9, for example, provide 10 discrete values.
- In practice, digital computers function more reliably if only two states are used. Because of the physical restriction of components, and because human logic tends to be binary (i.e., true/false, yes/no statements), digital components that are constrained to take discrete values are further constrained to take only two values and are said to be binary.

- Digital computers use the binary number system, which has two digits: 0 and 1. A binary digit is called a bit.
- Information is represented in digital computers in groups of bits.
- By using various coding techniques, groups of bits can be made to represent not only binary numbers but also other discrete symbols, such as decimal digits or letters of the alphabet.

- In contrast to the common decimal numbers that employ the base 10 system, binary numbers use a base 2 system with two digits: 0 and 1.
- For example, the binary number 1001011 represents a quantity that can be converted to a decimal number by multiplying each bit by the base 2 raised to an integer power as follows:

$$1 \times 2^6 + 0 \times 2^5 + 0 \times 2^4 + 1 \times 2^3 + 0 \times 2^2 + 1 \times 2^1 + 1 \times 2^0 = 75$$

- Hence, $(1001011)_2 = 75_{10}$

PROGRAM

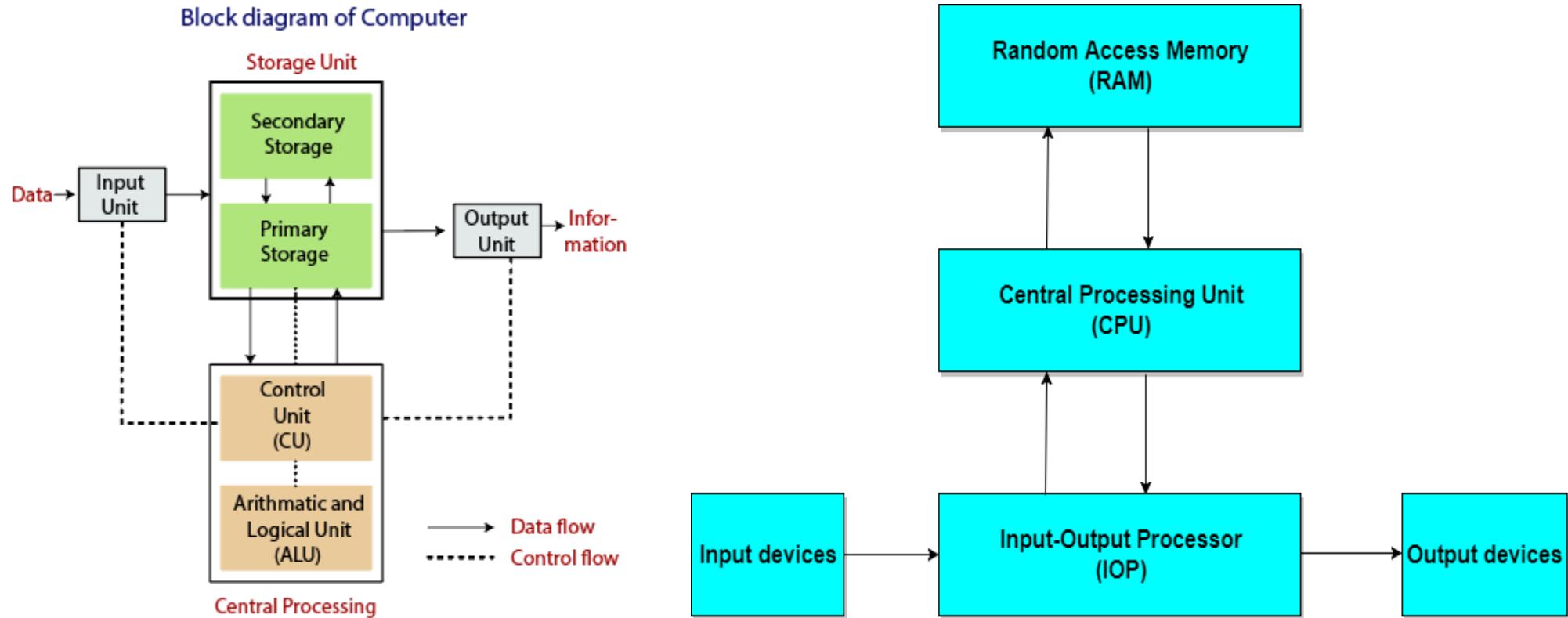
- A computer system is sometimes subdivided into two functional entities: hardware and software.
- The hardware of the computer consists of all the electronic components and electromechanical devices that comprise the physical entity of the device.
- Computer software consists of the instructions and data that the computer manipulates to perform various data-processing tasks.
- A sequence of instructions for the computer is called a program.
- The data that are manipulated by the program constitute the data base.

OPERATING SYSTEM

- The programs included in a systems software package are referred to as the operating system.
- They are distinguished from application programs written by the user for the purpose of solving particular problems.
- For example, a high-level language program written by a user to solve particular data-processing needs is an application program, but the compiler that translates the high-level language program to machine language is a system program.

COMPUTER HARDWARE

- The hardware of the computer is usually divided into three major parts:
- The **central processing unit (CPU)** contains an arithmetic and logic unit for manipulating data, a number of registers for storing data, and control circuits for fetching and executing instructions.
- The **memory** of a computer contains storage for instructions and data. It is called a random access memory (RAM) because the CPU can access any location in memory at random and retrieve the binary information within a fixed interval of time.
- The **input and output processor (IOP)** contains electronic circuits for communicating and controlling the transfer of information between the computer and the outside world. The input and output devices connected to the computer include keyboards, printers, terminals, magnetic disk drives, and other communication devices.



COMPUTER ORGANIZATION

- Computer organization is concerned with the way the hardware components operate and the way they are connected together to form the computer system.
- The various components are assumed to be in place and the task is to investigate the organizational structure to verify that the computer parts operate as intended.

COMPUTER DESIGN

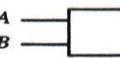
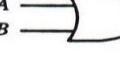
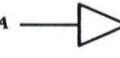
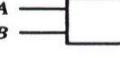
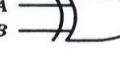
- Computer design is concerned with the hardware design of the computer.
- Once the computer specifications are formulated, it is the task of the designer to develop hardware for the system.
- Computer design is concerned with the determination of what hardware should be used and how the parts should be connected. This aspect of computer hardware is sometimes referred to as computer implementation.

COMPUTER ARCHITECTURE

- Computer architecture is concerned with the structure and behavior of the computer as seen by the user.
- It includes the information formats, the instruction set, and techniques for addressing memory.
- The architectural design of a computer system is concerned with the specifications of the various functional modules, such as processors and memories, and structuring them together into a computer system.

LOGIC GATES

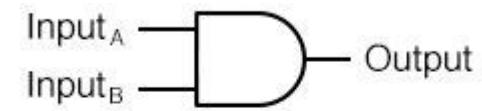
- Binary logic deals with binary variables and with operations that assume a logical meaning.
- It is used to describe, in algebraic or tabular form, the manipulation and processing of binary information.
- The manipulation of binary information is done by logic circuits called **gates**.
- Gates are blocks of hardware that produce signals of binary 1 or 0 when input logic requirements are satisfied.
- Each gate has a distinct graphic symbol and its operation can be described by means of an algebraic expression.
- The input-output relationship of the binary variables for each gate can be represented in tabular form by a **truth table**.

Name	Graphic symbol	Algebraic function	Truth table															
AND		$x = A \cdot B$ or $x = AB$	<table border="1"> <thead> <tr> <th>A</th><th>B</th><th>x</th></tr> </thead> <tbody> <tr> <td>0</td><td>0</td><td>0</td></tr> <tr> <td>0</td><td>1</td><td>0</td></tr> <tr> <td>1</td><td>0</td><td>0</td></tr> <tr> <td>1</td><td>1</td><td>1</td></tr> </tbody> </table>	A	B	x	0	0	0	0	1	0	1	0	0	1	1	1
A	B	x																
0	0	0																
0	1	0																
1	0	0																
1	1	1																
OR		$x = A + B$	<table border="1"> <thead> <tr> <th>A</th><th>B</th><th>x</th></tr> </thead> <tbody> <tr> <td>0</td><td>0</td><td>0</td></tr> <tr> <td>0</td><td>1</td><td>1</td></tr> <tr> <td>1</td><td>0</td><td>1</td></tr> <tr> <td>1</td><td>1</td><td>1</td></tr> </tbody> </table>	A	B	x	0	0	0	0	1	1	1	0	1	1	1	1
A	B	x																
0	0	0																
0	1	1																
1	0	1																
1	1	1																
Inverter		$x = A'$	<table border="1"> <thead> <tr> <th>A</th><th>x</th></tr> </thead> <tbody> <tr> <td>0</td><td>1</td></tr> <tr> <td>1</td><td>0</td></tr> </tbody> </table>	A	x	0	1	1	0									
A	x																	
0	1																	
1	0																	
Buffer		$x = A$	<table border="1"> <thead> <tr> <th>A</th><th>x</th></tr> </thead> <tbody> <tr> <td>0</td><td>0</td></tr> <tr> <td>1</td><td>1</td></tr> </tbody> </table>	A	x	0	0	1	1									
A	x																	
0	0																	
1	1																	
NAND		$x = (AB)'$	<table border="1"> <thead> <tr> <th>A</th><th>B</th><th>x</th></tr> </thead> <tbody> <tr> <td>0</td><td>0</td><td>1</td></tr> <tr> <td>0</td><td>1</td><td>1</td></tr> <tr> <td>1</td><td>0</td><td>1</td></tr> <tr> <td>1</td><td>1</td><td>0</td></tr> </tbody> </table>	A	B	x	0	0	1	0	1	1	1	0	1	1	1	0
A	B	x																
0	0	1																
0	1	1																
1	0	1																
1	1	0																
NOR		$x = (A + B)'$	<table border="1"> <thead> <tr> <th>A</th><th>B</th><th>x</th></tr> </thead> <tbody> <tr> <td>0</td><td>0</td><td>1</td></tr> <tr> <td>0</td><td>1</td><td>0</td></tr> <tr> <td>1</td><td>0</td><td>0</td></tr> <tr> <td>1</td><td>1</td><td>0</td></tr> </tbody> </table>	A	B	x	0	0	1	0	1	0	1	0	0	1	1	0
A	B	x																
0	0	1																
0	1	0																
1	0	0																
1	1	0																
Exclusive-OR (XOR)		$x = A \oplus B$ or $x = A'B + AB'$	<table border="1"> <thead> <tr> <th>A</th><th>B</th><th>x</th></tr> </thead> <tbody> <tr> <td>0</td><td>0</td><td>0</td></tr> <tr> <td>0</td><td>1</td><td>1</td></tr> <tr> <td>1</td><td>0</td><td>1</td></tr> <tr> <td>1</td><td>1</td><td>0</td></tr> </tbody> </table>	A	B	x	0	0	0	0	1	1	1	0	1	1	1	0
A	B	x																
0	0	0																
0	1	1																
1	0	1																
1	1	0																
Exclusive-NOR or equivalence		$x = (A \oplus B)'$ or $x = A'B' + AB$	<table border="1"> <thead> <tr> <th>A</th><th>B</th><th>x</th></tr> </thead> <tbody> <tr> <td>0</td><td>0</td><td>1</td></tr> <tr> <td>0</td><td>1</td><td>0</td></tr> <tr> <td>1</td><td>0</td><td>0</td></tr> <tr> <td>1</td><td>1</td><td>1</td></tr> </tbody> </table>	A	B	x	0	0	1	0	1	0	1	0	0	1	1	1
A	B	x																
0	0	1																
0	1	0																
1	0	0																
1	1	1																

AND Gate

- The AND gate produces the AND logic function: that is, the output is 1 if input A and input B are both equal to 1; otherwise, the output is 0.
- These conditions are also specified in the truth table for the AND gate. The table shows that output x is 1 only when both input A and input B are 1.
- The algebraic operation symbol of the AND function is the same as the multiplication symbol (.) of ordinary arithmetic.
- AND gates may have more than two inputs, and by definition, the output is 1 if and only if all inputs are 1.

2 - input AND gate

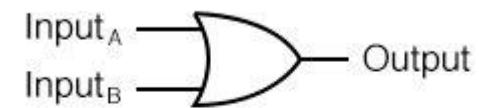


A	B	Output
0	0	0
0	1	0
1	0	0
1	1	1

OR Gate

- The OR gate produces the inclusive-OR function; that is, the output is 1 if input A or input B or both inputs are 1; otherwise, the output is 0.
- The algebraic symbol of the OR function is +, similar to arithmetic addition.
- OR gates may have more than two inputs, and by definition, the output is 1 if any input is 1.

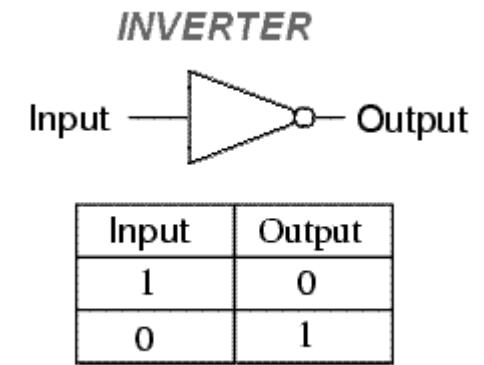
2 - input OR gate



A	B	Output
0	0	0
0	1	1
1	0	1
1	1	1

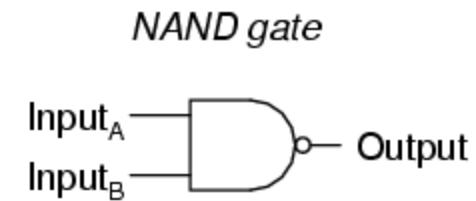
Inverter (NOT)

- The inverter circuit inverts the logic sense of a binary signal.
- It produces the NOT, or complement function.
- The algebraic symbol used for the logic complement is either a prime or a bar over the variable symbol.



NAND Gate

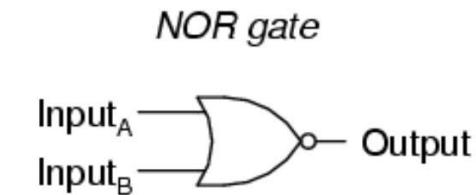
- This is a NOT-AND gate which is equal to an AND gate followed by a NOT gate.
- The outputs of all NAND gates are high if any of the inputs are low.
- The NAND function is the complement of the AND function, as indicated by the graphic symbol, which consists of an AND graphic symbol followed by a small circle.



A	B	Output
0	0	1
0	1	1
1	0	1
1	1	0

NOR Gate

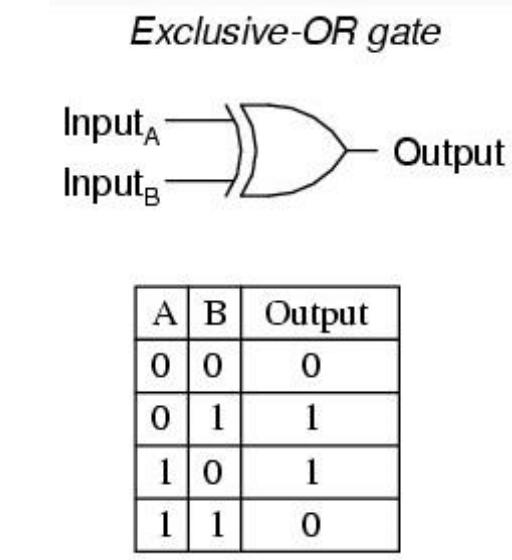
- This is a NOT-OR gate which is equal to an OR gate followed by a NOT gate.
- The outputs of all NOR gates are low if any of the inputs are high.
- The NOR gate is the complement of the OR gate and uses an OR graphic symbol followed by a small circle.



A	B	Output
0	0	1
0	1	0
1	0	0
1	1	0

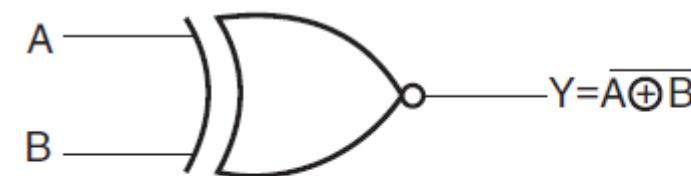
Exclusive-OR (XOR)

- It produces a high output if either, but not both, of its two inputs are high.
- The exclusive-OR gate has a graphic symbol similar to the OR gate except for the additional curved line on the input side.
- The output of this gate is 1 if any input is 1 but excludes the combination when both inputs are 1.



Exclusive-NOR (XNOR)

- The 'Exclusive-NOR' gate circuit does the opposite to the XOR gate.
- It will give a low output if either, but not both, of its two inputs are high.



$$Y = (\overline{A \oplus B}) = (A \cdot B + \overline{A} \cdot \overline{B})$$

A	B	Y
0	0	1
0	1	0
1	0	0
1	1	1

BOOLEAN ALGEBRA

- Boolean algebra is an algebra that deals with binary variables and logic operations.
- This forms the algebraic expression showing the operation of the logic circuit for each input variable either True or False that results in a logic “1” output.
- The three basic logic operations are AND, OR, and complement (NOT).
- The purpose of Boolean algebra is to facilitate the analysis and design of digital circuits. It provides a convenient tool to:
 - Express in algebraic form a truth table relationship between binary variables.
 - Express in algebraic form the input-output relationship of logic diagrams.
 - Find simpler circuits for the same function.

BOOLEAN FUNCTION

- Boolean function can be expressed algebraically with binary variables, the logic operation symbols, parentheses, and equal sign.
- For a given value of the variables, the Boolean function can be either 1 or 0.

E.g.,

$$F = x + y' z$$

The function F is equal to 1 if x is 1 or if both y' and z are equal to 1; F is equal to 0 otherwise. But saying that $y' = 1$ is equivalent to saying that $y = 0$ since y' is the complement of y. Therefore, we may say that F is equal to 1 if $x = 1$ or if $yz = 01$.

TRUTH TABLE

- The relationship between a function and its binary variables can be represented in a **truth table**.
- A truth table shows how the truth or falsity of a compound statement depends on the truth or falsity of the simple statements from which it's constructed.
- It defines the function of a logic gate by providing a concise list that shows all the output states in tabular form for each possible combination of input variable that the gate could encounter.
- To represent a function in a truth table we need a list of the 2^n combinations of the **n** binary variables.

LOGIC DIAGRAM

- Boolean function can be transformed from an algebraic expression into a logic diagram composed of AND, OR, and inverter/NOT gates.
- The logic diagram consists of gates and symbols that can directly replace an expression in Boolean arithmetic.
- This is a graphical representation of a logic circuit that shows the wiring and connections of each individual logic gate, represented by a specific graphical symbol that implements the logic circuit.

BOOLEAN ALGEBRA (REVISITED)

- A Boolean function specified by a truth table can be expressed algebraically in many different ways.
- **Boolean algebra** is the branch of algebra in which the values of the variables are the truth values true and false, usually denoted 1 and 0 respectively.
- By manipulating a Boolean expression according to Boolean algebra rules, one may obtain a simpler expression that will require fewer gates.
- It is used to analyze and simplify the digital (logic) circuits.

	AND Form	OR Form
Commutative Law	$A \cdot B = B \cdot A$	$A + B = B + A$
Associate Law	$(A \cdot B) \cdot C = A \cdot (B \cdot C)$	$(A + B) + C = A + (B + C)$
Distributive Law	$(A+B)+C = (A+C).(B+C)$	$(A + B) \cdot C = (A \cdot C) + (B \cdot C)$
Identity Law	$A \cdot 1 = A$	$A + 0 = A$
Zero and One Law	$A \cdot 0 = 0$	$A + 1 = 1$
Inverse Law	$A \cdot A' = 0$	$A + A' = 1$
Idempotent Law	$A \cdot A = A$	$A + A = A$
Absorption Law	$A(A+B) = A$	$A + A \cdot B = A$ $A + A'B = A+B$
DeMorgan's Law	$(A \cdot B)' = (A)' + (B)'$	$(A + B)' = (A)' \cdot (B)'$
Double Complement Law	$\overline{\overline{X}} = X$	

Source: Google Search

NUMBER SYSTEM

- The number system is a way to represent or express numbers.
- Based on the different symbol used to represent numbers, there are various types of number systems:
 - The decimal number system
 - The binary number system
 - The octal number system and
 - The hexadecimal number system
 - Binary Coded Decimal or BCD Numbering System

CONVERTING FROM BINARY TO DECIMAL

- decimal = $d_0 \times 2^0 + d_1 \times 2^1 + d_2 \times 2^2 + \dots$
- Example

binary number:	1	1	1	0	0	1
power of 2:	2^5	2^4	2^3	2^2	2^1	2^0

$$111001_2 = 1 \cdot 2^5 + 1 \cdot 2^4 + 1 \cdot 2^3 + 0 \cdot 2^2 + 0 \cdot 2^1 + 1 \cdot 2^0 = 57_{10}$$

CONVERTING FROM DECIMAL INTEGER TO BINARY

- Conversion steps:
 - Divide the number by 2.
 - Get the integer quotient for the next iteration.
 - Get the remainder for the binary digit.
 - Repeat the steps until the quotient is equal to 0.
- Example

Convert 13_{10} to binary :

Division by 2	Quotient	Remainder	Bit #
13/2	6	1	0
6/2	3	0	1
3/2	1	1	2
1/2	0	1	3

$$13_{10} = 1101_2$$

CONVERTING DECIMAL FRACTION TO BINARY

- Binary Fractions use the same weighting principle as decimal numbers except that each binary digit uses the base-2 numbering system.
- A decimal number representation of $(0.XY)_{10}$ can be converted into base of 2 and represented by $(0. a_1 , a_2 , a_3 , \text{etc.})_2$.
- The fraction number is multiplied by 2, the result of integer part is a_1 and fraction part multiply by 2, and then separate integer part from fraction, the integer part represents a_1 ; this process continues until the fraction becomes 0.

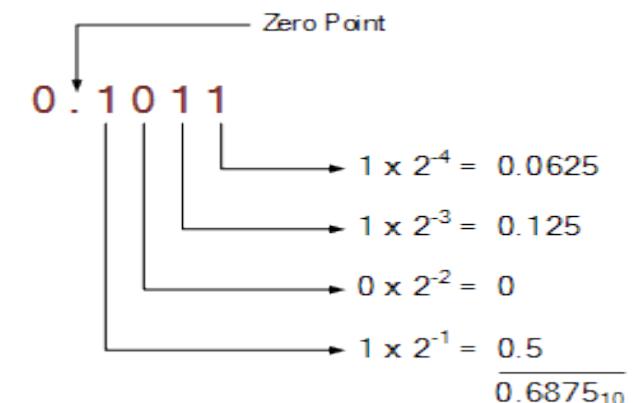
SOURCE: INTERNET

Example: $(0.625)_{10}$

	Integer	Fraction	Coefficient
$0.625 * 2 =$	1	. 25	$a_1 = 1$
$0.25 * 2 =$	0	. 5	$a_2 = 0$
$0.5 * 2 =$	1	. 0	$a_3 = 1$

Answer: $(0.625)_{10} = (0.a_1 a_2 a_3)_2 = (0.101)_2$

↑ ↑
MSB LSB

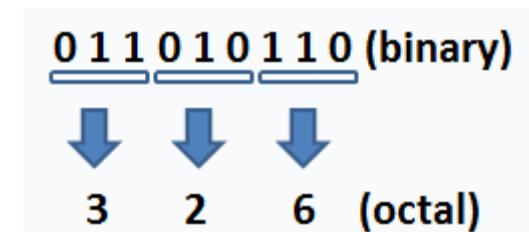
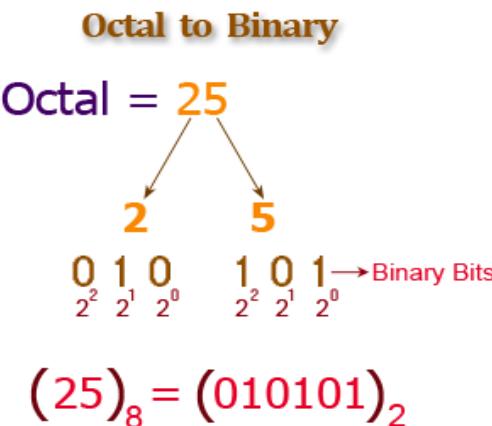


CONVERTING FROM OCTAL TO BINARY

- The octal numeral system, or oct for short, is the base-8 number system, and uses the digits 0 to 7.
- Octal numerals can be made from binary numerals by grouping consecutive binary digits into groups of three (starting from the right).

Decimal	Octal	Binary
0	0	0
1	1	1
2	2	10
3	3	11
4	4	100
5	5	101
6	6	110
7	7	111

SOURCE: INTERNET



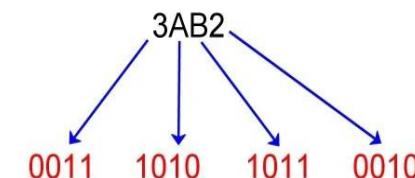
CONVERTING FROM HEX TO BINARY

- Hexadecimal (also base 16, or hex) is a positional numeral system with a radix, or base, of 16. It uses sixteen distinct symbols, most often the symbols 0–9 to represent values zero to nine, and A, B, C, D, E, F (or alternatively a, b, c, d, e, f) to represent values ten to fifteen.

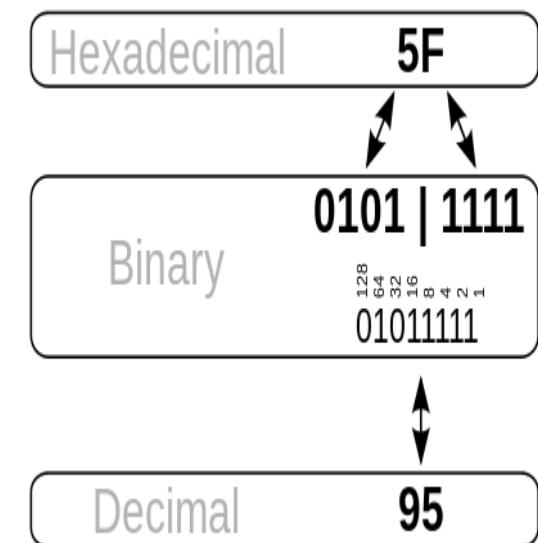
Decimal	Hexadecimal	Binary
0	0	0
1	1	1
2	2	10
3	3	11
4	4	100
5	5	101
6	6	110
7	7	111
8	8	1000
9	9	1001
10	A	1010
11	B	1011
12	C	1100
13	D	1101
14	E	1110
15	F	1111

SOURCE: INTERNET

Converting Hex to Binary



$$3AB2_{16} = 11101010110010_2$$



BINARY CODED DECIMAL (BCD)

- BCD or Binary Coded Decimal is that number system or code which has the binary numbers or digits to represent a decimal number where each digit is represented by a fixed number of binary bits, usually between four and eight.
- A decimal number contains 10 digits (0-9). Now the equivalent binary numbers can be found out of these 10 decimal numbers. In case of BCD the binary number formed by four binary digits, will be the equivalent code for the given decimal digits. In BCD we can use the binary number from 0000-1001 only, which are the decimal equivalent from 0-9 respectively.
- The BCD_{8421} code is so called because each of the four bits is given a ‘weighting’ according to its column value in the binary system. The least significant bit (lsb) has the weight or value 1, the next bit, going left, the value 2. The next bit has the value 4, and the most significant bit (msb) the value 8. E.g.,
 24_{10} in 8 bit binary would be 00011000 but in BCD_{8421} is 0010 0100.
 992_{10} in 16 bit binary would be 000000111100000 but in BCD_{8421} is 1001 1001 0010.

Decimal	BCD
0	0 0 0 0
1	0 0 0 1
2	0 0 1 0
3	0 0 1 1
4	0 1 0 0
5	0 1 0 1
6	0 1 1 0
7	0 1 1 1
8	1 0 0 0
9	1 0 0 1
-	1 0 1 0
-	1 0 1 1
-	1 1 0 0
-	1 1 0 1
-	1 1 1 0
-	1 1 1 1

Unused

Binary-Coded Decimal vs. Binary to Decimal Conversion

SOURCE: [HTTPS://REALPARS.COM/BCD/](https://realpars.com/bcd/)

Decimal Number	BCD	Binary
0	0000 0000 0000 0000	0000 0000 0000 0000
1	0000 0000 0000 0001	0000 0000 0000 0001
2	0000 0000 0000 0010	0000 0000 0000 0010
3	0000 0000 0000 0011	0000 0000 0000 0011
4	0000 0000 0000 0100	0000 0000 0000 0100
5	0000 0000 0000 0101	0000 0000 0000 0101
6	0000 0000 0000 0110	0000 0000 0000 0110
7	0000 0000 0000 0111	0000 0000 0000 0111
8	0000 0000 0000 1000	0000 0000 0000 1000
9	0000 0000 0000 1001	0000 0000 0000 1001
...		
9620	1001 0110 0010 0000	0010 0101 1001 0100
120	0001 0010 0000	0000 0000 0111 1000
→ 4568	0100 0101 0110 1000	0001 0001 1101 1000

REALPARS

Binary Base-2	Decimal Base-10	Hexa- Decimal Base-16	Octal Base-8	BCD Code
0000	0	0	0	0
0001	1	1	1	1
0010	2	2	2	2
0011	3	3	3	3
0100	4	4	4	4
0101	5	5	5	5
0110	6	6	6	6
0111	7	7	7	7
1000	8	8	10	8
1001	9	9	11	9
1010	10	A	12	---
1011	11	B	13	---
1100	12	C	14	---
1101	13	D	15	---
1110	14	E	16	---
1111	15	F	17	---

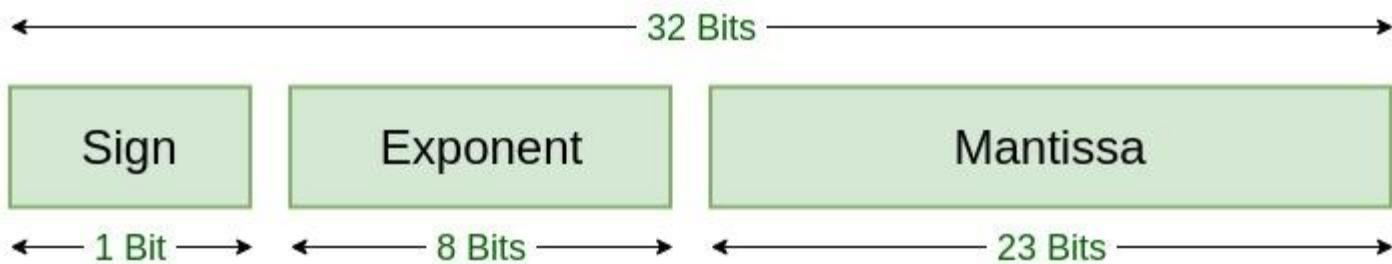
SOURCE: INTERNET

FLOATING POINT REPRESENTATION

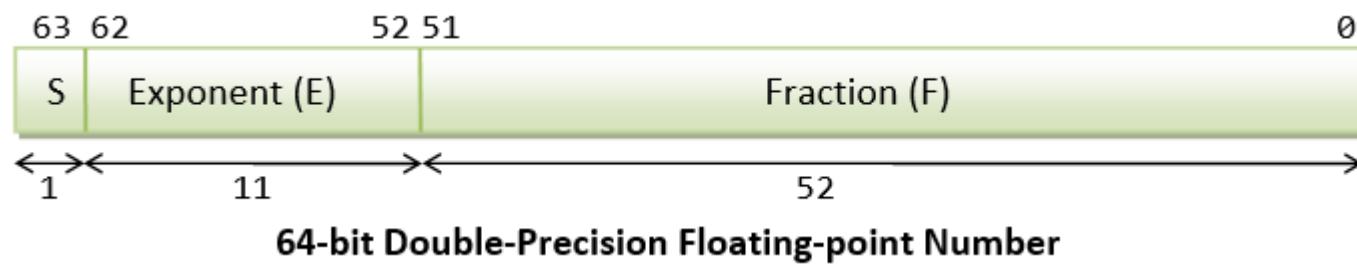
- Integers are whole numbers without fractional components. 1, 2, and 3 are integers, while 0.1, 2.2, and 3.0001 all have fractional components are called floating point numbers.
- The floating point unit performs floating point operations. Floating point numbers have a sign, a mantissa, and an exponent.
- The central processing unit (CPU) typically consists of an arithmetic logic unit (ALU), floating point unit (FLU/FPU), registers, control unit, and the cache memory.
- The arithmetic logic unit performs integer arithmetic operations such as addition, subtraction, and logic operations such as AND, OR, XOR, etc.

IEEE STANDARD FOR FLOATING POINT REPRESENTATION

- The Institute of Electrical and Electronics Engineers (IEEE) developed a standard to represent floating point numbers, referred to as IEEE 754.
- This standard defines a format for both single (32-bit) and double (64-bit) precision floating point numbers.
- Floating point numbers in single precision are represented by 32 bits while in order to increase the accuracy of a floating point number, IEEE 745 offers double precision represented by 64 bits.
- Decimal floating points are represented by $M \times 10^E$, where M is the signed mantissa (normalized mantissa) and E is the exponent (biased exponent).



Single Precision IEEE 754 Floating-Point Standard



- **Biased Exponent** is the exponent + 127 (01111111_2); therefore, the exponent is represented by a positive number.
- **Normalized Mantissa/mantissa** is represented by $1.M$, where M is called normalized mantissa; if $M = 00101$, then mantissa is 1.00101.
- Example:

Find normalized mantissa and biased exponent of $(111.0000111)_2$.

111.0000111 can be written in the form of $1.110000111 * 2^{10}$

Where

$$M = 110000111$$

$$\text{Biased exponent} = 10 + 01111111 = 10000001$$

The representation of 111.0000111 in single precision is

1bit	8 bits	23 bits
0	10000001	11000011100000000000000

- Represent 5.75 in IEEE 745 single precision.

$$-15.625 = (1111.101)_2$$

$$-1111.101 = -1.11101101 * 2^{11}$$

$$S = 1$$

$$\text{Normalized mantissa} = 0.11101101.$$

$$\text{Biased exponent} = 11 + 0111111 = 10000010.$$

IEEE745 single precision is

1 10000010 111011010000000000000000.

As examples, our 12-bit floating-point number with a binary representation of:

Sign Bit	Exponent					Mantissa						
1	0	0	0	0	1	0	0	0	0	0	0	1

converts to its decimal via $(-1)^1 \times 2^{1-7} \times 0.1000001 = -1 \times 2^{-6} \times 0.1000001 = -0.000001000001_2 = -(2^{-6} + 2^{-12})_{10} = -(1/64 + 1/4096) = -65/4096 = -0.015869140625$. While the floating-point number:

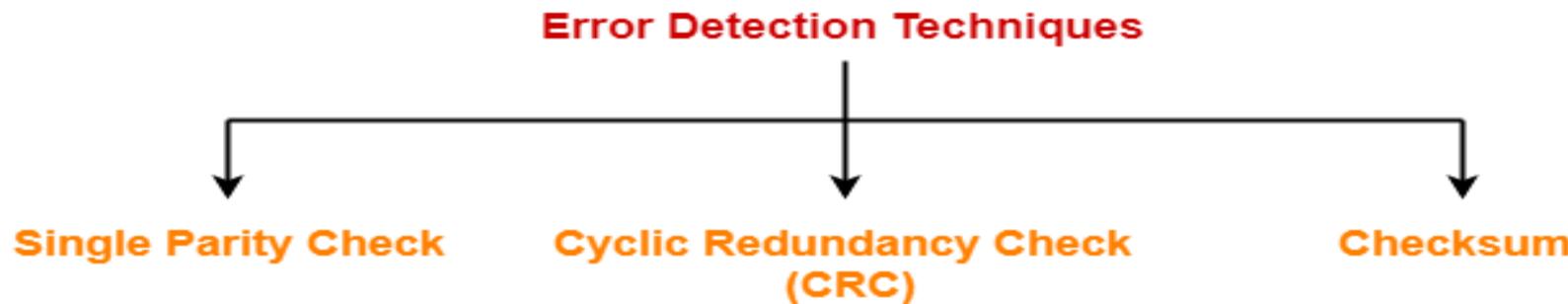
Sign Bit	Exponent					Mantissa						
0	0	1	1	0	1	0	1	0	1	0	0	1

converts to its decimal via $(-1)^0 \times 2^{6-7} \times 1.1010101 = (1 \times 2^{-1} \times 1.1010101)_2 = 0.11010101_2 = (2^{-1} + 2^{-2} + 2^{-4} + 2^{-6} + 2^{-8})_{10} = 1/2 + 1/4 + 1/16 + 1/64 + 1/256 = 213/256 = 0.83203125$.

ERROR DETECTION & CORRECTION CODES

ERROR

- The sequence of bits is called as “Data stream”.
- The change in position of single bit also leads to catastrophic (major) error in data output.
- The error detection and correction techniques are used to get the exact or approximate output.
- In a data sequence, if 1 is changed to zero or 0 is changed to 1, it is called “Bit error”.
- There are generally 3 types of errors occur in data transmission from transmitter to receiver. They are
 - **Single bit errors** (The change in one bit in the whole data sequence; occurs in parallel communication system)
 - **Multiple bit errors** (If there is change in two or more bits of data sequence of transmitter to receiver; occurs in both serial type and parallel type data communication networks)
 - **Burst errors** (The change of set of bits in data sequence; calculated in from the first bit change to last bit change; occurs in serial communication and they are difficult to solve)



PARITY (Vertical Redundancy Check (VRC))

A parity bit is used for error detection of information, since a bit or bits may be changed during the transmission of information from source to destination, a parity bit is an extra bit appended to the information. It represents whether the number of ones or zeroes is either even or odd in the original transmission and can alert the destination to a loss of information.

