# Lecture1
# Memory Hierarchy

The memory unit is an essential component in any digital computer since it is needed for storing programs and data. A very small computer with a limited application may be able to fulfill its intended task without the need of additional storage capacity. Most general-purpose computers would run more efficiently if they were equipped with additional storage beyond the capacity of the main memory. There is just not enough space in one memory unit to accommodate all the programs used in a typical computer. Moreover, most computer users accumulate and continue to accumulate large amounts of data-processing software. Not all accumulated information is needed by the processor at the same time. Therefore, it is more economical to use low-cost storage devices to serve as a backup for storing the information that is not currently used by the CPU.

The memory unit that communicates directly with the CPU is called the **main memory**. Devices that provide backup storage are called **auxiliary memory**. The most common auxiliary memory devices used in computer systems are magnetic disks and tapes. They are used for storing system programs, large data files, and other backup information. Only programs and data currently needed by the processor reside in main memory. All other information is stored in auxiliary memory and transferred to main memory when needed.

The total memory capacity of a computer can be visualized as being a hierarchy of components. The memory hierarchy system consists of 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. Figure (1) illustrates the components in a typical memory hierarchy. At the bottom of the hierarchy are the relatively slow magnetic tapes used to store removable files.

Next the magnetic disks used as backup storage. The main memory occupies a central position by being able to communicate directly with the CPU and with auxiliary memory devices through an I/O processor. When programs not residing in main memory are needed by the CPU, they are brought in from auxiliary memory. Programs not currently needed in main memory are transferred into auxiliary memory to provide space for currently used programs and data.

A special very-high-speed memory called a **cache** is sometimes used to increase the speed of processing by making current programs and data available to the CPU at a rapid rate. The cache memory is employed in computer systems to compensate for the speed differential between main memory access time and processor logic. CPU logic is usually faster than main memory access time, with the result that processing speed is limited primarily by the speed of main memory. A technique used to compensate for the mismatch in operating speeds is to employ an extremely fast, small cache between the CPU and main memory whose access time is close to processor logic clock cycle time.

The cache is used for storing segments of programs currently being executed in the CPU and temporary data frequently needed in the present calculations.
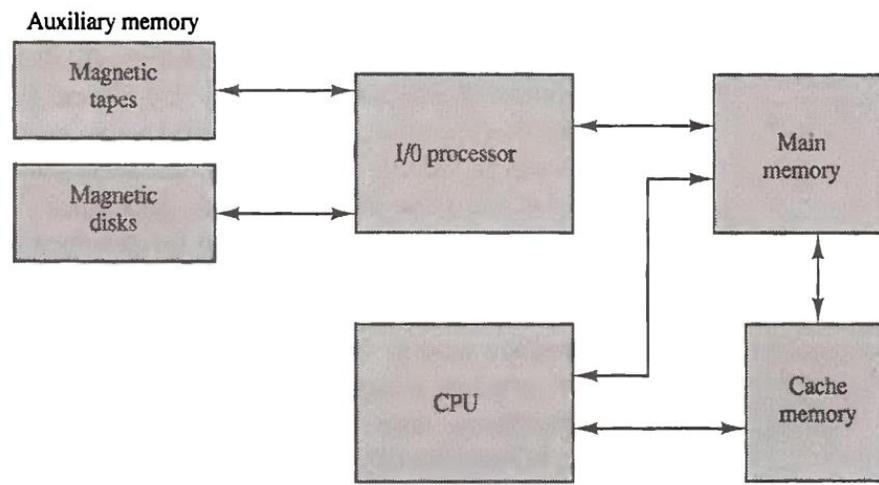


**Figure (1)**
**Memory hierarchy in computer system**

By making programs and data available at a rapid rate, it is possible to increase the performance rate of the computer.

While the I/O processor manages data transfers between auxiliary memory and main memory, the cache organization is concerned with the transfer of information between main memory and CPU. Thus, each is involved with a different level in the memory hierarchy system. The reason for having two or three levels of memory hierarchy is economics. As the storage capacity of the memory increases, the cost per bit for storing binary information decreases and the access time of the memory becomes longer. The auxiliary memory has a large storage capacity, is relatively inexpensive, but has low access speed compared to main memory. The cache memory is very small, relatively expensive, and has very high access speed. Thus as the memory access speed increases, so does its relative cost. The overall goal of using a memory hierarchy is to obtain the highest-possible average access speed while minimizing the total cost of the entire memory system. Auxiliary and cache memories are used for different purposes. The cache holds those parts of the program and data that are most heavily used, while the auxiliary memory holds those parts that are not presently used by the CPU. Moreover, the CPU has direct access to both cache and main memory but not to auxiliary memory. The transfer from auxiliary to main memory is usually done by means of direct memory access of large blocks of data. The typical access time ratio between cache and main memory is about 1 to 7.

**For example**, a typical cache memory may have an access time of 100 ns, while main memory access time may be 700 ns. Auxiliary memory average access time is usually 1000 times that of main memory. Block size in auxiliary memory typically ranges from 256 to 2048 words, while cache block size is typically from 1 to 16 words.

Many operating systems are designed to enable the CPU to process a number of independent programs concurrently. This concept, called **multiprogramming**, refers to the existence of two or more programs in different parts of the memory hierarchy at the same time. In this way it is possible to keep all parts of the computer busy by working with several programs in sequence.

**For example**, suppose that a program is being executed in the CPU and an I/O transfer is required. The CPU initiates the I/O processor to start executing the transfer. This leaves the CPU free to execute another program. In a multiprogramming system, when one program is waiting for input or output transfer, there is another program ready to utilize the CPU. With multiprogramming the need arises for running partial programs, for varying the amount of main memory in use by a given program, and for moving programs around the memory hierarchy. Computer programs are sometimes too long to be accommodated in the total space available in main memory. Moreover, a computer system uses many programs and all the programs cannot reside in main memory at all times. A program with its data normally resides in auxiliary memory. When the program or a segment of the program is to be executed, it is transferred to main memory to be executed by the CPU. Thus, one may think of auxiliary memory as containing the totality of information stored in a computer system. It is the task of the operating system to maintain in main memory a portion of this information that is currently active. The part of the computer system that supervises the flow of information between auxiliary memory and main memory is called the memory management system.
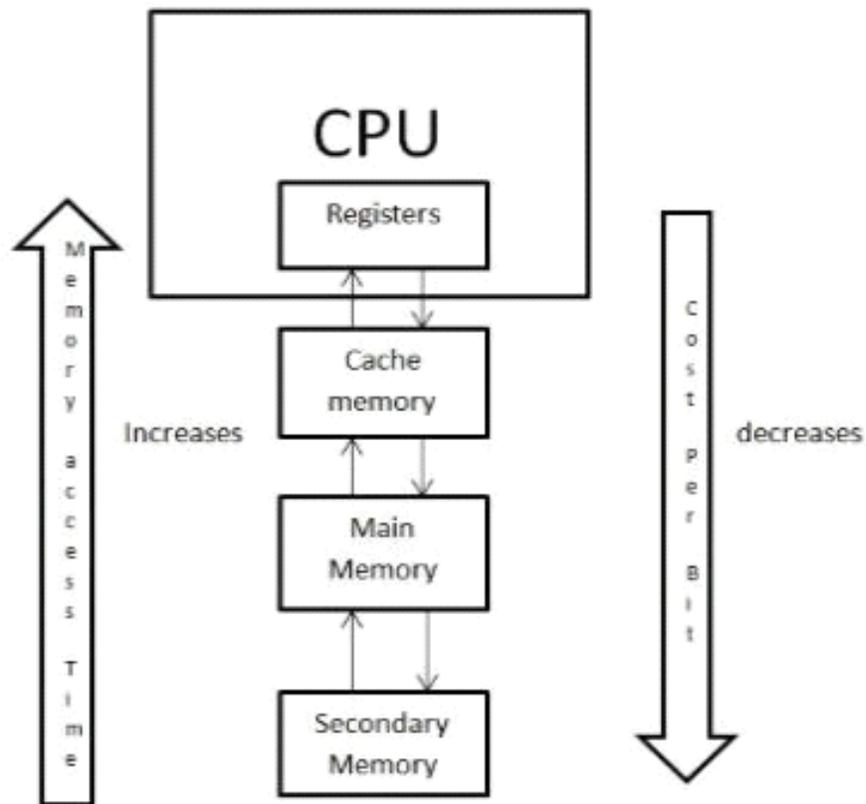
Figure 1

## Lecture2

## Main memory

The main memory is the central storage unit in a computer system. It is a relatively large and fast memory used to store programs and data during the computer operation. The principal technology used for the main memory is based on semiconductor integrated circuits. Integrated circuit RAM chips are available in two possible operating modes, static and dynamic. The static RAM consists essentially of internal flip-flops that store the binary information. The stored information remains valid as long as power is applied to the unit. The dynamic RAM stores the binary information in the form of electric charges that are applied to capacitors. The capacitors are provided inside the chip by MOS transistors. The stored charge on the capacitors tend to discharge with time and the capacitors must be periodically recharged by refreshing the dynamic memory. Refreshing is done by cycling through the words every few milliseconds to restore the decaying charge. The dynamic RAM offers reduced power consumption and larger

storage capacity in a single memory chip. The static RAM is easier to use and has shorter read and write cycles.

Most of the main memory in a general-purpose computer is made up of RAM integrated circuit chips, but a portion of the memory may be constructed with ROM chips. Originally, RAM was used to refer to a random-access memory, but now it is used to designate a read/write memory to distinguish it from a read-only memory, although ROM is also random access. RAM is used for storing the bulk of the programs and data that are subject to change. ROM is used for storing programs that are permanently resident in the computer and for tables of constants that do not change in value once the production of the computer is completed.

Among other things, the ROM portion of main memory is needed for storing an initial program called a **bootstrap loader**. The bootstrap loader is a program whose function is to start the computer software operating when power is turned on. Since RAM is volatile, its contents are destroyed when power is turned off. The contents of ROM remain unchanged after power is turned off and on again. The startup of a computer consists of turning the power on and starting the execution of an initial program. Thus when power is turned on, the hardware of the computer sets the program counter to the first address of the bootstrap loader. The bootstrap program loads a portion of the operating system from disk to main memory and control is then transferred to the operating system, which prepares the computer for general use. RAM and ROM chips are available in a variety of sizes. If the memory needed for the computer is larger than the capacity of one chip, it is necessary to combine a number of chips to form the required memory size. To demonstrate the chip interconnection,
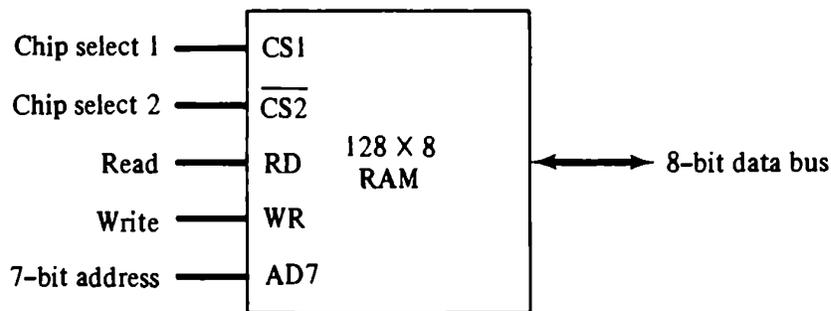
## RAM and ROM Chips

A RAM chip is better suited for communication with the CPU if it has one or more control inputs that select the chip only when needed. Another common feature is a bidirectional data bus that allows the transfer of data either from memory to CPU during a read operation or from CPU to memory during a write operation. A bidirectional bus can be constructed with three-state buffers. A three-state buffer output can be placed in one of three possible states: a signal equivalent to logic 1, a signal equivalent to logic 0, or a high- impedance state. The logic 1 and 0 are normal digital signals. The high- impedance state behaves like an open circuit, which means that the output does not carry a signal and has no logic significance.

The block diagram of a RAM chip is shown in Fig (2). The capacity of the memory is 128 words of eight bits (one byte) per word. This requires a 7-bit address and an 8-bit bidirectional data bus. The read and write inputs spear, the memory operation and the two chips select (CS) control inputs are ic: enabling the chip only when it is selected by the microprocessor. The availability of more than one control input to select the chip facilitates the decoding c: the address lines when multiple chips are used in the microcomputer. The read and write inputs are sometimes

combined into one line labeled R/W. When the chip is selected, the two binary states in this line specify the two operations of read or write.