- Even Parity

The extra bit (0 or 1) is chosen such that the number of ones becomes even.

- Odd Parity

The extra bit (0 or 1) is chosen such that the number of ones becomes odd.

3 bit data			Message with even parity		Message with odd parity	
A	B	C	Message	Parity	Message	Parity
0	0	0	000	0	000	1
0	0	1	001	1	001	0
0	1	0	010	1	010	0
0	1	1	011	0	011	1
1	0	0	100	1	100	0
1	0	1	101	0	101	1
1	1	0	110	0	110	1
1	1	1	111	1	111	0

Original Data	Even Parity	Odd Parity
00000000	0	1
01011011	1	0
01010101	0	1
11111111	0	1
10000000	1	0
01001001	1	0

Cyclic Redundancy Check (CRC)

- CRC is commonly used to detect accidental changes to data transmitted via telecommunications networks and storage devices.
- A cyclic code is a linear (n, k) block code with the property that every cyclic shift of a codeword results in another code word. Here k indicates the length of the message at transmitter (the number of information bits). n is the total length of the message after adding check bits. (actual data and the check bits). n, k is the number of check bits. The codes used for cyclic redundancy check there by error detection are known as CRC codes (Cyclic redundancy check codes). Cyclic redundancy-check codes are shortened cyclic codes.
- CRC involves binary division of the data bits being sent by a predetermined divisor agreed upon by the communicating system. The divisor is generated using polynomials. So, CRC is also called polynomial code checksum.

$$\begin{array}{r}
 1011 \overline{)1101000} \\
 \quad\quad\quad 1111 \leftarrow k \\
 \underline{1011} \\
 \quad\quad\quad 1100 \\
 \underline{1011} \\
 \quad\quad\quad 1110 \\
 \underline{1011} \\
 \quad\quad\quad 1010 \\
 \underline{1011} \\
 \quad\quad\quad 001 \leftarrow \text{remainder}
 \end{array}$$

101101110000
1101

011001110000
1101

000011110000
1101

000000100000
1101

000000010100
1101

000000001110
1101

000000000011

Senders Side

$$\begin{array}{r}
 111101 \\
 1101 \quad | \quad \boxed{100100000} \\
 \oplus \quad 1101 \\
 \hline
 01000 \\
 \oplus 1101 \\
 \hline
 01010 \\
 \oplus 1101 \\
 \hline
 001100 \\
 \oplus 1101 \\
 \hline
 \underline{\underline{0001}}
 \end{array}$$

CRC bits → 001

$$\begin{aligned}
 \text{Transmitted bits} &= \text{Original Message} + \text{CRC bits} \\
 &= 100100000 + 001 \\
 &= 100100001
 \end{aligned}$$

⊕ represents bitwise XOR

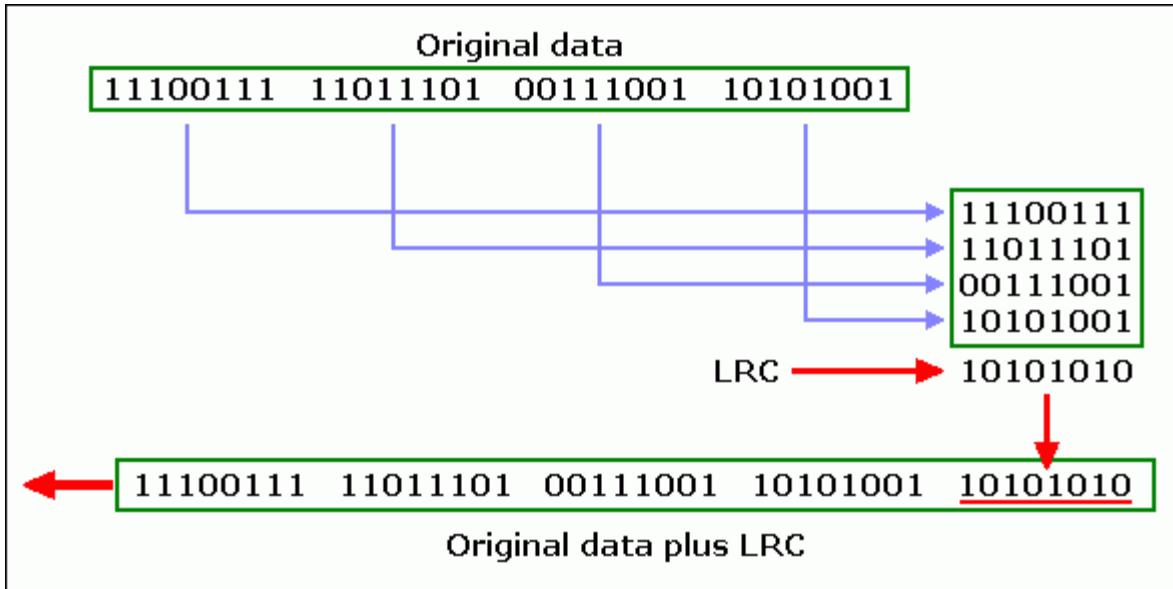
Receivers Side

$$\begin{array}{r}
 111101 \\
 1101 \quad | \quad \boxed{100100001} \\
 \oplus \quad 1101 \\
 \hline
 01000 \\
 \oplus 1101 \\
 \hline
 01010 \\
 \oplus 1101 \\
 \hline
 01110 \\
 \oplus 1101 \\
 \hline
 001101 \\
 \oplus 1101 \\
 \hline
 \underline{\underline{0000}}
 \end{array}$$

Remainder is
zero, So data is
accepted

Longitudinal Redundancy Check (LRC)

- In longitudinal redundancy method, a BLOCK of bits are arranged in a table format (in rows and columns) and we will calculate the parity bit for each column separately.
- The set of these parity bits are also sent along with our original data bits.
- Longitudinal redundancy check is a bit by bit parity computation, as we calculate the parity of each column individually.
- LRC increases the likelihood of detecting burst error.
- However, if two bits in one data unit are damaged and two bits in exactly the same positions in another data unit are also damaged, the LRC checker will not detect an error.



101000/1 001100/1 11011101 11100111
10101010 (LRC)

Calculate the LRC for Data Received

101000/1
001100/1
11011101
11100111

→ LRC Calculated by Receiver 10101010
→ Compare with LRC Received 10101010

Checksum

- A checksum number is appended to the packet sequence so that the sum of data plus checksum is zero.
- When received, the packet sequence may be added, along with the checksum, by a local microprocessor. If the sum is nonzero, an error has occurred.
- The checksum method includes parity bits, check digits and longitudinal redundancy check (LRC).

example: add two 16-bit integers

	1	1	1	0	0	1	1	0	0	1	1	0
	1	1	0	1	0	1	0	1	0	1	0	1
<hr/>												
wraparound	1	1	0	1	1	1	0	1	1	1	0	1
sum	1	0	1	1	1	0	1	1	1	1	1	0
checksum	0	1	0	0	0	1	0	0	1	0	0	0

Note: when adding numbers, a carryout from the most significant bit needs to be added to the result

If $k = 4$, and $n = 8$ then

$$\begin{array}{r}
 k=4, \quad n=8 \\
 10110011 \\
 10101011 \\
 \hline
 01011110 \\
 1 \\
 \hline
 01011111 \\
 01011010 \\
 \hline
 10111001 \\
 11010101 \\
 \hline
 10001110 \\
 1 \\
 \hline
 10001111 \\
 01110000
 \end{array}$$

Sum : 10001111
Checksum 01110000

At sender side

$$\begin{array}{r}
 10110011 \\
 10101011 \\
 \hline
 01011110 \\
 1 \\
 \hline
 01011111 \\
 01011010 \\
 \hline
 10111001 \\
 11010101 \\
 \hline
 10001110 \\
 1 \\
 \hline
 10001111 \\
 01110000
 \end{array}$$

Sum: 11111111

Complement = 00000000
Conclusion = Accept data

At receiver side

Hamming Codes

- Hamming code is a set of error-correction codes that can be used to **detect and correct the errors** that can occur when the data is moved or stored from the sender to the receiver.
- Redundant bits: Redundant bits are extra binary bits that are generated and added to the information-carrying bits of data transfer to ensure that no bits were lost during the data transfer.

The number of redundant bits can be calculated using the following formula:

$$2^r \geq m + r + 1 \text{ where, } r = \text{redundant bit, } m = \text{data bit}$$

- Parity Bits: A parity bit is a bit appended to a data of binary bits to ensure that the total number of 1's in the data is even or odd.

The key to the Hamming Code is the use of extra parity bits to allow the identification of a single error. Create the code word as follows:

1. Mark all bit positions that are powers of two as parity bits. (positions 1, 2, 4, 8, 16, 32, 64, etc.)
2. All other bit positions are for the data to be encoded. (positions 3, 5, 6, 7, 9, 10, 11, 12, 13, 14, 15, 17, etc.)
3. Each parity bit calculates the parity for some of the bits in the code word. The position of the parity bit determines the sequence of bits that it alternately checks and skips.

Position 1: check 1 bit, skip 1 bit, check 1 bit, skip 1 bit, etc. (1,3,5,7,9,11,13,15,...)

Position 2: check 2 bits, skip 2 bits, check 2 bits, skip 2 bits, etc. (2,3,6,7,10,11,14,15,...)

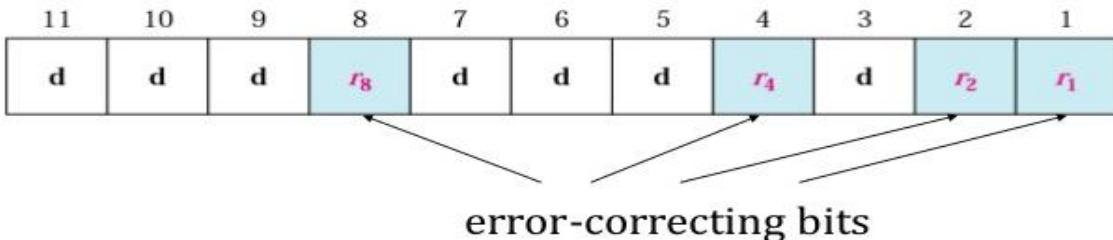
Position 4: check 4 bits, skip 4 bits, check 4 bits, skip 4 bits, etc.
(4,5,6,7,12,13,14,15,20,21,22,23,...)

Position 8: check 8 bits, skip 8 bits, check 8 bits, skip 8 bits, etc. (8-15,24-31,40-47,...)

Position 16: check 16 bits, skip 16 bits, check 16 bits, skip 16 bits, etc. (16-31,48-63,80-95,...)

Position 32: check 32 bits, skip 32 bits, check 32 bits, skip 32 bits, etc. (32-63,96-127,160-191,...)
etc.

4. Set a parity bit to 1 if the total number of ones in the positions it checks is odd. Set a parity bit to 0 if the total number of ones in the positions it checks is even.



Example of redundancy bit calculation

d	d	d	r8	d	d	d	r4	d	r2	r1

Input Data 1001101

Adding r1	1	0	0	r8	1	1	0	r4	1	r2	1
	1	0	0	r8	1	1	0	r4	1	r2	1
Adding r2	1	0	0	r8	1	1	0	r4	1	0	1
Adding r4	1	0	0	r8	1	1	0	0	1	0	1

Adding r8	1	0	0	1	1	1	0	0	1	0	1
	1	0	0	1	1	1	0	0	1	0	1

Output data:10011100101

example

Error detecting using hamming code
error

11	10	9	8	7	6	5	4	3	2	1
1	0	0	1	0	1	0	0	1	0	1

If no.1,s
even 0If no.1,s is
odd 1

r1=1

11	10	9	8	7	6	5	4	3	2	1
1	0	0	1	0	1	0	0	1	0	1

8 4 2 1

0 1 1 1

7

11	10	9	8	7	6	5	4	3	2	1
1	0	0	1	0	1	0	0	1	0	1

It mean the 7
bit is
corrupted

11	10	9	8	7	6	5	4	3	2	1
1	0	0	1	0	1	0	0	1	0	1

r8=0

DUALITY PRINCIPLE

- Dual:

The dual of a Boolean expression is the expression one obtains by interchanging addition and multiplication and interchanging 0's and 1's. The dual of the function F is denoted F^d .

- Duality Principle:

This principle states that any algebraic equality derived from these axioms will still be valid whenever the OR and AND operators, and identity elements 0 and 1, have been interchanged. i.e. changing every OR into AND and vice versa, and every 0 into 1 and vice versa, i.e.,

If F and G are Boolean functions such that $F = G$, then $F^d = G^d$.

Example:

The dual of $xy' + x'z = (x + y') \cdot (x' + z)$.

CHAPTER-2

A BRIEF HISTORY OF COMPUTERS

- The chapter gives a overview of the evolution of computer technology from early digital computers to the latest microprocessors.
- The generation of computers is classified as:
 - The First Generation: Vacuum Tubes
 - The Second Generation: Transistors
 - The Third Generation: Integrated Circuits
 - Later Generations
- Each new generation is characterized by greater processing performance, larger memory capacity, and smaller size than the previous one.

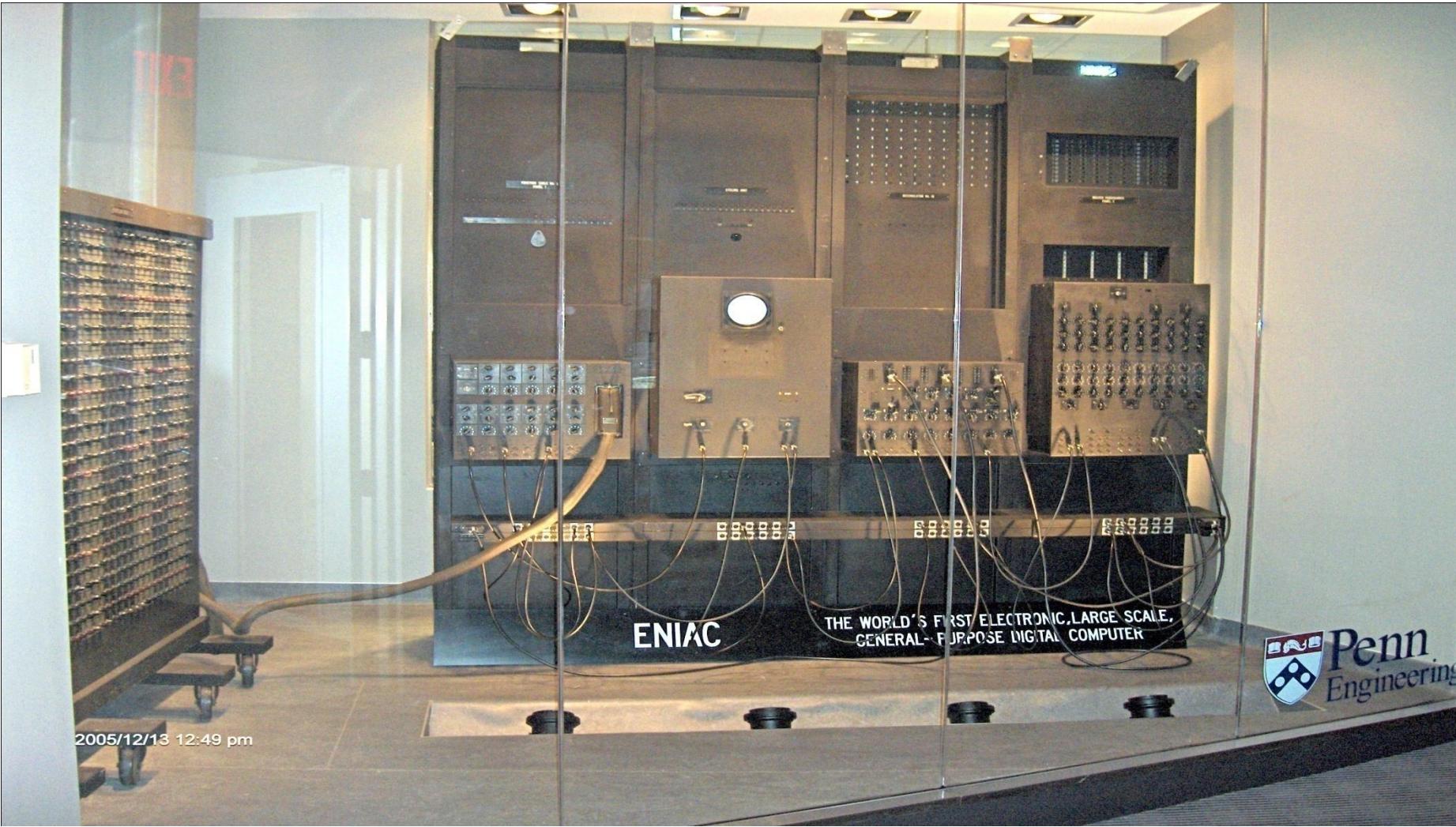
Generation	Approximate Dates	Technology	Typical Speed (operations per second)
1	1946–1957	Vacuum tube	40,000
2	1958–1964	Transistor	200,000
3	1965–1971	Small- and medium-scale integration	1,000,000
4	1972–1977	Large-scale integration	10,000,000
5	1978–1991	Very-large-scale integration	100,000,000
6	1991–	Ultra-large-scale integration	1,000,000,000

THE FIRST GENERATION: VACUUM TUBES

- ENIAC

The ENIAC (Electronic Numerical Integrator And Computer), designed and constructed at the University of Pennsylvania, was the world's first general purpose electronic digital computer. The ENIAC was a decimal rather than a binary machine. That is, numbers were represented in decimal form, and arithmetic was performed in the decimal system. Its memory consisted of 20 accumulators, each capable of holding a 10-digit decimal number. A ring of 10 vacuum tubes represented each digit. At any time, only one vacuum tube was in the ON state, representing one of the 10 digits. The major drawback of the ENIAC was that it had to be programmed manually by setting switches and plugging and unplugging cables.

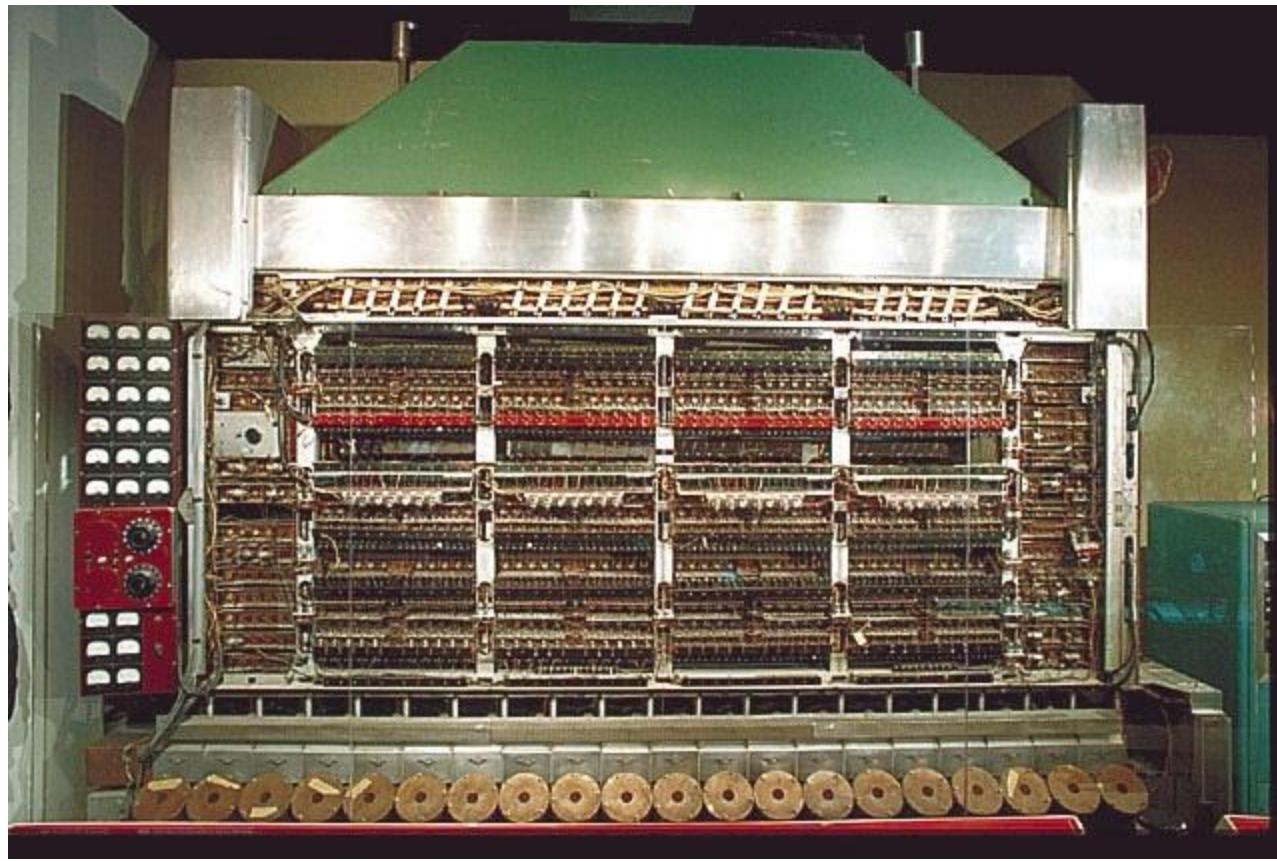
SOURCE: INTERNET



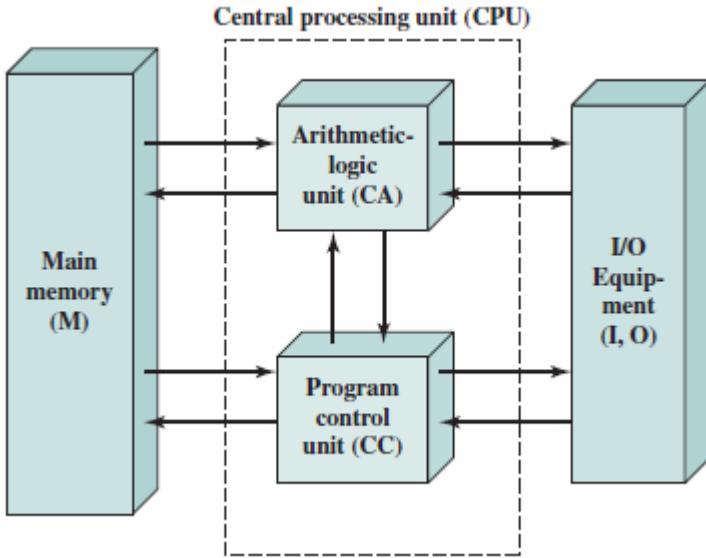
- THE VON NEUMANN MACHINE

The task of entering and altering programs for the ENIAC was extremely tedious. But suppose a program could be represented in a form suitable for storing in memory alongside the data. Then, a computer could get its instructions by reading them from memory, and a program could be set or altered by setting the values of a portion of memory. This idea, known as the stored-program concept, is usually attributed to the ENIAC designers, most notably the mathematician John von Neumann, who was a consultant on the ENIAC project.

In 1946, von Neumann and his colleagues began the design of a new stored program computer, referred to as the IAS computer, at the Princeton Institute for Advanced Studies.



SOURCE: INTERNET

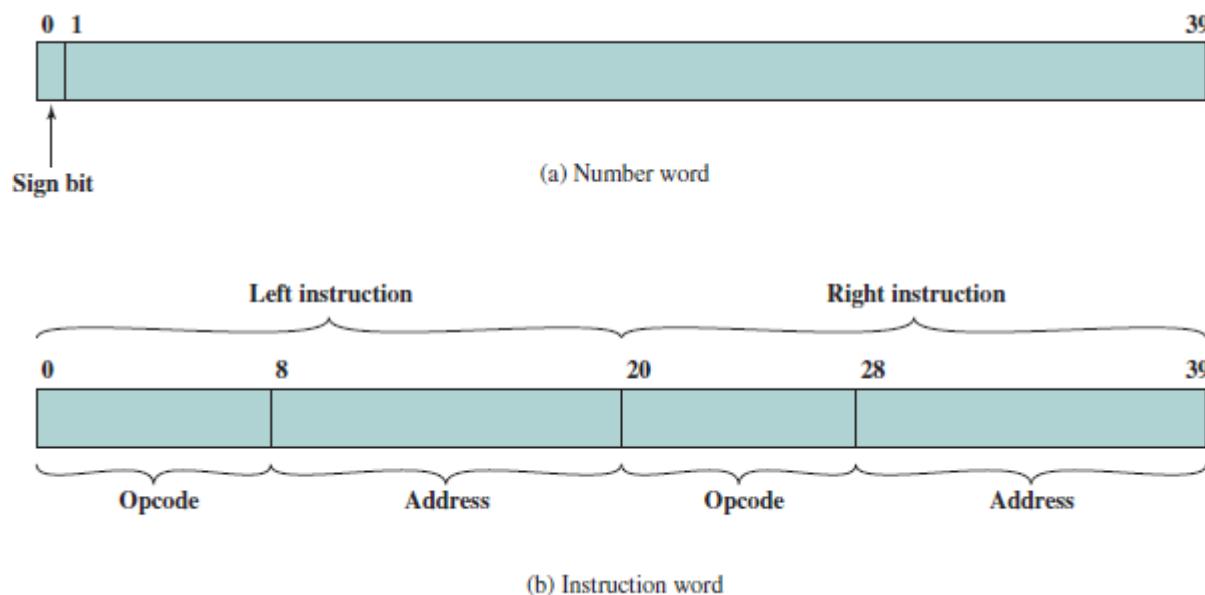


Structure of the IAS Computer

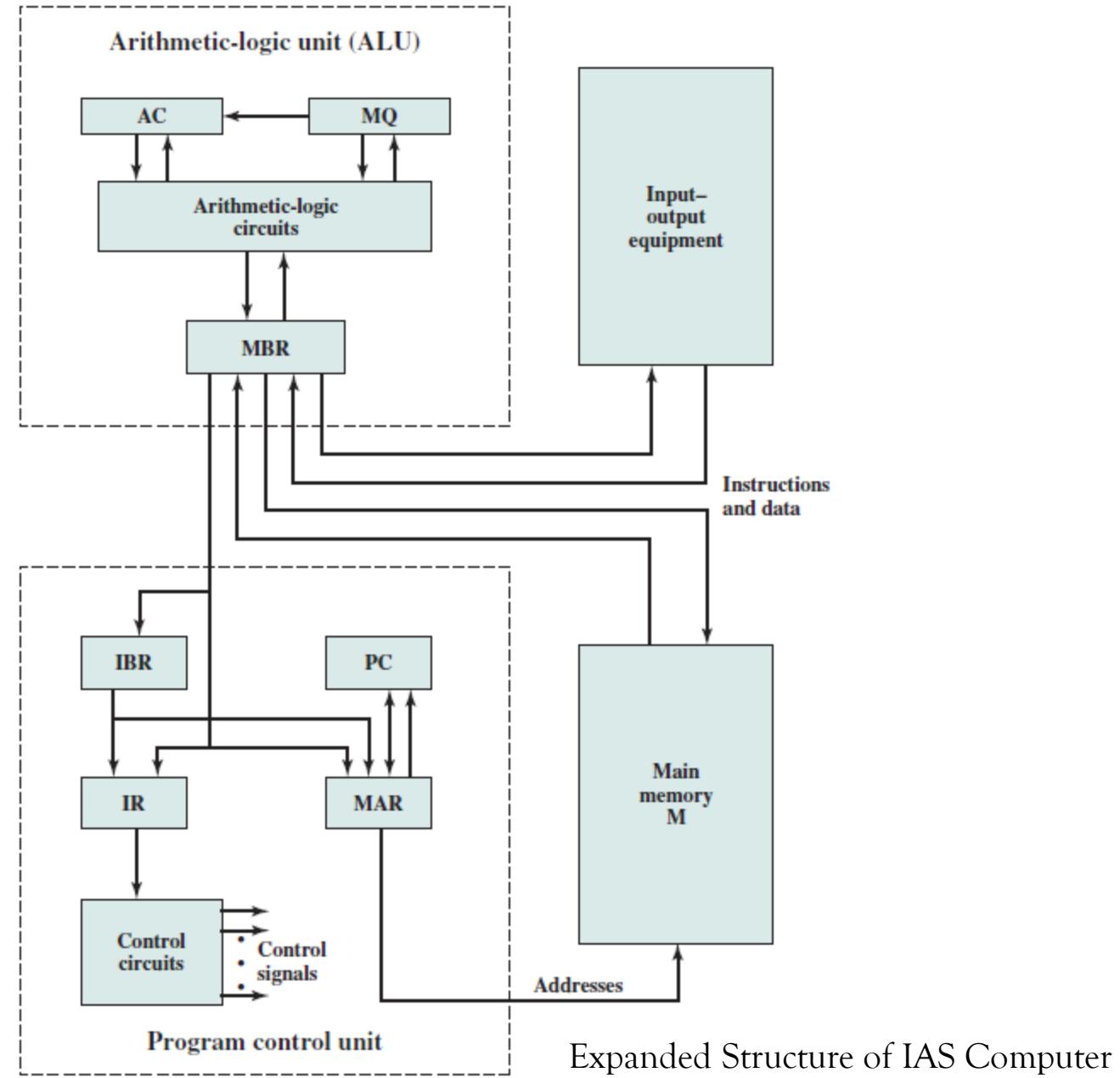
- A **main memory**, which stores both data and instructions.
- An **arithmetic and logic unit** (ALU) capable of operating on binary data.
- A **control unit**, which interprets the instructions in memory and causes them to be executed.
- **Input/output (I/O)** equipment operated by the control unit.

With rare exceptions, all of today's computers have this same general structure and function and are thus referred to as **von Neumann machines**.

- The **memory** of the IAS consists of 1000 storage locations, called words, of 40 binary digits (bits) each. Both data and instructions are stored there. Numbers are represented in binary form, and each instruction is a binary code. Each number is represented by a sign bit and a 39-bit value. A word may also contain two 20-bit instructions, with each instruction consisting of an 8-bit operation code (opcode) specifying the operation to be performed and a 12-bit address designating one of the words in memory (numbered from 0 to 999).



- The **control unit** operates the IAS by fetching instructions from memory and executing them one at a time.
- The control unit and the ALU contain storage locations, called registers, defined as follows:
 - **Memory buffer register (MBR):** Contains a word to be stored in memory or sent to the I/O unit, or is used to receive a word from memory or from the I/O unit.
 - **Memory address register (MAR):** Specifies the address in memory of the word to be written from or read into the MBR.
 - **Instruction register (IR):** Contains the 8-bit opcode instruction being executed.
 - **Instruction buffer register (IBR):** Employed to hold temporarily the righthand instruction from a word in memory.
 - **Program counter (PC):** Contains the address of the next instruction pair to be fetched from memory.
 - **Accumulator (AC) and multiplier quotient (MQ):** Employed to hold temporarily operands and results of ALU operations. For example, the result of multiplying two 40-bit numbers is an 80-bit number; the most significant 40 bits are stored in the AC and the least significant in the MQ.



- The IAS operates by repetitively performing an **instruction cycle**. Each instruction cycle consists of two sub cycles.
 - During the **fetch cycle**, the opcode of the next instruction is loaded into the IR and the address portion is loaded into the MAR. This instruction may be taken from the IBR, or it can be obtained from memory by loading a word into the MBR, and then down to the IBR, IR, and MAR.
 - Once the opcode is in the IR, the **execute cycle** is performed. Control circuitry interprets the opcode and executes the instruction by sending out the appropriate control signals to cause data to be moved or an operation to be performed by the ALU.

- The IAS computer had a total of **21 instructions** which can be grouped as follows:
 - **Data transfer:** Move data between memory and ALU registers or between two ALU registers.
 - **Unconditional branch:** Normally, the control unit executes instructions in sequence from memory. This sequence can be changed by a branch instruction, which facilitates repetitive operations.
 - **Conditional branch:** The branch can be made dependent on a condition, thus allowing decision points.
 - **Arithmetic:** Operations performed by the ALU.
 - **Address modify:** Permits addresses to be computed in the ALU and then inserted into instructions stored in memory. This allows a program considerable addressing flexibility.

Instruction Type	Opcode	Symbolic Representation	Description
Data transfer	00001010	LOAD MQ	Transfer contents of register MQ to the accumulator AC
	00001001	LOAD MQ,M(X)	Transfer contents of memory location X to MQ
	00100001	STOR M(X)	Transfer contents of accumulator to memory location X
	00000001	LOAD M(X)	Transfer M(X) to the accumulator
	00000010	LOAD -M(X)	Transfer -M(X) to the accumulator
	00000011	LOAD M(X)	Transfer absolute value of M(X) to the accumulator
	00000100	LOAD - M(X)	Transfer - M(X) to the accumulator
Unconditional branch	00001101	JUMP M(X,0:19)	Take next instruction from left half of M(X)
	00001110	JUMP M(X,20:39)	Take next instruction from right half of M(X)
Conditional branch	00001111	JUMP + M(X,0:19)	If number in the accumulator is nonnegative, take next instruction from left half of M(X)
	00010000	JUMP + M(X,20:39)	If number in the accumulator is nonnegative, take next instruction from right half of M(X)
Arithmetic	00000101	ADD M(X)	Add M(X) to AC; put the result in AC
	00000111	ADD M(X)	Add M(X) to AC; put the result in AC
	00000110	SUB M(X)	Subtract M(X) from AC; put the result in AC
	00001000	SUB M(X)	Subtract M(X) from AC; put the remainder in AC
	00001011	MUL M(X)	Multiply M(X) by MQ; put most significant bits of result in AC, put least significant bits in MQ
	00001100	DIV M(X)	Divide AC by M(X); put the quotient in MQ and the remainder in AC
	00010100	LSH	Multiply accumulator by 2; that is, shift left one bit position
	00010101	RSH	Divide accumulator by 2; that is, shift right one position
	00010010	STOR M(X,8:19)	Replace left address field at M(X) by 12 rightmost bits of AC
Address modify	00010011	STOR M(X,28:39)	Replace right address field at M(X) by 12 rightmost bits of AC

The IAS Instruction Set

- COMMERCIAL COMPUTERS

The UNIVAC I was the first successful commercial computer. It was intended for both scientific and commercial applications.

The UNIVAC II, which had greater memory capacity and higher performance than the UNIVAC I, was delivered in the late 1950s and illustrates several trends that have remained characteristic of the computer industry.

First, advances in technology allow companies to continue to build larger, more powerful computers. Second, each company tries to make its new machines backward compatible with the older machines.

The UNIVAC division also began development of the 1100 series of computers, which was to be its major source of revenue. This series illustrates a distinction that existed at one time. The first model, the UNIVAC 1103, and its successors for many years were primarily intended for scientific applications, involving long and complex calculations.

IBM, then the major manufacturer of punched-card processing equipment, delivered its first electronic stored-program computer, the 701, in 1953. The 701 was intended primarily for scientific applications. In 1955, IBM introduced the companion 702 product, which had a number of hardware features that suited it to business applications. These were the first of a long series of 700/7000 computers that established IBM as the overwhelmingly dominant computer manufacturer.



UNIVAC I



UNIVAC II



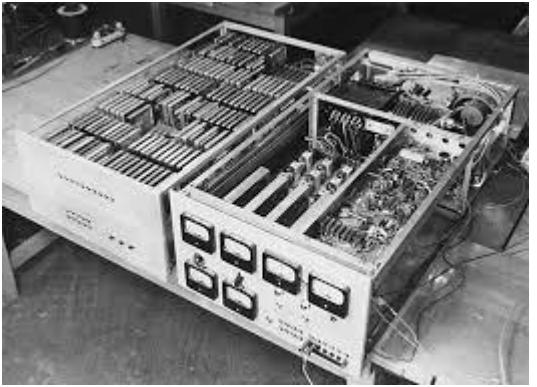
UNIVAC 1103



IBM 701

THE SECOND GENERATION: TRANSISTORS

- The first major change in the electronic computer came with the replacement of the vacuum tube by the transistor. The transistor is smaller, cheaper, and dissipates less heat than a vacuum tube but can be used in the same way as a vacuum tube to construct computers.
- Unlike the vacuum tube, which requires wires, metal plates, a glass capsule, and a vacuum, **the transistor is a solid-state device**, made from silicon.
- *The use of the transistor defines the second generation of computers.* It has become widely accepted to classify computers into generations based on the fundamental hardware technology employed.



- The second generation saw the introduction of more complex arithmetic and logic units and control units, the use of high-level programming languages, and the provision of *system software* with the computer.
- In broad terms, system software provided the ability to load programs, move data to peripherals, and libraries to perform common computations, similar to what modern OSes like Windows and Linux do.
- The second generation is noteworthy also for the appearance of the Digital Equipment Corporation (DEC), who delivered its first computer, the (Programmed Data Processor) PDP-1. This computer and this company began the minicomputer phenomenon that became prominent in the third generation.
- From the introduction of the 700 series in 1952 to the introduction of the last member of the 7000 series in 1964, IBM 7094 which underwent an evolution that is typical of computer products. Successive members of the product line show increased performance, increased capacity, and/or lower cost.