The function table listed in Fig. (2) Specifies the operation of the RAM chip. The unit is in operation only when CS1 = 1 and CS2 = 0. The bar on top of the second select variable indicates that this input is enabled when it is equal to 0. If the chip select inputs are not enabled, or if they are enabled but the read or write inputs are not enabled, the memory is inhibited, and its data bus is in a high-impedance state. When CS1 = 1 and CS2 = 0, the memory can be placed in a write or read mode. When the WR input is enabled, the memory stores a byte from the data bus into a location specified by the address input lines. When the RD input is enabled, the content of the selected byte is placed into the data bus. The RD and WR signals control the memory operation as well as the bus buffers associated with the bidirectional data bus.
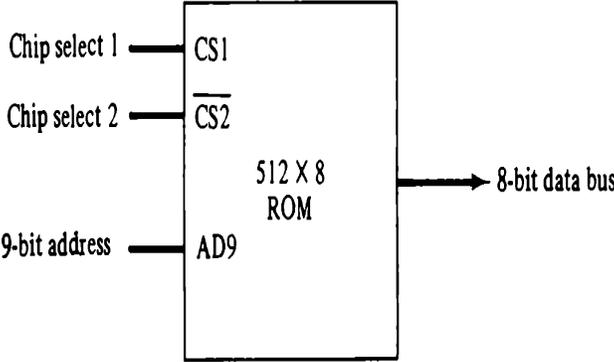


(a) Block diagram

| CS1 | $\overline{CS2}$ | RD | WR | Memory function | State of data bus |
|---|---|---|---|---|---|
| 0 | 0 | × | × | Inhibit | High-impedance |
| 0 | 1 | × | × | Inhibit | High-impedance |
| 1 | 0 | 0 | 0 | Inhibit | High-impedance |
| 1 | 0 | 0 | 1 | Write | Input data to RAM |
| 1 | 0 | 1 | × | Read | Output data from RAM |
| 1 | 1 | × | × | Inhibit | High-impedance |

(b) Function table

**Fig (2) Typical RAM chip**

A ROM chip is organized externally in a similar manner. However, since a ROM can only read, the data bus can only be in an output mode. The block diagram of a ROM chip is shown in Fig(3). For the same-size chip, it is possible to have more bits of ROM than of RAM, because the internal binary cells in ROM occupy less space than in RAM. For this reason, the diagram Specifies a 512-byte ROM, while the RAM has only 128 bytes.

The nine address lines in the ROM chip specify any one of the 512 bytes stored in it. The two chip select inputs must be CS1 = 1 and CS2 = 0 for the unit to operate. Otherwise, the data bus is in a high-impedance state. There is no need for a read or write control because the unit can only read.



**Typical ROM chip**

**Lecture3**
**Memory Address Map**

The designer of a computer system must calculate the amount of memory required for the particular application and assign it to either RAM or ROM. The interconnection between memory and processor is then established from knowledge of the size of memory needed and the type of RAM and ROM chips available. The addressing of memory can be established by means of a table that specifies the memory address assigned to each chip. The table, called a memory address map, is a pictorial representation of assigned address space for each chip in the system.

To demonstrate with a particular example, assume that a computer system needs 512 bytes of RAM and 512 bytes of ROM. The RAM and ROM chips to be used are specified in Figs. (2) and (3). The memory address map for this configuration is shown in Table 1 below. The component column specifies whether a RAM or a ROM chip is used. The hexadecimal address column assigns a range of hexadecimal equivalent addresses for each chip. The address bus lines are listed in the third column. Although there are 16 lines in the address bus, the table shows only 10 lines because the other 6 are not used in this example and are assumed to be zero. The small x's under the address bus lines designate those lines that must be connected to the address inputs in each chip. The RAM chips have 128 bytes and need seven address lines. The ROM chip has 512 bytes and needs 9 address lines. The x's are always assigned to the low-order bus lines: lines 1 through 7 for the RAM and lines 1 through 9 for the ROM. It is now necessary to distinguish between four RAM chips by assigning to each a different address. For this particular example we choose bus lines 8 and 9 to represent four distinct binary combinations. Note that any other pair of unused bus lines can be chosen for this purpose. The table clearly shows that the nine low-order bus lines constitute a memory space for RAM equal to 29 = 512 bytes. The distinction between a RAM and ROM address is done with another bus line. Here we choose line 10 for this purpose. When line 10 is 0, the CPU selects a RAM, and when this line is equal to 1, it selects the ROM. The equivalent hexadecimal address for each chip is obtained from the information under the address bus assignment. The address bus lines are subdivided into groups of four bits each so that each group can be represented with a hexadecimal digit. The first hexadecimal digit represents lines 13 to 16 and is always 0. The next hexadecimal digit represents lines 9 to 12, but lines 11 and 12 are always 0. The range of hexadecimal addresses for each component is determined from the x's associated with it. These x's represent a binary number that can range from an all-0's to an all-l's value.

## TABLE 12-1 Memory Address Map for Microprocomputer

| Component | Hexadecimal address | Address bus | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | 10 | 9 | 8 | 7 | 6 | 5 | 4 | 3 | 2 | 1 |
| RAM 1 | 0000–007F | 0 | 0 | 0 | x | x | x | x | x | x | x |
| RAM 2 | 0080–00FF | 0 | 0 | 1 | x | x | x | x | x | x | x |
| RAM 3 | 0100–017F | 0 | 1 | 0 | x | x | x | x | x | x | x |
| RAM 4 | 0180–01FF | 0 | 1 | 1 | x | x | x | x | x | x | x |
| ROM | 0200–03FF | 1 | x | x | x | x | x | x | x | x | x |

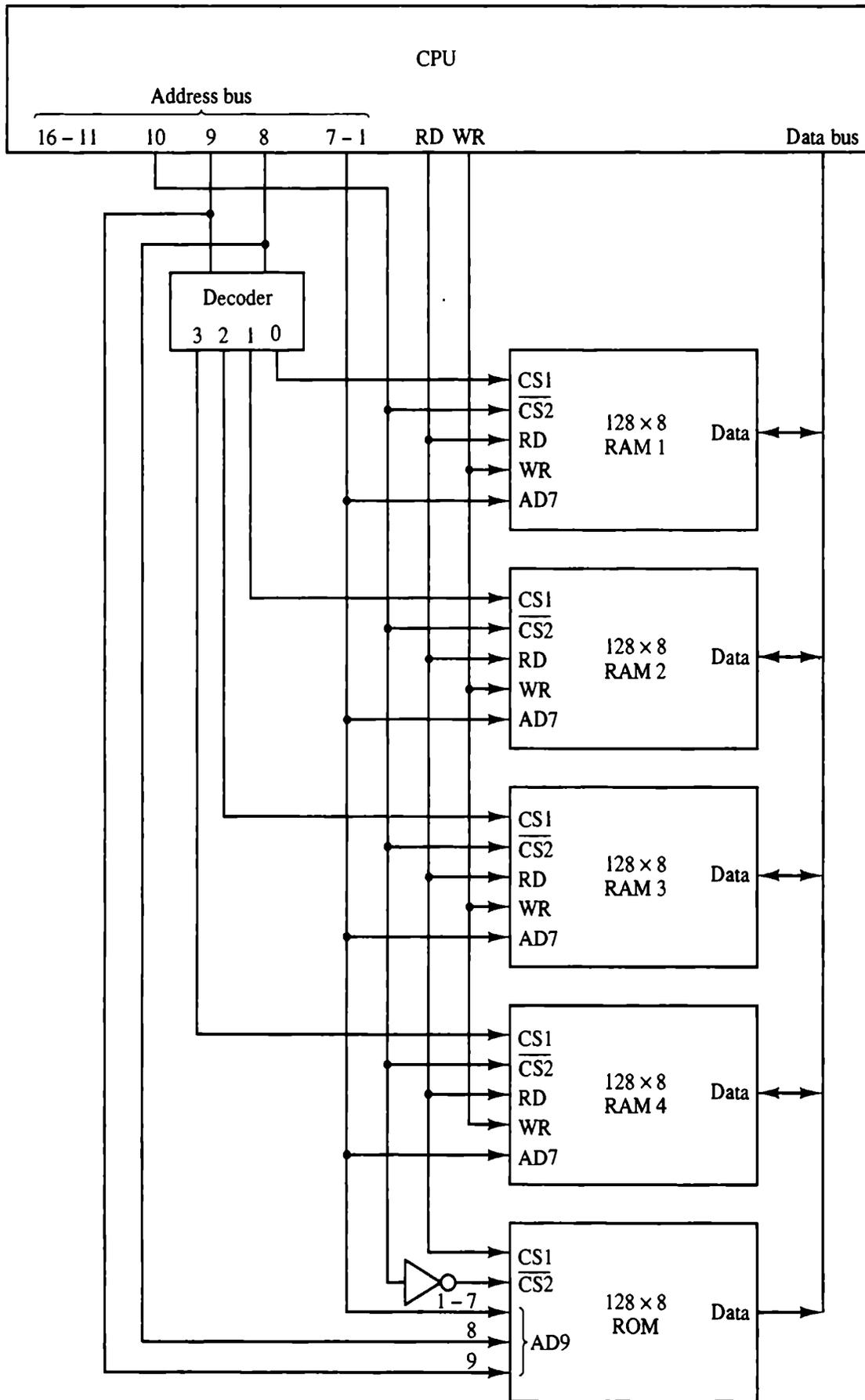Table 1: Memory address map for microcomputer

**Fig (4) Memory connection to CPU**

## Lecture4
## 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*. The reason for this property may be understood considering that a typical computer program flows in a straight-line fashion with program loops and subroutine calls encountered frequently. When a program loop is executed, the CPU repeatedly refers to the set of instructions in memory that constitute the loop. Every time a given subroutine is called, its set of instructions is fetched from memory. Thus loops and subroutines tend to localize the references to memory for fetching instructions. To a lesser degree, memory references to data also tend to be localized. **Table-lookup** procedures repeatedly refer to that portion in memory where the table is stored. Iterative procedures refer to common memory locations and array of numbers are confined within a local portion of memory. The result of all these observations is the locality of reference property, which states that over a short interval of time, the addresses generated by a typical program refer to a few localized areas of memory repeatedly, while the remainder of memory is accessed relatively infrequently.

If the active portions of the program and data are placed in a fast small memory, the average memory access time can be reduced, thus reducing the total execution time of the program. Such a fast small memory is referred to as a **cache memory**. It is placed between the CPU and main memory as illustrated in Fig. (1). the cache memory access time is less than the access time of main memory by a factor of 5 to 10. The cache is the fastest component in the memory hierarchy and approaches the speed of CPU components.

The fundamental idea of cache organization is that by keeping the most frequently accessed instructions and data in the fast cache memory, the average memory access time will approach the access time of the cache. Although the cache is only a small fraction of the size of main memory, a large fraction of memory requests will be found in the fast cache memory because of the locality of reference property of programs.

The basic operation of the cache is as follows:

When the CPU needs to access memory, the cache is examined. If the word is found in the cache, it is read from the fast memory. If the word addressed by the CPU is not found in the cache, the main memory is accessed to read the word. A block of words containing the one just accessed is then transferred from main memory to cache memory. The block size may vary from one word (the one just accessed) to about 16 words adjacent to the one just accessed. In this

manner, some data are transferred to cache so that future references to memory find the required words in the fast cache memory.

The performance of cache memory is frequently measured in terms of a quantity called *hit ratio*. When the CPU refers to memory and finds the word in cache, it is said to produce a **hit**. If the word is not found in cache, it is in main memory and it counts as **a miss**. The ratio of the number of hits divided by the total CPU references to memory (hits plus misses) is the hit ratio. The hit ratio is best measured experimentally by running representative programs in the computer and measuring the number of hits and misses during a given interval of time. Hit ratios of 0.9 and higher have been reported. This high ratio verifies the validity of the locality of reference property.

The average memory access time of a computer system can be improved considerably by use of a cache. If the hit ratio is high enough so that most of the time the CPU accesses the cache instead of main memory, the average access time is closer to the access time of the fast cache memory. For example, a computer with cache access time of 100 ns, a main memory access time of 1000 ns, and a hit ratio of 0.9 produces an average access time of 200 ns. This is a considerable improvement over a similar computer without a cache memory, whose access time is 1000 ns. The basic characteristic of cache memory is its fast access time. Therefore, very little or no time must be wasted when searching for words in the cache.