SOURCE: INTERNET



DEC PDP-1

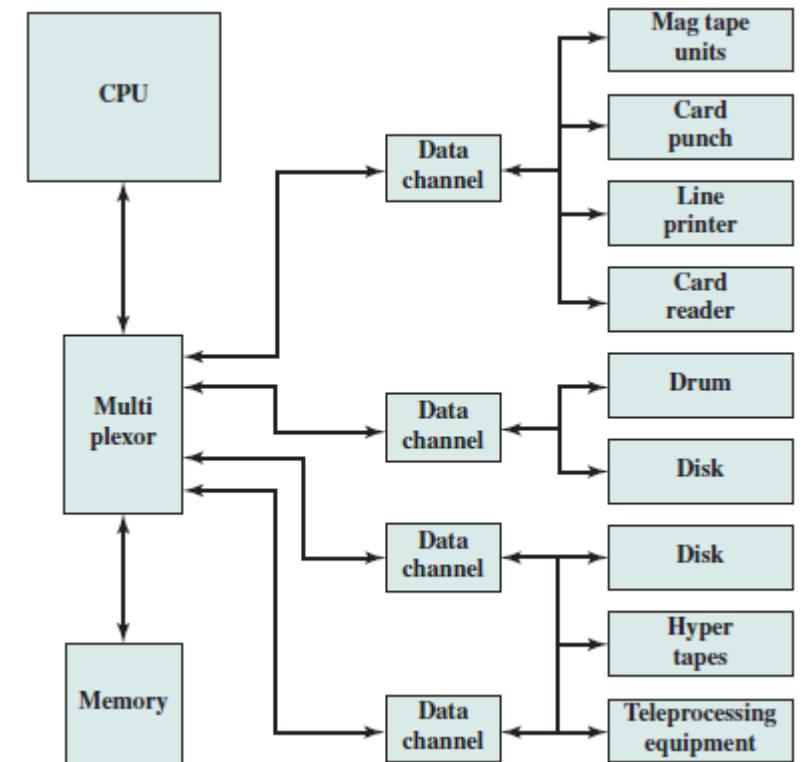


IBM 7094

Model Number	First Delivery	CPU Tech-nology	Memory Tech-nology	Cycle Time (μs)	Memory Size (K)	Number of Opcodes	Number of Index Registers	Hardwired Floating-Point	I/O Overlap (Channels)	Instruc-tion Fetch Overlap	Speed (relative to 701)
701	1952	Vacuum tubes	Electrostatic tubes	30	2–4	24	0	no	no	no	1
704	1955	Vacuum tubes	Core	12	4–32	80	3	yes	no	no	2.5
709	1958	Vacuum tubes	Core	12	32	140	3	yes	yes	no	4
7090	1960	Transistor	Core	2.18	32	169	3	yes	yes	no	25
7094 I	1962	Transistor	Core	2	32	185	7	yes (double precision)	yes	yes	30
7094 II	1964	Transistor	Core	1.4	32	185	7	yes (double precision)	yes	yes	50

Example members of the IBM 700/7000 Series

- The most important of these is the use of **data channels**. A data channel is an independent I/O module with its own processor and instruction set. In a computer system with such devices, the CPU does not execute detailed I/O instructions. Such instructions are stored in a main memory to be executed by a special-purpose processor in the data channel itself.
- The CPU initiates an I/O transfer by sending a control signal to the data channel, instructing it to execute a sequence of instructions in memory. The data channel performs its task independently of the CPU and signals the CPU when the operation is complete. This arrangement relieves the CPU of a considerable processing burden.
- Another new feature is the **multiplexor**, which is the central termination point for data channels, the CPU, and memory. The multiplexor schedules access to the memory from the CPU and data channels, allowing these devices to act independently



An IBM 7094 Configuration

THE THIRD GENERATION: INTEGRATED CIRCUITS

- A single, self-contained transistor is called a **discrete component**. Throughout the 1950s and early 1960s, electronic equipment was composed largely of discrete components—transistors, resistors, capacitors, and so on.
- In 1958 came the achievement that revolutionized electronics and started the era of microelectronics: the invention of the **integrated circuit**.
- It is the integrated circuit that defines the third generation of computers.

- MICROELECTRONICS

Microelectronics means, literally, “small electronics.” The basic elements of a digital computer, as we know, must perform storage, movement, processing, and control functions. Only two fundamental types of components are required: **gates** and **memory cells**.

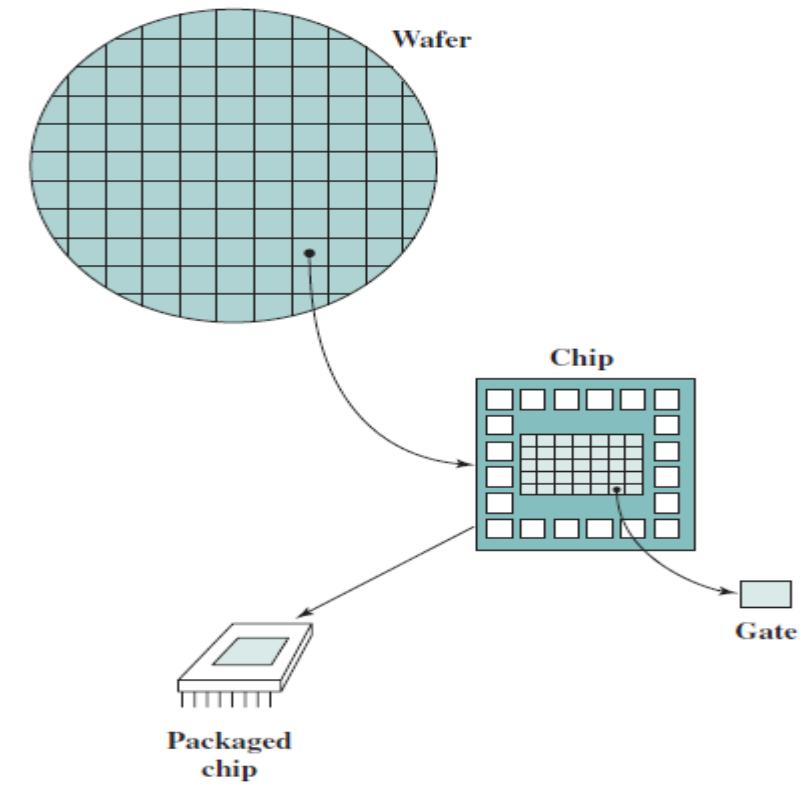
A **gate** is a device that implements a simple Boolean or logical function, such as IF A AND B ARE TRUE THEN C IS TRUE (AND gate). Such devices are called gates because they control data flow in much the same way that canal gates control the flow of water.

The **memory cell** is a device that can store one bit of data; that is, the device can be in one of two stable states at any time. By interconnecting large numbers of these fundamental devices, we can construct a computer.

- Four basic functions of a computer:
 - **Data storage:** Provided by memory cells.
 - **Data processing:** Provided by gates.
 - **Data movement:** The paths among components are used to move data from memory to memory and from memory through gates to memory.
 - **Control:** The paths among components can carry control signals. For example, a gate will have one or two data inputs plus a control signal input that activates the gate. When the control signal is ON, the gate performs its function on the data inputs and produces a data output. Similarly, the memory cell will store the bit that is on its input lead when the WRITE control signal is ON and will place the bit that is in the cell on its output lead when the READ control signal is ON.
- Thus, a computer consists of gates, memory cells, and interconnections among these elements. The gates and memory cells are, in turn, constructed of simple digital electronic components.

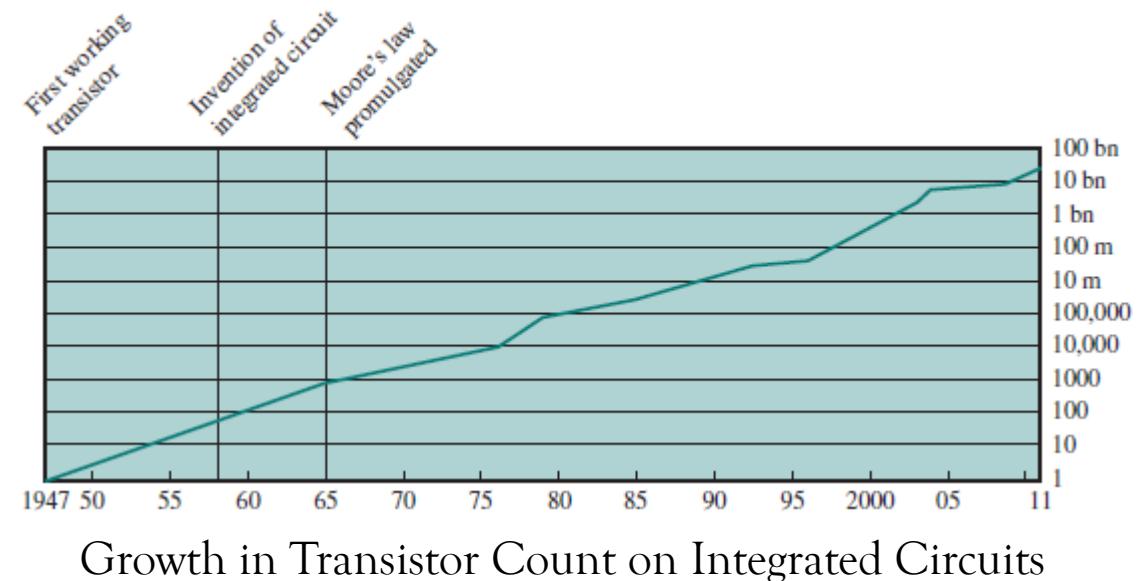
- The integrated circuit exploits the fact that such components as transistors, resistors, and conductors can be fabricated from a semiconductor such as silicon.
- It is merely an extension of the solid-state art to fabricate an entire circuit in a tiny piece of silicon rather than assemble discrete components made from separate pieces of silicon into the same circuit.

- A thin wafer of silicon is divided into a matrix of small areas, each a few millimeters square.
- The identical circuit pattern is fabricated in each area, and the wafer is broken up into chips.
- Each chip consists of many gates and/or memory cells plus a number of input and output attachment points.
- This chip is then packaged in housing that protects it and provides pins for attachment to devices beyond the chip.
- A number of these packages can then be interconnected on a printed circuit board to produce larger and more complex circuits.



Relationship among Wafer, Chip, and Gate

- Initially, only a few gates or memory cells could be reliably manufactured and packaged together. These early integrated circuits are referred to as **small scale integration (SSI)**.
- As time went on, it became possible to pack more and more components on the same chip.
- Gordon Moore, cofounder of Intel, observed that the number of transistors that could be put on a single chip was doubling every year and correctly predicted that this pace would continue into the near future.



Growth in Transistor Count on Integrated Circuits

- The consequences of Moore's law are profound:
 - a. The cost of a chip has remained virtually unchanged during this period of rapid growth in density. This means that the cost of computer logic and memory circuitry has fallen at a dramatic rate.
 - b. Because logic and memory elements are placed closer together on more densely packed chips, the electrical path length is shortened, increasing operating speed.
 - c. The computer becomes smaller, making it more convenient to place in a variety of environments.
 - d. There is a reduction in power and cooling requirements.
 - e. The interconnections on the integrated circuit are much more reliable than solder connections. With more circuitry on each chip, there are fewer interchip connections.

- IBM SYSTEM/360

In 1964, IBM announced the System/360, a new family of computer products. The 360 product line was incompatible with older IBM machines. Thus, the transition to the 360 would be difficult for the current customer base. Hence, it break out of some of the constraints of the 7000 architecture and targeted to produce a system capable of evolving with the new integrated circuit technology.

With some modifications and extensions, the architecture of the 360 remains to this day the architecture of IBM's mainframe computers.

The System/360 was the industry's first planned family of computers. The concept of a family of compatible computers was both novel and extremely successful.

The System/360 not only dictated the future course of IBM but also had a profound impact on the entire industry. Many of its features have become standard on other large computers.

- The characteristics of a family are as follows:
 - **Similar or identical instruction set:** In many cases, the exact same set of machine instructions is supported on all members of the family. Thus, a program that executes on one machine will also execute on any other. In some cases, the lower end of the family has an instruction set that is a subset of that of the top end of the family. This means that programs can move up but not down.
 - **Similar or identical operating system:** The same basic operating system is available for all family members. In some cases, additional features are added to the higher-end members.
 - **Increasing speed:** The rate of instruction execution increases in going from lower to higher family members.
 - **Increasing number of I/O ports:** The number of I/O ports increases in going from lower to higher family members.
 - **Increasing memory size:** The size of main memory increases in going from lower to higher family members.
 - **Increasing cost:** At a given point in time, the cost of a system increases in going from lower to higher family members.

Characteristic	Model 30	Model 40	Model 50	Model 65	Model 75
Maximum memory size (bytes)	64K	256K	256K	512K	512K
Data rate from memory (Mbytes/s)	0.5	0.8	2.0	8.0	16.0
Processor cycle time (μ s)	1.0	0.625	0.5	0.25	0.2
Relative speed	1	3.5	10	21	50
Maximum number of data channels	3	3	4	6	6
Maximum data rate on one channel (Kbytes/s)	250	400	800	1250	1250

Key Characteristics of the System/360 Family

- DEC PDP-8

In the same year that IBM shipped its first System/360, another momentous first shipment occurred: PDP-8 from Digital Equipment Corporation (DEC). At a time when the average computer required an airconditioned room, the PDP-8 (dubbed a minicomputer by the industry, after the miniskirt of the day) was small enough that it could be placed on top of a lab bench or be built into other equipment. It could not do everything the mainframe could, and was cheap enough for each lab technician to have one. In contrast, the System/360 series of mainframe computers introduced just a few months before cost hundreds of thousands of dollars.

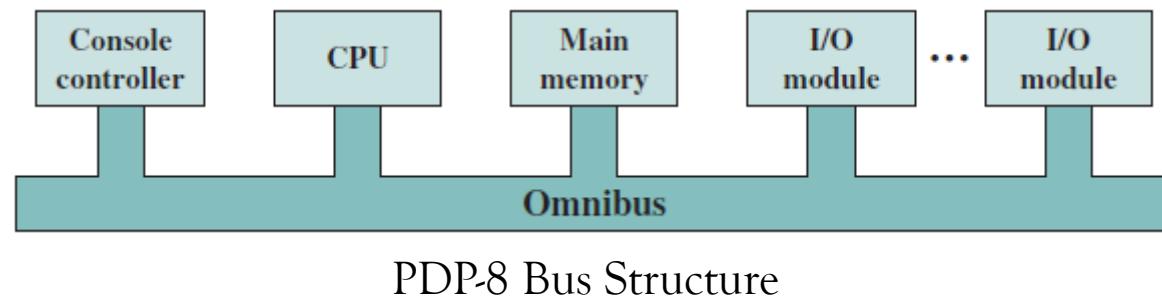
The low cost and small size of the PDP-8 enabled another manufacturer to purchase a PDP-8 and integrate it into a total system for resale. These other manufacturers came to be known as original equipment manufacturers (OEMs), and the OEM market became and remains a major segment of the computer marketplace.

Model	First Shipped	Cost of Processor + 4K 12-bit Words of Memory (\$1000s)	Data Rate from Memory (words/µs)	Volume (cubic feet)	Innovations and Improvements
PDP-8	4/65	16.2	1.26	8.0	Automatic wire-wrapping production
PDP-8/5	9/66	8.79	0.08	3.2	Serial instruction implementation
PDP-8/1	4/68	11.6	1.34	8.0	Medium-scale integrated circuits
PDP-8/L	11/68	7.0	1.26	2.0	Smaller cabinet
PDP-8/E	3/71	4.99	1.52	2.2	Omnibus
PDP-8/M	6/72	3.69	1.52	1.8	Half-size cabinet with fewer slots than 8/E
PDP-8/A	1/75	2.6	1.34	1.2	Semiconductor memory; floating-point processor

Evolution of the PDP-8

In contrast to the central-switched architecture used by IBM on its 700/7000 and 360 systems, later models of the PDP-8 used a structure that is now virtually universal for microcomputers: the bus structure.

The PDP-8 bus, called the Omnibus, consists of 96 separate signal paths, used to carry control, address, and data signals. Because all system components share a common set of signal paths, their use can be controlled by the CPU. This architecture is highly flexible, allowing modules to be plugged into the bus to create various configurations.



LATER GENERATIONS

- Beyond the third generation there is less general agreement on defining generations of computers.
- There have been a number of later generations, based on advances in integrated circuit technology.
- With the introduction of large-scale integration (LSI), more than 1000 components can be placed on a single integrated circuit chip. Very-large-scale integration (VLSI) achieved more than 10,000 components per chip, while current ultra-large-scale integration (ULSI) chips can contain more than one billion components.

- SEMICONDUCTOR MEMORY

The first application of integrated circuit technology to computers was construction of the processor (the control unit and the arithmetic and logic unit) out of integrated circuit chips. But it was also found that this same technology could be used to construct memories.

In the 1950s and 1960s, most computer memory was constructed from tiny rings of ferromagnetic material, each about a sixteenth of an inch in diameter.

These rings were strung up on grids of fine wires suspended on small screens inside the computer. Magnetized one way, a ring (called a core) represented a one; magnetized the other way, it stood for a zero. Magnetic-core memory was rather fast; it took as little as a millionth of a second to read a bit stored in memory. But it was expensive, bulky, and used destructive readout: The simple act of reading a core erased the data stored in it. It was therefore necessary to install circuits to restore the data as soon as it had been extracted.

Then, in 1970, Fairchild produced the first relatively capacious semiconductor memory. This chip, about the size of a single core, could hold 256 bits of memory. It was nondestructive and much faster than core. It took only 70 billionths of a second to read a bit. However, the cost per bit was higher than for that of core.

In 1974, a seminal event occurred: The price per bit of semiconductor memory dropped below the price per bit of core memory. Following this, there has been a continuing and rapid decline in memory cost accompanied by a corresponding increase in physical memory density. This has led the way to smaller, faster machines with memory sizes of larger and more expensive machines from just a few years earlier.

Since 1970, semiconductor memory has been through 13 generations: 1K, 4K, 16K, 64K, 256K, 1M, 4M, 16M, 64M, 256M, 1G, 4G, and, as of this writing, 16 Gbits on a single chip ($1K = 2^{10}$, $1M = 2^{20}$, $1G = 2^{30}$). Each generation has provided four times the storage density of the previous generation, accompanied by declining cost per bit and declining access time.

- MICROPROCESSORS

As time went on, more and more elements were placed on each chip, so that fewer and fewer chips were needed to construct a single computer processor.

A breakthrough was achieved in 1971, when Intel developed its 4004. The 4004 was the first chip to contain all of the components of a CPU on a single chip: The microprocessor was born.

The 4004 can add two 4-bit numbers and can multiply only by repeated addition. By today's standards, the 4004 is hopelessly primitive, but it marked the beginning of a continuing evolution of microprocessor capability and power.

This evolution can be seen most easily in the number of bits that the processor deals with at a time.

- **Data bus width:** the number of bits of data that can be brought into or sent out of the processor at a time.
- Number of bits in the accumulator or in the set of general-purpose registers.

The next major step in the evolution of the microprocessor was the introduction in 1972 of the Intel 8008. This was the first 8-bit microprocessor and was almost twice as complex as the 4004. This was the first general-purpose microprocessor, designed to be the CPU of a general-purpose microcomputer. The 8080, however, is faster, has a richer instruction set, and has a large addressing capability.

About the same time, 16-bit microprocessors began to be developed. One of these was the 8086. The next step in this trend occurred in 1981, when both Bell Labs and Hewlett-Packard developed 32-bit, single-chip microprocessors. Intel introduced its own 32-bit microprocessor, the 80386, in 1985.

(a) 1970s Processors					
	4004	8008	8080	8086	8088
Introduced	1971	1972	1974	1978	1979
Clock speeds	108 kHz	108 kHz	2 MHz	5 MHz, 8 MHz, 10 MHz	5 MHz, 8 MHz
Bus width	4 bits	8 bits	8 bits	16 bits	8 bits
Number of transistors	2300	3500	6000	29,000	29,000
Feature size (μm)	10		6	3	6
Addressable memory	640 Bytes	16 kB	64 kB	1 MB	1 MB

	80286	386TM DX	386TM SX	486TM DX CPU
Introduced	1982	1985	1988	1989
Clock speeds	6 MHz–12.5 MHz	16 MHz–33 MHz	16 MHz–33 MHz	25 MHz–50 MHz
Bus width	16 bits	32 bits	16 bits	32 bits
Number of transistors	134,000	275,000	275,000	1.2 million
Feature size (μm)	1.5	1	1	0.8–1
Addressable memory	16 MB	4 GB	16 MB	4 GB
Virtual memory	1 GB	64 TB	64 TB	64 TB
Cache	—	—	—	8 kB

(c) 1990s Processors

	486TM SX	Pentium	Pentium Pro	Pentium II
Introduced	1991	1993	1995	1997
Clock speeds	16 MHz–33 MHz	60 MHz–166 MHz	150 MHz–200 MHz	200 MHz–300 MHz
Bus width	32 bits	32 bits	64 bits	64 bits
Number of transistors	1.185 million	3.1 million	5.5 million	7.5 million
Feature size (μm)	1	0.8	0.6	0.35
Addressable memory	4 GB	4 GB	64 GB	64 GB
Virtual memory	64 TB	64 TB	64 TB	64 TB
Cache	8 kB	8 kB	512 kB L1 and 1 MB L2	512 kB L2

(d) Recent Processors

	Pentium III	Pentium 4	Core 2 Duo	Core i7 EE 990
Introduced	1999	2000	2006	2011
Clock speeds	450–660 MHz	1.3–1.8 GHz	1.06–1.2 GHz	3.5 GHz
Bus width	64 bits	64 bits	64 bits	64 bits
Number of transistors	9.5 million	42 million	167 million	1170 million
Feature size (nm)	250	180	65	32
Addressable memory	64 GB	64 GB	64 GB	64 GB
Virtual memory	64 TB	64 TB	64 TB	64 TB
Cache	512 kB L2	256 kB L2	2 MB L2	1.5 MB L2/12 MB L3

Evolution of Intel Microprocessors

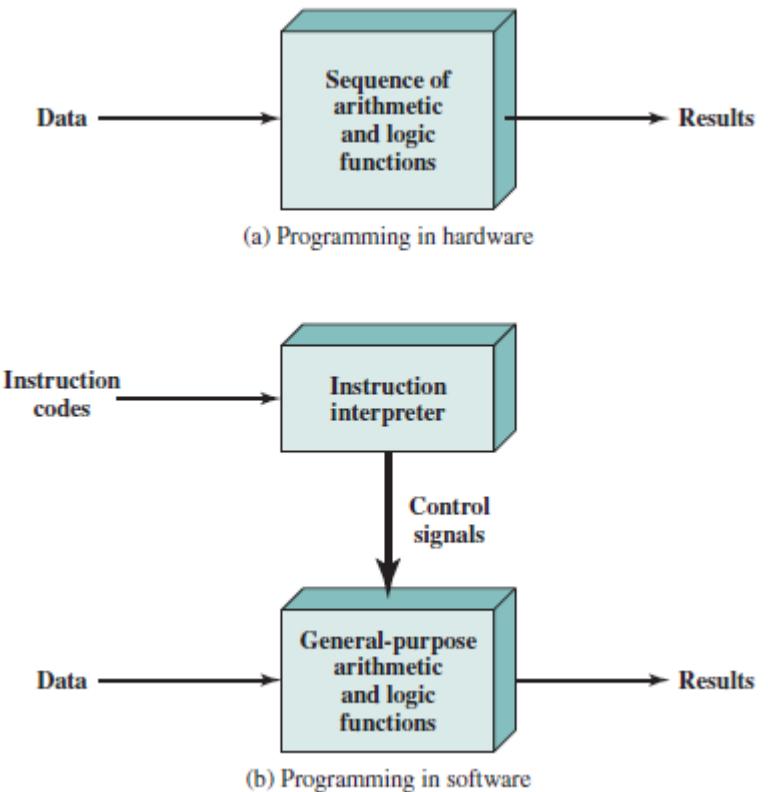
COMPUTER TYPES

The different types of computers are:

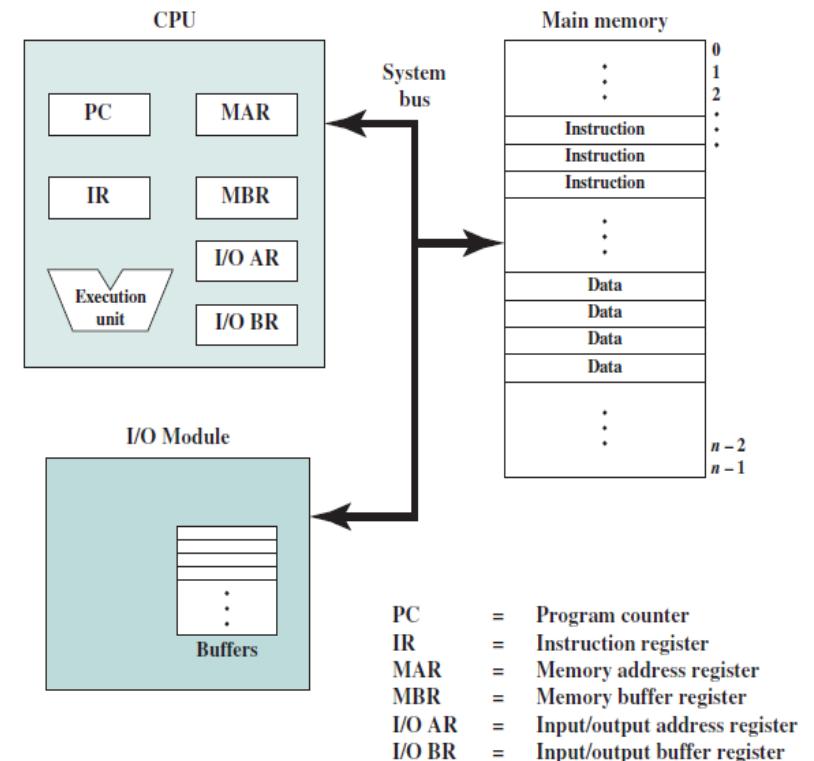
- **Personal computers:** This is the most common type found in homes, schools, business offices, etc. It is the most common type of desktop computers with processing and storage units along with various input and output devices.
- **Notebook computers:** These are compact and portable versions of PC.
- **Work stations:** These have high resolution input/output (I/O) graphics capability, but with same dimensions as that of desktop computer. These are used in engineering applications of interactive design work.
- **Enterprise systems:** These are used for business data processing in medium to large corporations that require much more computing power and storage capacity than work stations. Internet associated with servers have become a dominant worldwide source of all types of information.
- **Super computers:** These are used for large scale numerical calculations required in the applications like weather forecasting etc.

- A computer consists of CPU (central processing unit), memory, and I/O components, with one or more modules of each type.
- These components are interconnected in some fashion to achieve the basic function of the computer, which is to execute programs.
- Thus, at a top level, we can characterize a computer system by describing
 - (1) the external behavior of each component, that is, the data and control signals that it exchanges with other components and
 - (2) the interconnection structure and the controls required to manage the use of the interconnection structure.

- A computer design referred to as the von Neumann architecture is based on three key concepts:
 - Data and instructions are stored in a single read-write memory.
 - The contents of this memory are addressable by location, without regard to the type of data contained there.
 - Execution occurs in a sequential fashion (unless explicitly modified) from one instruction to the next.
- There is a small set of basic logic components that can be combined in various ways to store binary data and perform arithmetic and logical operations on that data. The process of connecting the various components in the desired configuration as a form of programming is termed a *hardwired program* where the resulting “program” is in the form of hardware.
- Instead of rewiring the hardware for each new program, all we need to do is provide a new sequence of codes. Each code is, in effect, an instruction, and part of the hardware interprets each instruction and generates control signals. To distinguish this new method of programming, a sequence of codes or instructions is called *software*.

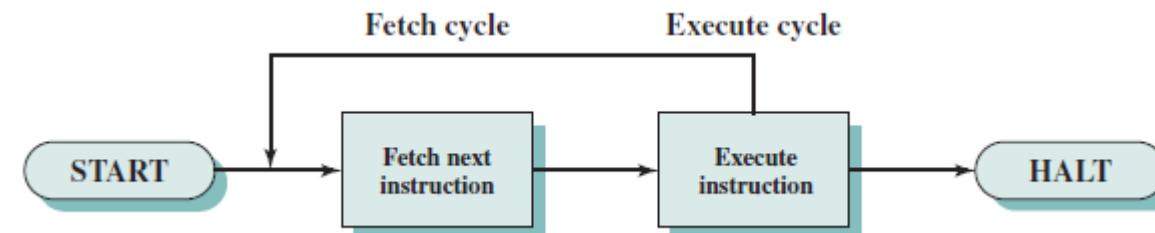


- Data and instructions must be put into the system. For this we need some sort of input module. This module contains basic components for accepting data and instructions in some form and converting them into an internal form of signals usable by the system. A means of reporting results is needed, and this is in the form of an output module. Taken together, these are referred to as ***I/O components***.
- There must be a place to store temporarily both instructions and data. That module is called ***memory, or main memory***, to distinguish it from external storage or peripheral devices.
- The CPU exchanges data with memory. For this purpose, it typically makes use of ***two internal (to the CPU) registers***: a ***memory address register (MAR)***, which specifies the address in memory for the next read or write, and a ***memory buffer register (MBR)***, which contains the data to be written into memory or receives the data read from memory. Similarly, an ***I/O address register (I/OAR)*** specifies a particular I/O device. An ***I/O buffer (I/OBR)*** register is used for the exchange of data between an I/O module and the CPU.



COMPUTER FUNCTION

- The basic function performed by a computer is execution of a program, which consists of a set of instructions stored in memory.
- The processor does the actual work by executing instructions specified in the program.
- In its simplest form, instruction processing consists of two steps:
 - The processor reads (fetches) instructions from memory one at a time and executes each instruction.
 - Program execution consists of repeating the process of instruction fetch and instruction execution.
- The processing required for a single instruction is called an **instruction cycle** consisting of two steps are referred to as **the fetch cycle** and **the execute cycle**.

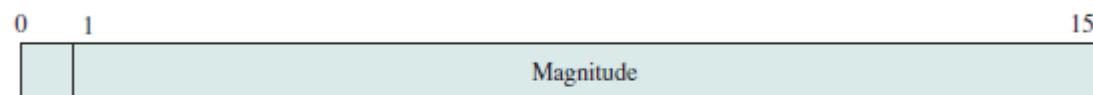


- At the beginning of each instruction cycle, the processor **fetches** an instruction from memory. In a typical processor, a register called the program counter (PC) holds the address of the instruction to be fetched next. Unless told otherwise, the processor always increments the PC after each instruction fetch so that it will fetch the next instruction in sequence (i.e., the instruction located at the next higher memory address).
- The fetched instruction is loaded into a register in the processor known as the *instruction register* (IR). The instruction contains bits that specify the action the processor is to take. The processor interprets the instruction and performs the required action.

- In general, these actions fall into four categories:
 - **Processor-memory:** Data may be transferred from processor to memory or from memory to processor.
 - **Processor-I/O:** Data may be transferred to or from a peripheral device by transferring between the processor and an I/O module.
 - **Data processing:** The processor may perform some arithmetic or logic operation on data.
 - **Control:** An instruction may specify that the sequence of execution be altered.
- The **execution** cycle for a particular instruction may involve more than one reference to memory. Also, instead of memory references, an instruction may specify an I/O operation.



(a) Instruction format



(b) Integer format

Program counter (PC) = Address of instruction

Instruction register (IR) = Instruction being executed

Accumulator (AC) = Temporary storage

(c) Internal CPU registers

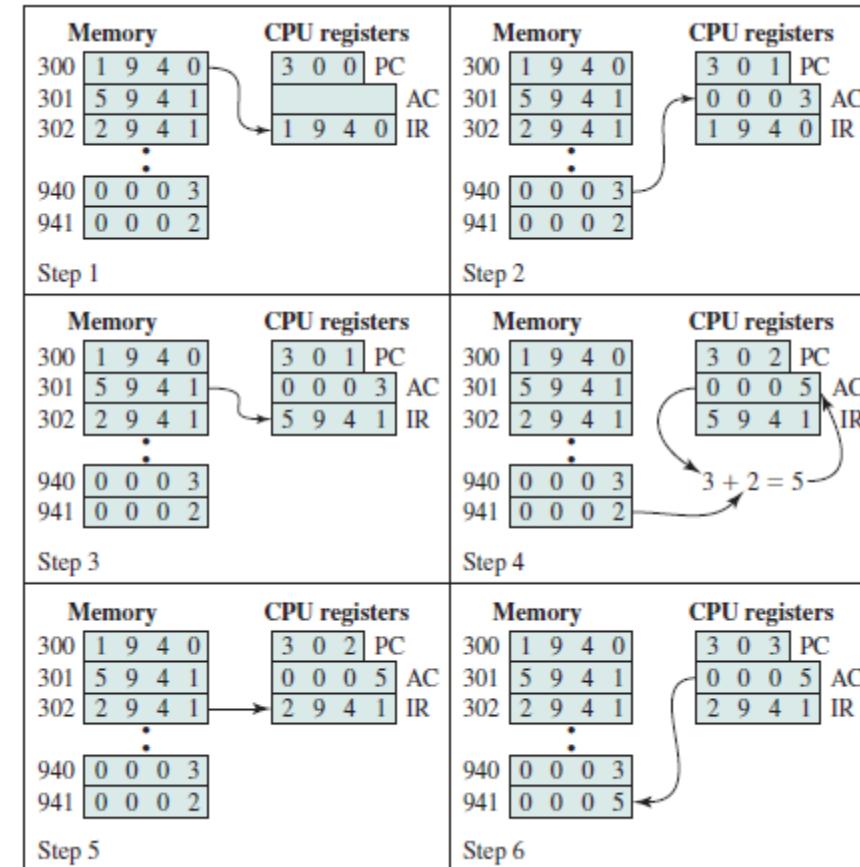
0001 = Load AC from memory

0010 = Store AC to memory

0101 = Add to AC from memory

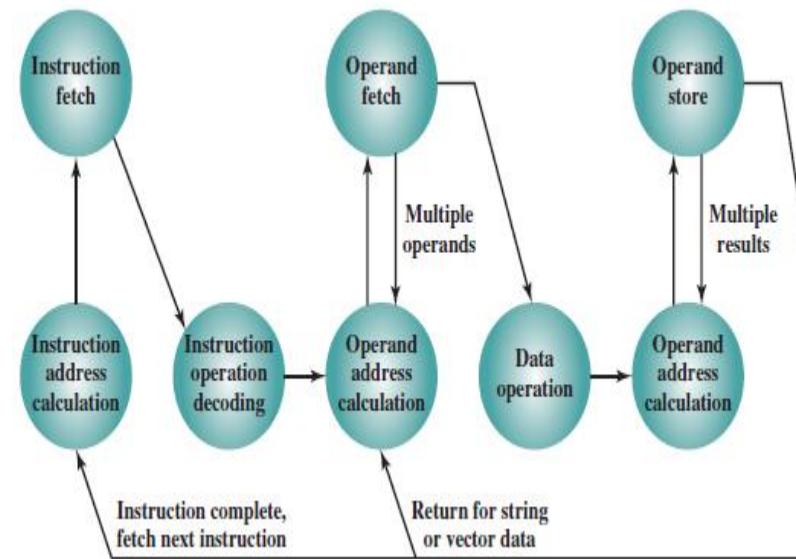
(d) Partial list of opcodes

Characteristics of a Hypothetical Machine



Example of Program Execution
(contents of memory and registers in hexadecimal)

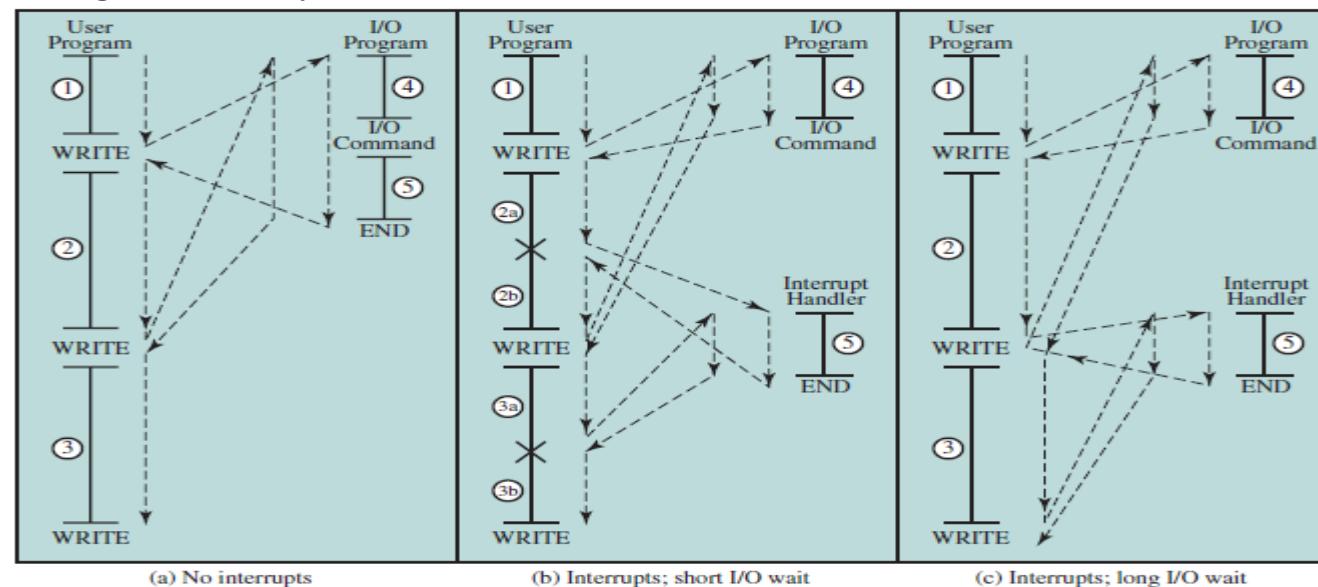
- **Instruction address calculation (iac):** Determine the address of the next instruction to be executed.
- **Instruction fetch (if):** Read instruction from its memory location into the processor.
- **Instruction operation decoding (iod):** Analyze instruction to determine type of operation to be performed and operand(s) to be used.
- **Operand address calculation (oac):** If the operation involves reference to an operand in memory or available via I/O, then determine the address of the operand.
- **Operand fetch (of):** Fetch the operand from memory or read it in from I/O.
- **Data operation (do):** Perform the operation indicated in the instruction.
- **Operand store (os):** Write the result into memory or out to I/O.