The transformation of data from main memory to cache memory is referred to as a **mapping process.** Three types of mapping procedures are of practical interest when considering the organization of cache memory:

1. Associative mapping
2. Direct mapping
3. Set-associative mapping

To help in the discussion of these three mapping procedures we will use a specific example of a memory organization as shown in Fig (5). The main memory can store 32K words of 12 bits each. The cache is capable of storing 512 of these words at any given time. For every word stored in cache, there is a duplicate copy in main memory. The CPU communicates with both memories. It first sends a 15-bit address to cache. If there is a hit, the CPU accepts the 12-bit data from cache. If there is a miss, the CPU reads the word from main memory and the word is then transferred to cache.
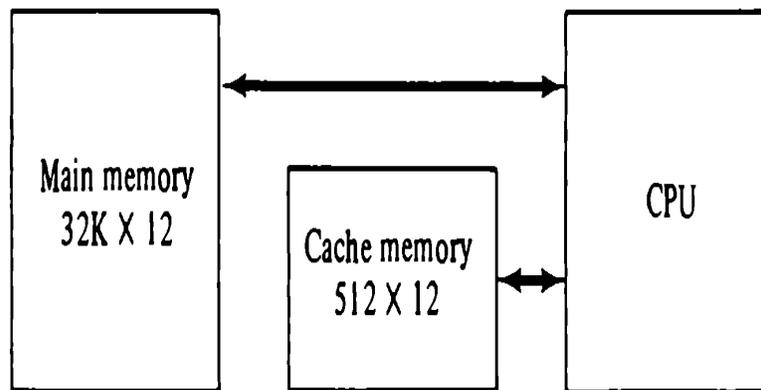
**Figure (5)**
**Example of cache memory.**

## Lecture5
## Associative Mapping

The fastest and most flexible cache organization uses an associative memory. This organization is illustrated in Fig. (6). The associative memory stores both the address and content (data) of the memory word. This permits any location in cache to store any word from main memory. The diagram shows three words presently stored in the cache. The address value of 15 bits is shown as a five-digit octal number and its corresponding 12-bit word is shown as a four-digit octal number. A CPU address of 15 bits is placed in the argument register and the associative memory is searched for a matching address. If the address is found, the corresponding 12-bit data is read and sent to the CPU. If no match occurs, the main memory is accessed for the word. The address-data pair is then transferred to the associative cache memory. If the cache is full, an address-data pair must be displaced to make room for a pair that is needed and not presently in the cache. The decision as to what pair is replaced is determined from the replacement algorithm that the designer chooses for the cache. A simple procedure is to replace cells of the cache in round-robin order whenever a new word is requested from main memory. This constitutes a first-in first-out (FIFO) replacement policy.
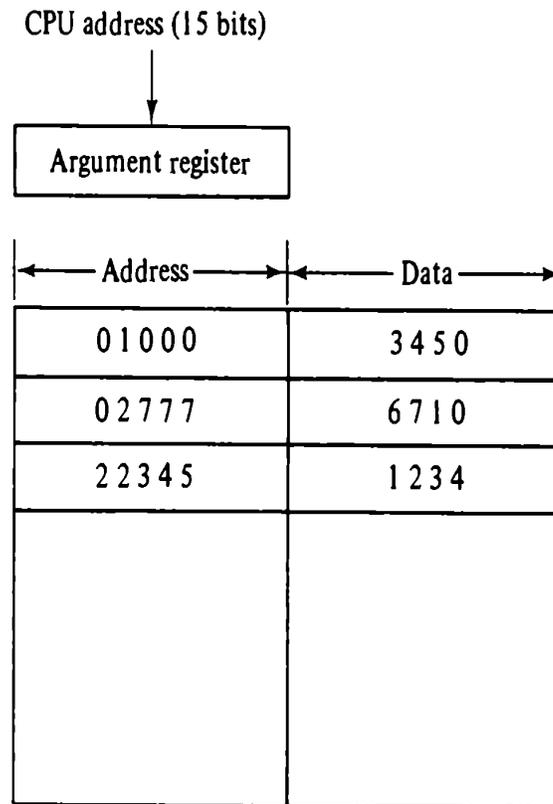
CPU address (15 bits)

Argument register

| ← Address → | ← Data → |
|---|---|
| 0 1 0 0 0 | 3 4 5 0 |
| 0 2 7 7 7 | 6 7 1 0 |
| 2 2 3 4 5 | 1 2 3 4 |
|  |  |

**Fig (6)**
**Associative mapping cache (all numbers in octal).**

## Direct Mapping

Associative memories are expensive compared to random-access memories because of the added logic associated with each cell. The possibility of using a random-access memory for the cache is investigated in Fig. (7). The CPU address of 15 bits is divided into two fields. The nine least significant bits constitute the index field and the remaining six bits from the tag field. The figure shows that main memory needs an address that includes both the tag and the index bits. The number of bits in the index field is equal to the number of address bits required to access the cache memory.

In the general case, there are $2^k$ words in cache memory and $2^n$ words in main memory. The n-bit memory address is divided into two fields: k bits for the index field and n - k bits for the tag field. The direct mapping cache organization uses the w-bit address to access the main memory and the k-bit index to access the cache. The internal organization of the words in the cache memory is as shown in Fig. (8-b). Each word in cache consists of the data word and its

associated tag. When a new word is first brought into the cache, the tag bits are stored alongside the data bits. When the CPU generates a memory request, the index field is used for the address to access the cache.
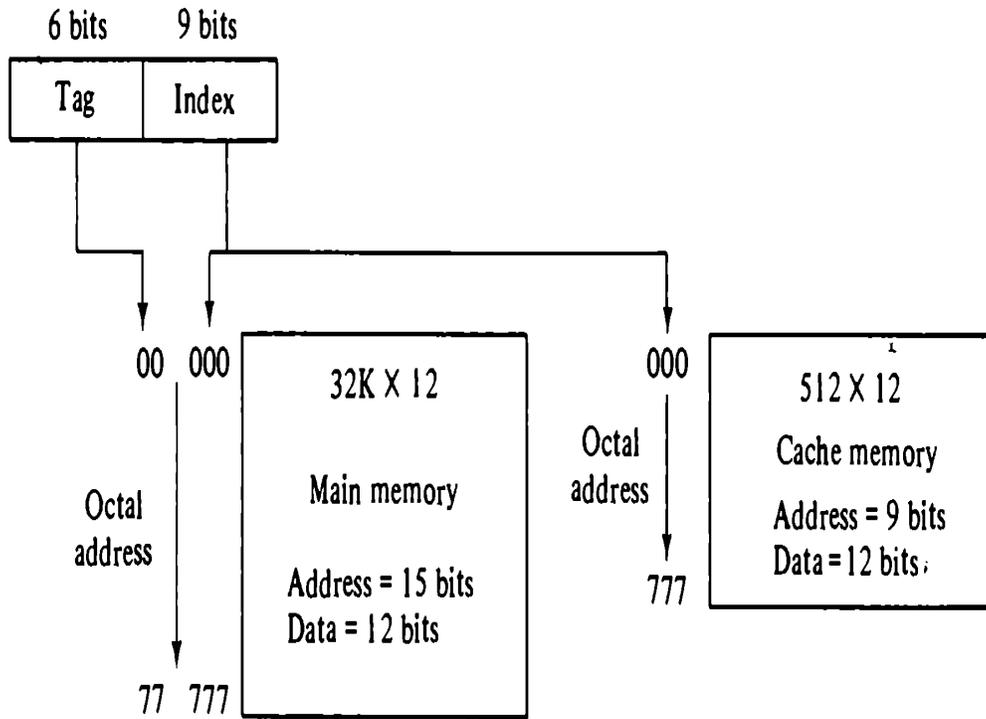


**Fig (7)**

**Addressing relationships between main and cache memories.**

| Memory address | Memory data |
|---|---|
| 00000 | 1 2 2 0 |
| | |
| 00777 | 2 3 4 0 |
| 01000 | 3 4 5 0 |
| | |
| 01777 | 4 5 6 0 |
| 02000 | 5 6 7 0 |
| | |
| 02777 | 6 7 1 0 |

(a)  Main memory

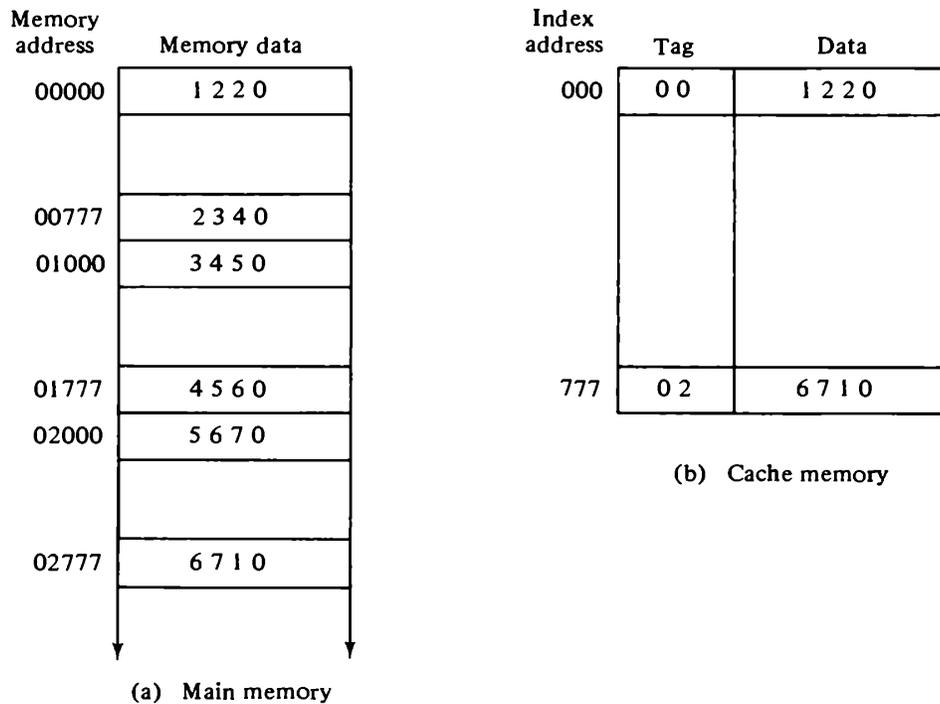| Index address | Tag | Data |
|---|---|---|
| 000 | 0 0 | 1 2 2 0 |
| | | |
| 777 | 0 2 | 6 7 1 0 |

(b)  Cache memory

**Fig (8)**
**Direct mapping cache organization.**

The tag field of the CPU address is compared with the tag in the word read from the cache. If the two tags match, there is a **hit** and the desired data word is in cache. If there is no match, there is a **miss** and the required word is read from main memory. It is then stored in the cache together with the new tag, replacing the previous value. The disadvantage of direct mapping is that the hit ratio can drop considerably if two or more words whose addresses have the same index but different tags are accessed repeatedly. However, this possibility is minimized by the fact that such words are relatively far apart in the address range (multiples of 512 locations in this example.)