Instruction Cycle State Diagram

- INTERRUPTS

Virtually all computers provide a mechanism by which other modules (I/O, memory) may interrupt the normal processing of the processor. **Interrupts** are provided primarily as a way to improve processing efficiency.

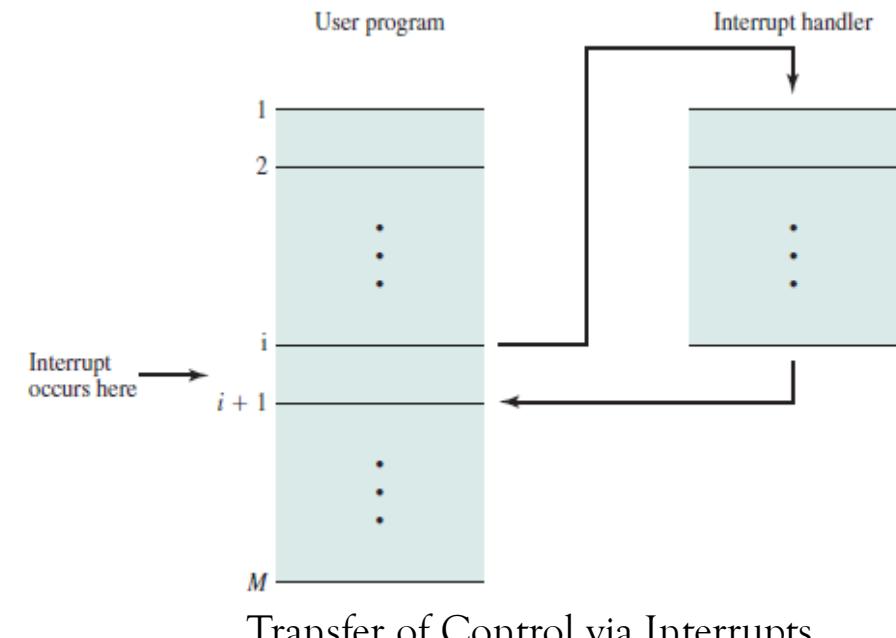


X = interrupt occurs during course of execution of user program

Program Flow of Control without and with Interrupts

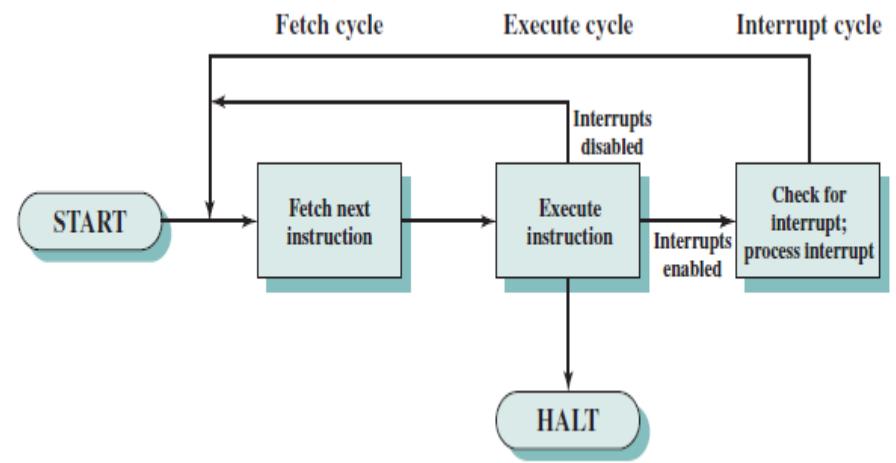
- INTERRUPTS AND THE INSTRUCTION CYCLE

With interrupts, the processor can be engaged in executing other instructions while an I/O operation is in progress. When the external device becomes ready to be serviced—that is, when it is ready to accept more data from the processor—the I/O module for that external device sends an interrupt request signal to the processor. The processor responds by suspending operation of the current program, branching off to a program to service that particular I/O device, known as an **interrupt handler**, and resuming the original execution after the device is serviced.



Transfer of Control via Interrupts

- To accommodate interrupts, an interrupt cycle is added to the instruction cycle.
- In the interrupt cycle, the processor checks to see if any interrupts have occurred, indicated by the presence of an interrupt signal.
- If no interrupts are pending, the processor proceeds to the fetch cycle and fetches the next instruction of the current program.
- If an interrupt is pending, the processor does the following:
 - It suspends execution of the current program being executed and saves its context. This means saving the address of the next instruction to be executed (current contents of the program counter) and any other data relevant to the processor's current activity.
 - It sets the program counter to the starting address of an interrupt handler routine.



- MULTIPLE INTERRUPTS

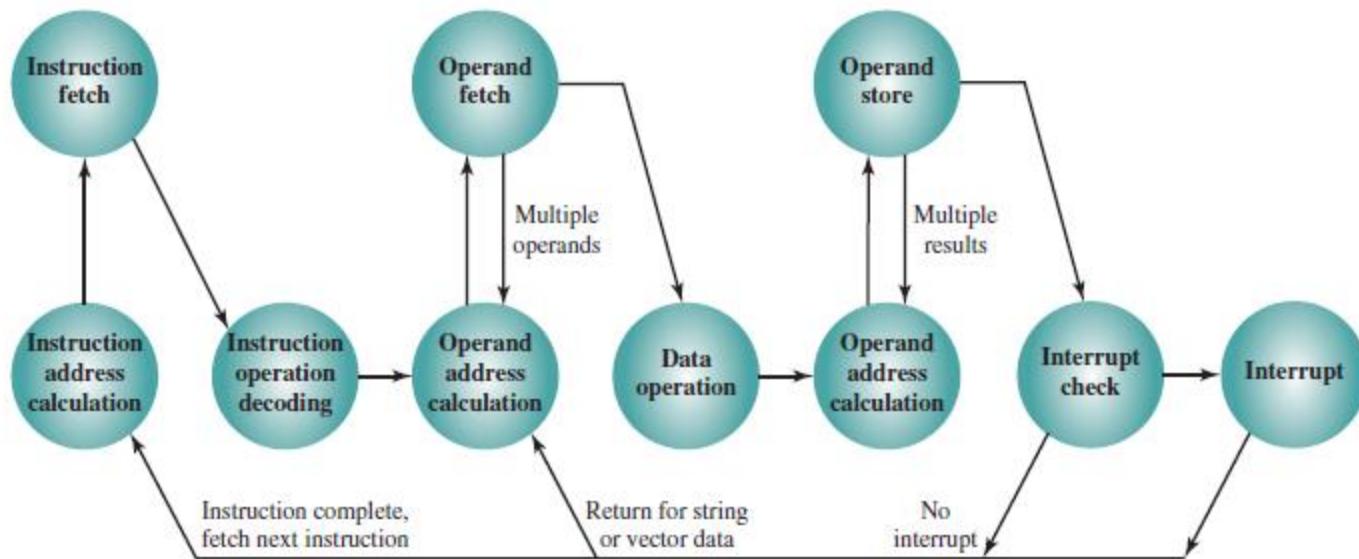
For some scenarios in computer, multiple interrupts can occur. Two approaches can be taken to dealing with multiple interrupts.

1. The first is to disable interrupts while an interrupt is being processed. A **disabled interrupt** simply means that the processor can and will ignore that interrupt request signal.

If an interrupt occurs during this time, it generally remains pending and will be checked by the processor after the processor has enabled interrupts. Thus, when a user program is executing and an interrupt occurs, interrupts are disabled immediately. After the interrupt handler routine completes, interrupts are enabled before resuming the user program, and the processor checks to see if additional interrupts have occurred. This approach is nice and simple, as interrupts are handled in strict sequential order.

The drawback to the preceding approach is that it does not take into account relative priority or time-critical needs.

2. A second approach is to define priorities for interrupts and to allow an interrupt of higher priority to cause a lower-priority interrupt handler to be itself interrupted.



Instruction Cycle State Diagram, with Interrupts

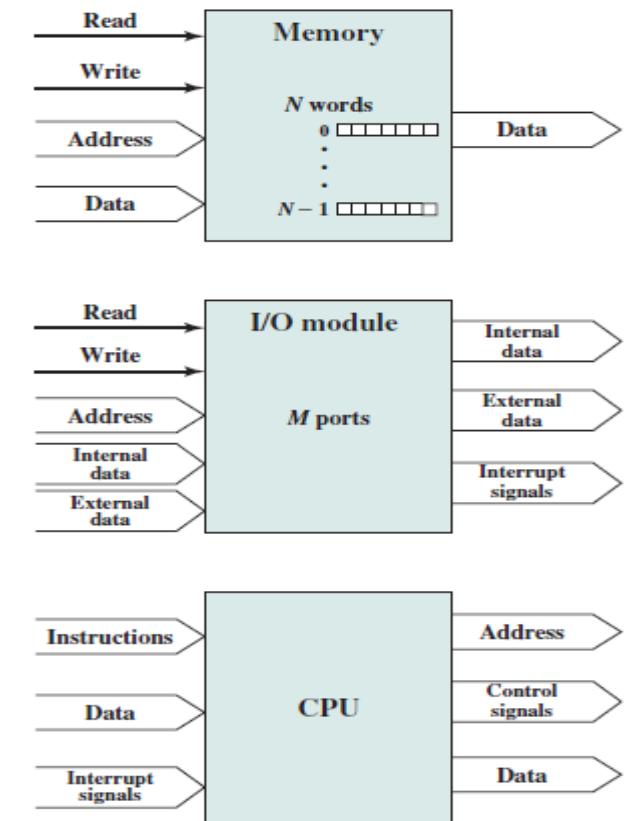
- I/O Function

An I/O module (e.g., a disk controller) can exchange data directly with the processor. Just as the processor can initiate a read or write with memory, designating the address of a specific location, the processor can also read data from or write data to an I/O module. In this latter case, the processor identifies a specific device that is controlled by a particular I/O module. Thus, an instruction sequence is similar in form of program execution to that with I/O instructions rather than memory-referencing instructions.

In some cases, it is desirable to allow I/O exchanges to occur directly with memory. In such a case, the processor grants to an I/O module the authority to read from or write to memory, so that the I/O-memory transfer can occur without tying up the processor. During such a transfer, the I/O module issues read or write commands to memory, relieving the processor of responsibility for the exchange. This operation is known as **direct memory access (DMA)**.

INTERCONNECTION STRUCTURES

- A computer consists of a set of components or modules of three basic types (*processor, memory, I/O*) that communicate with each other.
- The collection of paths connecting the various modules is called the *interconnection structure*. The design of this structure will depend on the exchanges that must be made among modules.
- The interconnection structure must support the following types of transfers:
 - **Memory to processor:** The processor reads an instruction or a unit of data from memory.
 - **Processor to memory:** The processor writes a unit of data to memory.
 - **I/O to processor:** The processor reads data from an I/O device via an I/O module.
 - **Processor to I/O:** The processor sends data to the I/O device.
 - **I/O to or from memory:** For these two cases, an I/O module is allowed to exchange data directly with memory, without going through the processor, using direct memory access.

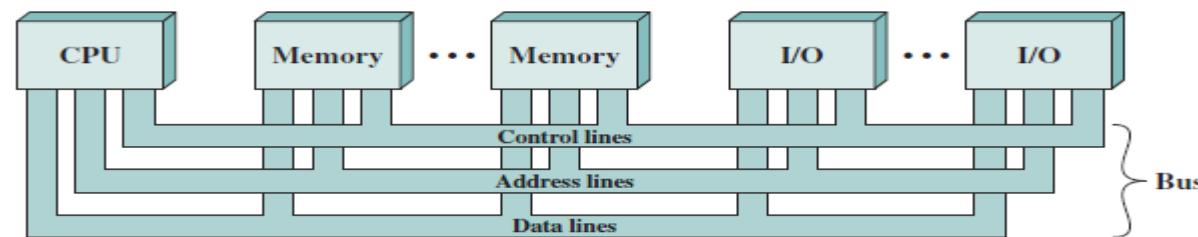


BUS INTERCONNECTION

- A bus is a communication pathway connecting two or more devices.
- A key characteristic of a bus is that it is a shared transmission medium where multiple devices connect to the bus, and a signal transmitted by any one device is available for reception by all other devices attached to the bus.
- If two devices transmit during the same time period, their signals will overlap and become garbled. Thus, only one device at a time can successfully transmit.
- Computer systems contain a number of different buses that provide pathways between components at various levels of the computer system hierarchy.
- A bus that connects major computer components (processor, memory, I/O) is called a **system bus**.
- The most common computer interconnection structures are based on the use of one or more system buses.

Bus Structure

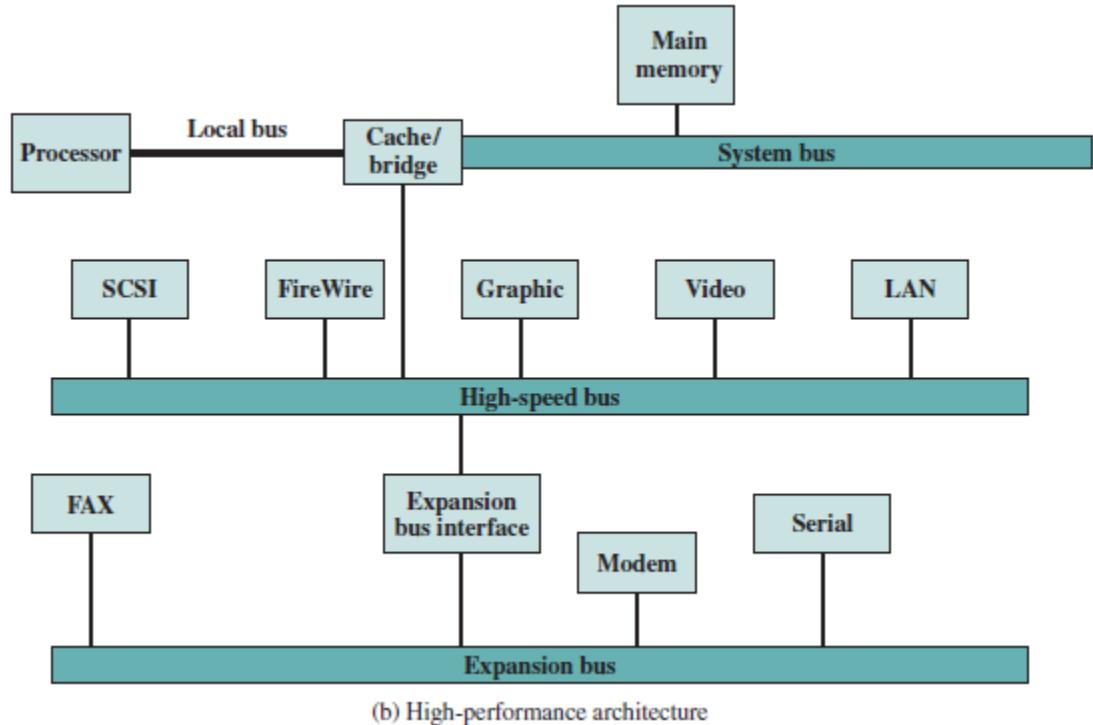
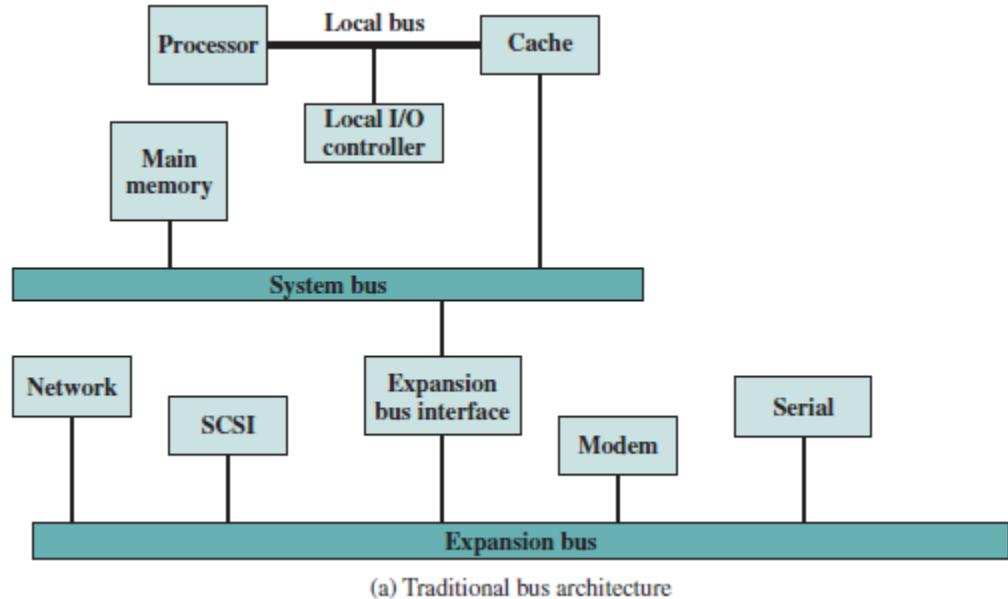
- A system bus consists, typically, of from about fifty to hundreds of separate lines.
- Each line is assigned a particular meaning or function.
- Although there are many different bus designs, on any bus the lines can be classified into three functional groups:
 - The **data lines** provide a path for moving data among system modules. These lines, collectively, are called the **data bus**. The data bus may consist of 32, 64, 128, or even more separate lines, the number of lines being referred to as the width of the data bus.
 - The **address lines** are used to designate the source or destination of the data on the data bus. The width of the **address bus** determines the maximum possible memory capacity of the system.
 - The **control lines** are used to control the access to and the use of the data and address lines. The data and address lines are shared by all components, their use is being controlled by **control bus**.



- Control signals transmit both *command* and *timing information* among system modules.
- Timing signals indicate the validity of data and address information. Command signals specify operations to be performed.
- Typical control lines include:
 - **Memory write:** causes data on the bus to be written into the addressed location
 - **Memory read:** causes data from the addressed location to be placed on the bus
 - **I/O write:** causes data on the bus to be output to the addressed I/O port
 - **I/O read:** causes data from the addressed I/O port to be placed on the bus
 - **Transfer ACK:** indicates that data have been accepted from or placed on the bus
 - **Bus request:** indicates that a module needs to gain control of the bus
 - **Bus grant:** indicates that a requesting module has been granted control of the bus
 - **Interrupt request:** indicates that an interrupt is pending
 - **Interrupt ACK:** acknowledges that the pending interrupt has been recognized
 - **Clock:** is used to synchronize operations
 - **Reset:** initializes all modules.

Multiple-bus Hierarchies

- If a great number of devices are connected to the bus, performance will suffer.
 - In general, the more devices attached to the bus, the greater the bus length and hence the greater the propagation delay.
 - The bus may become a bottleneck as the aggregate data transfer demand approaches the capacity of the bus.
- The traditional bus architecture is reasonably efficient but begins to break down as higher and higher performance is seen in the I/O devices. In response to these growing demands, a common approach taken by industry is to build a highspeed bus that is closely integrated with the rest of the system, requiring only a bridge between the processor's bus and the high-speed bus. This arrangement is sometimes known as a **mezzanine architecture**.
- The advantage of this arrangement is that the high-speed bus brings high demand devices into closer integration with the processor and at the same time is independent of the processor. Thus, differences in processor and high-speed bus speeds and signal line definitions are tolerated. Changes in processor architecture do not affect the high-speed bus, and vice versa.



Elements of Bus Design

- **Bus Types**

Bus lines can be separated into two generic types: *dedicated* and *multiplexed*.

Physical dedication refers to the use of multiple buses, each of which connects only a subset of modules.

- **Method of Arbitration**

In all but the simplest systems, more than one module may need control of the bus where only one unit at a time can successfully transmit over the bus to which some method of arbitration is needed.

The various methods can be roughly classified as being either centralized arbitration or distributed arbitration. In a centralized scheme, a single hardware device, referred to as a bus controller or arbiter, is responsible for allocating time on the bus. In a distributed scheme, there is no central controller. Rather, each module contains access control logic and the modules act together to share the bus.

Type	Bus Width
Dedicated	Address
Multiplexed	Data
Method of Arbitration	Data Transfer Type
Centralized	Read
Distributed	Write
Timing	Read-modify-write
Synchronous	Read-after-write
Asynchronous	Block

With both methods of arbitration, the purpose is to designate one device, either the processor or an I/O module, as master. The master may then initiate a data transfer (e.g., read or write) with some other device, which acts as slave for this particular exchange.

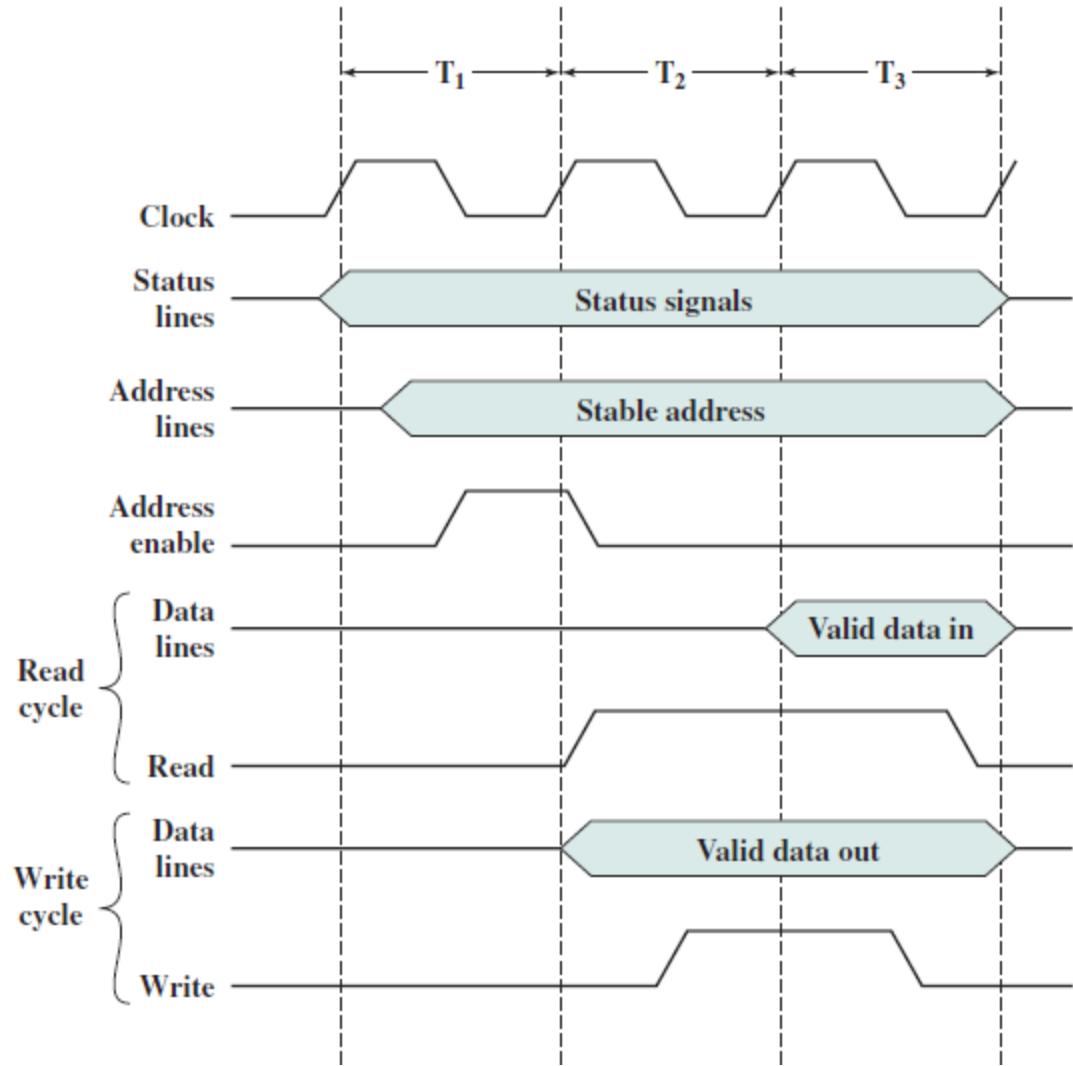
- **Timing**

Timing refers to the way in which events are coordinated on the bus. Buses use either synchronous timing or asynchronous timing.

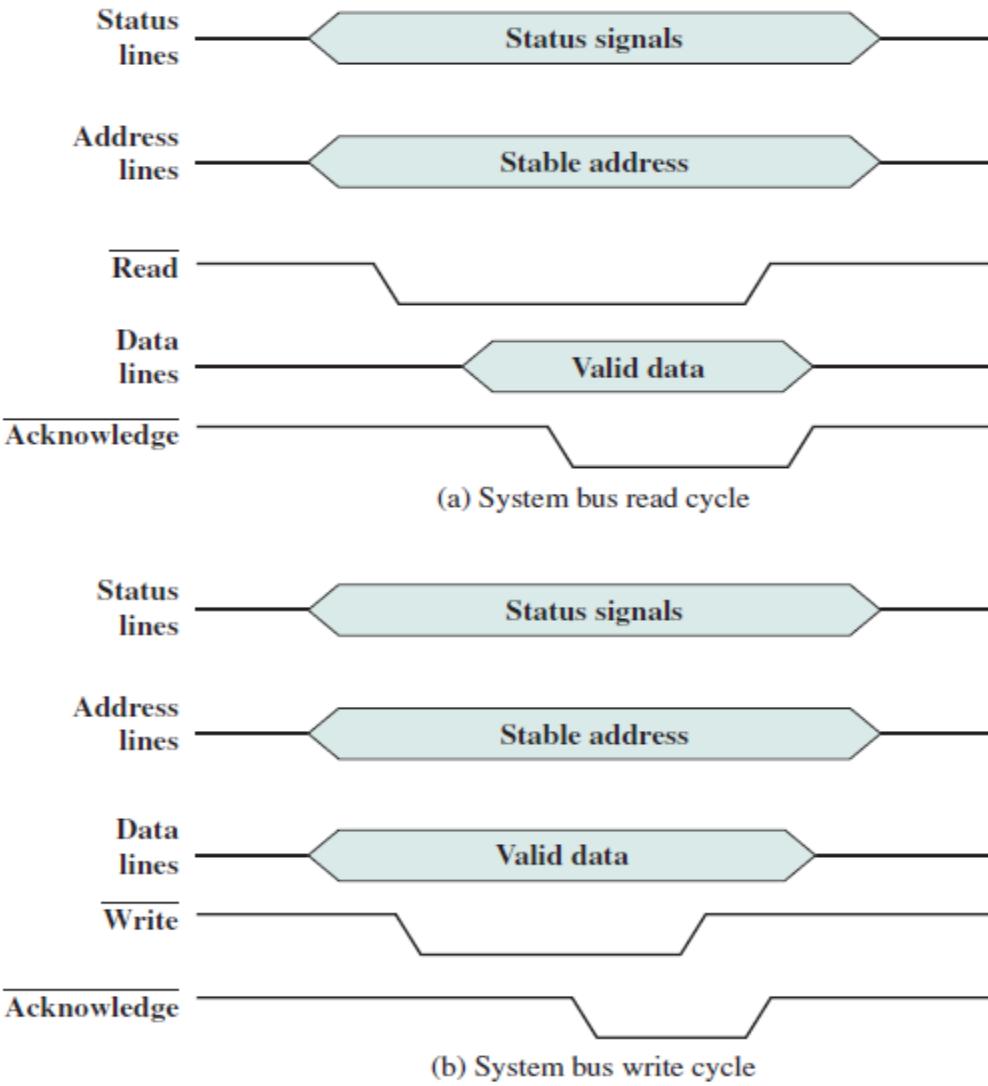
With synchronous timing, the occurrence of events on the bus is determined by a clock.

A single 1-0 transmission is referred to as a clock cycle or bus cycle and defines a time slot.

With asynchronous timing, the occurrence of one event on a bus follows and depends on the occurrence of a previous event.



Timing of Synchronous Bus Operations



Timing of Asynchronous Bus Operations

CHAPTER-3: CENTRAL PROCESSING UNIT

ADDITION AND SUBTRACTION OF SIGNED NUMBERS

When an integer binary number is positive, the sign is represented by 0 and the magnitude by a positive binary number. When the number is negative, the sign is represented by 1 but the rest of the number may be represented in one of three possible ways:

1. Signed-magnitude representation
2. Signed-1's complement representation
3. Signed 2's complement representation

Example, three different ways to represent - 14

In signed-magnitude representation 1 0001110

In signed-1's complement representation 1 11 10001

In signed-2's complement representation 1 11 10010

- The signed-magnitude system is used in ordinary arithmetic but is awkward when employed in computer arithmetic. Therefore, the signed-complement is normally used.
- The 1's complement imposes difficulties because it has two representations of 0 (+ 0 and - 0).
- The 1's complement is useful as a logical operation since the change of 1 to 0 or 0 to 1 is equivalent to a logical complement operation.

Arithmetic Addition

- The addition of two numbers in the signed-magnitude system follows the rules of ordinary arithmetic.
- If the signs are the same, we add the two magnitudes and give the sum the common sign.
- If the signs are different, we subtract the smaller magnitude from the larger and give the result the sign of the larger magnitude. For example, $(+25) + (-37) = - (37 - 25) = -12$
- This is a process that requires the comparison of the signs and the magnitudes and then performing either addition or subtraction.

- By contrast, the rule for adding numbers in the signed-2's complement system does not require a comparison or subtraction, only addition and complementation.
- The procedure is very simple and can be stated as follows: Add the two numbers, including their sign bits, and discard any carry out of the sign (leftmost) bit position.
- Note that negative numbers must initially be in 2' s complement and that if the sum obtained after the addition is negative, it is in 2's complement form.

$$\begin{array}{r} +6 \quad 00000110 \\ +13 \quad 00001101 \\ \hline +19 \quad 00010011 \end{array} \qquad \begin{array}{r} -6 \quad 11111010 \\ +13 \quad 00001101 \\ \hline +7 \quad 00000111 \end{array}$$
$$\begin{array}{r} +6 \quad 00000110 \\ -13 \quad 11110011 \\ \hline -7 \quad 11111001 \end{array} \qquad \begin{array}{r} -6 \quad 11111010 \\ -13 \quad 11110011 \\ \hline -19 \quad 11101101 \end{array}$$

Arithmetic Subtraction

- Subtraction of two signed binary numbers when negative numbers are in 2's complement form is very simple and can be stated as follows: Take the 2's complement of the subtrahend (including the sign bit) and add it to the minuend (including the sign bit). A carry out of the sign bit position is discarded.
- This procedure stems from the fact that a subtraction operation can be changed to an addition operation if the sign of the subtrahend is changed. This is demonstrated by the following relationship:

$$(\pm A) - (+ B) = (\pm A) + (-B)$$

$$(\pm A) - (c B) = (\pm A) + (+ B)$$

Consider the subtraction of $(-6) - (-13) = +7$. In binary with eight bits this is written as $11111010 - 11110011$. The subtraction is changed to addition by taking the 2's complement of the subtrahend (-13) to give $(+13)$. In binary this is $11111010 + 00001101 = 100000111$. Removing the end carry, we obtain the correct answer $00000111 (+ 7)$.

It is worth noting that binary numbers in the signed-2's complement system are added and subtracted by the same basic addition and subtraction rules as unsigned numbers. Therefore, computers need only one common hardware circuit to handle both types of arithmetic. The user or programmer must interpret the results of such addition or subtraction differently depending on whether it is assumed that the numbers are signed or unsigned.

Overflow

- When two numbers of n digits each are added and the sum occupies $n + 1$ digits, we say that an **overflow** occurred.
- An overflow is a problem in digital computers because the width of registers is finite. A result that contains $n + 1$ bits cannot be accommodated in a register with a standard length of n bits.
- For this reason, many computers detect the occurrence of an overflow, and when it occurs, a corresponding flip-flop is set which can then be checked by the user.
- If two numbers are added, and they are both positive or both negative, then overflow occurs if and only if the result has the opposite sign.
- An *overflow condition* can be detected by observing the carry into the sign bit position and the carry out of the sign bit position. If these two carries are not equal, an overflow condition is produced.