To see how the direct-mapping organization operates, consider the numerical example shown in Fig. (8). the word at address zero is presently stored in the cache (index = 000, tag = 00, data = 1220). Suppose that the CPU now wants to access the word at address 02000. The index address is 000, so it is used to access the cache. The two tags are then compared. The cache tag is 00 but the address tag is 02, which does not produce a match. Therefore, the main memory is accessed and the data word 5670 is transferred to the CPU. The cache word at index address 000 is then replaced with a tag of 02 and data of 5670. The direct-mapping example just described uses a block size of one word. The same organization but using a block size of 8 words is shown in   Fig. (9)
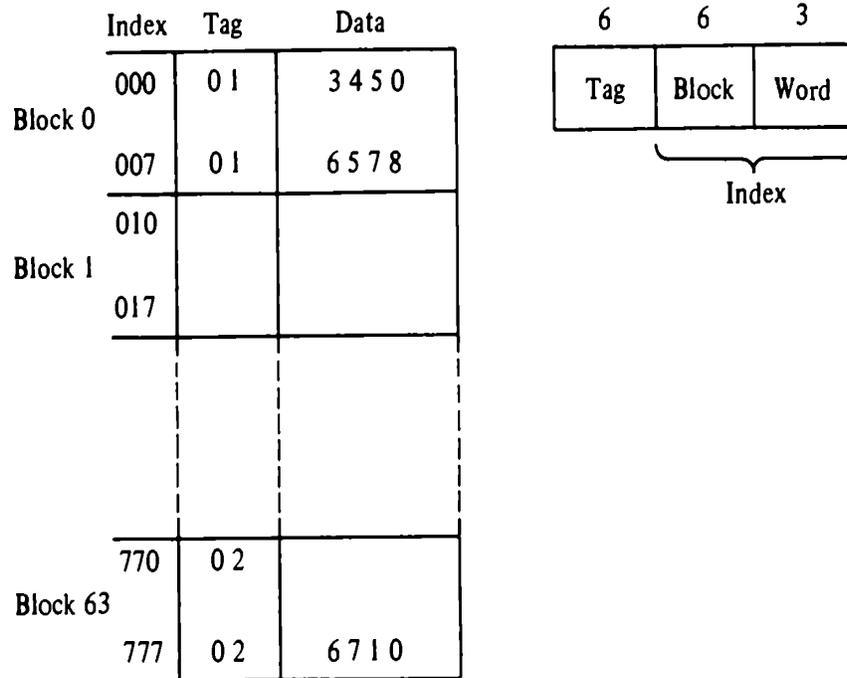
| Index | Tag | Data |
|-------|-----|------|
| 000 | 0 1 | 3 4 5 0 |
| 007 | 0 1 | 6 5 7 8 |
| 010 | | |
| 017 | | |
| ⋮ | ⋮ | ⋮ |
| 770 | 0 2 | |
| 777 | 0 2 | 6 7 1 0 |

Block 0 (000–007), Block 1 (010–017), Block 63 (770–777)

| 6 | 6 | 3 |
|-----|-------|------|
| Tag | Block | Word |

Index = Block + Word

**Fig (9)**
**Direct mapping cache with block size of 8 words.**

The index field is now divided into two parts: the block field and the word field. In a 512-word cache there are 64 blocks of 8 words each, since 64 x 8 = 512. The block number is specified with a 6-bit field and the word within the block is specified with a 3-bit field. The tag field stored within the cache is common to all eight words of the same block. Every time a miss occurs, an entire block of eight words must be transferred from main memory to cache memory.

Although this takes extra time, the hit ratio will most likely improve with a larger block size because of the sequential nature of computer programs.

## Lecture6

## Set Associative Mapping

It was mentioned previously that the disadvantage of direct mapping is that two words with the same index in their address but with different tag values cannot reside in cache memory at the same time. A third type of cache organization, called set-associative mapping, is an improvement over the direct- mapping organization in that each word of cache can store two or more words of memory under the same index address. Each data word is stored together with

its tag and the number of tag-data items in one word of cache is said to form a set. An example of a set-associative cache organization for a set size of two is shown in Fig. (10). Each index address refers to two data words and their associated tags. Each tag requires six bits, and each data word has 12 bits, so the word length is 2(6 + 12) = 36 bits. An index address of nine bits can accommodate 512 words. Thus, the size of cache memory is 512 x 36. It can accommodate 1024 words of main memory since each word of cache contains two data words. In general, a set-associative cache of set size k will accommodate k words of main memory in each word of cache.

| Index | Tag | Data | Tag | Data |
|-------|-----|------|-----|------|
| 000 | 0 1 | 3 4 5 0 | 0 2 | 5 6 7 0 |
| | | | | |
| 777 | 0 2 | 6 7 1 0 | 0 0 | 2 3 4 0 |

**Fig (10)**
**Two-way set-associative mapping cache.**

The octal numbers listed in Fig. (10) are with reference to the main memory contents illustrated in Fig. (8-a). The words stored at addresses 01000 and 02000 of main memory are stored in cache memory at index address 000. Similarly, the words at addresses 02777 and 00777 are stored in cache at index address 777. When the CPU generates a memory request, the index value of the address is used to access the cache. The tag field of the CPU address is then compared with both tags in the cache to determine if a match occurs. The comparison logic is done by an associative search of the tags in the set similar to an associative memory search: thus, the name "set-associative". The hit ratio will improve as the set size increases because more words with the same index, but different tags can reside in cache. However, an increase in

the set size increases the number of bits in words of cache and requires more complex comparison logic.

When a miss occurs in a set-associative cache and the set is full, it is necessary to replace one of the tag-data items with a new value. The most common replacement algorithms used are **random replacement, first-in, first- out (FIFO), and least recently used (LRU)**. With the random replacement policy, the control chooses one tag-data item for replacement at random. The FIFO procedure selects for replacement the item that has been in the set the longest. The LRU algorithm selects for replacement the item that has been least recently used by the CPU. Both FIFO and LRU can be implemented by adding a few extra bits in each word of cache.

## Lecture7
## Writing into Cache

An important aspect of cache organization is concerned with memory write requests. When the CPU finds a word in cache during a read operation, the main memory is not involved in the transfer. However, if the operation is a write, there are two ways that the system can precede.

**The first procedure** is simplest and most commonly used procedure is to update main memory with every memory write operation, with cache memory being updated in parallel if it contains the word at the specified address. This is called the **write-through method**. This method has the advantage that main memory always contains the same data as the cache. This characteristic is important in systems with direct memory access transfers. It ensures that the data residing in main memory are valid at all times so that an I/O device communicating through DMA would receive the most recent updated data.

**The second procedure** is called the **write-back method**. In this method only the cache location is updated during a write operation. The location is then marked by a flag so that later when the word is removed from the cache it is copied into main memory. The reason for the write-back method is that during the time a word resides in the cache, it may be updated several times; however, as long as the word remains in the cache, it does not matter whether the copy in main memory is out of date, since requests from the word are filled from the cache. It is only when the word is displaced from the cache that an accurate copy need be rewritten into main memory. Analytical results indicate that the number of memory writes in a typical program ranges between 10 and 30 percent of the total references to memory.

## Cache Initialization

One more aspect of cache organization that must be taken into consideration is the problem of initialization. The cache is initialized when power is applied to the computer or when the main memory is loaded with a complete set of programs from auxiliary memory. After initialization the cache is considered to be empty, but in effect it contains some non-valid data. It is

customary to include with each word in cache a valid bit to indicate whether or not the word contains valid data.

The cache is initialized by clearing all the valid bits to 0. The valid bit of a particular cache word is set to 1 the first time this word is loaded from main memory and stays set unless the cache has to be initialized again. The introduction of the valid bit means that a word in cache is not replaced by another word unless the valid bit is set to 1 and a mismatch of tags occurs. If the valid bit happens to be 0, the new word automatically replaces the invalid data. Thus, the initialization condition has the effect of forcing misses from the cache until it fills with valid data.

## Virtual Memory

In a memory hierarchy system, programs and data are first stored in auxiliary memory. Portions of a program or data are brought into main memory as they are needed by the CPU. Virtual memory is a concept used in some large computer systems that permit the user to construct programs as though a large memory space were available, equal to the totality of auxiliary memory. Each address that is referenced by the CPU goes through an address mapping from the so-called virtual address to a physical address in main memory. Virtual memory is used to give programmers the illusion that they have a very large memory at their disposal, even though the computer actually has a relatively small main memory. A virtual memory system provides a mechanism for translating program-generated addresses into correct main memory locations.

This is done dynamically, while programs are being executed in the CPU. The translation or mapping is handled automatically by the hardware by means of a mapping table.

## Lecture8
## Address Space and Memory Space

An address used by a programmer will be called a virtual address, and the set of such addresses the address space. An address in main memory is called a location or physical address. The set of such locations is called the memory space. Thus, the address space is the set of addresses generated by programs as they reference instructions and data; the memory space consists of the actual main memory locations directly addressable for processing. In most computers the address and memory spaces are identical. The address space is allowed to be larger than the memory space in computers with virtual memory.

As an illustration, consider a computer with a main-memory capacity of 32K words (K = 1024). Fifteen bits are needed to specify a physical address in memory since $32K = 2^{15}$. Suppose that the computer has available auxiliary memory for storing $2^{20} = 1024K$ words. Thus, auxiliary memory has a capacity for storing information equivalent to the capacity of 32 main memories.

Denoting the address space by N and the memory space by M, we then have for this example N = 1024K and M = 32K.

In a multiprogramming computer system, programs and data are transferred to and from auxiliary memory and main memory based on demands imposed by the CPU. Suppose that program 1 is currently being executed in the CPU.

Program 1 and a portion of its associated data are moved from auxiliary memory into main memory as shown in Fig. (11). Portions of programs and data need not be in contiguous locations in memory since information is being moved in and out, and empty spaces may be available in scattered locations in memory.

In a virtual memory system, programmers are told that they have the total address space at their disposal. Moreover, the address field of the instruction code has a sufficient number of bits to specify all virtual addresses. In our example, the address field of an instruction code will consist of 20 bits, but physical memory addresses must be specified with only 15 bits. Thus, CPU will reference instructions and data with a 20-bit address, but the information at this address must be taken from physical memory because access to auxiliary storage for individual words will be prohibitively long. (Remember that for efficient transfers, auxiliary storage moves an entire record to the main memory.) A table is then needed, as shown in Fig. (11), to map a virtual address of 20 bits to a physical address of 15 bits. The mapping is a dynamic operation, which means that every address is translated immediately as a word is referenced by CPU.
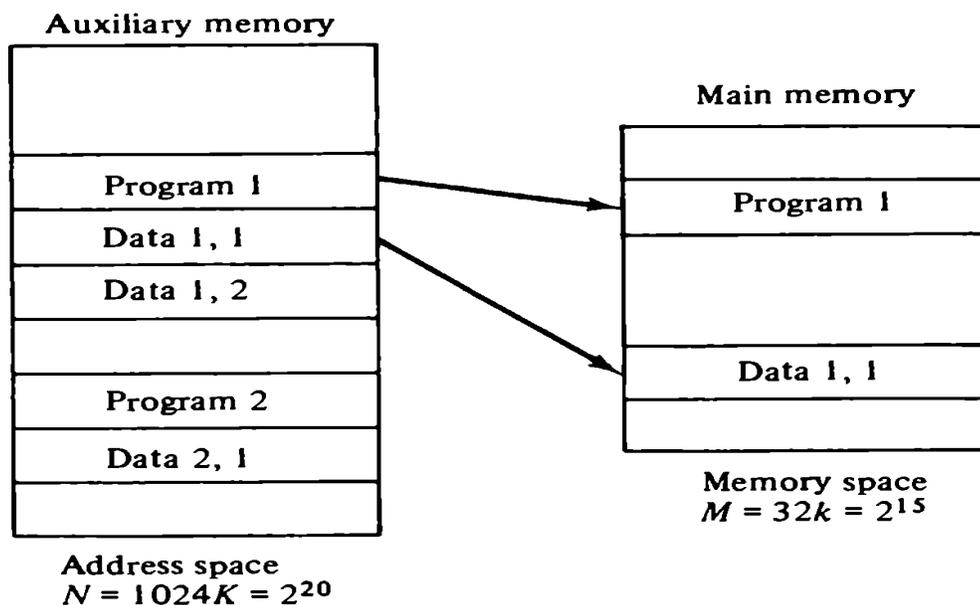


**Fig (11)**
**Relation between address and memory space in a virtual memory system.**

The mapping table may be stored in a separate memory as shown in Fig. (12), or in main memory. In the first case, an additional memory unit is required as well as one extra memory access time. In the second case, the table takes space from main memory and two accesses to memory are required with the program running at half speed. A third alternative is to use an associative memory.
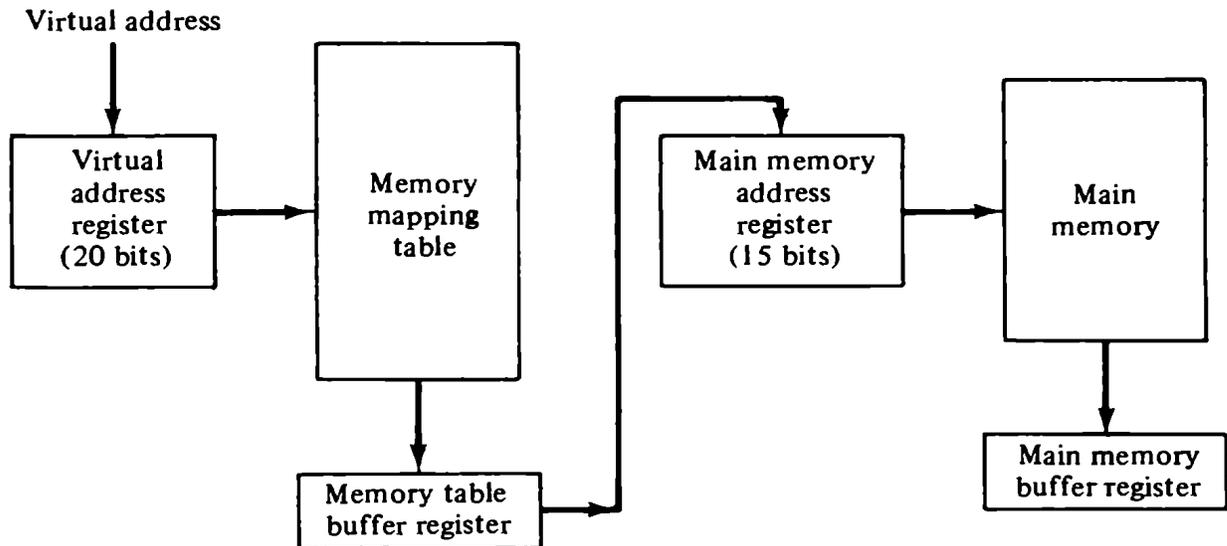


**Fig (12)**
**Memory table for mapping a virtual address.**

## Lecture9
## Address Mapping Using Pages

The table implementation of the address mapping is simplified if the information in the address space and the memory space are each divided into groups of fixed size. The physical memory is broken down into groups of equal size called blocks, which may range from 64 to 4096 words each. The term page refers to groups of address space of the same size. For example, if a page or block consists of 1K words, then, using the previous example, address space is divided into 1024 pages and main memory is divided into 32 blocks. Although both a page and a block are split into groups of IK words, a page refers to the organization of address space, while a block refers to the organization of memory space. The programs are also considered to be split into pages.