Example: Two signed binary numbers, +70 and +80, are stored in two eight-bit registers. The range of numbers that each register can accommodate is from binary +127 to binary -128. Since the sum of the two numbers is +150, it exceeds the capacity of an eight-bit register.

$$\begin{array}{r}
 \text{carries: } 0 \ 1 \\
 +70 \quad 0 \ 1000110 \\
 +80 \quad 0 \ 1010000 \\
 \hline
 +150 \quad 1 \ 0010110
 \end{array}$$

$$\begin{array}{r}
 \text{carries: } 1 \ 0 \\
 -70 \quad 1 \ 0111010 \\
 -80 \quad 1 \ 0110000 \\
 \hline
 -150 \quad 0 \ 1101010
 \end{array}$$

$\begin{array}{r} 1001 = -7 \\ +0101 = 5 \\ \hline 1110 = -2 \end{array}$	$\begin{array}{r} 1100 = -4 \\ +0100 = 4 \\ \hline 10000 = 0 \end{array}$
(a) $(-7) + (+5)$	(b) $(-4) + (+4)$
$\begin{array}{r} 0011 = 3 \\ +0100 = 4 \\ \hline 0111 = 7 \end{array}$	$\begin{array}{r} 1100 = -4 \\ +1111 = -1 \\ \hline 11011 = -5 \end{array}$
(c) $(+3) + (+4)$	(d) $(-4) + (-1)$
$\begin{array}{r} 0101 = 5 \\ +0100 = 4 \\ \hline 1001 = \text{Overflow} \end{array}$	$\begin{array}{r} 1001 = -7 \\ +1010 = -6 \\ \hline 10011 = \text{Overflow} \end{array}$
(e) $(+5) + (+4)$	(f) $(-7) + (-6)$

Addition of Numbers in Twos Complement Representation

$\begin{array}{r} 0010 = 2 \\ +1001 = -7 \\ \hline 1011 = -5 \end{array}$	$\begin{array}{r} 0101 = 5 \\ +1110 = -2 \\ \hline 10011 = 3 \end{array}$
(a) $M = 2 = 0010$ $S = 7 = 0111$ $-S = 1001$	(b) $M = 5 = 0101$ $S = 2 = 0010$ $-S = 1110$
$\begin{array}{r} 1011 = -5 \\ +1110 = -2 \\ \hline 11001 = -7 \end{array}$	$\begin{array}{r} 0101 = 5 \\ +0010 = 2 \\ \hline 0111 = 7 \end{array}$
(c) $M = -5 = 1011$ $S = 2 = 0010$ $-S = 1110$	(d) $M = 5 = 0101$ $S = -2 = 1110$ $-S = 0010$
$\begin{array}{r} 0111 = 7 \\ +0111 = 7 \\ \hline 1110 = \text{Overflow} \end{array}$	$\begin{array}{r} 1010 = -6 \\ +1100 = -4 \\ \hline 10110 = \text{Overflow} \end{array}$
(e) $M = 7 = 0111$ $S = -7 = 1001$ $-S = 0111$	(f) $M = -6 = 1010$ $S = 4 = 0100$ $-S = 1100$

Subtraction of Numbers in Twos Complement Representation ($M - S$)

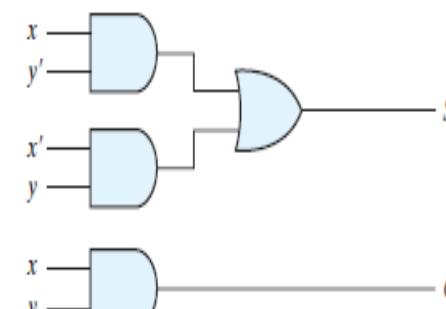
Half Adder

- The half adder circuit needs two binary inputs and two binary outputs. The input variables designate the augend and addend bits; the output variables produce the sum and carry. We assign symbols x and y to the two inputs and S (for sum) and C (for carry) to the outputs. The C output is 1 only when both inputs are 1. The S output represents the least significant bit of the sum.
- The simplified sum-of-products expressions are:

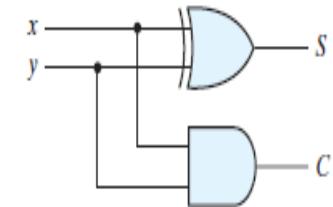
$$S = x'y + xy'$$

$$C = xy$$

		Half Adder	
x	y	C	S
0	0	0	0
0	1	0	1
1	0	0	1
1	1	1	0



(a) $S = xy' + x'y$
 $C = xy$

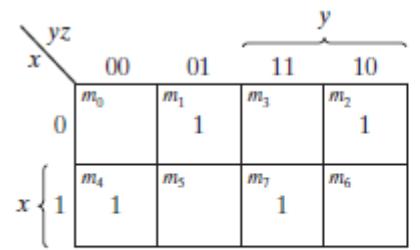
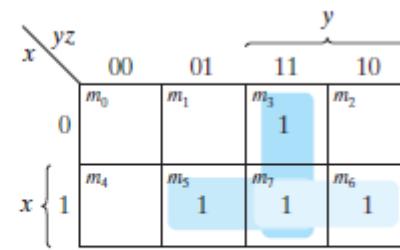
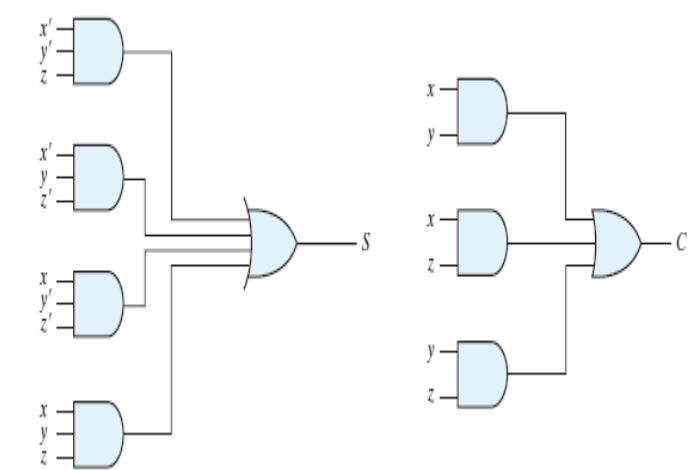


(b) $S = x \oplus y$
 $C = xy$

Full Adder

- Addition of n-bit binary numbers requires the use of a full adder, and the process of addition proceeds on a bit-by-bit basis, right to left, beginning with the least significant bit. A full adder is a combinational circuit that forms the arithmetic sum of three bits. It consists of three inputs and two outputs.
- The simplified expressions are: $S = x'y'z + x'yz' + xy'z' + xyz$
 $C = xy + xz + yz$

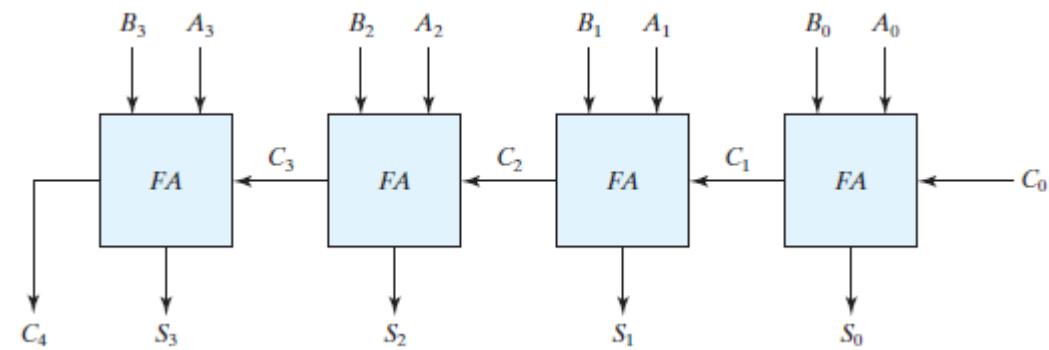
Full Adder			C	S
x	y	z		
0	0	0	0	0
0	0	1	0	1
0	1	0	0	1
0	1	1	1	0
1	0	0	0	1
1	0	1	1	0
1	1	0	1	0
1	1	1	1	1

(a) $S = x'y'z + x'yz' + xy'z' + xyz$ (b) $C = xy + xz + yz$ 

Binary Adder

- A binary adder is a digital circuit that produces the arithmetic sum of two binary numbers.
- It can be constructed with full adders connected in cascade, with the output carry from each full adder connected to the input carry of the next full adder in the chain.
- Addition of n-bit numbers requires a chain of n full adders or a chain of one-half adder and n-1 full adders.
- Example, 4-bit Ripple Carry Adder

Subscript <i>i</i>:	3	2	1	0	
Input carry	0	1	1	0	C_i
Augend	1	0	1	1	A_i
Addend	0	0	1	1	B_i
Sum	1	1	1	0	S_i
Output carry	0	0	1	1	C_{i+1}



Carry Propagation

- As in any combinational circuit, the signal must propagate through the gates before the correct output sum is available in the output terminals. The total propagation time is equal to the propagation delay of a typical gate, times the number of gate levels in the circuit. The longest propagation delay time in an adder is the time it takes the carry to propagate through the full adders.
- The number of gate levels for the carry propagation can be found from the circuit of the full adder.
- The carry propagation time is an important attribute of the adder because it limits the speed with which two numbers are added.
- Since all other arithmetic operations are implemented by successive additions, the time consumed during the addition process is critical. An obvious solution for reducing the carry propagation delay time is to employ faster gates with reduced delays. However, physical circuits have a limit to their capability.
- Another solution is to increase the complexity of the equipment in such a way that the carry delay time is reduced. There are several techniques for reducing the carry propagation time in a parallel adder. The most widely used technique employs the principle of *carry lookahead logic*.

Consider the circuit of the full adder with new binary variables:

$$P_i = A_i \oplus B_i$$

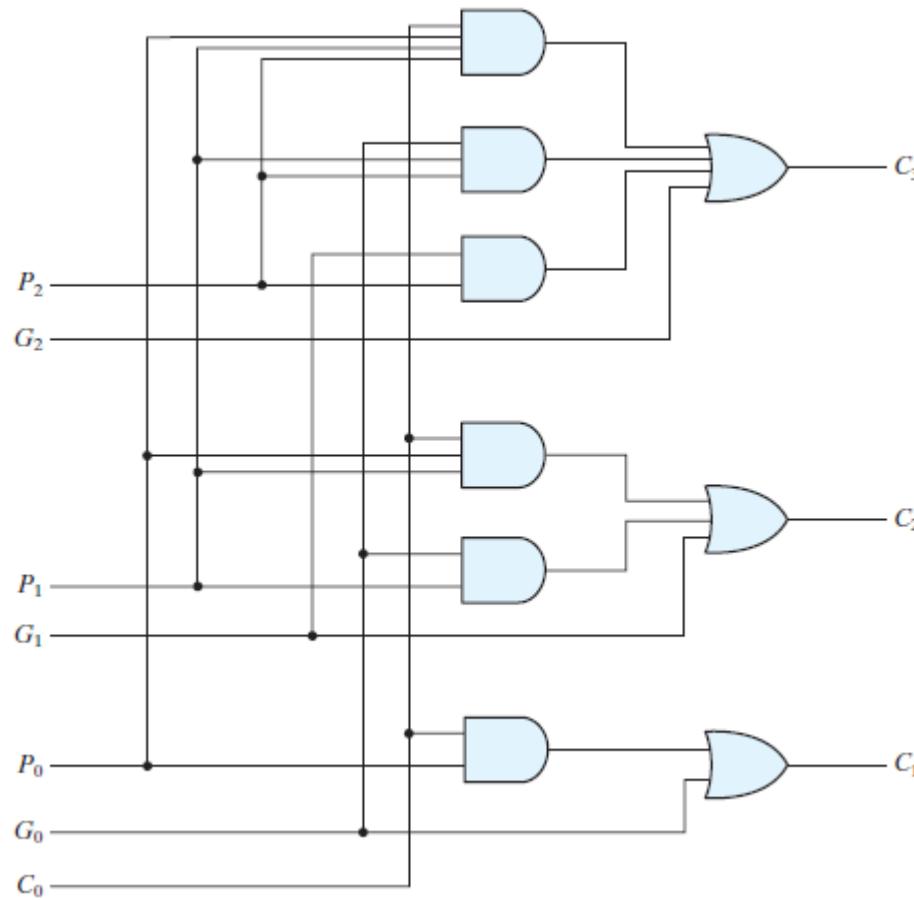
$$G_i = A_i B_i$$

the output sum and carry can respectively be expressed as

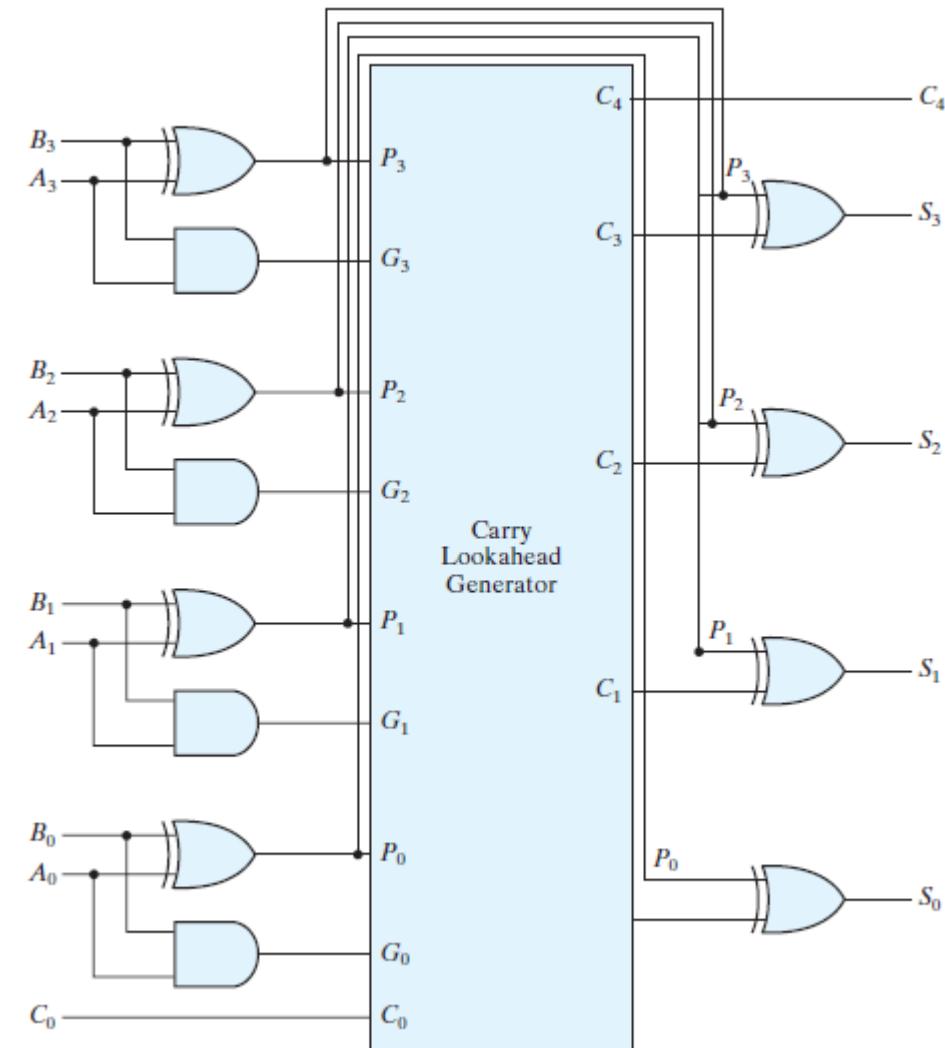
$$S_i = P_i \oplus C_i$$

$$C_{i+1} = G_i + P_i C_i$$

G_i is called a *carry generate*, and it produces a carry of 1 when both A_i and B_i are 1, regardless of the input carry C_i . P_i is called a *carry propagate*, because it determines whether a carry into stage i will propagate into stage $i + 1$ (i.e., whether an assertion of C_i will propagate to an assertion of C_{i+1}).



Logic diagram of carry lookahead generator



Four-bit adder with carry lookahead

KARNAUGH MAP (K-MAP)

- Karnaugh maps are used to simplify a Boolean function without using Boolean algebra theorems.
- A K-map is also another way to represent the truth table of a function.
- K-maps are made of cells where each cell represents a minterm. Cells marked with a one will be the minterms used for the sum of the minterms representation of a function.
- Conversely, cells marked with a zero will be used for the product of the maxterms representation.
- **Adjacent Cell**

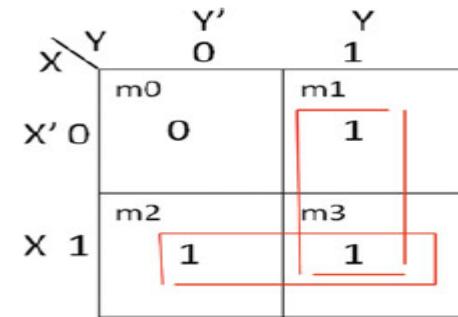
Two cells are adjacent if they differ on only one variable. The cells $X'Y'$ and $X'Y$ are adjacent because their only difference is Y' and Y . Adjacent cells can be combined in order to simplify a K-map's function.

Simplify the following function:

$$F(X, Y) = X'Y + XY' + XY$$

or

$$F(X, Y) = m_1 + m_2 + m_3$$



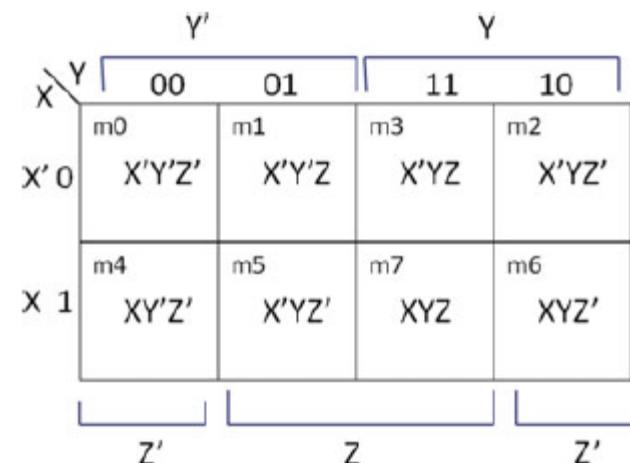
The cells m_2 and m_3 are adjacent, so they can be combined. Likewise, the cells m_1 and m_3 can be combined. By reading the map, you will have the simplified function.

Cells m_2 and m_3 are the entire row X , and cells m_1 and m_3 are the entire column Y , with the other cell being zero. Therefore,

$$F(X, Y) = X + Y$$

THREE-VARIABLE MAP

- A three-variable K-map contains eight cells, and each cell represents a minterm.
- Adjacent cells can be grouped together in a K-map; in a K-map it can combine 2 cells, 4 cells, 8 cells, and 16 cells.



		Y'		Y	
		00	01	11	10
X	Y	m0	m1	m3	m2
X' 0	0	1	1	1	1
X 1	1	0	0	0	0

		Y'		Y	
		00	01	11	10
X	Y	m0	m1	m3	m2
X' 0	0	1	0	0	1
X 1	1	1	0	0	1

		Y'		Y	
		00	01	11	10
X	Y	m0	m1	m3	m2
X' 0	0	0	0	1	1
X 1	1	1	1	1	1

		Y'		Y	
		00	01	11	10
X	Y	m0	m1	m3	m2
X' 0	0	1	1	0	0
X 1	1	0	0	1	1

		Y'		Y	
		00	01	11	10
X	Y	m0	m1	m3	m2
X' 0	0	1	1	1	1
X 1	1	1	1	1	1

(a) All adjacent cells in columns Z' and columns Y are one.

$$F(X, Y, Z) = Z' + Y'$$

(b) The cells in row X' , columns Y' are adjacent, as are the cells in row X , columns Y .

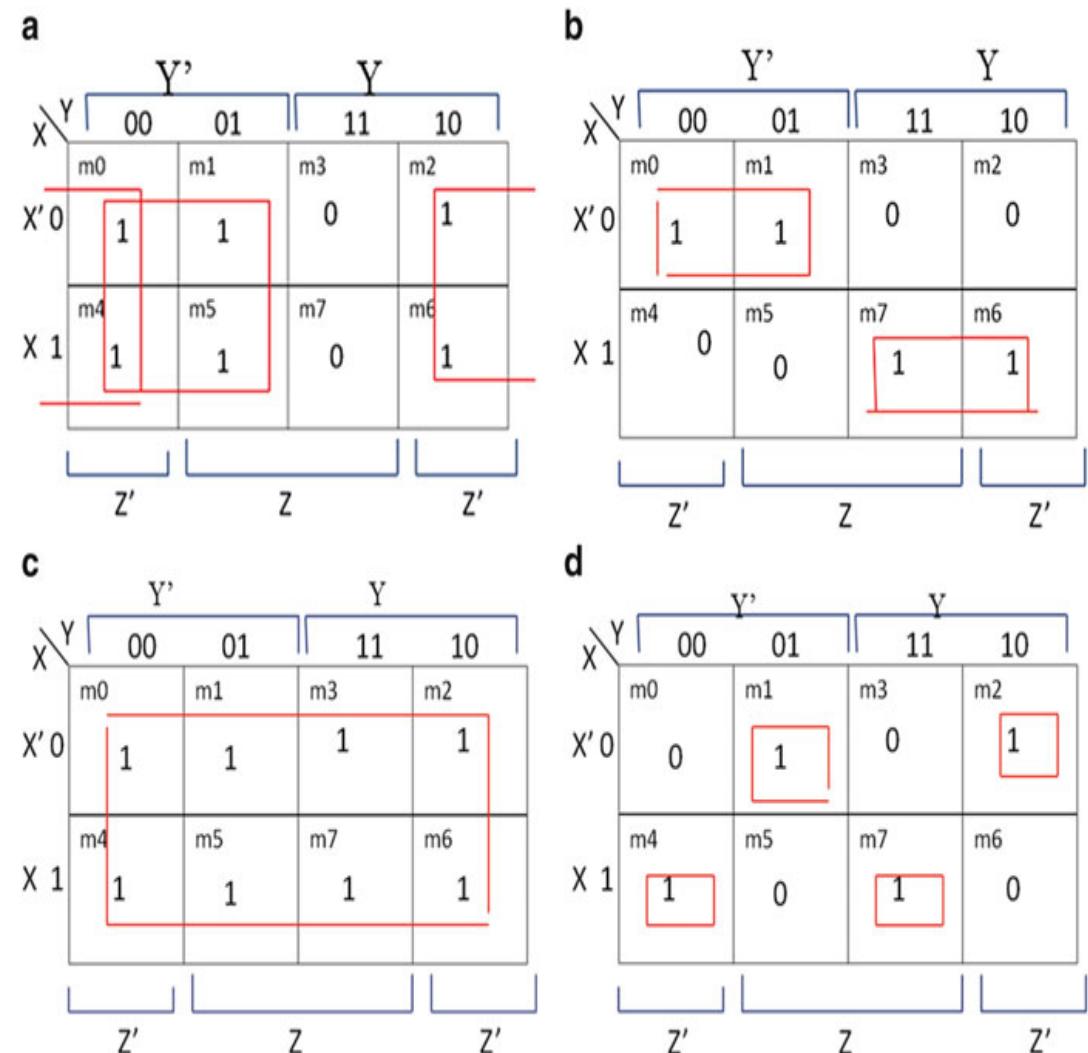
$$F(X, Y, Z) = X'Y' + XY$$

(c) All cells are one, so the function is always equal to one.

$$F(X, Y, Z) = 1$$

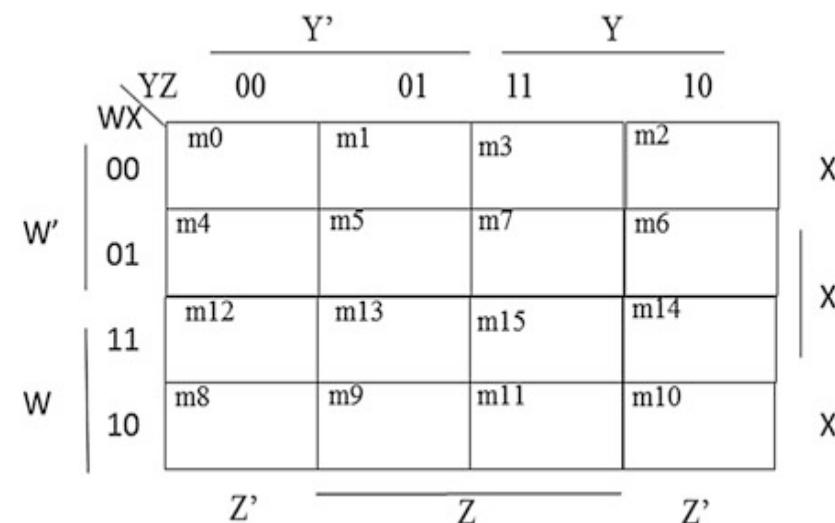
(d) Without adjacent cells to simplify the terms, the function equals the ones.

$$F(X, Y, Z) = X'Y'Z + X'YZ' + XY'Z' + XYZ$$

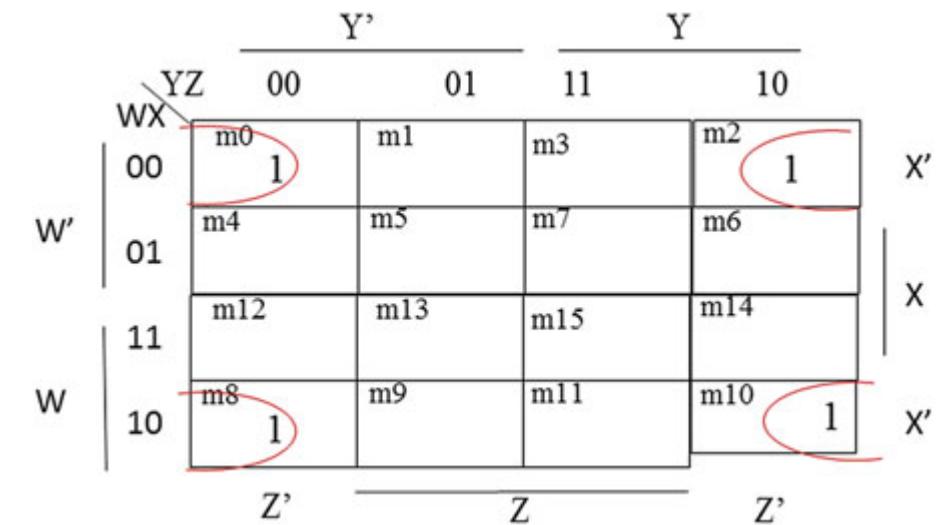


FOUR-VARIABLE MAP

- Four-variable K-maps contain 16 cells.



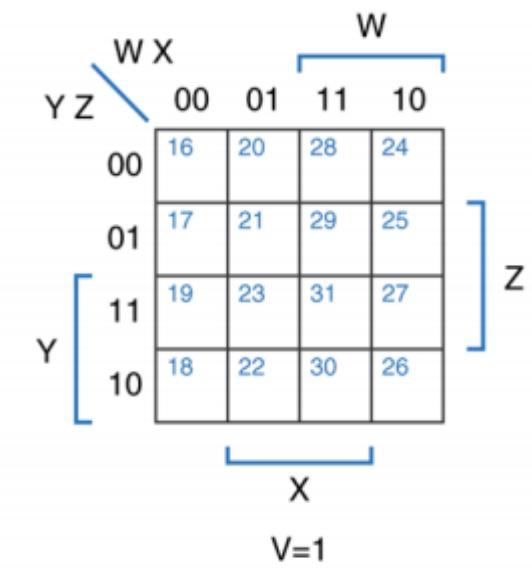
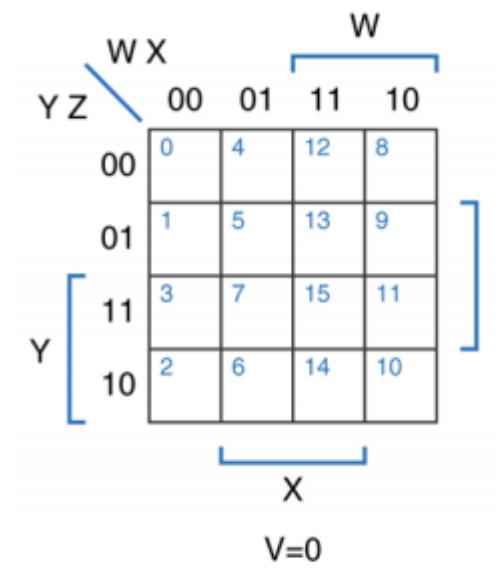
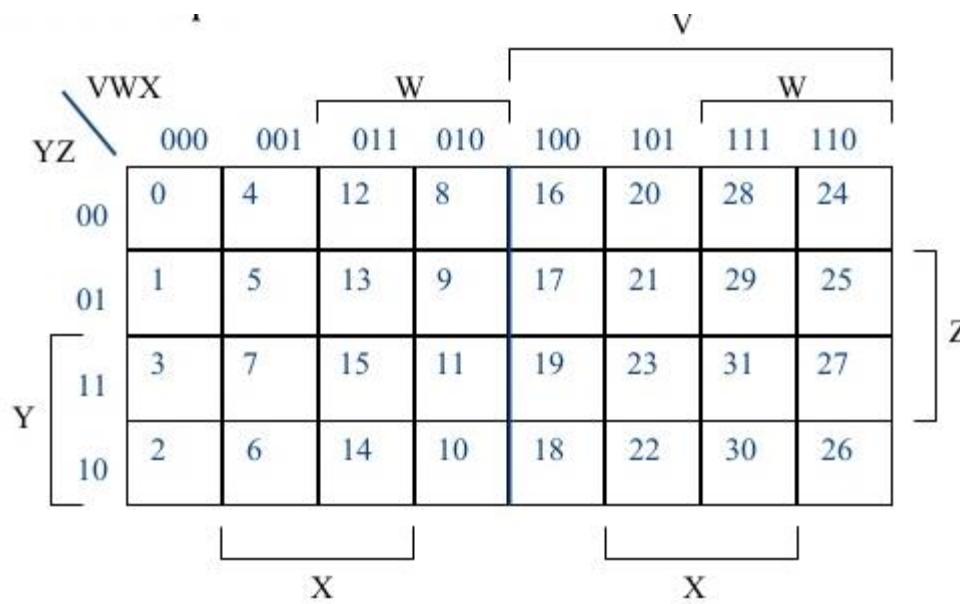
Example, K-map for $F(W,X,Y,Z) = m_0 + m_2 + m_8 + m_{10}$



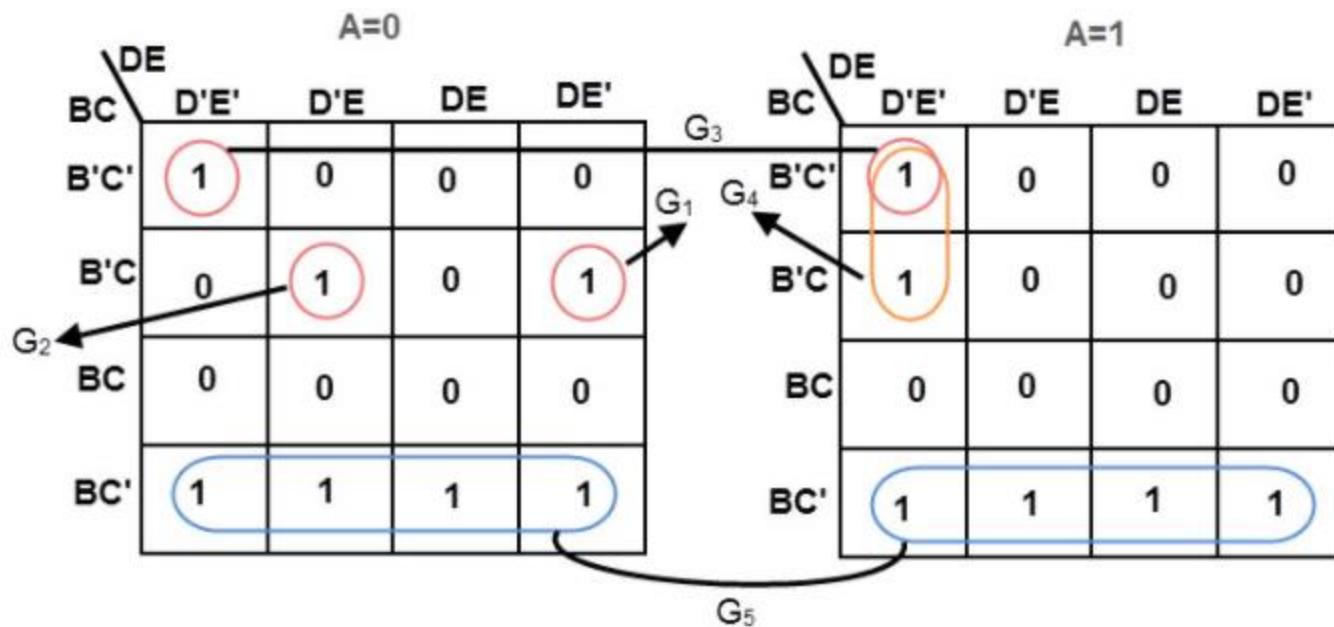
The simplified function is $F(W,X,Y,Z) = X'Z'$.

FIVE-VARIABLE MAP

- A 5-variable Boolean function can have a maximum of 32 minterms. Hence, 5-variable K-Map contains 32 cells.
- A 5-variable K-map is formed using two connected 4-variable map.



$$f(A, B, C, D, E) = \sum m(0, 5, 6, 8, 9, 10, 11, 16, 20, 42, 25, 26, 27)$$



SUM OF PRODUCTS (SOP) AND PRODUCT OF SUMS (POS)

- The sum of products of a function is its simplified sum of minterms while product of sums of a function is its simplified product of maxterms.

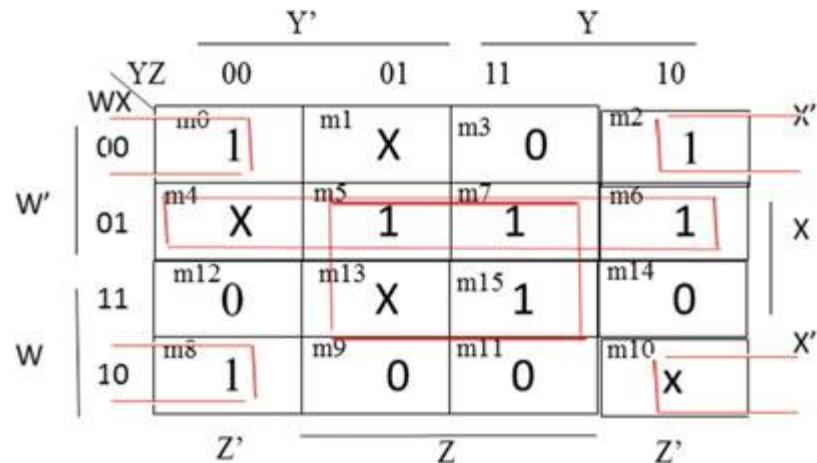
Example,

Consider function $F1(X,Y,Z) = XY' + YZ' + XZ$ which is represented by the sum of products.

Consider function $F2(X,Y,Z) = (X + Y)(Y + Z')(X + Z)$ which is represented by the product of sums.

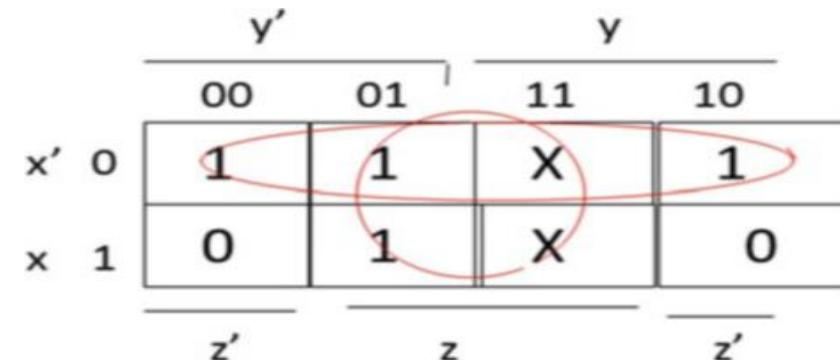
DON'T CARE CONDITIONS

- In a truth table, if certain combinations of the input variables are impossible, they are considered don't care conditions.
- These conditions are where the output of the function does not matter.
- A truth table or K-map cell marked with a “X” or “d” is a don't care term, and output will not be affected whether it is a one or zero.
- The don't care can be used to expand the adjacency of cells in a K-map to further simplify a function, since their output does not matter.
- Since a don't care can output either a zero or one, we can assume it is a one in order to expand a grouping of adjacent cells.



K-map with don't care minterms

$$F(W,X,Y,Z) = XZ + X'Z' + XW'$$



K-map for $F(X,Y,Z) = m_0 + m_1 + m_2 + m_5$ and $D(X,Y,Z) = m_3 + m_7$
 Simplified: $F(X,Y,Z) = X' + Z$.

When minterms for function F are don't care terms, the don't care function D is equal to the sum of the don't care minterm(s).

SIGNED OPERAND MULTIPLICATION

- Multiplication of binary numbers is performed in the same way as multiplication of decimal numbers.
- The multiplicand is multiplied by each bit of the multiplier, starting from the least significant bit. Each such multiplication forms a partial product.
- Successive partial products are shifted one position to the left. The final product is obtained from the sum of the partial products.

For instance, the multiplicand bits are B1 and B0, the multiplier bits are A1 and A0, and the product is C3C2C1C0. The first partial product is formed by multiplying B1B0 by A0.

The multiplication of two bits such as A0 and B0 produces a 1 if both bits are 1; otherwise, it produces a 0. This is identical to an AND operation.

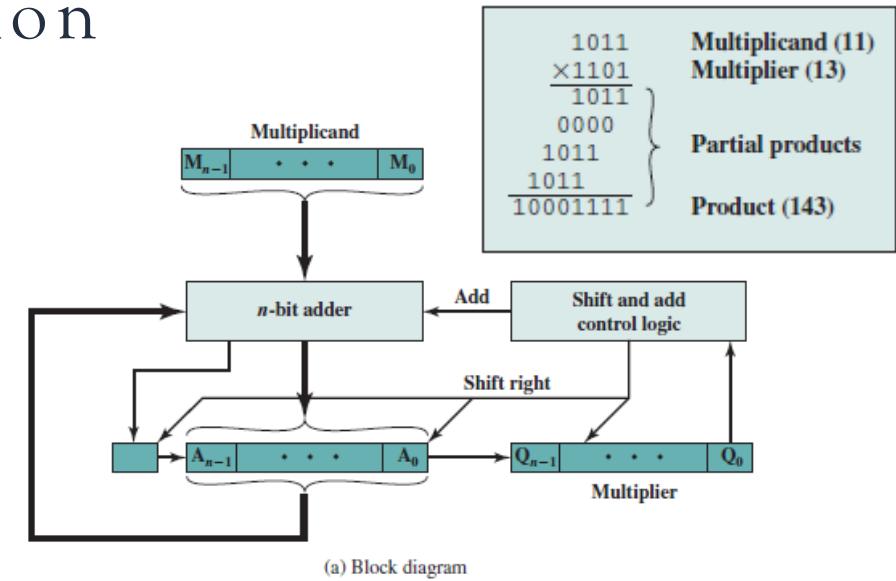
Therefore, the partial product can be implemented with AND gates. The second partial product is formed by multiplying B1B0 by A1 and shifting one position to the left.

	B_1	B_0		
	A_1	A_0		
	$\underline{A_0B_1}$	A_0B_0		
	A_1B_1	A_1B_0		
	$\underline{C_3}$	C_2	C_1	C_0

		1	1	1	0	1	1
	x	1	1	1	1	0	0
		1	1	1	1	0	1
		1	1	1	0	0	1
		1	1	1	0	1	1
		1	1	0	0	1	1
		0	1
		0	1

Twos Complement Multiplication

- If the two numbers are considered to be unsigned integers, then suppose we multiplied 11 (1011) by 13 (1101) to get 143 (10001111). If we interpret these as twos complement numbers, we have -5 (1011) times -3 (1101) equals -113 (10001111). Unfortunately, this simple scheme will not work for multiplication.
- Straightforward multiplication will not work if both the multiplicand and multiplier are negative. In fact, it will not work if either the multiplicand or the multiplier is negative.
- From the block diagram, the multiplier and multiplicand are loaded into two registers (Q and M). A third register, the A register, is also needed and is initially set to 0. There is also a 1-bit C register, initialized to 0, which holds a potential carry bit resulting from addition.



(a) Block diagram

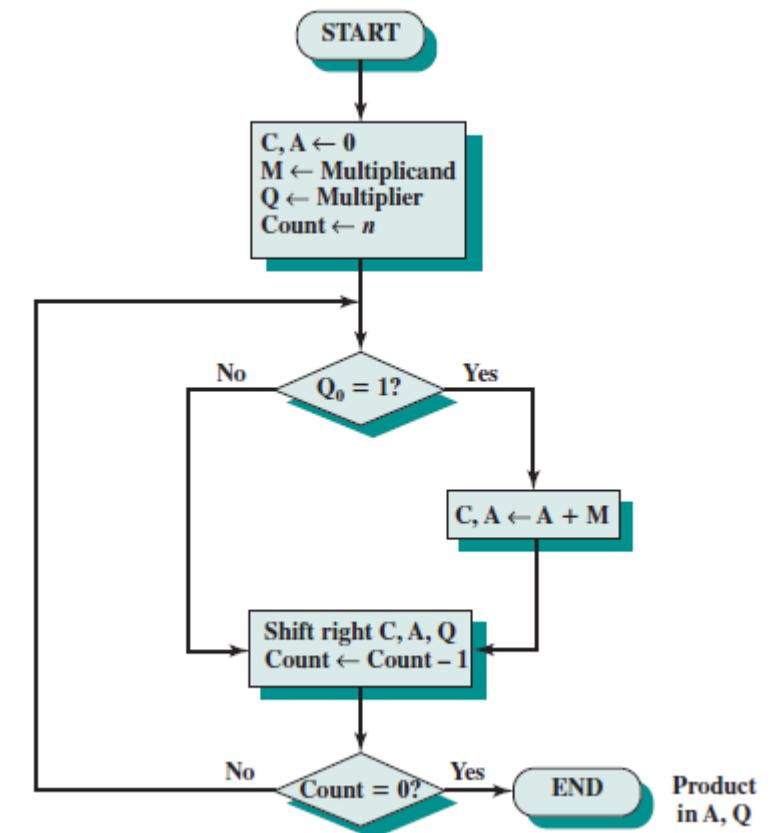
C	A	Q	M		Initial values
0	0000	1101	1011		
0	1011	1101	1011	Add	{ First cycle
0	0101	1110	1011	Shift	
0	0010	1111	1011	Shift	{ Second cycle
0	1101	1111	1011	Add	{ Third cycle
0	0110	1111	1011	Shift	
1	0001	1111	1011	Add	{ Fourth cycle
0	1000	1111	1011	shift	

(b) Example from Figure 10.7 (product in A, Q)

- Recall that any unsigned binary number can be expressed as a sum of powers of 2. Thus, $1101 = 1 * 2^3 + 1 * 2^2 + 0 * 2^1 + 1 * 2^0 = 2^3 + 2^2 + 2^0$
- Further, the multiplication of a binary number by 2^n is accomplished by shifting that number to the left n bits.
- Hence, the multiplication process is re-casted to make the generation of partial products by multiplication explicit.
- The only difference is that it recognizes that the partial products should be viewed as $2n$ -bit numbers generated from the n -bit multiplicand.
- Thus, as an unsigned integer, the 4-bit multiplicand 1011 is stored in an 8-bit word as 00001011. Each partial product (other than that for 20) consists of this number shifted to the left, with the unoccupied positions on the right filled with zeros (e.g., a shift to the left of two places yields 00101100).

$\begin{array}{r} 1011 \\ \times 1101 \\ \hline 1011 \\ 0000 \\ 1011 \\ \hline 1011 \end{array}$	Multiplicand (11) Multiplier (13)
Partial products Product (143)	

$\begin{array}{r} 1011 \\ \times 1101 \\ \hline \end{array}$	$\begin{array}{ll} 00001011 & 1011 \times 1 \times 2^0 \\ 00000000 & 1011 \times 0 \times 2^1 \\ 00101100 & 1011 \times 1 \times 2^2 \\ 01011000 & 1011 \times 1 \times 2^3 \\ \hline \end{array}$
--	--



- The problem i.e. straightforward multiplication will not work if the multiplicand is negative is that each contribution of the negative multiplicand as a partial product must be a negative number on a $2n$ -bit field; the sign bits of the partial products must line up.
- Considering the same example, if these are treated as unsigned integers, the multiplication of $9 * 3 = 27$ proceeds simply. However, if 1001 is interpreted as the twos complement value -7, then each partial product must be a negative twos complement number of $2n$ (8) bits.
- Note that this is accomplished by padding out each partial product to the left with binary 1s.
- Again, if the multiplier is negative, straightforward multiplication also will not work. The reason is that the bits of the multiplier no longer correspond to the shifts or multiplications that must take place.
- Example, the 4-bit decimal number -3 is written 1101 in twos complement. If we simply took partial products based on each bit position, we would have the following correspondence: $1101 = -(1 * 2^3 + 1 * 2^2 + 0 * 2^1 + 1 * 2^0) = -(2^3 + 2^2 + 2^0)$
- In fact, what is desired is $-(2^1 + 2^0)$. So this multiplier cannot be used directly in the manner we have been describing.

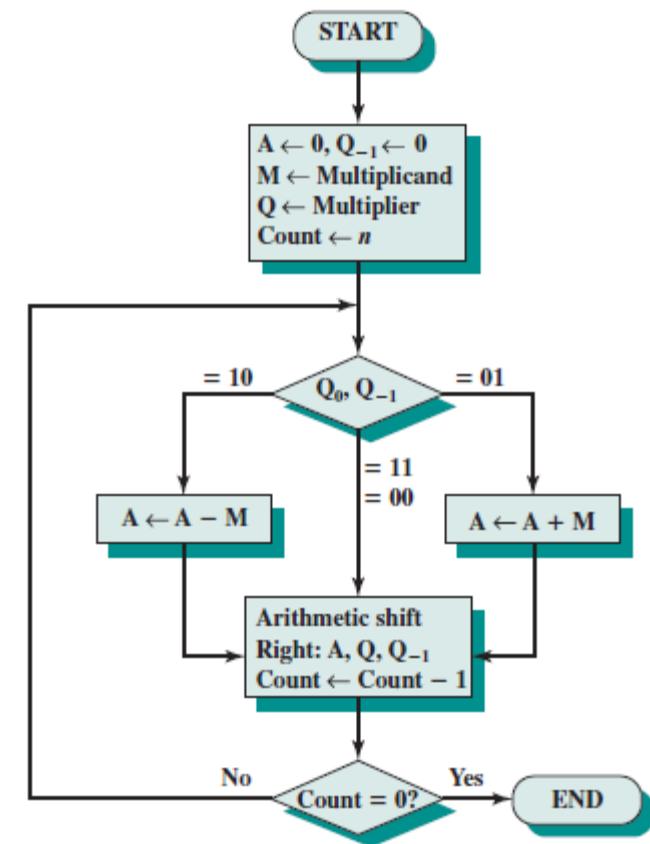
$ \begin{array}{r} 1001 \quad (9) \\ \times 0011 \quad (3) \\ \hline 00001001 \quad 1001 \times 2^0 \\ 00010010 \quad 1001 \times 2^1 \\ \hline 00011011 \quad (27) \end{array} $	$ \begin{array}{r} 1001 \quad (-7) \\ \times 0011 \quad (3) \\ \hline 11111001 \quad (-7) \times 2^0 = (-7) \\ 11110010 \quad (-7) \times 2^1 = (-14) \\ \hline 11101011 \quad (-21) \end{array} $
---	--

(a) Unsigned integers

(b) Twos complement integers

Booth's Algorithm

- To overcome the dilemma of binary multiplication, one of the solution would be to convert both multiplier and multiplicand to positive numbers, perform the multiplication, and then take the twos complement of the result if and only if the sign of the two original numbers differed. Implementers have preferred to use techniques that do not require this final transformation step.
- One of the most common of these is Booth's algorithm. This algorithm also has the benefit of speeding up the multiplication process, relative to a more straightforward approach.



PROCESSOR ORGANIZATION

- 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:
 - The register set stores intermediate data used during the execution of the instructions.
 - The arithmetic logic unit (ALU) performs the required micro-operations 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.

- A processor designed should be able to do:

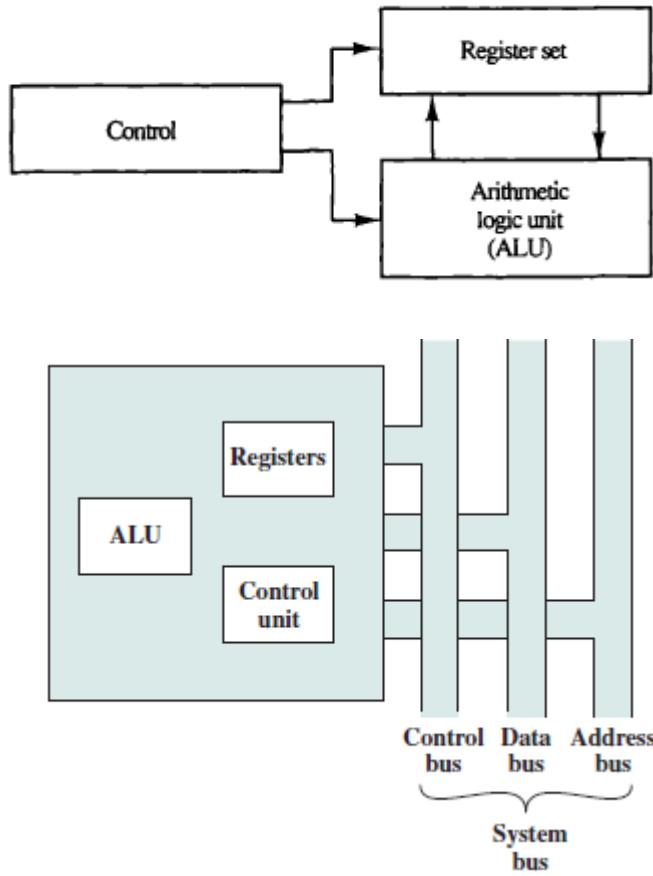
Fetch instruction: The processor reads an instruction from memory (register, cache, main memory).

Interpret instruction: The instruction is decoded to determine what action is required.

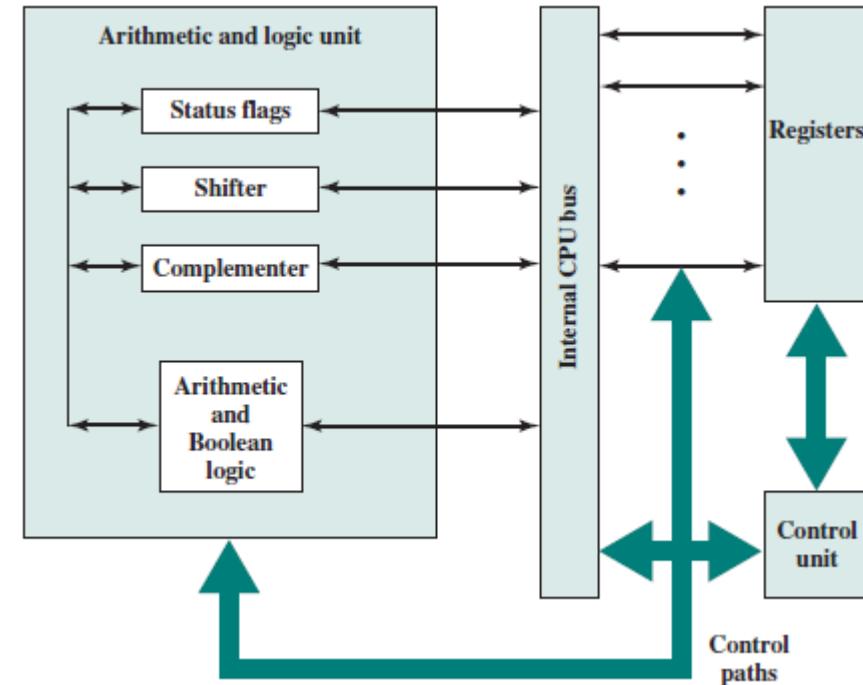
Fetch data: The execution of an instruction may require reading data from memory or an I/O module.

Process data: The execution of an instruction may require performing some arithmetic or logical operation on data.

Write data: The results of an execution may require writing data to memory or an I/O module.



The CPU with the System Bus



Internal Structure of the CPU

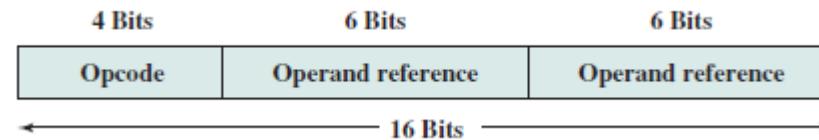
GENERAL REGISTER ORGANIZATION

- Within the processor, there is a set of registers that function as a level of memory above main memory and cache in the hierarchy. The registers in the processor perform two roles:
 - **User-visible registers:** Enable the machine- or assembly language programmer to minimize main memory references by optimizing use of registers. A user-visible register is one that may be referenced by means of the machine language that the processor executes. We can characterize these in the following categories:
 - *General purpose:* can be assigned to a variety of functions by the programmer and can contain the operand for any opcode.
 - *Data:* may be used only to hold data and cannot be employed in the calculation of an operand address.
 - *Address:* may themselves be somewhat general purpose, or they may be devoted to a particular addressing mode.
 - *Condition codes (flags):* are bits set by the processor hardware as the result of operations.
 - **Control and status registers:** Used by the control unit to control the operation of the processor and by privileged, operating system programs to control the execution of programs.

CONTROL UNIT

INSTRUCTION(S)

- Within the computer, each instruction is represented by a sequence of bits.
- The instruction is divided into fields, corresponding to the constituent elements of the instruction.



- The operation of the processor is determined by the instructions it executes, referred to as *machine instructions* or *computer instructions*.
- The collection of different instructions that the processor can execute is referred to as the processor's instruction set.

Instruction Types

- **Data processing:** Arithmetic and logic instructions
- **Data storage:** Movement of data into or out of register and or memory locations
- **Data movement:** I/O instructions
- **Control:** Test and branch instructions

Arithmetic instructions provide computational capabilities for processing numeric data.

Logic (Boolean) instructions operate on the bits of a word as bits rather than as numbers; thus, they provide capabilities for processing any other type of data the user may wish to employ. These operations are performed primarily on data in processor registers. Therefore, there must be *memory instructions* for moving data between memory and the registers.

I/O instructions are needed to transfer programs and data into memory and the results of computations back out to the user.

Test instructions are used to test the value of a data word or the status of a computation.

Branch instructions are then used to branch to a different set of instructions depending on the decision made.

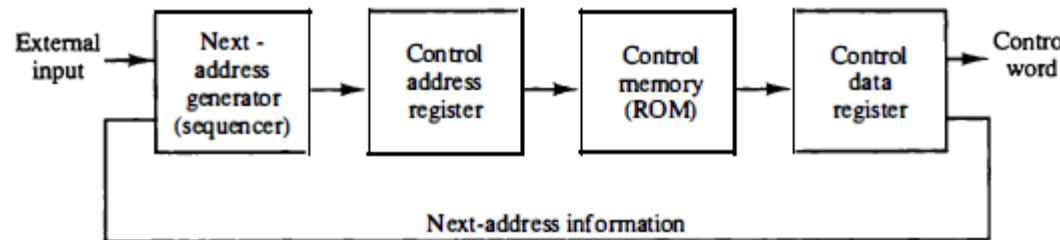
Instruction Formats

- An instruction format defines the layout of the bits of an instruction, in terms of its constituent fields.
- An instruction format must include an opcode and, implicitly or explicitly, zero or more operands. Each explicit operand is referenced using one of the addressing modes.
- For most instruction sets, more than one instruction format is used.
- 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.

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.

- 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:
 - **Single accumulator organization** (Example, ADD X: The ADD instruction in this case results in the operation $AC \leftarrow AC + M[X]$. AC is the accumulator register and M [X] symbolizes the memory word located at address X).
 - **General register organization** (Example, ADD R1, R2, R3: denotes the operation $R1 \leftarrow R2 + R3$; If $R1=R2$ then ADD R1 , R2 viz. $R1 \leftarrow R1 + R2$; MOV R1 , R2: denotes the transfer $R1 \leftarrow R2$ (or $R2 \leftarrow R1$, depending on the particular computer); ADD R1 , X: would specify the operation $R1 \leftarrow R1 + M[X]$. It has two address fields, one for register R1 and the other for the memory address X).
 - **Stack organization** (Example, PUSH X: will push the word at address X to the top of the stack).



- Operations specified by computer instructions are executed on some data stored in memory or processor registers.
- Operands residing in memory are specified by their memory address.
- 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.

Example, a CPU with 16 registers R_0 through R_{15} will have a register address field of four bits. The binary number 0101, for example, will designate register R_5 .

Hence based on these information, instruction formats are categorized as:

1. Three-Address Instructions
2. Two-Address Instructions
3. One-Address Instructions
4. Zero-Address Instructions
5. RISC (Reduced Instruction Set Computer) Instructions

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 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.

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.

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 register and assume that the AC contains the result of all operations.

Zero-Address Instructions:

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.

RISC:

Early computers had small and simple instruction sets, forced mainly by the need to minimize the hardware used to implement them. As digital hardware became cheaper with the advent of integrated circuits, computer instructions tended to increase both in number and complexity.

A computer with a large number of instructions is classified as a *complex instruction set computer, abbreviated CISC*.

A number of computer designers recommended that computers use fewer instructions with simple constructs so they can be executed much faster within the CPU without having to use memory as often. This type of computer is classified as a *reduced instruction set computer or RISC*.

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

RISC vs. CISC

SOURCE: INTERNET

CISC	RISC
Emphasis on hardware	Emphasis on software
Multiple instruction sizes and formats	Instructions of same set with few formats
Less registers	Uses more registers
More addressing modes	Fewer addressing modes
Extensive use of microprogramming	Complexity in compiler
Instructions take a varying amount of cycle time	Instructions take one cycle time
Pipelining is difficult	Pipelining is easy

RISC	CISC
Multiple register sets, often consisting of more than 256 registers	Single register set, typically 6 to 16 registers total
Three register operands allowed per instruction (e.g., add R1, R2, R3)	One or two register operands allowed per instruction (e.g., add R1, R2)
Parameter passing through efficient on-chip register windows	Parameter passing through inefficient off-chip memory
Single-cycle instructions (except for load and store)	Multiple-cycle instructions
Hardwired control	Microprogrammed control
Highly pipelined	Less pipelined
Simple instructions that are few in number	Many complex instructions
Fixed length instructions	Variable length instructions
Complexity in compiler	Complexity in microcode
Only load and store instructions can access memory	Many instructions can access memory
Few addressing modes	Many addressing modes

Execution Of A Complete Instruction

- Consider the instruction Add (R3),R1 which adds the contents of a memory-location pointed by R3 to register R1.
- Executing this instruction requires the following actions:
 - 1) Fetch the instruction.
 - 2) Fetch the first operand.
 - 3) Perform the addition.
 - 4) Load the result into R1.
- Control sequence for execution of this instruction is as follows:
 - 1) PCout, MARin, Read, Select4, Add, Zin
 - 2) Zout, PCin, Yin, WMFC
 - 3) MDRout, IRin
 - 4) R3out, MARin, Read
 - 5) R1out, Yin, WMFC
 - 6) MDRout, SelectY, Add, Zin
 - 7) Zout, R1in, End

- Instruction execution proceeds as follows:

Step1 → The instruction-fetch operation is initiated by loading contents of PC into MAR & sending a Read request to memory. The Select signal is set to Select4, which causes the Mux to select constant 4. This value is added to operand at input B (PC's content), and the result is stored in Z

Step2 → Updated value in Z is moved to PC.

Step3 → Fetched instruction is moved into MDR and then to IR.

Step4 → Contents of R3 are loaded into MAR & a memory read signal is issued.

Step5 → Contents of R1 are transferred to Y to prepare for addition.

Step6 → When Read operation is completed, memory-operand is available in MDR, and the addition is performed.

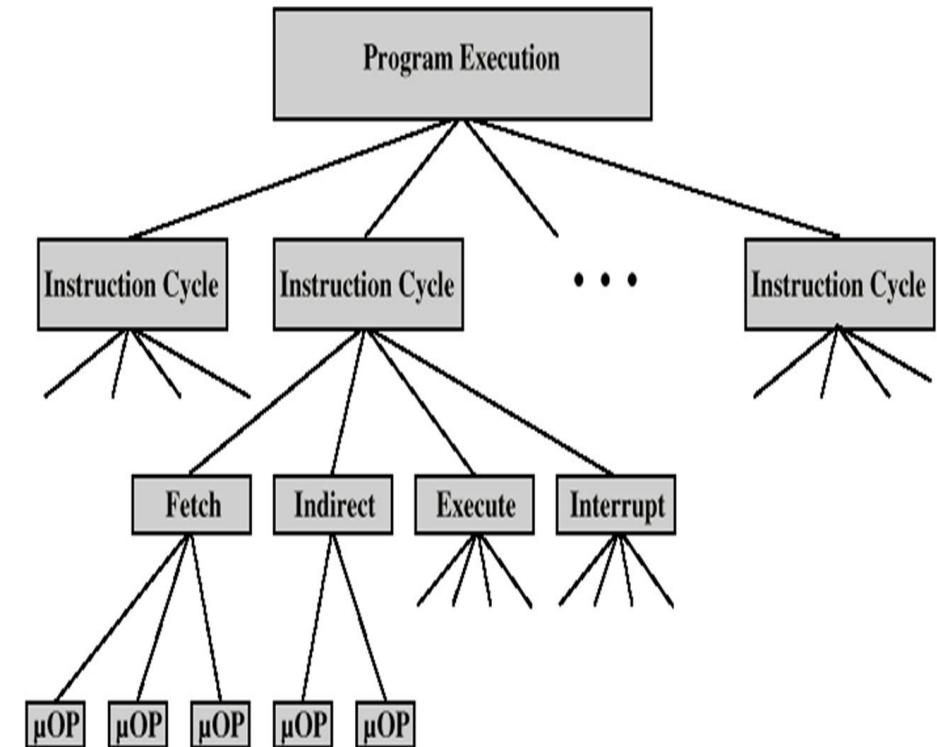
Step7 → Sum is stored in Z, then transferred to R1. The End signal causes a new instruction fetch cycle to begin by returning to step1

Micro-Operation:

The execution of a program consists of the sequential execution of instructions. Each instruction is executed during an instruction cycle made up of shorter sub-cycles (e.g., fetch, indirect, execute, interrupt). The performance of each sub-cycle involves one or more shorter operations, that is, micro-operations.

Micro-operations are the functional, or atomic, operations of a processor.

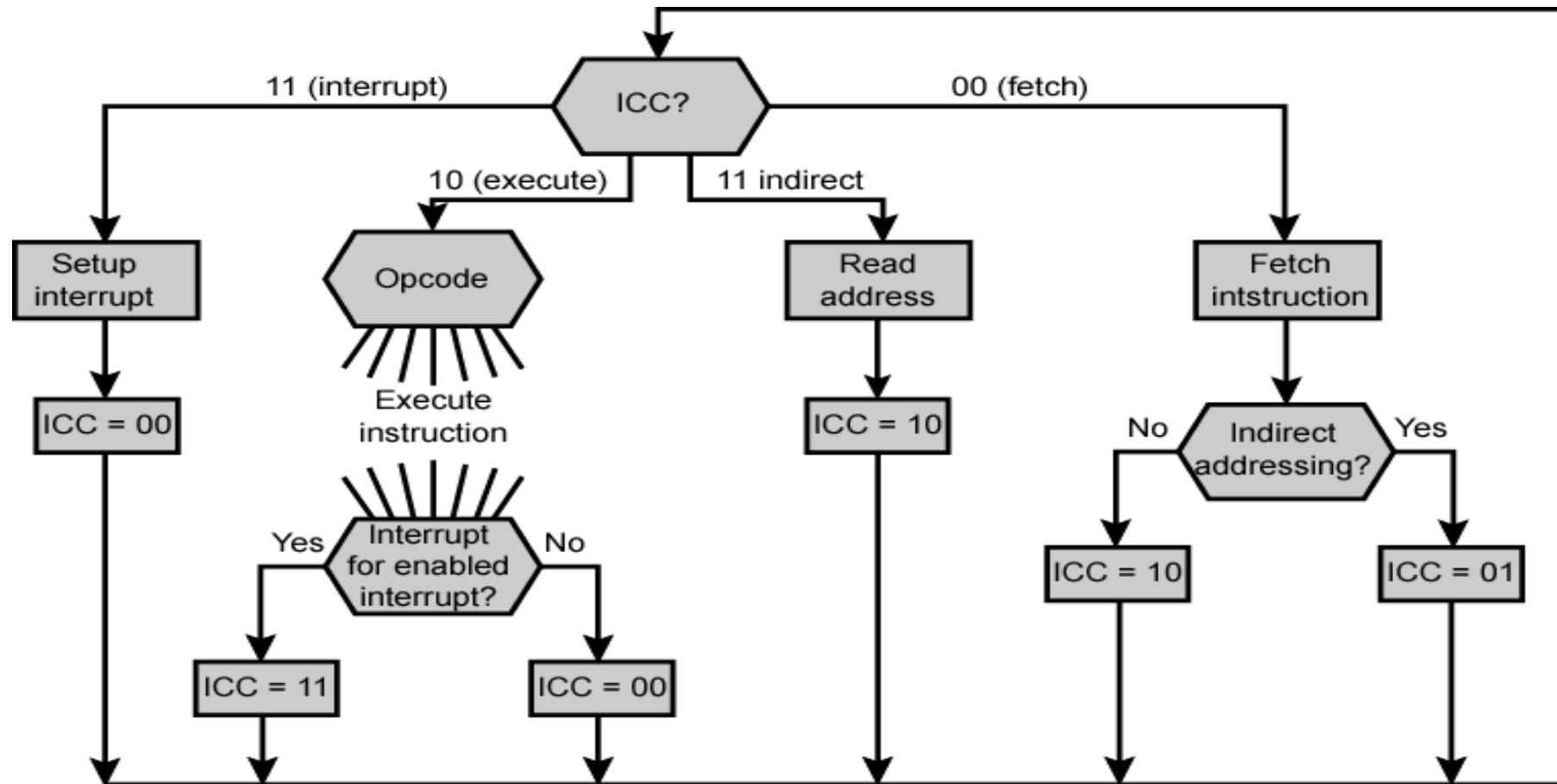
1. The Fetch Cycle (4 registers: MAR, MBR, PC, IR)
2. The Indirect Cycle
3. The Interrupt Cycle
4. The Execute Cycle
5. The Instruction Cycle



Instruction Cycle

- Each phase decomposed into sequence of elementary micro-operations
- E.g. fetch, indirect, and interrupt cycles
- Execute cycle
 - One sequence of micro-operations for each opcode
- Need to tie sequences together
- Assume new 2-bit register
 - Instruction cycle code (ICC) designates which part of cycle processor is in
 - 00: Fetch
 - 01: Indirect
 - 10: Execute
 - 11: Interrupt

Flowchart For Instruction Cycle



Types Of Micro-operation

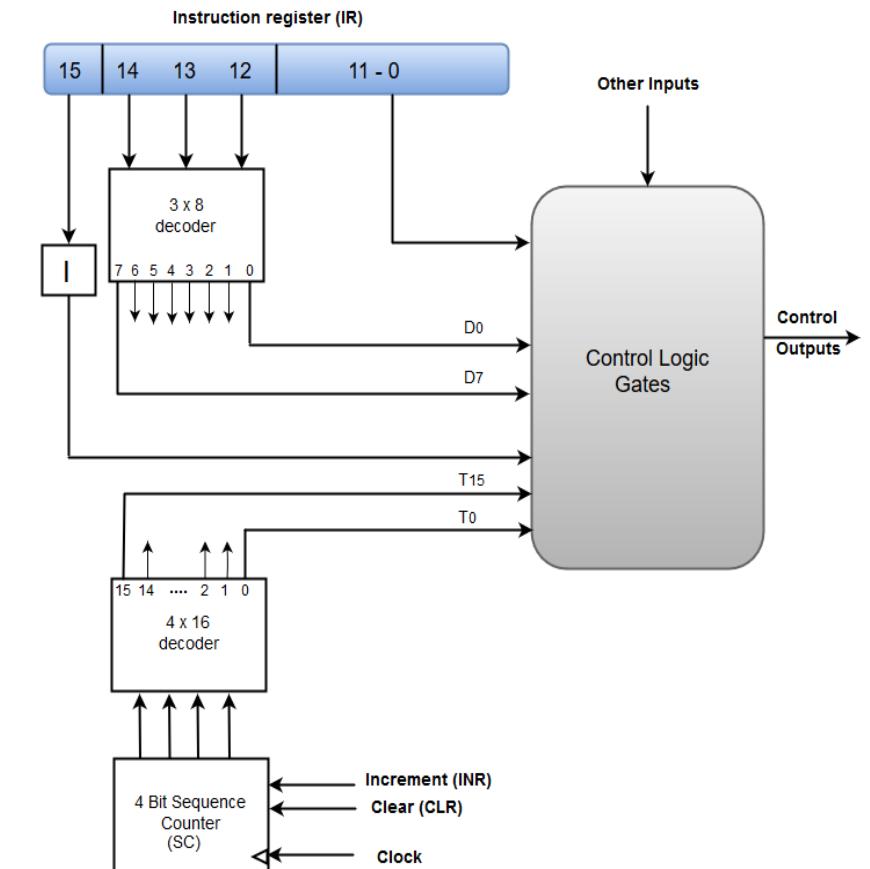
- Transfer data between registers
- Transfer data from register to external
- Transfer data from external to register
- Perform arithmetic or logical operations

HARDWIRE AND MICRO PROGRAMMED CONTROL

The Hardwired Control organization involves the control logic to be implemented with gates, flip-flops, decoders, and other digital circuits.

- A Hard-wired Control consists of two decoders, a sequence counter, and a number of logic gates.
- An instruction fetched from the memory unit is placed in the instruction register (IR).
- The component of an instruction register includes; I bit, the operation code, and bits 0 through 11.
- The operation code in bits 12 through 14 are coded with a 3 x 8 decoder.
- The outputs of the decoder are designated by the symbols D_0 through D_7 .
- The operation code at bit 15 is transferred to a flip-flop designated by the symbol I.
- The operation codes from bits 0 through 11 are applied to the control logic gates.
- The Sequence counter (SC) can count in binary from 0 through 15.

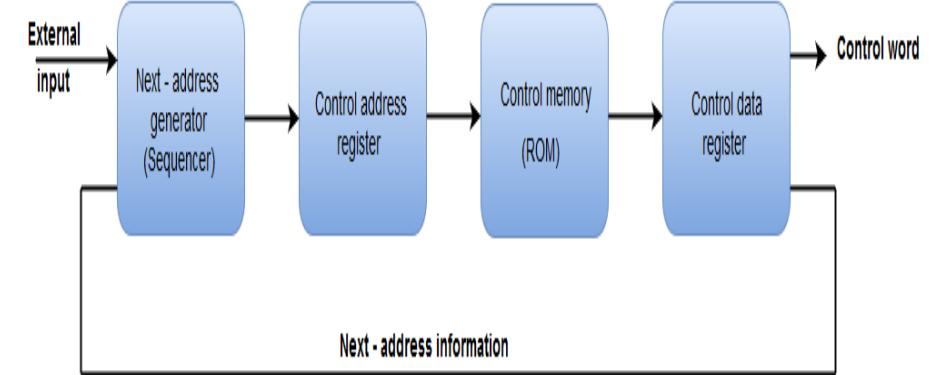
Control Unit of a Basic Computer:



The Microprogrammed Control organization is implemented by using the programming approach. In Microprogrammed Control, the micro-operations are performed by executing a program consisting of micro-instructions.

- The Control memory address register specifies the address of the micro-instruction.
- The Control memory is assumed to be a ROM, within which all control information is permanently stored.
- The control register holds the microinstruction fetched from the memory.
- The micro-instruction contains a control word that specifies one or more micro-operations for the data processor.
- While the micro-operations are being executed, the next address is computed in the next address generator circuit and then transferred into the control address register to read the next microinstruction.
- The next address generator is often referred to as a micro-program sequencer, as it determines the address sequence that is read from control memory.

Microprogrammed Control Unit of a Basic Computer:



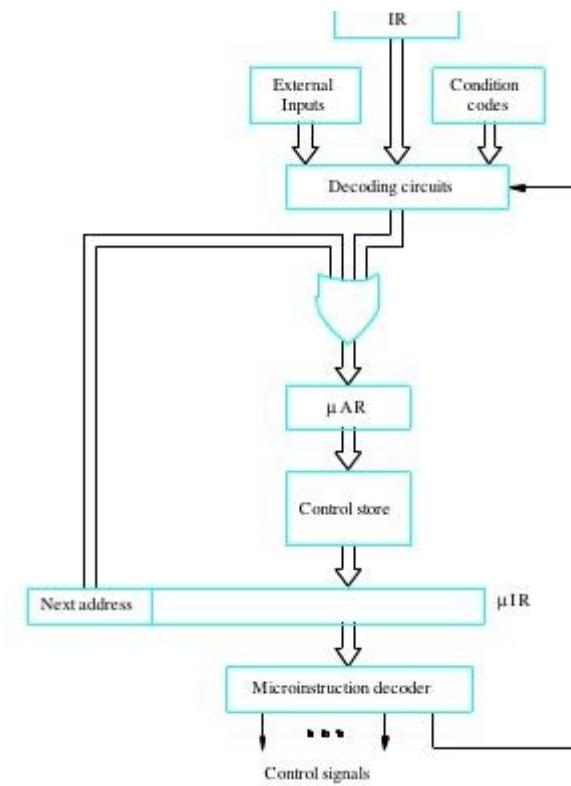
Microprogramming

- The control unit is responsible for initiating the sequence of microoperations that comprise instructions.
 - When these control signals are generated by hardware, the control unit is **hardwired**.
 - When these control signals originate in data stored in a special unit and constitute a program on the small scale, the control unit is **microprogrammed**.
- In a microprogrammed processor, the control unit consists of:
 - **Control address register (CAR)** – contains the address of the next microinstruction to be executed.
 - **Control data register (CDR)** – contains the microinstruction to be executed.
 - **The sequencer** – determines the next address from within control memory
 - **Control memory** – where microinstructions are stored.

- A micro-program control unit can be viewed as consisting of two parts:
 - The control memory that stores the microinstructions.
 - Sequencing circuit that controls the generation of the next address.
- A *micro-program sequencer* attached to a control memory inputs certain bits of the microinstruction, from which it determines the next address for control memory. A typical sequencer provides the following address-sequencing capabilities:
 - Increment the present address for control memory.
 - Branches to an address as specified by the address field of the micro instruction.
 - Branches to a given address if a specified status bit is equal to 1.
 - Transfer control to a new address as specified by an external source (Instruction Register).
 - Has a facility for subroutine calls and returns.
- Depending on the current microinstruction condition flags, and the contents of the instruction register, a control memory address must be generated for the next micro instruction.

Microinstructions With Next-address Field

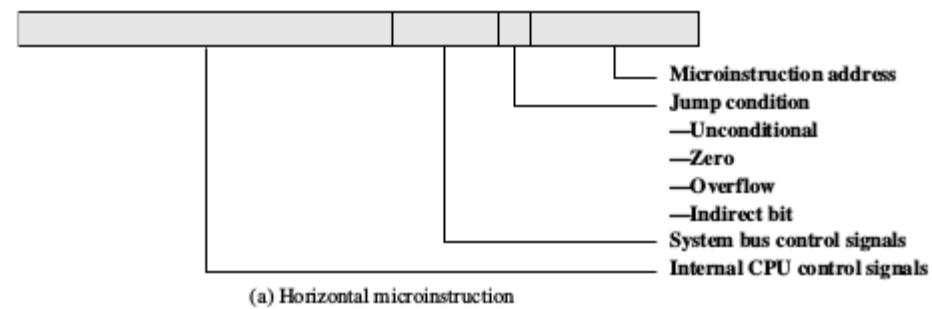
- The microprogram requires several branch microinstructions, which perform no useful operation in the datapath.
- A powerful alternative approach is to include an address field as a part of every microinstruction to indicate the location of the next microinstruction to be fetched.
- Pros: separate branch microinstructions are virtually eliminated; few limitations in assigning addresses to microinstructions.
- Cons: additional bits for the address field (around 1/6)



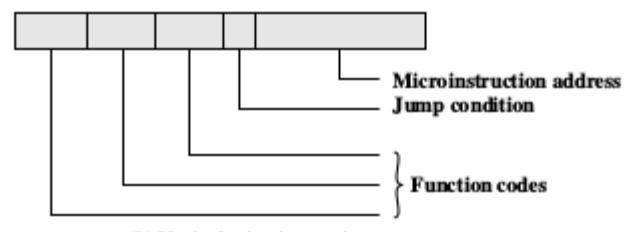
Prefetching Microinstructions

- One drawback of microprogrammed control is that it leads to a slower operating speed because of the time it takes to fetch microinstructions from the control store.
- Faster operation is achieved if the next microinstruction is prefetched while the current one is being executed. In this way , the execution time can be overlapped with the fetch time.
- *Prefetching microinstructions* presents some organizational difficulties . Sometimes the status flags and the result of the currently executed microinstruction are needed to determine the address of the next microinstruction. In these cases , the fetch must be repeated with the correct address, which requires more complex hardware. However , the disadvantages are minor, and the prefetching technique is often used.