Portions of programs are moved from auxiliary memory to main memory in records equal to the size of a page. The term "page frame" is sometimes used to denote a block.

Consider a computer with an address space of 8K and a memory space of 4K. If we split each into groups of IK words we obtain eight pages and four blocks as shown in Fig. (13). At any given time, up to four pages of address space may reside in main memory in any one of the four blocks. The mapping from address space to memory space is facilitated if each virtual address is considered to be represented by two numbers: **a page number address** and **a line within the page.** In a computer with 2P words per page, p bits are used to specify a line address and the remaining high-order bits of the virtual address specify the page number. In the example of Fig. (13), a virtual address has 13 bits. Since each page consists of $2^{10} = 1024$ words, the high-order three bits of a virtual address will specify one of the eight pages and the low-order 10 bits give the line address within the page. Note that the line address in address space and memory space is the same; the only mapping required is from a page number to a block number.

The organization of the memory mapping table in a paged system is shown in Fig (14). The memory-page table consists of eight words, one for each page. The address in the page table denotes the page number and the content of the word gives the block number where that page is stored in main memory. The table shows that pages 1, 2, 5, and 6 are now available in main memory in blocks 3, 0, 1, and 2, respectively. A presence bit in each location indicates whether the page has been transferred from auxiliary memory into main memory. A 0 in the presence bit indicates that this page is not available in main memory. The CPU references a word in memory with a virtual address of 13 bits. The three high-order bits of the virtual address specify a page number and also an address for the memory-page table. The content of the word in the memory page table at the page number address is read out into the memory table buffer register. If the presence bit is a 1, the block number thus read is transferred to the two high-order bits of the main memory address register.

The line number from the virtual address is transferred into the 10 low-order bits of the memory address register. A read signal to main memory transfers the content of the word to the main memory buffer register ready to be used by the CPU. If the presence bit in the word read from the page table is 0, it signifies that the content of the word referenced by the virtual address does not reside in main memory.

| Address space |
| --- |
| Page 0 |
| Page 1 |
| Page 2 |
| Page 3 |
| Page 4 |
| Page 5 |
| Page 6 |
| Page 7 |

Address space
$N = 8K = 2^{13}$

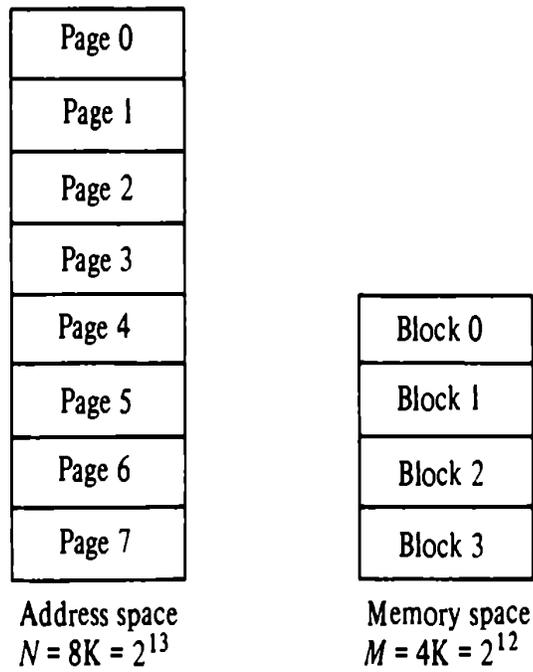| Memory space |
| --- |
| Block 0 |
| Block 1 |
| Block 2 |
| Block 3 |

Memory space
$M = 4K = 2^{12}$

**Fig (13)**
**Address space and memory space split into groups of IK words.**

A call to the operating system is then generated to fetch the required page from auxiliary memory and place it into main memory before resuming computation.
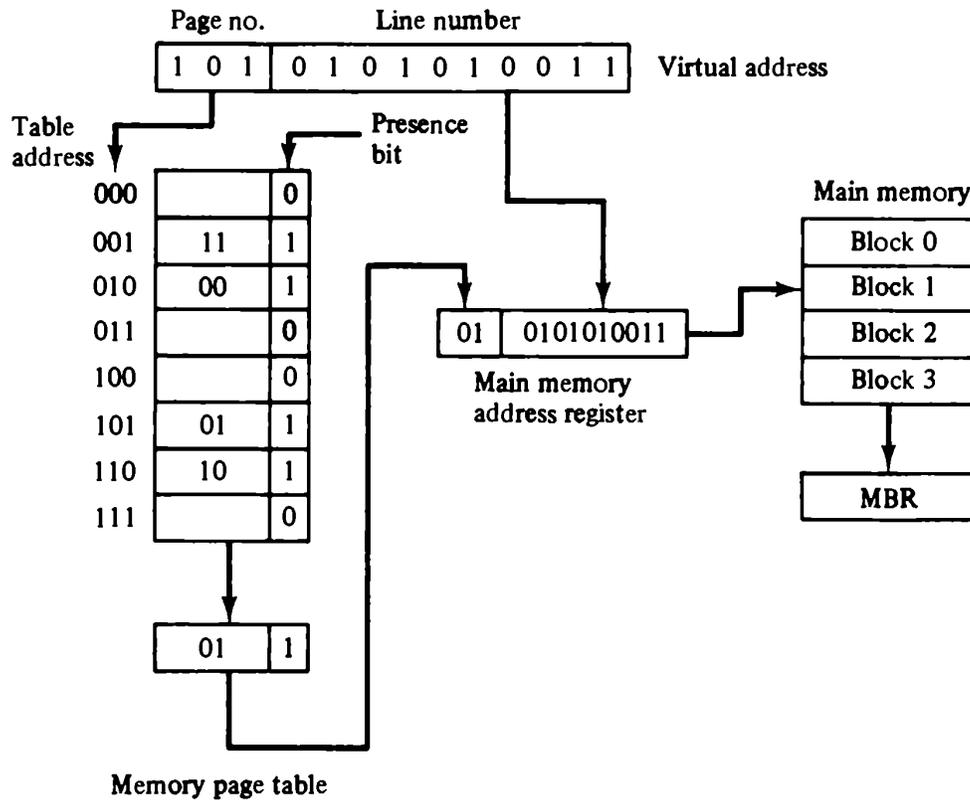
**Fig (14)**
**Memory table in a paged system.**

## Lecture10

## Associative Memory Page Table

A random-access memory page table is inefficient with respect to storage utilization. In the example of Fig. (14) We observe that eight words of memory are needed; one for each page, but at least four words will always be marked empty because main memory cannot accommodate more than four blocks. In general, a system with n pages and m blocks would require a memory-page table of n locations of which up to m blocks will be marked with block numbers and all others will be empty. As a second numerical example, consider an address space of 1024K words and memory space of 32K words. If each page or block contains IK words, the number of pages is 1024 and the number of blocks 32. The capacity of the memory-page table must be 1024 words and only 32 locations may have a presence bit equal to 1. At any given time, at least 992 locations will be empty and not in use.

A more efficient way to organize the page table would be to construct it with a number of words equal to the number of blocks in main memory. In this way the size of the memory is reduced, and each location is fully utilized. This method can be implemented by means of an associative

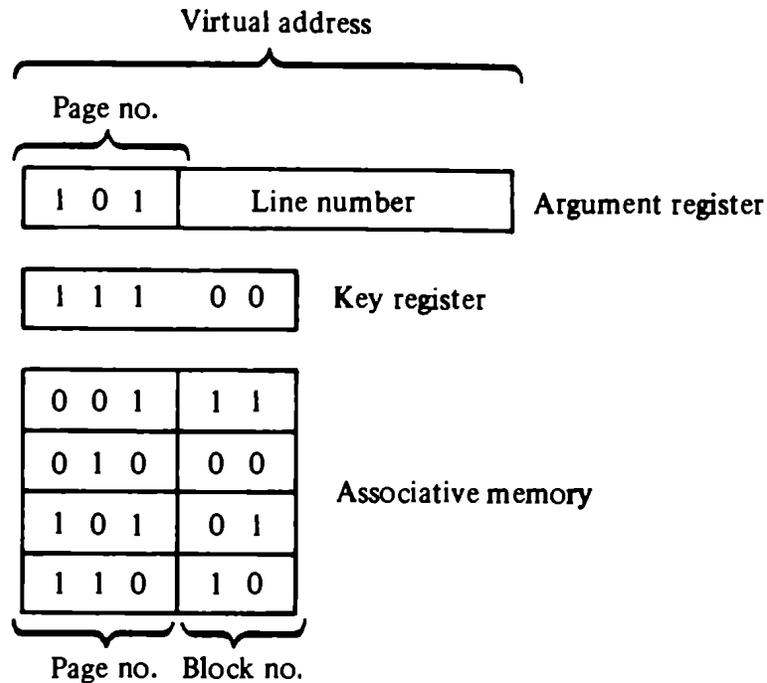memory with each word in memory containing a page number together with its corresponding block number.



**Fig (15)**
**An associative memory page table.**

The page field in each word is compared with the page number in the virtual address. If a match occurs, the word is read from memory and its corresponding block number is extracted.

Consider again the case of eight pages and four blocks as in the example of Fig. (14). We replace the random access memory-page table with an associative memory of four words as shown in Fig(15). Each entry in the associative memory array consists of two fields. The first three bits specify a field for storing the page number. The last two bits constitute a field for storing the block number. The virtual address is placed in the argument register. The page number bits in the argument register are compared with all page numbers in the page field of the associative memory. If the page number is found, the 5-bit word is read out from memory. The corresponding block number, being in the same word, is transferred to the main memory address register. If no match occurs, a call to the operating system is generated to bring the required page from auxiliary memory.

## Lecture11

## Page Replacement

A virtual memory system is a combination of hardware and software techniques. The memory management software system handles all the software operations for the efficient utilization of memory space. It must decide (1) which page in main memory ought to be removed to make room for a new page, (2) when a new page is to be transferred from auxiliary memory to main memory, and (3) where the page is to be placed in main memory. The hardware mapping mechanism and the memory management software together constitute the architecture of a virtual memory.