Vertical organization	Horizontal organization
Highly encoded schemes that use compact codes to specify only a small number of control functions in each microinstruction are referred to as a vertical organization	The minimally encoded scheme in which many resources can be controlled with a single microinstruction is called a horizontal organization
This approach results in considerably slower operating speeds because more microinstructions are needed to perform the desired control functions	This approach is useful when a higher operating speed is desired and when the machine structure allows parallel use of resources



(a) Horizontal microinstruction

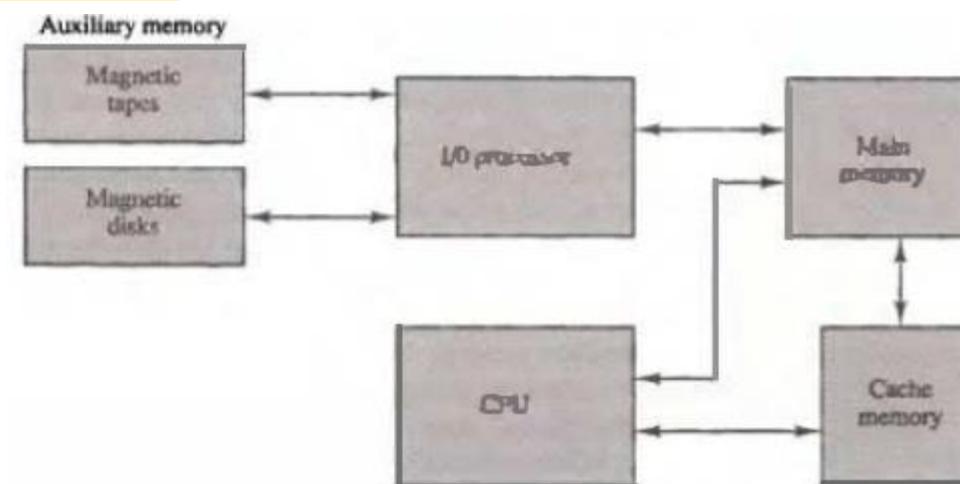


(b) Vertical microinstruction

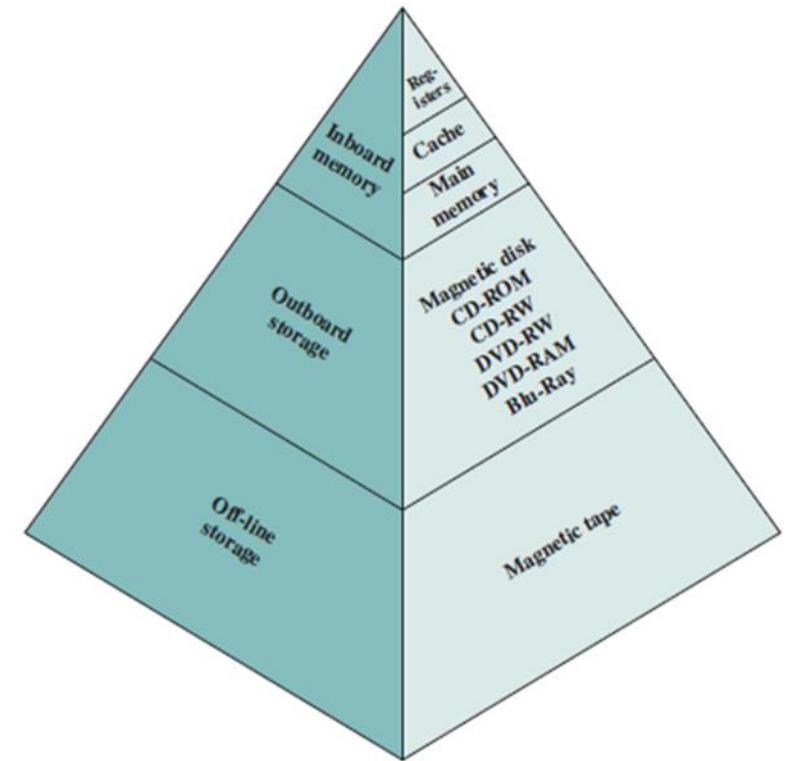
MEMORY ORGANIZATION

- The memory unit is an essential component in any digital computer since it is needed for storing programs and data.
- Most general-purpose computers would run more efficiently if they were equipped with additional storage beyond the capacity of the main memory.
- 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.
- Main memory is the principal internal memory system of the computer.

- The total memory capacity of a computer can be visualized as being a hierarchy of components.
- The memory hierarchy system consists of all 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.



- As one goes down the hierarchy, the following occur:
 - Decreasing cost per bit
 - Increasing capacity
 - Increasing access time
 - Decreasing frequency of access of the memory by the processor
- Thus, smaller, more expensive, faster memories are supplemented by larger, cheaper, slower memories.



The three forms of memory just described are, typically, volatile and employ semiconductor technology.

The use of three levels exploits the fact that semiconductor memory comes in a variety of types, which differ in speed and cost.

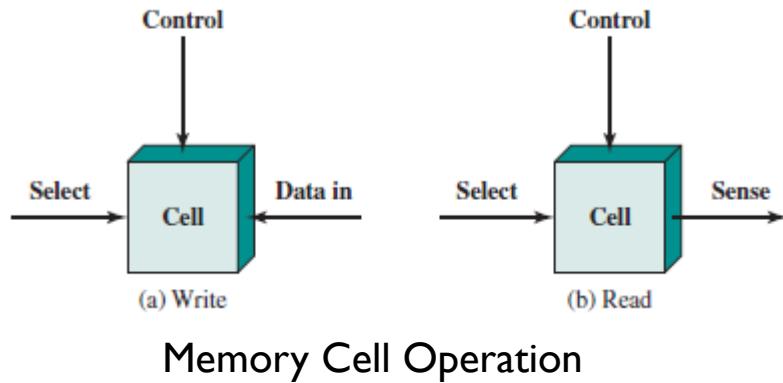
Data are stored more permanently on external mass storage devices, of which the most common are hard disk and removable media, such as removable magnetic disk, tape, and optical storage.

External, nonvolatile memory is also referred to as secondary memory or auxiliary memory. These are used to store program and data files and are usually visible to the programmer only in terms of files and records, as opposed to individual bytes or words.

Disk is also used to provide an extension to main memory known as virtual memory.

Semiconductor Main Memory

- In earlier computers, the most common form of random-access storage for computer main memory employed an array of doughnut-shaped ferromagnetic loops referred to as cores. Hence, main memory was often referred to as core.
- The advent of, and advantages of, microelectronics has long since vanquished the magnetic core memory. Today, the use of semiconductor chips for main memory is almost universal.
- The basic element of a semiconductor memory is the memory cell. Although a variety of electronic technologies are used, all semiconductor memory cells share certain properties:
 - They exhibit two stable (or semi-stable) states, which can be used to represent binary 1 and 0.
 - They are capable of being written into (at least once), to set the state.
 - They are capable of being read to sense the state.



Memory Type	Category	Erasure	Write Mechanism	Volatility
Random-access memory (RAM)	Read-write memory	Electrically, byte-level	Electrically	Volatile
Read-only memory (ROM)	Read-only memory	Not possible	Masks	Nonvolatile
Programmable ROM (PROM)				
Erasable PROM (EPROM)	Read-mostly memory	UV light, chip-level	Electrically	Nonvolatile
Electrically Erasable PROM (EEPROM)		Electrically, byte-level		
Flash memory		Electrically, block-level		

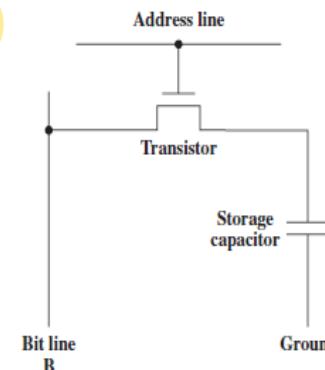
Semiconductor Memory Types

DRAM and SRAM

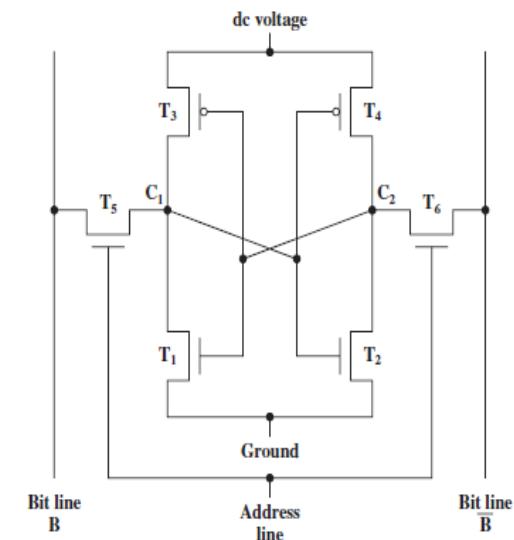
- The most common type of semiconductor memory is referred to as random-access memory (RAM).
- One distinguishing characteristic of memory that is designated as RAM is that it is possible both to read data from the memory and to write new data into the memory easily and rapidly.
- A **dynamic RAM (DRAM)** is made with cells that store data as charge on capacitors. The presence or absence of charge in a capacitor is interpreted as a binary 1 or 0.
- Because capacitors have a natural tendency to discharge, dynamic RAMs require periodic charge refreshing to maintain data storage. The term *dynamic* refers to this tendency of the stored charge to leak away, even with power continuously applied.
- In contrast, a **static RAM (SRAM)** is a digital device that uses the same logic elements used in the processor. In a SRAM, binary values are stored using traditional flip-flop logic-gate configurations.
- A static RAM will hold its data as long as power is supplied to it.

SRAM versus DRAM

Both static and dynamic RAMs are volatile; that is, power must be continuously supplied to the memory to preserve the bit values. A dynamic memory cell is simpler and smaller than a static memory cell. Thus, DRAM is more dense (smaller cells = more cells per unit area) and less expensive than a corresponding SRAM. On the other hand, a DRAM requires the supporting refresh circuitry. For larger memories, the fixed cost of the refresh circuitry is more than compensated for by the smaller variable cost of DRAM cells. Thus, DRAMs tend to be favored for large memory requirements. A final point is that SRAMs are somewhat faster than DRAMs. Because of these relative characteristics, SRAM is used for cache memory (both on and off chip), and DRAM is used for main memory.



(a) Dynamic RAM (DRAM) cell



(b) Static RAM (SRAM) cell

ROM

- A **read-only memory (ROM)** contains a permanent pattern of data that cannot be changed. A ROM is nonvolatile; that is, no power source is required to maintain the bit values in memory. While it is possible to read a ROM, it is not possible to write new data into it.
- When only a small number of ROMs with a particular memory content is needed, a less expensive alternative is the **programmable ROM (PROM)**. Like the ROM, the PROM is nonvolatile and may be written into only once.
- The optically **erasable programmable read-only memory (EPROM)** is read and written electrically, as with PROM. However, before a write operation, all the storage cells must be erased to the same initial state by exposure of the packaged chip to ultraviolet radiation.
- A more attractive form of read-mostly memory is **electrically erasable programmable read-only memory (EEPROM)**. This is a read-mostly memory that can be written into at any time without erasing prior contents; only the byte or bytes addressed are updated. The write operation takes considerably longer than the read operation, on the order of several hundred microseconds per byte.
- Another form of semiconductor memory is **flash memory** (so named because of the speed with which it can be reprogrammed)

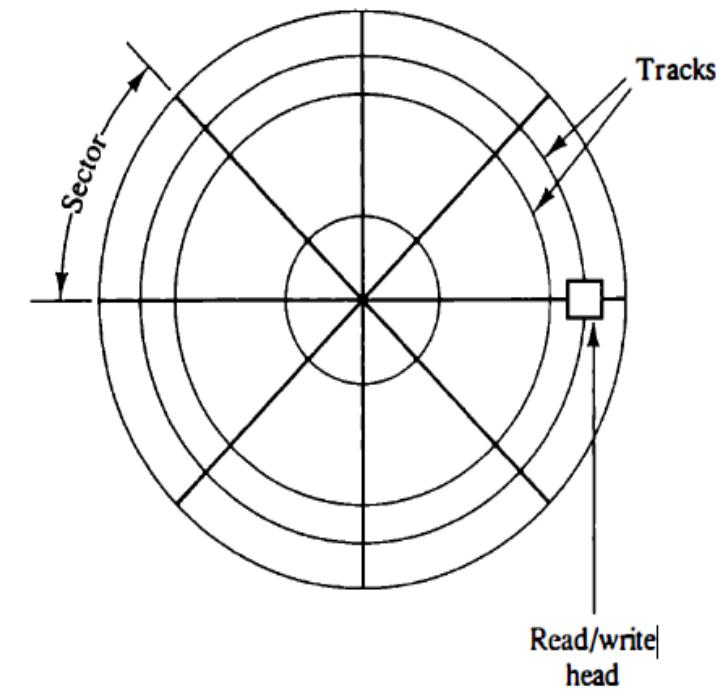
Like EEPROM, flash memory uses an electrical erasing technology. An entire flash memory can be erased in one or a few seconds, which is much faster than EPROM. In addition, it is possible to erase just blocks of memory rather than an entire chip. Flash memory gets its name because the microchip is organized so that a section of memory cells are erased in a single action or “flash.” However, flash memory does not provide byte-level erasure. Like EPROM, flash memory uses only one transistor per bit, and so achieves the high density (compared with EEPROM) of EPROM.

Auxiliary Memory

- The most common auxiliary memory devices used in computer systems are magnetic disks and tapes.
- To understand fully the physical mechanism of auxiliary memory devices one must have a knowledge of magnetics, electronics, and electromechanical systems.
- 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 started 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.

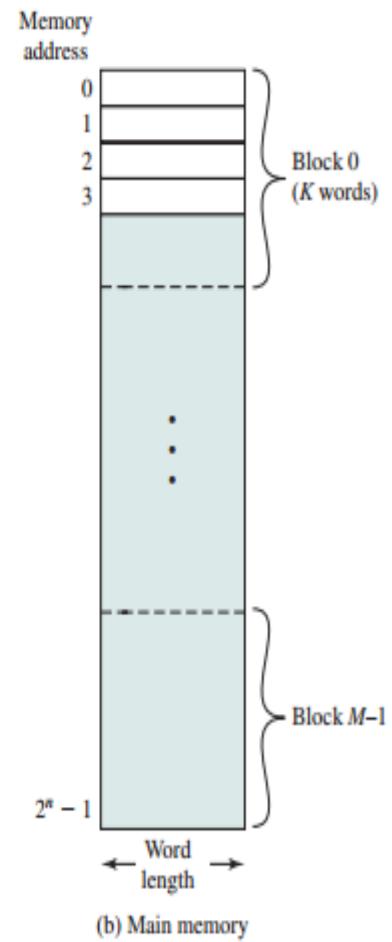
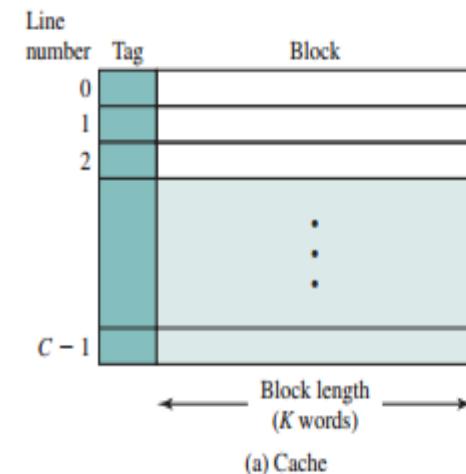
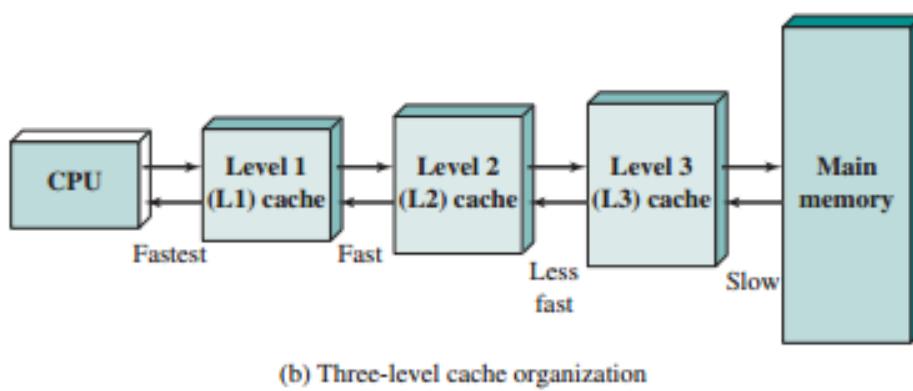
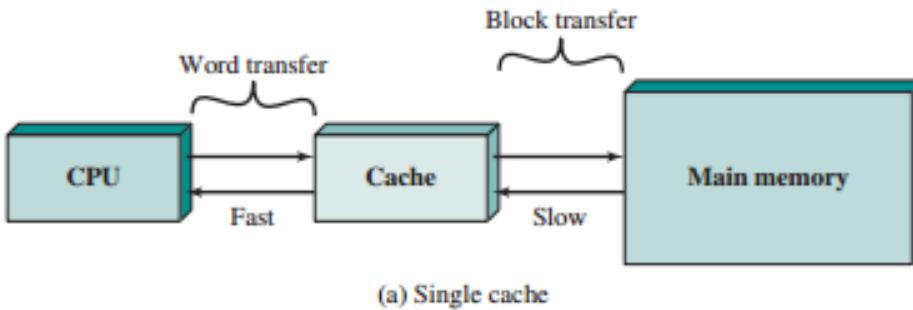
- A **magnetic disk** is a circular plate constructed of metal or plastic coated with magnetized material. Often both sides of the disk are used and several disks may be stacked on one spindle with read/write heads available on each surface. All disks rotate together at high speed and are not stopped or started for 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*.
- 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.



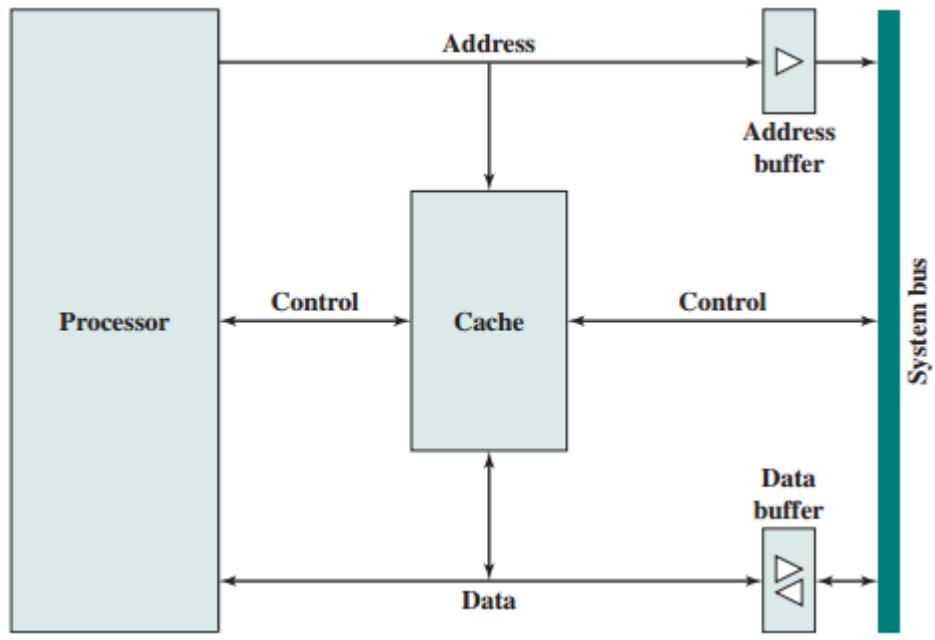
- A **magnetic tape** transport consists of the electrical, mechanical, and electronic components to provide the parts and control mechanism for a magnetic-tape unit. The tape itself is a strip of plastic coated with a magnetic recording medium. Bits are recorded as magnetic spots on the tape along several 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 and the number of characters in the record. Records may be of fixed or variable length.

Cache Memory

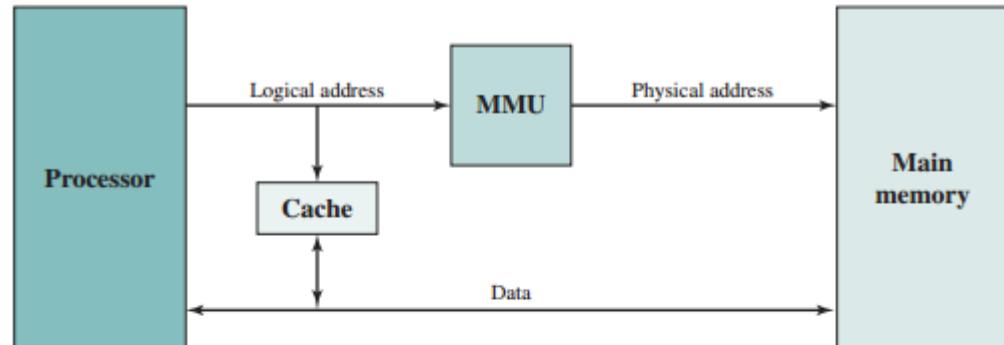
- Analysis of a large number of typical programs has shown that the references to memory at any given interval of time tend to be confined within a few localized areas in memory. This phenomenon is known as the **property of locality of reference**.
- A special very-high speed memory called a Cache is 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.
- 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.
- 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 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.



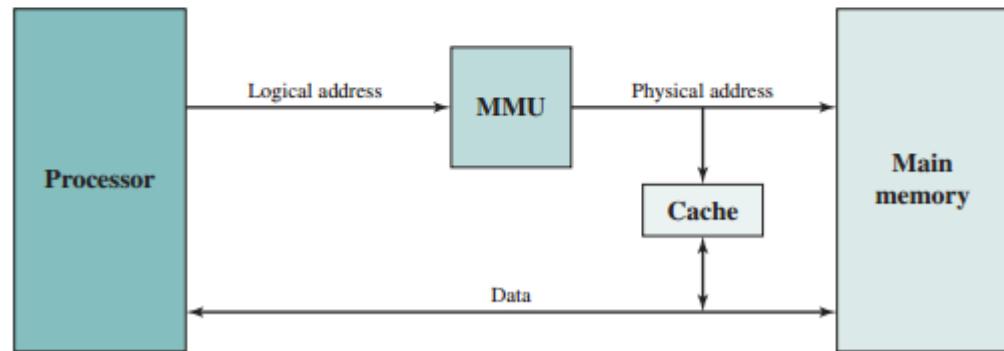
Typical Cache Organization



Cache Addresses	Write Policy
Logical	Write through
Physical	Write back
Cache Size	Line Size
Mapping Function	Number of Caches
Direct	Single or two level
Associative	Unified or split
Set associative	
Replacement Algorithm	
Least recently used (LRU)	
First in first out (FIFO)	
Least frequently used (LFU)	
Random	



(a) Logical cache



(b) Physical cache

Three Related Concepts For Internal Memory

- **Word:** The “natural” unit of organization of memory. The size of a word is typically equal to the number of bits used to represent an integer and to the instruction length. Unfortunately, there are many exceptions.
- **Addressable units:** In some systems, the addressable unit is the word. However, many systems allow addressing at the byte level. In any case, the relationship between the length in bits A of an address and the number N of addressable units is $2^A = N$.
- **Unit of transfer:** For main memory, this is the number of bits read out of or written into memory at a time. The unit of transfer need not equal a word or an addressable unit. For external memory, data are often transferred in much larger units than a word, and these are referred to as blocks.

Method of Accessing Units of Data

- **Sequential access:** Memory is organized into units of data, called records. Access must be made in a specific linear sequence. Stored addressing information is used to separate records and assist in the retrieval process. A shared read-write mechanism is used, and this must be moved from its current location to the desired location, passing and rejecting each intermediate record.
- **Direct access:** As with sequential access, direct access involves a shared read-write mechanism. However, individual blocks or records have a unique address based on physical location. Access is accomplished by direct access to reach a general vicinity plus sequential searching, counting, or waiting to reach the final location. Again, access time is variable.
- **Random access:** Each addressable location in memory has a unique, physically wired-in addressing mechanism. The time to access a given location is independent of the sequence of prior accesses and is constant. Thus, any location can be selected at random and directly addressed and accessed. Main memory and some cache systems are random access.
- **Associative:** This is a random access type of memory that enables one to make a comparison of desired bit locations within a word for a specified match, and to do this for all words simultaneously. Thus, a word is retrieved based on a portion of its contents rather than its address. As with ordinary random-access memory, each location has its own addressing mechanism, and retrieval time is constant independent of location or prior access patterns. Cache memories may employ associative access.

Performance Parameters

- **Access time (latency):** For random-access memory, this is the time it takes to perform a read or write operation, that is, the time from the instant that an address is presented to the memory to the instant that data have been stored or made available for use. For non-random-access memory, access time is the time it takes to position the read-write mechanism at the desired location.
- **Memory cycle time:** This concept is primarily applied to random-access memory and consists of the access time plus any additional time required before a second access can commence. This additional time may be required for transients to die out on signal lines or to regenerate data if they are read destructively. Note that memory cycle time is concerned with the system bus, not the processor.
- **Transfer rate:** This is the rate at which data can be transferred into or out of a memory unit. For random-access memory, it is equal to $1/(cycle\ time)$.

For non-random-access memory, the following relationship holds:

$$T_n = T_A + \frac{n}{R}$$

where

T_n = Average time to read or write n bits

T_A = Average access time

n = Number of bits

R = Transfer rate, in bits per second (bps)

- 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 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 transformation of data from main memory to cache memory is referred to as a **mapping process**.

Cache versus Auxiliary Memory

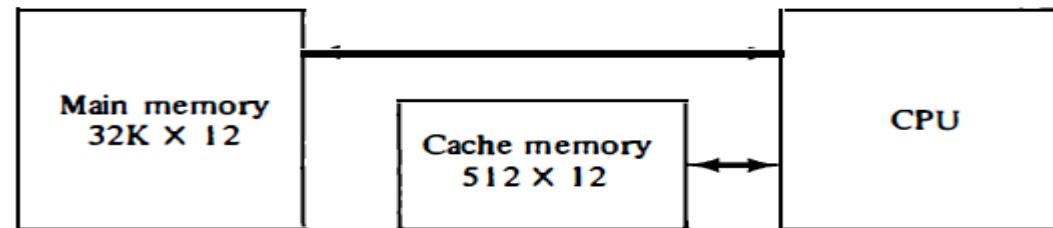
Auxiliary and cache memories are used for different purposes. The cache holds those parts of the program and data that are most heavily used, while the auxiliary memory holds those parts that are not presently used by the CPU.

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 a 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 from 26 to 20 words, while cache block size is typically from 1 to 16 words.

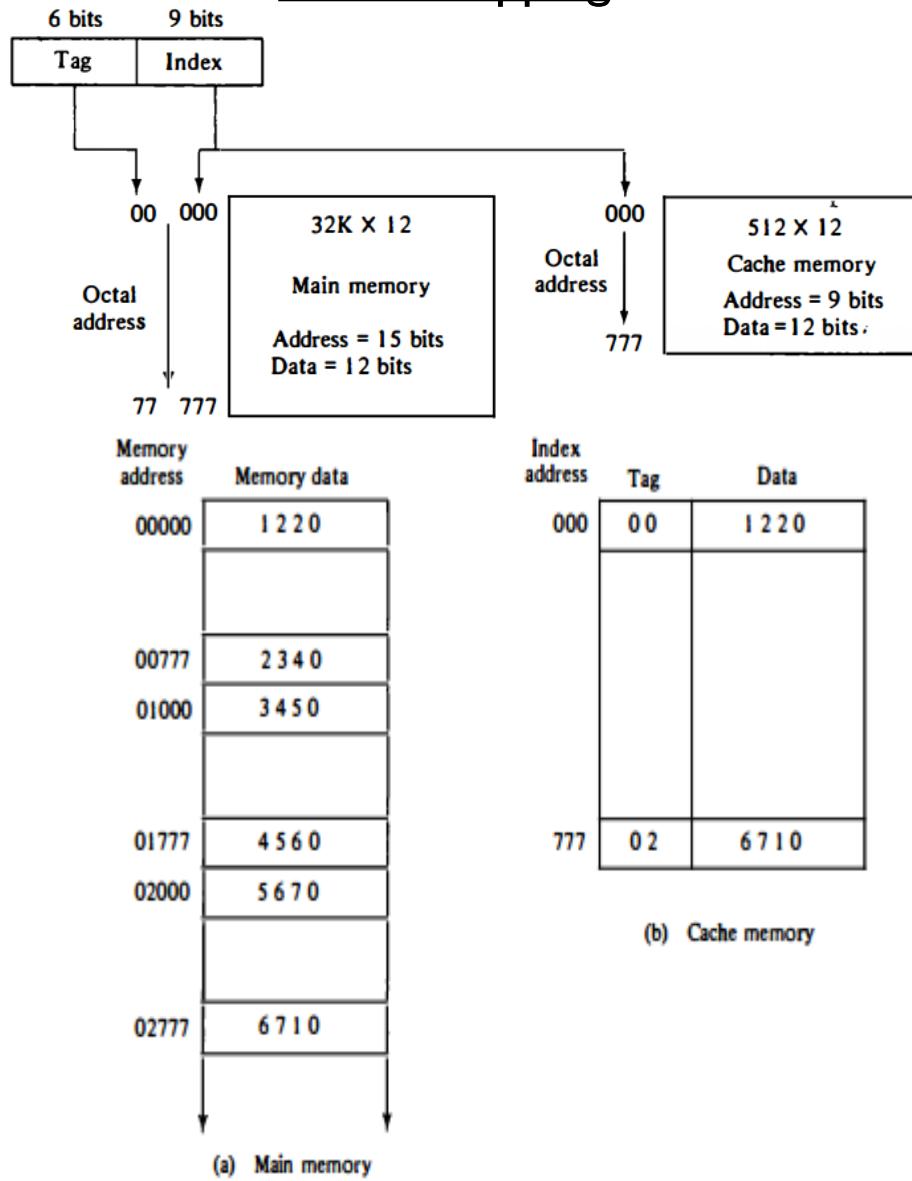
- Three types of mapping procedures are of practical interest when considering the organization of cache memory:
 1. Associative mapping
 2. Direct mapping
 3. Set-associative mapping
- The main memory can store 32K words of 12 bits each. The cache is capable of storing 512 of these words at any given 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.



- For all three cases, the example includes the following elements:
 - The cache can hold 64 Kbytes.
 - Data are transferred between main memory and the cache in blocks of 4 bytes each.
This means that the cache is organized as $16K = 2^{14}$ lines of 4 bytes each.
 - The main memory consists of 16 Mbytes, with each byte directly addressable by a 24-bit address ($2^{24} = 16M$). Thus, for mapping purposes, we can consider main memory to consist of 4M blocks of 4 bytes each.

Direct Mapping

SOURCE: MORRIS MANO, THIRD EDITION



The simplest technique, known as direct mapping, maps each block of main memory into only one possible cache line. The mapping is expressed as

$$i = j \text{ modulo } m$$

where

i = cache line number

j = main memory block number

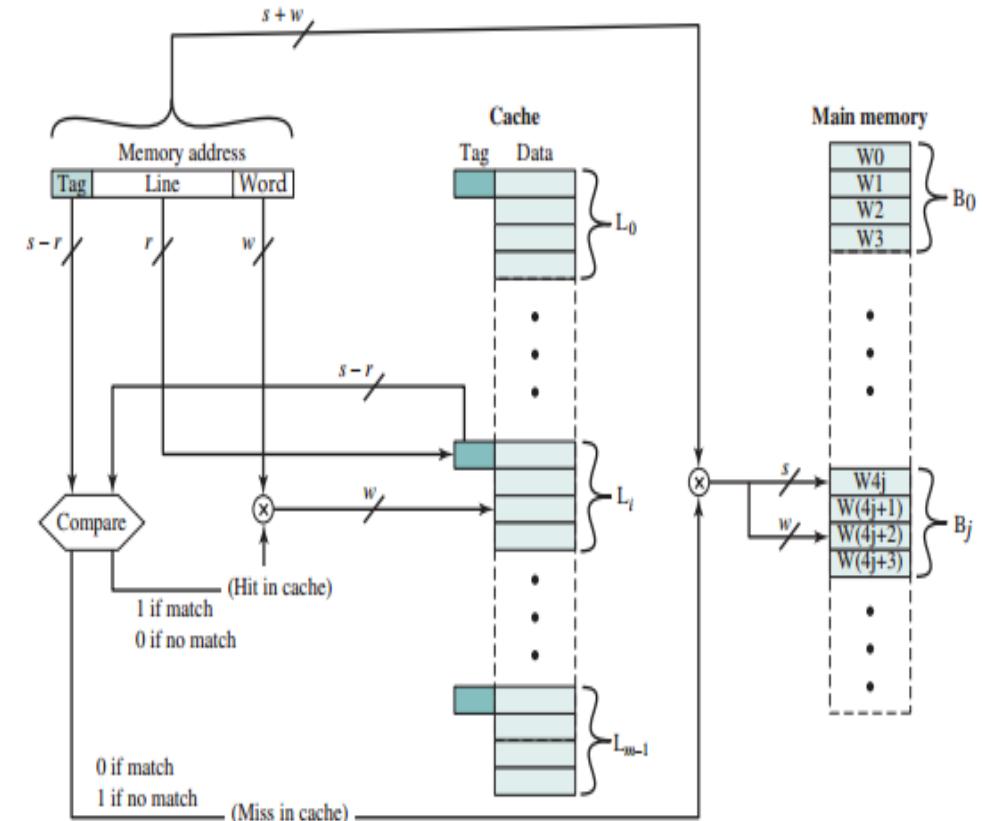
m = number of lines in the cache

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.

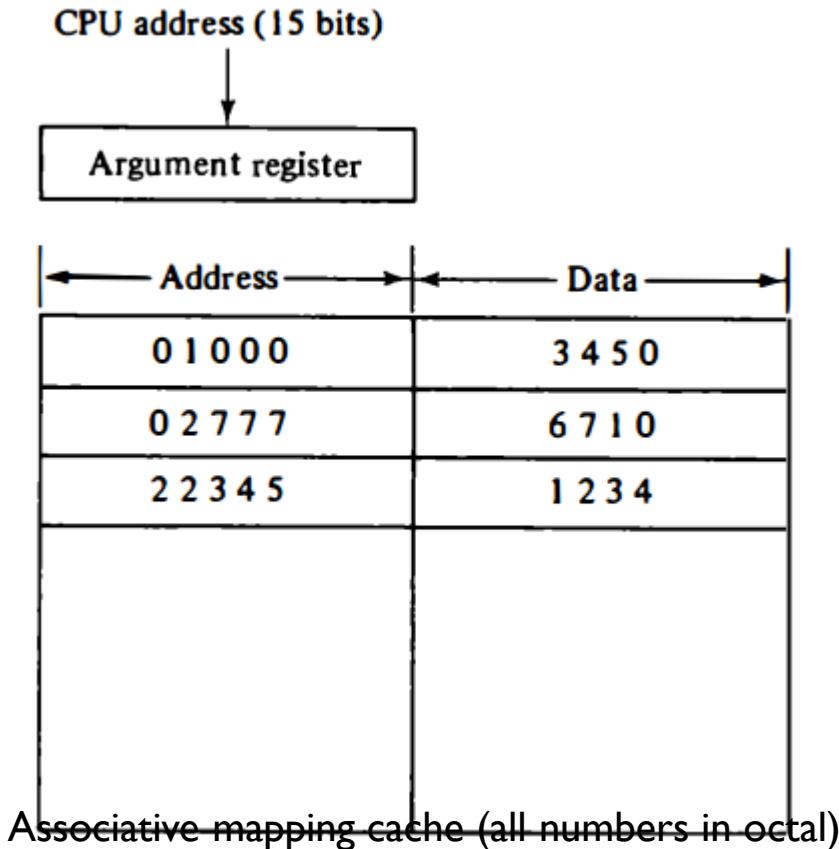
The disadvantage of direct mapping is that the hit ratio can drop considerably if two or more words whose addresses have the same index but different tags are accessed repeatedly. However, this possibility is made by the fact that such words are relatively far apart in the address range.