When a program starts execution, one or more pages are transferred into main memory and the page table is set to indicate their position. The program is executed from main memory until it attempts to reference a page that is still in auxiliary memory. This condition is called page fault. When page fault occurs, the execution of the present program is suspended until the required page is brought into main memory. Since loading a page from auxiliary memory to main memory is basically an I/O operation, the operating system assigns this task to the I/O processor. In the meantime, control is transferred to the next program in memory that is waiting to be processed in the CPU. Later, when the memory block has been assigned and the transfer completed, the original program can resume its operation. When a page fault occurs in a virtual memory system, it signifies that the page referenced by the CPU is not in main memory. A new page is then transferred from auxiliary memory to main memory. If main memory is full, it would be necessary to remove a page from a memory block to make room for the new page. The policy for choosing pages to remove is determined from the replacement algorithm that is used. The goal of a replacement policy is to try to remove the page least likely to be referenced in the immediate future. Two of the most common replacement algorithms used are the first-in, first-out (FIFO) and the least recently used (LRU). The FIFO algorithm selects for replacement the page that has been in memory the longest time. Each time a page is loaded into memory, its identification number is pushed into a FIFO stack. FIFO will be full whenever memory has no more empty blocks. When a new page must be loaded, the page least recently brought in is removed. The page to be removed is easily determined because its identification number is at the top of the FIFO stack. The FIFO replacement policy has the advantage of being easy to implement. It has the disadvantage that under certain circumstances pages are removed and loaded from memory too frequently. The LRU policy is more difficult to implement but has been more attractive on the assumption that the least recently used page is a better candidate for removal than the least recently loaded page as in FIFO. The LRU algorithm can be implemented by associating a counter with every page that is in main memory. When a page is referenced, its

associated counter is set to zero. At fixed intervals of time, the counters associated with all pages presently in memory are incremented by 1. The least recently used page is the page with the highest count. The counters are often called ==aging registers,== as their count indicates their age, that is, how long ago their associated pages have been referenced.

## Memory Management Hardware

In a multiprogramming environment where many programs reside in memory it becomes necessary to move programs and data around the memory, to vary the amount of memory in use by a given program, and to prevent a program from changing other programs. The demands on computer memory brought about by multiprogramming have created the need for a memory management system. ==A memory management system is== a collection of hardware and software procedures for managing the various programs residing in memory. The memory management software is part of an overall operating system available in many computers. Here we are concerned with the hardware unit associated with the memory management system.

The basic ==components of a memory management== unit are:

1. A facilityألية for dynamic storage relocation that maps logical memory references into physical memory addresses
2. A provisionبند for sharing common programs stored in memory by different users
3. Protection of information against unauthorized access between users and preventing users from changing operating system functions

The fixed page size used in the virtual memory system causes certain difficulties with respect to program size and the logical structure of programs. It is more convenient to divide programs and segment data into logical parts called segments. A ==segment is a set== of logically related instructions or data elements associated with a given name. Segments may be generated by the programmer or by the operating system. Examples of segments are a subroutine, an array of data, a table of symbols, or a user's program.

The sharing of common programs is an integral part of a multiprogramming system. Other system programs residing in memory are also shared by all users in a multiprogramming system without having to produce multiple copies. The third issue in multiprogramming is protecting one program from unwanted interaction with another. An example of unwanted interaction is one user's unauthorized copying of another user's program. Another aspect of protection is concerned with preventing the occasionalعرضي  user from performing operating system functions and thereby interrupting the orderly sequence of operations in a computer installation. The secrecy of certain programs must be kept from unauthorized personnel to prevent abuses in the confidential activities of an organization.
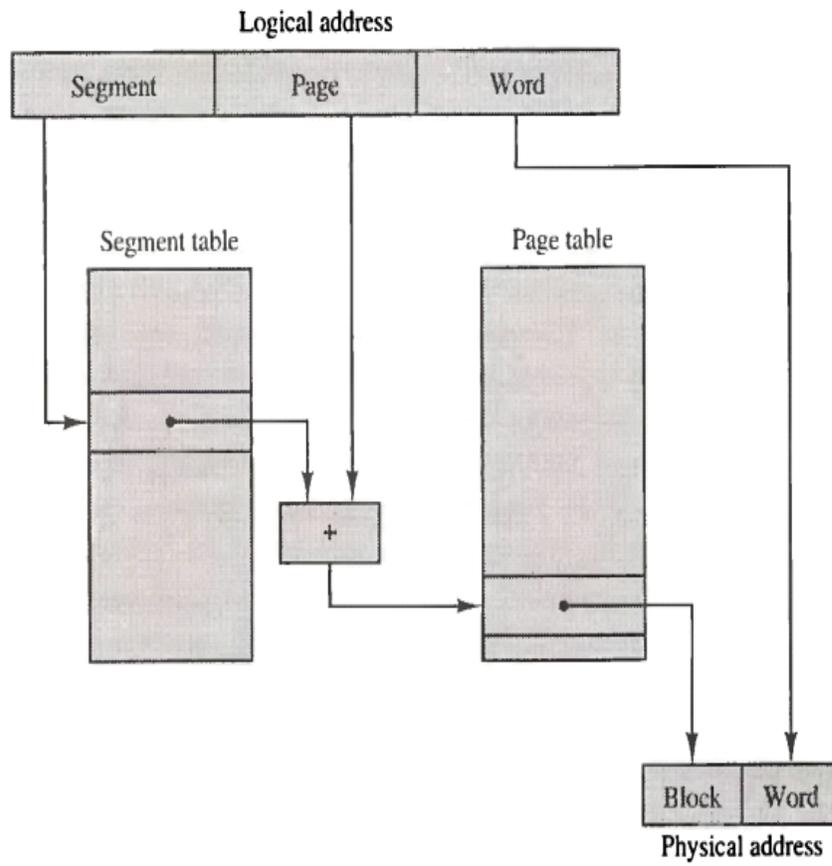
The address generated by a segmented program is called a **logical address**. This is similar to a virtual address except that logical address space is associated with variable-length segments rather than fixed-length pages. The logical address may be larger than the physical memory address as in virtual memory, but it may also be equal, and sometimes even smaller than the length of the physical memory address. In addition to relocation information, each segment has protection information associated with it. Shared programs are placed in a unique segment in each user's logical address space so that a single physical copy can be shared. The function of the memory management unit is to map logical addresses into physical addresses similar to the virtual memory mapping concept.
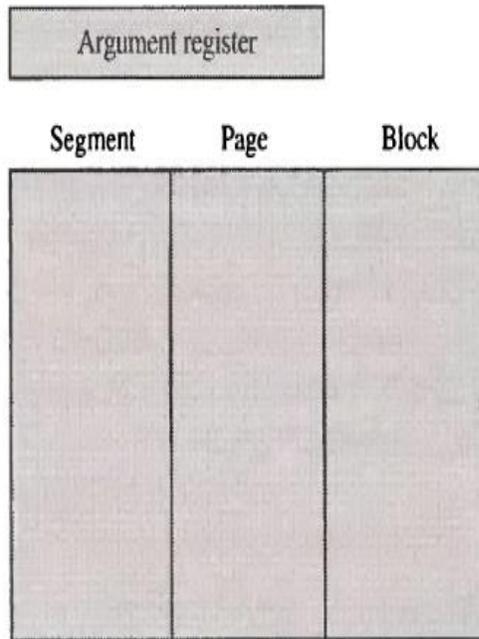
## Lecture12
## Segmented 'Page Mapping

It was already mentioned that the property of logical space is that it uses variable-length segments. The length of each segment is allowed to grow and contract according to the needs of the program being executed. One way of specifying the length of a segment is by associating with it a number of equal-size pages. To see how this is done, consider the logical address shown in Fig. (16). the logical address is partitioned into three fields. **The segment field** specifies a segment number. **The page field** specifies the page within the segment and **the word field** gives the specific word within the page. A page field of k bits can specify up to $2^k$ pages. A segment number may be associated with just one page or with as many as $2^k$ pages. Thus, the length of a segment would vary according to the number of pages that are assigned to it. The mapping of the logical address into a physical address is done by means of two tables, as shown in Fig. (16-a). The segment number of the logical address specifies the address for the segment table. The entry in the segment table is a pointer address for a page table base. The page table base is added to the page number given in the logical address. The sum produces a pointer address to an entry in the page table. The value found in the page table provides the block number in physical memory. The concatenation of the block field with the word field produces the final physical mapped address. The two mapping tables may be stored in two separate small memories or in main memory. In either case, a memory reference from the CPU will require three accesses to memory: one from the segment table, one from the page table, and the third from main memory. This would slow the system significantly when compared to a conventional system that requires only one reference to memory. To avoid this speed penalty, a fast associative memory is used to hold the most recently referenced table entries. (This type of memory is sometimes called a translation look aside buffer, abbreviated TLB.) The first time a given block is referenced, its value together with the corresponding segment and page numbers are entered into the associative memory as shown in Fig. (16-b). Thus, the mapping process is first attempted by associative search with the given segment and page numbers. If it succeeds, the mapping delay is only that of the associative memory. If no match occurs, the slower table

mapping of Fig. (16-a). is used and the result transformed into the associative memory for future reference.



Logical address

| Segment | Page | Word |

(a) Logical to physical address mapping

(b) Associative memory translation look-aside buffer (TLB)

**Fig (16)**
**Mapping in segmented-page memory management unit.**

## Lecture13

### Numerical Example

A numerical example may clarify the operation of the memory management unit. Consider the 20-bit logical address specified in Fig. (17-a). The 4-bit segment number specifies one of 16 possible segments. The 8-bit page number can specify up to 256 pages, and the 8-bit word field implies a page size of 256 words. This configuration allows each segment to have any number of pages up to 256. The smallest possible segment will have one page or 256 words. The largest possible segment will have 256 pages, for a total of 256 x 256 = 64K words.

The physical memory shown in Fig. (17-b). consists of 220 words of 32 bits each. The 20-bit address is divided into two fields: a 12-bit block number and an 8-bit word number. Thus, physical memory is divided into 4096 blocks of 256 words each. A page in a logical address has a corresponding block in physical memory. Note that both the logical and physical address has 20 bits.

In the absence of a memory management unit, the 20-bit address from the CPU can be used to access physical memory directly. Consider a program loaded into memory that requires five pages. The operating system may assign to this program segment 6 and pages 0 through 4, as shown in Fig. (18-a). The total logical address range for the program is from hexadecimal 60000

to 604FF. When the program is loaded into physical memory, it is distributed among five blocks in physical memory where the operating system finds empty spaces. The correspondence between each memory block and logical page number is then entered in a table as shown in Fig (18-b). The information from this table is entered in the segment and page tables as shown in Fig (18-a).

Now consider the specific logical address given in Fig. (18). The 20-bit address is listed as a five-digit hexadecimal number. It refers to word number 7E of page 2 in segment 6. The base of segment 6 in the page table is at address 35. Segment 6 has associated with it five pages, as shown in the page table at addresses 35 through 39. Page 2 of segment 6 is at address 35 + 2 = 37. The physical memory block is found in the page table to be 019. Word 7E in block 19 gives the 20-bit physical address 0197E. Note that page 0 of segment 6 maps into block 12 and page 1 maps into block 0. The associative memory in Fig (18-b). shows that pages 2 and 4 of segment 6 have been referenced previously and therefore their corresponding block numbers are stored in the associative memory.

From this example it should be evident that the memory management system can assign any number of pages to each segment. Each logical page can be mapped into any block in physical memory. Pages can move to different blocks in memory depending on memory space requirements. The only updating required is the change of the block number in the page table. Segments can grow or shrink independently without affecting each other. Different segments can use the same block of memory if it is required to share a program by many users. For example, block number 12 in physical memory can be assigned a second logical address F0000 through F00FF. This specifies segment number 15 and page 0, which maps to block 12 as shown in     Fig (18-a).

**Figure 12-23**   Example of logical and physical memory address assignment.

Hexadecimal
address            Page number

| Hex address | Page number |
|---|---|
| 60000 | Page 0 |
| 60100 | Page 1 |
| 60200 | Page 2 |
| 60300 | Page 3 |
| 60400 / 604FF | Page 4 |

| Segment | Page | Block |
|---|---|---|
| 6 | 00 | 012 |
| 6 | 01 | 000 |
| 6 | 02 | 019 |
| 6 | 03 | 053 |
| 6 | 04 | A61 |

(a)   Logical address assignment

(b)   Segment-page versus
       memory block assignment

Fig (17) Example of logical and physical memory address assignment

**Logical address (in haxadecimal)**

| 6 | 02 | 7E |
|---|----|----|

| Segment table | | Page table | | Physical memory | |
|---|---|---|---|---|---|

Segment table:

| 0 | |
|---|---|
| 6 | 35 |
| | |
| F | A3 |

Page table:

| 00 | |
|----|-----|
| 35 | 012 |
| 36 | 000 |
| 37 | 019 |
| 38 | 053 |
| 39 | A61 |
| A3 | 012 |

Physical memory:

| 00000 | Block 0 |
|-------|---------|
| 000FF | |
| 01200 | Block 12 |
| 012FF | |
| 01900 | 32-bit word |
| 0197E | |
| 019FF | |

**(a) Segment and page table mapping**

| Segment | Page | Block |
|---------|------|-------|
| 6 | 02 | 019 |
| 6 | 04 | A61 |
| | | |
| | | |

**(b) Associative memory (TLB)**

**Figure 12-24** Logical to physical memory mapping example (all numbers are in hexadecimal).

# Lecture14

## Memory Protection

Memory protection can be assigned to the physical address or the logical address. The protection of memory through the physical address can be done by assigning to each block in memory a number of protection bits that indicate the type of access allowed to its corresponding block. Every time a page is moved from one block to another it would be necessary to update the block protection bits. A much better place to apply protection is in the logical address space rather than the physical address space. This can be done by including protection information within the segment table or segment register of the memory

management hardware. The content of each entry in the segment table or a segment register is called a ==descriptor==. A typical descriptor would contain, in addition to a base address field, one or two additional fields for protection purposes. A typical format for a segment descriptor is shown in Fig. (19). The base address field gives the base of the page table address ==in a segmented-page== organization or the block base address ==in a segment register== organization. This is the address used in mapping from a logical to the physical address. The ==length field== gives the segment size by specifying the maximum number of pages assigned to the segment. The length field is compared against the page number in the logical address. A size violation occurs if the page number falls outside the segment length boundary. ==Thus, a given program and its data cannot access memory not assigned to it by the operating system.==

The protection held in a segment descriptor specifies the access rights available to a particular segment. In a segmented-page organization, each entry in the page table may have its own protection field to describe the access rights of each page. The protection information is set into the descriptor by the master control program of the operating system. Some of the access rights of interest that are used for protecting the programs residing in memory are: مهم

1. Full read and writes privileges
2. Read only (write protection)
3. Execute only (program protection)
4. System only (operating system protection)

Full read and write privileges are given to a program when it is executing its own instructions. Write protection is useful for sharing system programs such as utility programs and other library routines. These system programs are stored in an area of memory where they can be shared by many users. They can be read by all programs, but no writing is allowed. This protects them from being changed by other programs. The execute-only condition protects programs from being copied. It restricts the segment to be referenced only during the instruction fetch phase but

not during the execute phase. Thus, it allows the users to execute the segment program instructions but prevents them from reading the instructions as data for the purpose of copying their content. Portions of the operating system will reside in memory at any given time. These system programs must be protected by making them inaccessible to unauthorized users. The operating system protection condition is placed in the descriptors of all operating system programs to prevent the occasional user from accessing operating system segments.

| Base address | Length | Protection |
|---|---|---|

**Fig (19) Format of a typical segment descriptor**

## Input-Output Interface and Peripherals

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

1. Peripherals are electromechanical and electromagnetic devices, and their manner of operation is different from the operation of the CPU and memory, which are electronic devices. Therefore, a conversion of signal values may be required.
2. The data transfer rate of peripherals is usually slower than the transfer rate of the CPU, and consequently, a synchronization mechanism may be needed.
3. Data codes and formats in peripherals differ from the word format in the CPU and memory.
4. The operating modes of peripherals are different from each other, and each must be controlled so as not to disturb the operation of other peripherals connected to the CPU.

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

## I/O Bus and Interface Modules

A typical communication link between the processor and several peripherals s shown in Fig. (20). The I/O bus consists of data lines, address lines, and control lines. The magnetic disk, printer, and terminal are employed in practically any general-purpose computer. Magnetic tape is used in some computers for backup storage. It also synchronizes the data flow and supervises the transfer between peripheral and processor. Each peripheral has its own controller that operates the particular electromechanical device. For example, the printer controller controls the paper motion, the print timing, and the selection of printing characters. A controller may be housed separately or may be physically integrated with the peripheral.
The I/O bus from the processor is attached to all peripheral interfaces. To communicate with a particular device, the processor places a device address on the address lines. Each interface attached to the I/O bus contains an address decoder that monitors the address lines. When the interface detects its own address, it activates the path between the bus lines and the device

that it controls. All peripherals whose address does not correspond to the address in the bus are disabled by their interface.

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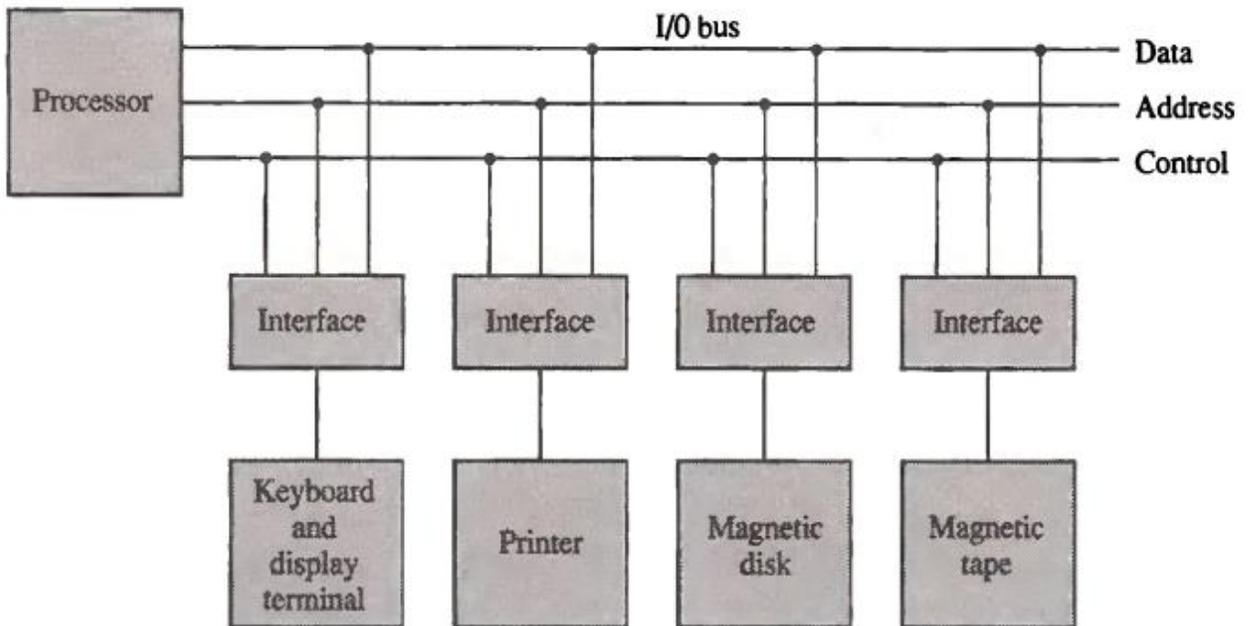


**Fig (20)**
**Connection of I/O bus to input-output devices.**

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

A control command is issued to activate the peripheral and to inform it what to do. For example, a magnetic tape unit may be instructed to backspace the tape by one record, to rewind the tape, or to start the tape moving in the forward direction. The particular control command issued depends on the peripheral, and each peripheral receives its own distinguished sequence of control commands, depending on its mode of operation. A status command is used to test various status conditions in the interface and the peripheral. For example, the computer may wish to check the status of the peripheral before a transfer is initiated. During the transfer,

one or more errors may occur which are detected by the interface. These errors are designated by setting bits in a status register that the processor can read at certain intervals.

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

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

## Lecture15
## Modes of Transfer

Binary information received from an external device is usually stored in memory for later processing. Information transferred from the central computer into an external device originates in the memory unit. The CPU merely executes the I/O instructions and may accept the data temporarily, but the ultimate source or destination is the memory unit. Data transfer between the central computer and I/O devices may be handled in a variety of modes. Some modes use the CPU as an intermediate path; others transfer the data directly to and from the memory unit. Data transfer to and from peripherals may be handled in one of three possible modes:

1. Programmed I/O
2. Interrupt-initiated I/O
3. Direct memory access (DMA)

Programmed I/O operations are the result of I/O instructions written in the computer program. Each data item transfer is initiated by an instruction in the program. Usually, the transfer is to and from a CPU register and peripheral. Other instructions are needed to transfer the data to and from CPU and memory. Transferring data under program control requires constant monitoring of the peripheral by the CPU. Once a data transfer is initiated, the CPU is required to monitor the interface to see when a transfer can again be made.

It is up to the programmed instructions executed in the CPU to keep close tabs on everything that is taking place in the interface unit and the I/O device. In the programmed I/O method, the CPU stays in a program loop until the I/O unit indicates that it is ready for data transfer. This is a

time-consuming process since it keeps the processor busy needlessly. It can be avoided by using an interrupt facility and special commands to inform the interface to issue an interrupt request signal when the data are available from the device. In the meantime, the CPU can proceed to execute another program. The interface meanwhile keeps monitoring the device. When the interface determines that the device is ready for data transfer, it generates an interrupt request to the computer. Upon detecting the external interrupt signal, the CPU momentarily stops the task it is processing, branches to a service program to process the I/O transfer, and then returns to the task it was originally performing. Transfer of data under programmed I/O is between CPU and peripheral. In direct memory access (DMA), the interface transfers data into and out of the memory unit through the memory bus. The CPU initiates the transfer by supplying the interface with the starting address and the number of words needed to be transferred and then proceeds to execute other tasks. When the transfer is made, the DMA requests memory cycles through the memory bus. When the request is granted by the memory controller, the DMA transfers the data directly into memory. The CPU merely delays its memory access operation to allow the direct memory I/O transfer. Since peripheral speed is usually slower than processor speed, I/O-memory transfers are infrequent compared to processor access to memory. Many computers combine the interface logic with the requirements for direct memory access into one unit and call it an I/O processor (IOP). The IOP can handle many peripherals through a DMA and interrupt facility. In such a system, the computer is divided into three separate modules: the memory unit, the CPU, and the IOP.

## Example of Programmed I/O

In the programmed I/O method, the I/O device does not have direct access to memory. A transfer from an I/O device to memory requires the execution of several instructions by the CPU, including an input instruction to transfer the data from the device to the CPU and a store instruction to transfer the data from the CPU to memory. Other instructions may be needed to verify that the data are available from the device and to count the numbers of words transferred.

An example of data transfer from an I/O device through an interface into the CPU is shown in Fig (21). The device transfers bytes of data one at a time as they are available. When a byte of data is available, the device places it in the I/O bus and enables its data valid line. The interface accepts the byte into its data register and enables the data accepted line.
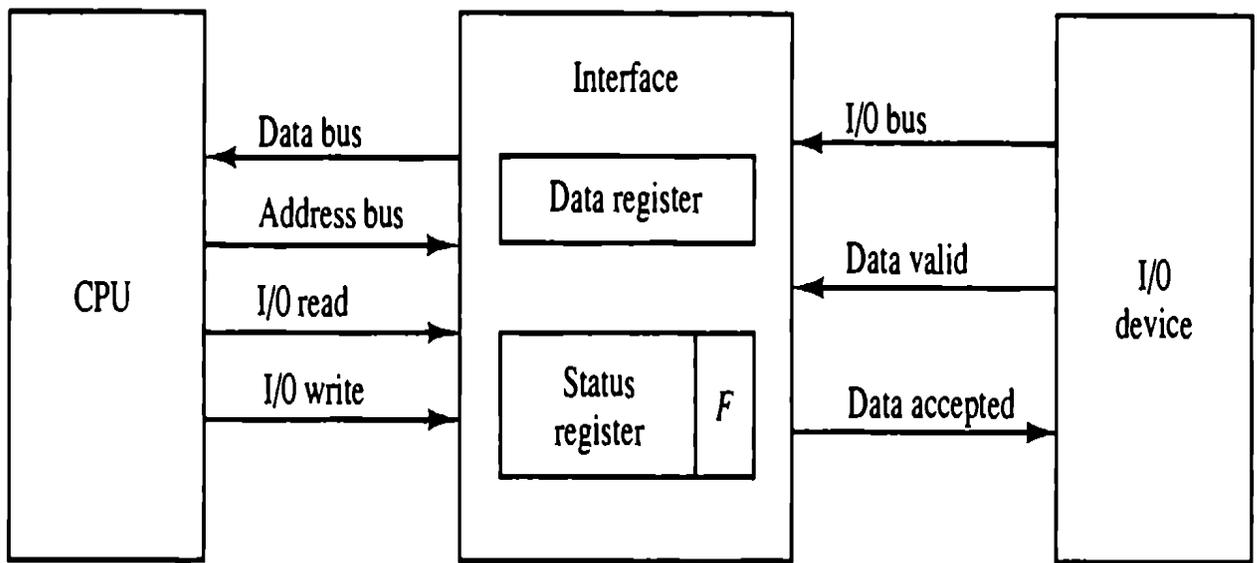
The interface sets a bit in the status register that we will refer to as an F or "flag" bit. The device can now disable the data valid line, but it will not transfer another byte until the data accepted line is disabled by the interface. This is according to the handshaking procedure established in Fig (21).