For purposes of cache access, each main memory address can be viewed as consisting of three fields. The least significant w bits identify a unique word or byte within a block of main memory; in most contemporary machines, the address is at the byte level. The remaining s bits specify one of the 2^s blocks of main memory. The cache logic interprets these s bits as a tag of $s - r$ bits (most significant portion) and a line field of r bits. This latter field identifies one of the $m = 2^r$ lines of the cache. To summarize,

- Address length = $(s + w)$ bits
- Number of addressable units = 2^{s+w} words or bytes
- Block size = line size = 2^w words or bytes
- Number of blocks in main memory = $\frac{2^{s+w}}{2^w} = 2^s$
- Number of lines in cache = $m = 2^r$
- Size of cache = 2^{r+w} words or bytes
- Size of tag = $(s - r)$ bits

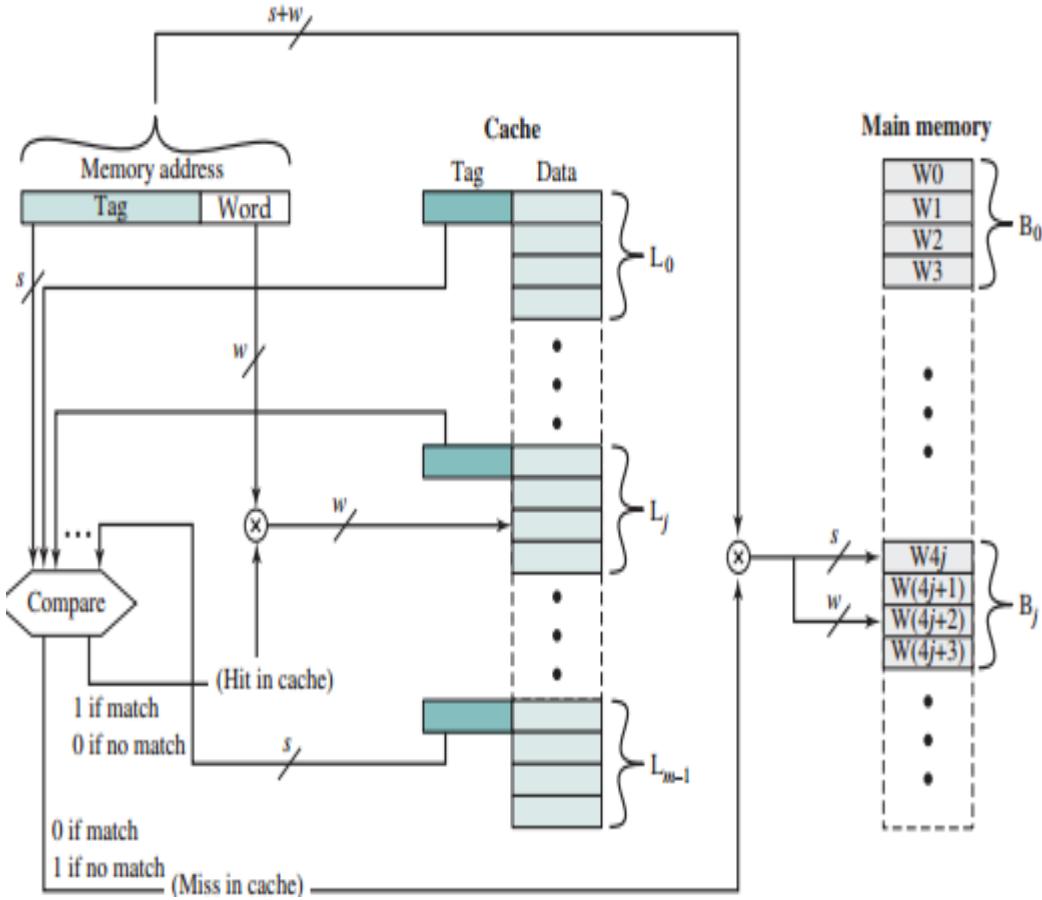


Associative Mapping



The fastest and most flexible cache organization uses an associative memory. The associative memory stores both the address ad content (data) of the memory word. This permits any location in cache to store any word from main memory.

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.



Note that no field in the address corresponds to the line number, so that the number of lines in the cache is not determined by the address format. To summarize,

- Address length = $(s + w)$ bits
- Number of addressable units = 2^{s+w} words or bytes
- Block size = line size = 2^w words or bytes
- Number of blocks in main memory = $\frac{2^{s+w}}{2^w} = 2^s$
- Number of lines in cache = undetermined
- Size of tag = s bits

Set-associative mapping is a compromise that exhibits the strengths of both the direct and associative approaches while reducing their disadvantages.

In this case, the cache consists of a number sets, each of which consists of a number of lines. The relationships are

$$\begin{aligned} m &= v * k \\ i &= j \text{ modulo } v \end{aligned}$$

where

i = cache set number

j = main memory block number

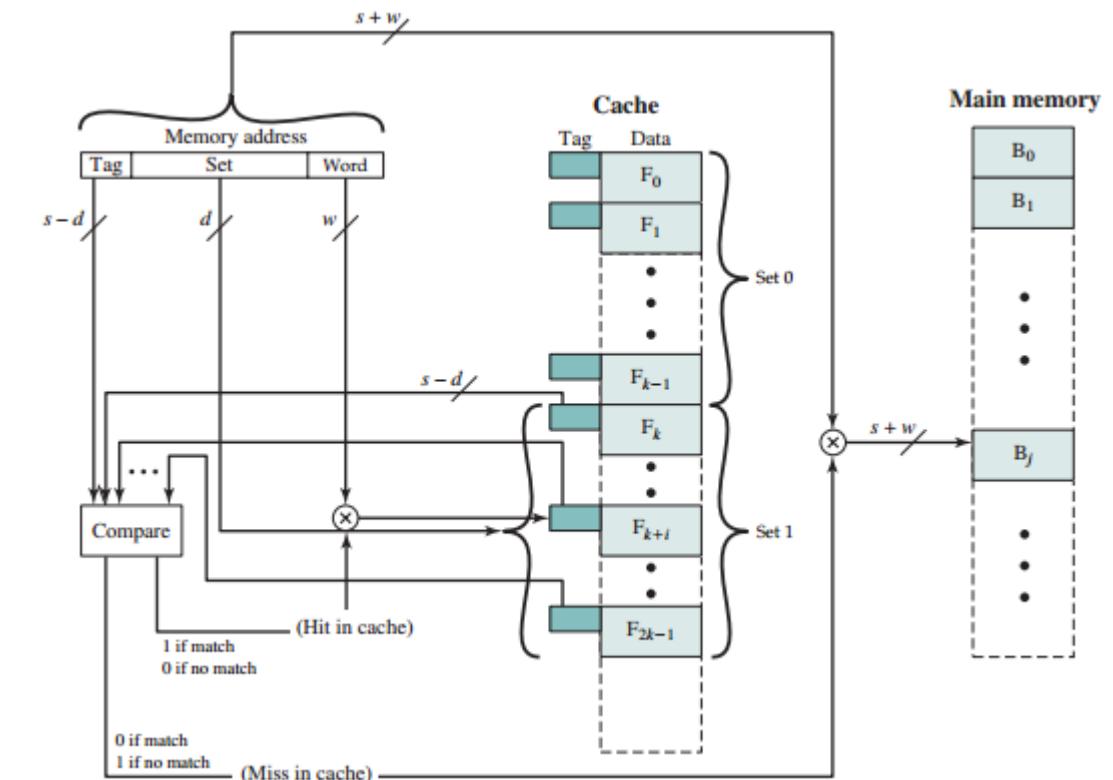
m = number of lines in the cache

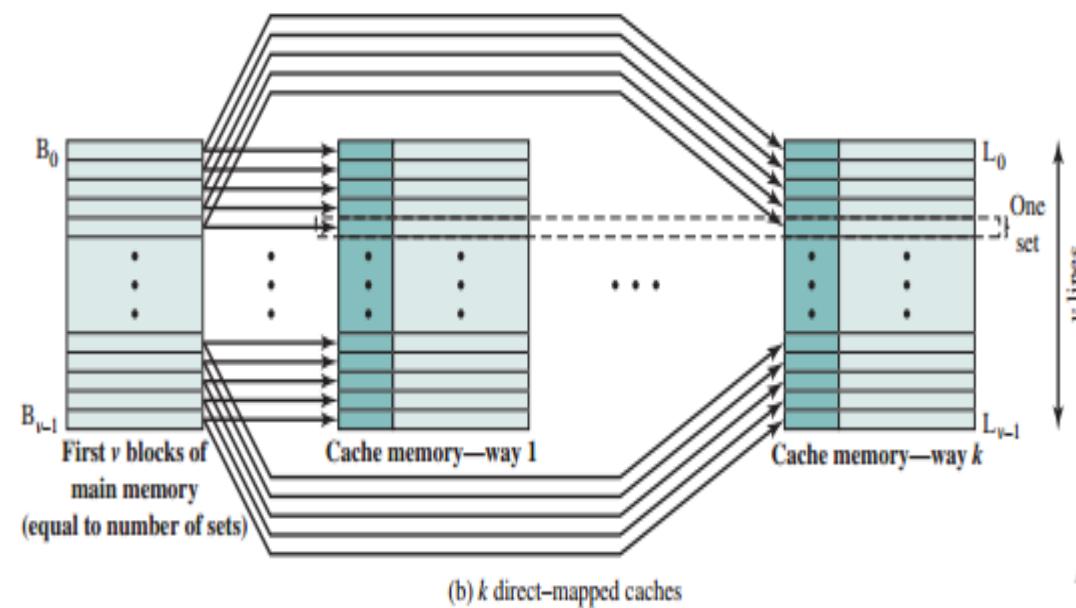
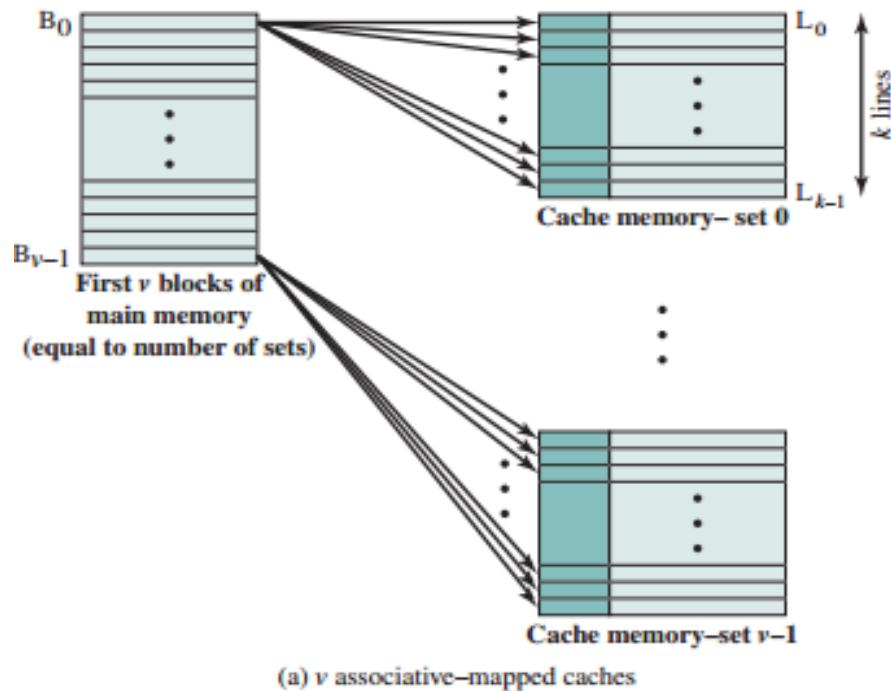
v = number of sets

k = number of lines in each set

This is referred to as k -way set-associative mapping.

- Block size = line size = 2^w words or bytes
- Number of blocks in main memory = $\frac{2^{s+w}}{2^w} = 2^s$
- Number of lines in set = k
- Number of sets = $v = 2^d$
- Number of lines in cache = $m = kv = k \times 2^d$
- Size of cache = $k \times 2^{d+w}$ words or bytes
- Size of tag = $(s - d)$ bits





Mapping from Main Memory to Cache: k -Way Set Associative

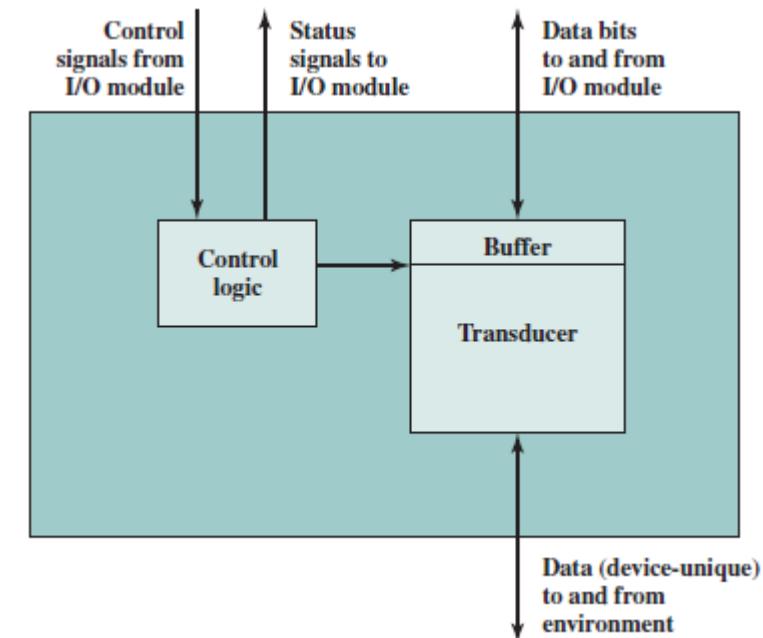
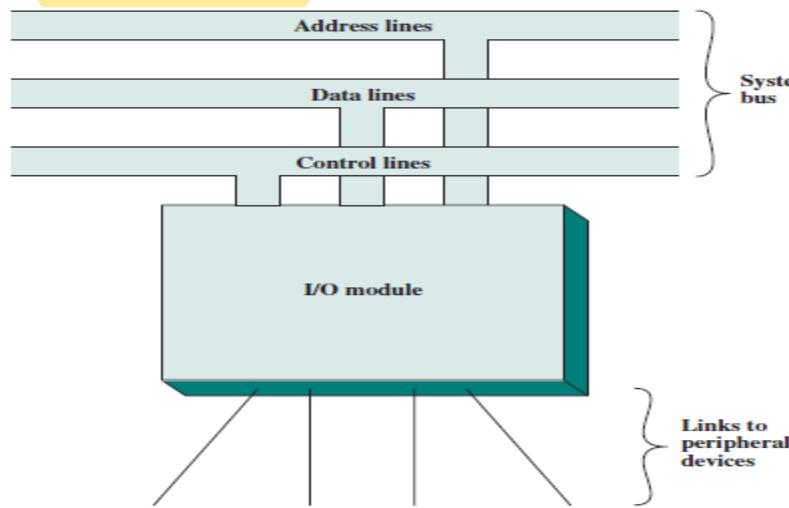
INPUT/OUTPUT

External Devices

Human readable: Suitable for communicating with the computer user

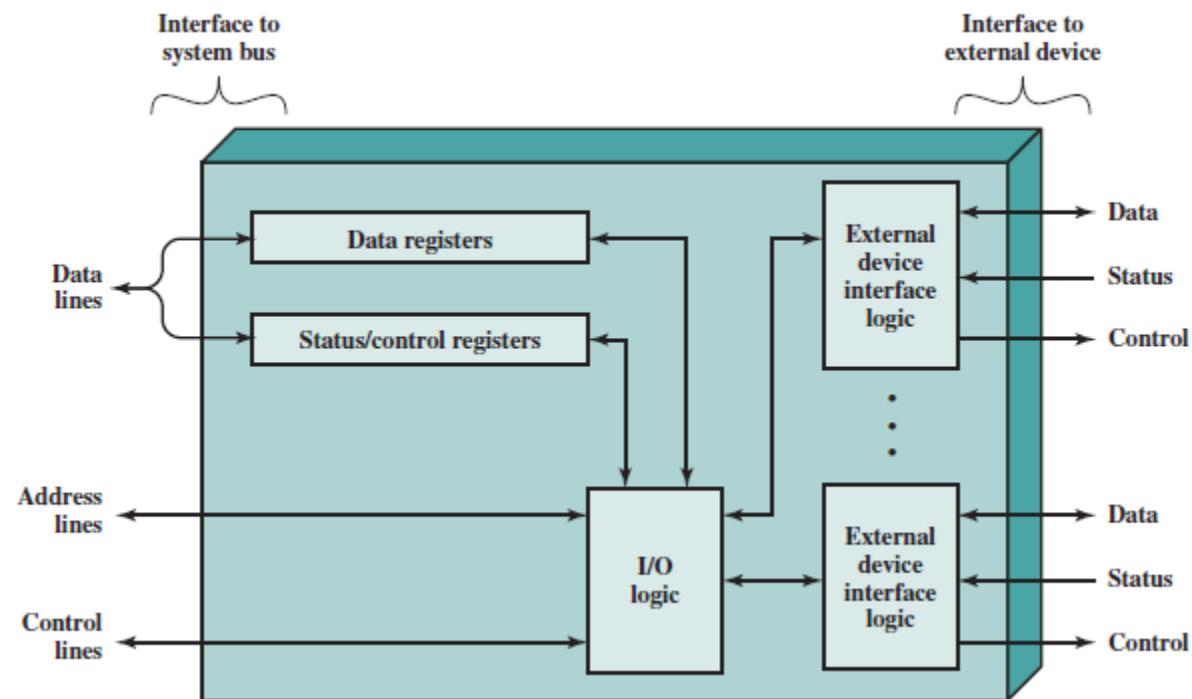
Machine readable: Suitable for communicating with equipment

Communication: Suitable for communicating with remote devices



I/O Modules

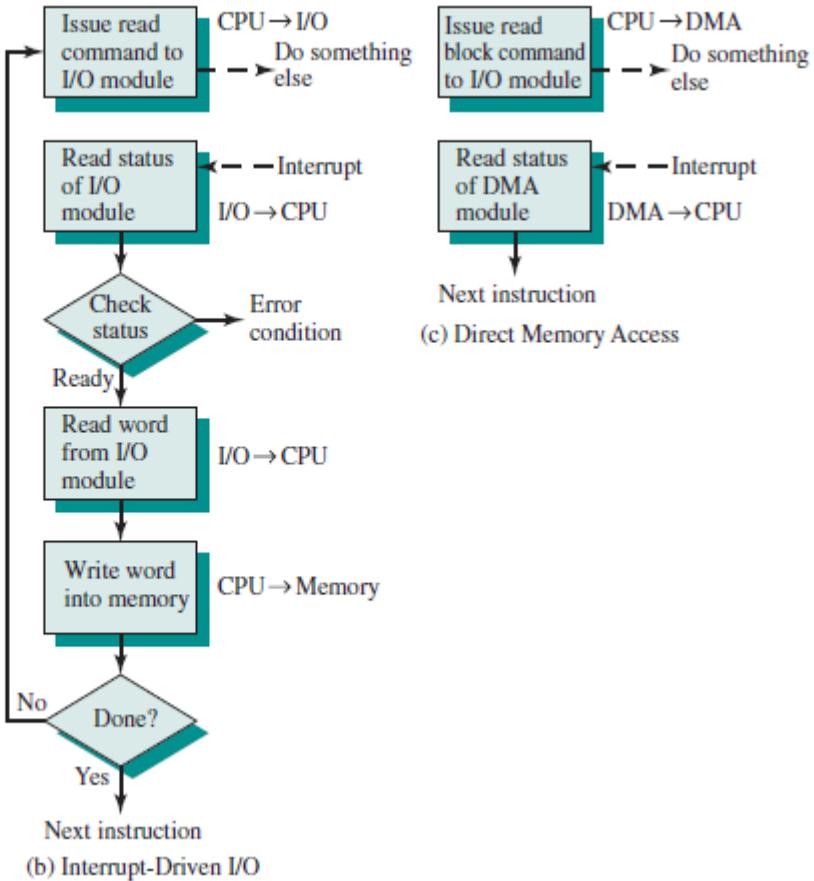
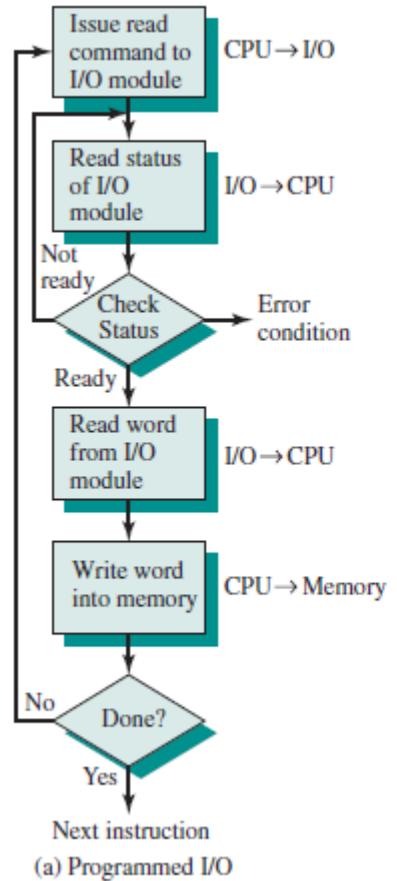
- Control and timing
- Processor communication
- Device communication
- Data buffering
- Error detection



Programmed I/O

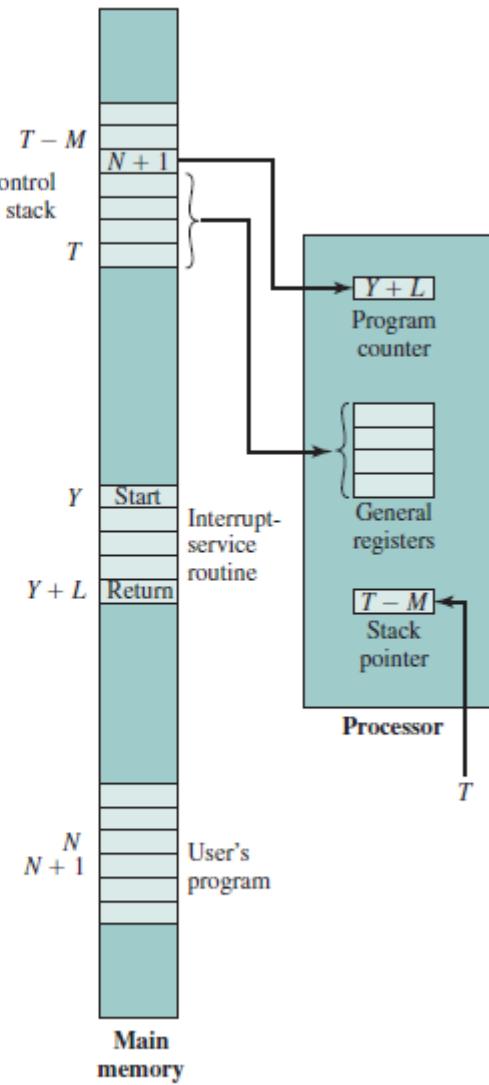
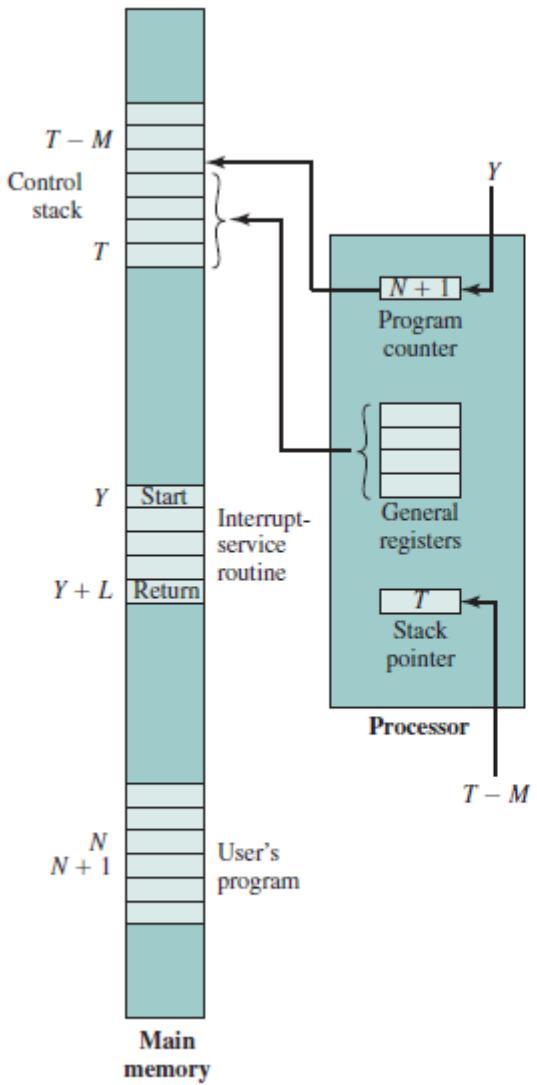
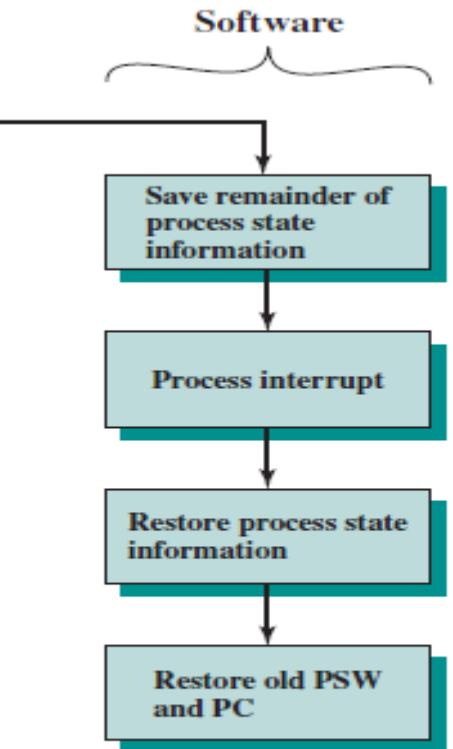
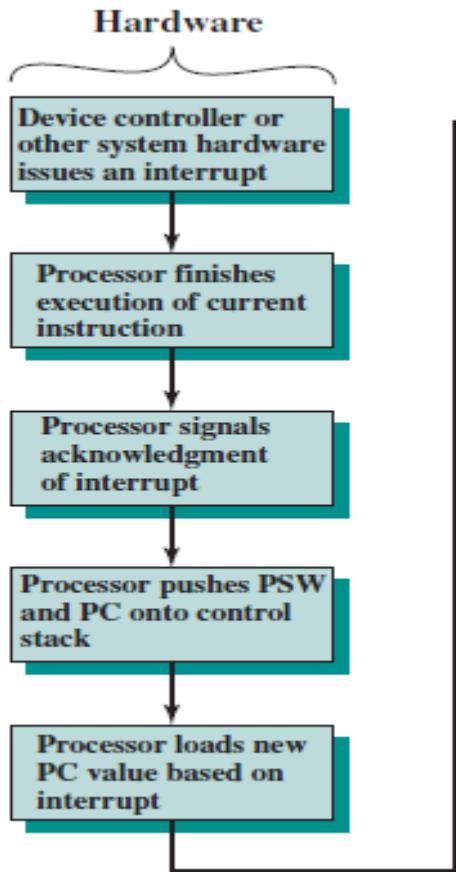
- With *programmed I/O*, data are exchanged between the processor and the I/O module.
- I/O Commands
 - Control: Used to activate a peripheral and tell it what to do.
 - Test: Used to test various status conditions associated with an I/O module and its peripherals.
 - Read: Causes the I/O module to obtain an item of data from the peripheral and place it in an internal buffer
 - Write: Causes the I/O module to take an item of data (byte or word) from the data bus and subsequently transmit that data item to the peripheral.

	No Interrupts	Use of Interrupts
I/O-to-memory transfer through processor	Programmed I/O	Interrupt-driven I/O
Direct I/O-to-memory transfer		Direct memory access (DMA)



Interrupt-driven I/O

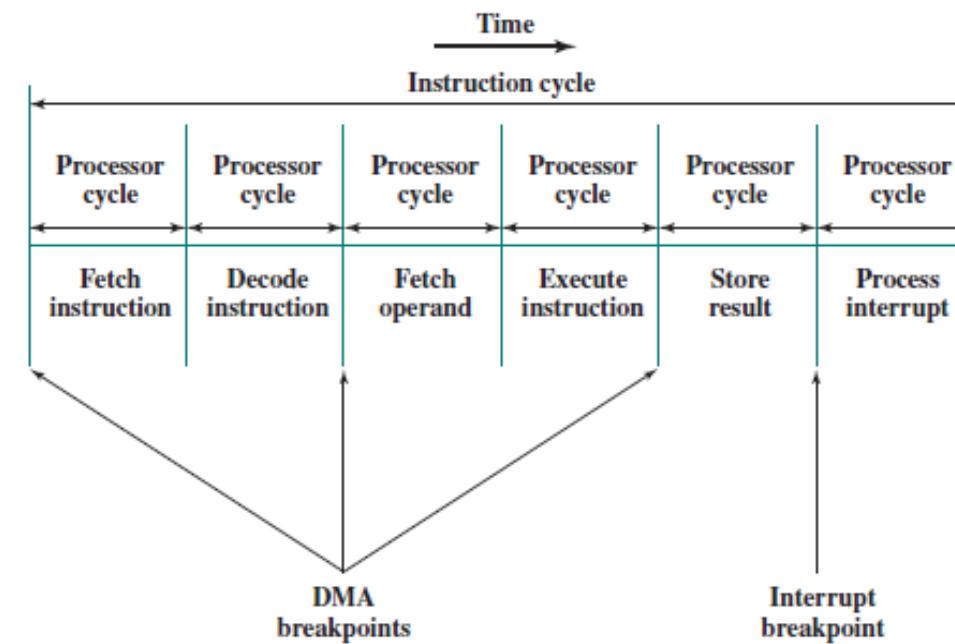
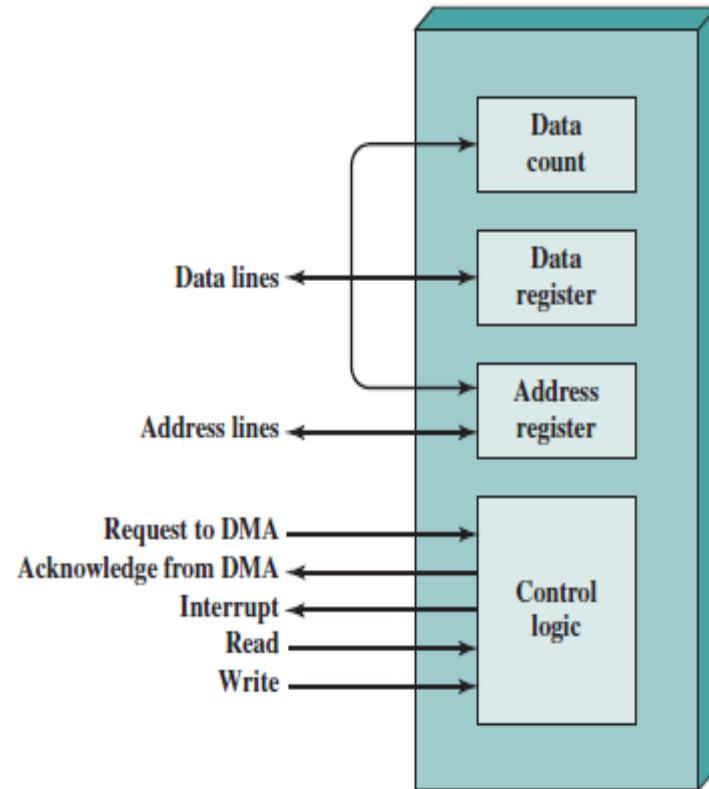
- The problem with programmed I/O is that the processor has to wait a long time for the I/O module of concern to be ready for either reception or transmission of data. The processor, while waiting, must repeatedly interrogate the status of the I/O module. As a result, the level of the performance of the entire system is severely degraded.
- An alternative is for the processor to issue an I/O command to a module and then go on to do some other useful work. The I/O module will then interrupt the processor to request service when it is ready to exchange data with the processor. The processor then executes the data transfer, as before, and then resumes its former processing.

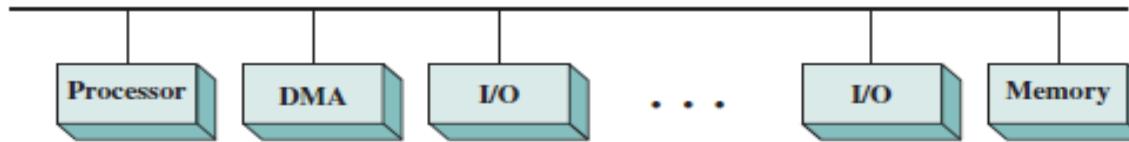


PSW=Program Status Word
PC= Program Counter

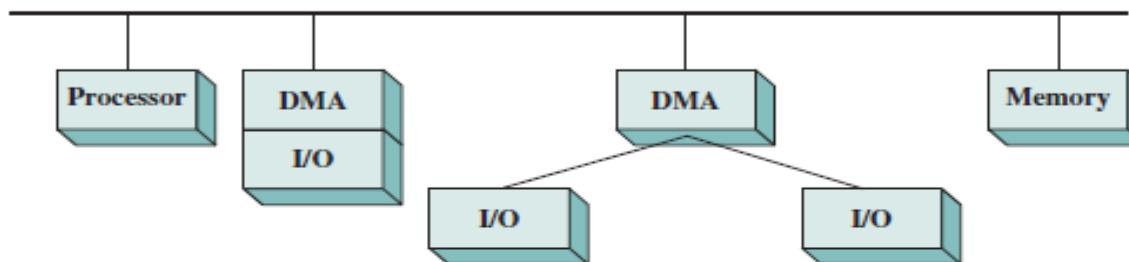
Direct Memory Access (DMA)

- Interrupt-driven I/O, though more efficient than simple programmed I/O, still requires the active intervention of the processor to transfer data between memory and an I/O module, and any data transfer must traverse a path through the processor. Thus, both these forms of I/O suffer from two inherent drawbacks:
 - The I/O transfer rate is limited by the speed with which the processor can test and service a device.
 - The processor is tied up in managing an I/O transfer; a number of instructions must be executed for each I/O transfer.
- When large volumes of data are to be moved, a more efficient technique is required: direct memory access (DMA).

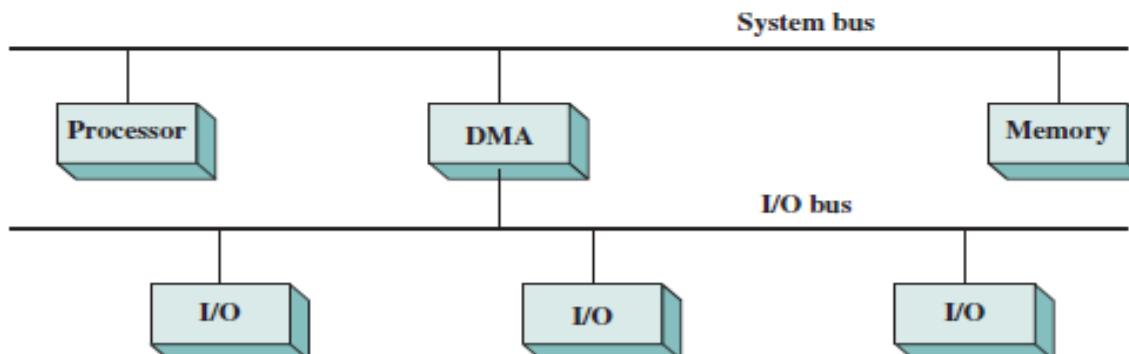




(a) Single-bus, detached DMA



(b) Single-bus, integrated DMA-I/O



(c) I/O bus

SERIAL COMMUNICATION

- A data communication processor is an I/O processor that distributes and collects data from many remote terminals connected through telephone and other communication lines. It is a specialized I/O processor designed to communicate directly with data communication networks.
- With the use of a data communication processor, the computer can service fragments of each network demand in an interspersed manner and thus have the apparent behavior of serving many users at once. In this way the computer is able to operate efficiently in a time-sharing environment.

- The most striking difference between an I/O processor and a data communication processor is in the way the processor communicates with the I/O devices. An I/O processor communicates with the peripherals through a common I/O bus that is comprised of many data and control lines. All peripherals share the common bus and use it to transfer information to and from the I/O processor. A data communication processor communicates with each terminal through a single pair of wires. Both data and control information are trans-character received.
- Another procedure used in asynchronous terminals involving a human operator is to echo the character. The character transmitted from the keyboard to the computer is recognized by the processor and retransmitted to the terminal printer

- In synchronous transmission, where an entire block of characters is transmitted, each character has a parity bit for the receiver to check. After the entire block is sent, the transmitter sends one more character that constitutes a parity over the length of the message. This character is called a longitudinal redundancy check (LRC). Another method used for checking errors in transmission is the cyclic redundancy check (CRC).
- A communication line may be connected to a synchronous or asynchronous interface, depending on the transmission method of the remote terminal. An asynchronous interface receives serial data with start and stop bits in each character. Synchronous transmission does not use start-stop bits to frame characters and therefore makes more efficient use of the communication link

- Contrary to asynchronous transmission, where each character can be sent separately with its own start and stop bits, synchronous transmission must send a continuous message in order to maintain synchronism. The message consists of a group of bits transmitted sequentially as a block of data. The entire block is transmitted with special control characters at the beginning and end of the block. The control characters at the beginning of the block supply the information needed to separate the incoming bits into individual characters.

- Data can be transmitted between two points in three different modes: simplex, half-duplex, or full-duplex.
- A half-duplex transmission system is one that is capable of transmitting in both directions but data can be transmitted in only one direction at a time.
- A full-duplex transmission can send and receive data in both directions simultaneously.
- The communication lines, modems, and other equipment used in the transmission of information between two or more stations is called a data link.
- The orderly transfer of information in a data link is accomplished by means of a protocol .