A program is written for the computer to check the flag in the status register to determine if a byte has been placed in the data register by the I/O device. This is done by reading the status

register into a CPU register and checking the value of the flag bit. If the flag is equal to 1, the CPU reads the data from the data register. The flag bit is then cleared to 0 by either the CPU or the interface, depending on how the interface circuits are designed. Once the flag is cleared, the interface disables the data accepted line and the device can then transfer the next data byte.

A flowchart of the program that must be written for the CPU is shown in Fig. (22). It is assumed that the device is sending a sequence of bytes that must be stored in memory. The transfer of each byte requires three instructions:

l: Read the status register.
2. Check the status of the flag bit and branch to step 1 if not set or to step
3 if set.
3. Read the data register.

Each byte is read into a CPU register and then transferred to memory with a store instruction. A common I/O programming task is to transfer a block of words from an I/O device and store them in a memory buffer.



$F$ = Flag bit

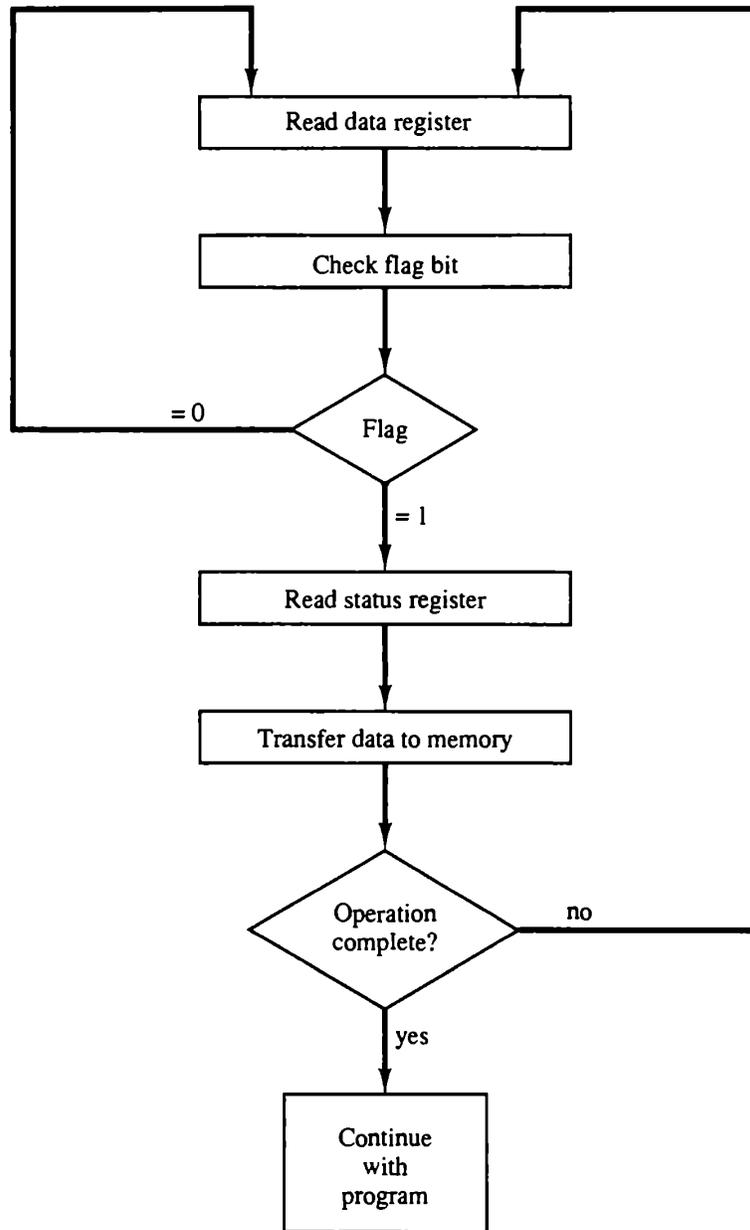**Fig (21)**
**Data transfer from I/O device to CPU.**

**Fig (22)**
**Flowchart for CPU program to input data.**

The programmed I/O method is particularly useful in small low-speed computers or in systems that are dedicated to monitor a device continuously. The difference in information transfer rate between the CPU and the I/O device makes this type of transfer inefficient. To see why this is inefficient, consider a typical computer that can execute the two instructions that read the status register and check the flag in 1. Assume that the input device transfers its data at an average rate of 100 bytes per second. This is equivalent to one byte every 10,000. This means

that the CPU will check the flag 10,000 times between each transfer. The CPU is wasting time while checking the flag instead of doing some other useful processing tasks.

## Lecture16

## Direct Memory Access (DMA)

The transfer of data between a fast storage device such as magnetic disk and memory is often limited by the speed of the CPU. Removing the CPU from the path and letting the peripheral device manage the memory buses directly.

The data formats of peripheral devices differ from memory and CPU data formats. The IOP must structure data words from many different sources. For example, it may be necessary to take four bytes from an input device and pack them into one 32-bit word before the transfer to memory. Data are gathered in the IOP at the device rate and bit capacity while the CPU is executing its own program. After the input data are assembled into a memory word, they are transferred from IOP directly into memory by "stealing" one memory cycle from the CPU. Similarly, an output word transferred from memory to the IOP is directed from the IOP to the output device at the device rate and bit capacity. The communication between the IOP and the devices attached to it is similar to the program control method of transfer. Communication with memory is similar to the direct memory access method. The way by which the CPU and IOP communicate depends on the level of sophistication included in the system. In very-large-scale computers, each processor is independent of all others and any one processor can initiate an operation. In most computer systems, the CPU is the master while the IOP is a slave processor. The CPU is assigned the task of initiating all operations, but I/O instructions are executed in the IOP. CPU instructions provide operations to start an I/O transfer and also to test I/O status conditions needed for making decisions on various I/O activities. The IOP, in turn, typically asks for CPU attention by means of an interrupt. It also responds to CPU requests by placing a status word in a prescribed location in memory to be examined later by a CPU program. When an I/O operation is desired, the CPU informs the IOP where to find the I/O program and then leaves the transfer details to the IOP.

Instructions that are read from memory by an IOP are sometimes called commands, to distinguish them from instructions that are read by the CPU. Otherwise, an instruction and a command have similar functions. Commands are prepared by experienced programmers and are stored in memory. The command words constitute the program for the IOP. The CPU informs the IOP where to find the commands in memory when it is time to execute the I/O program.

## CPU-IOP Communication

The communication between CPU and IOP may take different forms, depending on the particular computer considered. In most cases the memory unit acts as a message center where each processor leaves information for the other. To appreciate the operation of a typical IOP, we will illustrate by a specific example the method by which the CPU and IOP communicate. This is a simplified example that omits many operating details in order to provide an overview of basic concepts. The sequence of operations may be carried out as shown in the flowchart of Fig (23). The CPU sends an instruction to test the IOP path.

CPU operations            IOP operations

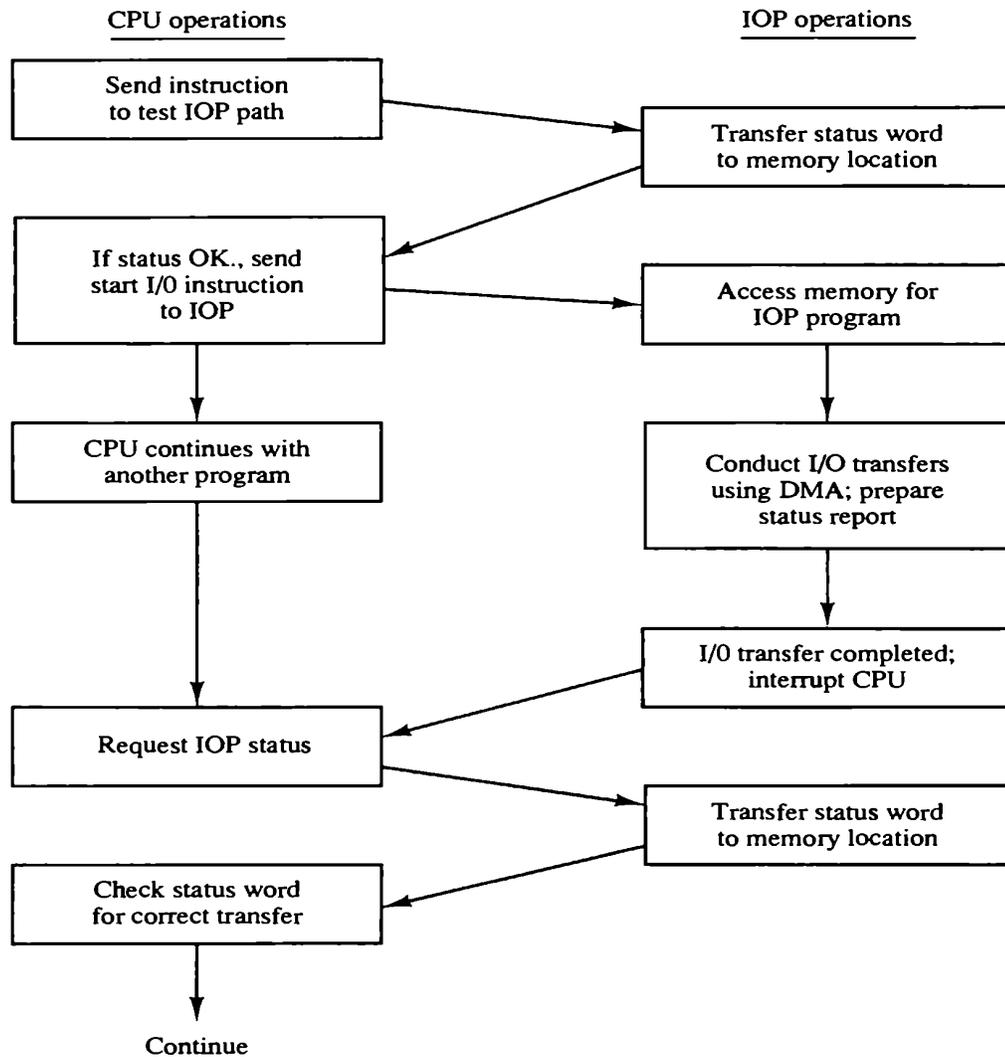| Send instruction to test IOP path |
| Transfer status word to memory location |
| If status OK., send start I/0 instruction to IOP |
| Access memory for IOP program |
| CPU continues with another program |
| Conduct I/O transfers using DMA; prepare status report |
| I/0 transfer completed; interrupt CPU |
| Request IOP status |
| Transfer status word to memory location |
| Check status word for correct transfer |
| Continue |

**Fig (23)**
**CPU-IOP communication.**

The IOP responds by inserting a status word in memory for the CPU to check. The bits of the status word indicate the condition of the IOP and I/O device, such as IOP overload condition,

device busy with another transfer, or device ready for I/O transfer. The CPU refers to the status word in memory to decide what to do next. If all is in order, the CPU sends the instruction to start I/O transfer. The memory address received with this instruction tells the IOP where to find its program. The CPU can now continue with another program while the IOP is busy with the I/O program. Both programs refer to memory by means of DMA transfer. When the IOP terminates the execution of its program, it sends an interrupt request to the CPU. The CPU responds to the interrupt by issuing an instruction to read the status from the IOP. The IOP responds by placing the contents of its status report into a specified memory location. The status word indicates whether the transfer has been completed or if any errors occurred during the transfer. From inspection of the bits in the status word, the CPU determines if the I/O operation was completed satisfactorily without errors.

The IOP takes care of all data transfers between several I/O units and the memory while the CPU is processing another program. The IOP and CPU are competing for the use of memory, so the number of devices that can be in operation is limited by the access time of the memory. It is not possible to saturate the memory by I/O devices in most systems, as the speed of most devices is much slower than the CPU. However, some very fast units, such as magnetic disks, can use an appreciable number of the available memory cycles. In that case, the speed of the CPU may deteriorate because it will often have to wait for the IOP to conduct memory transfers.

## CPU BASICS

A typical CPU has three major components: (1) register set, (2) arithmetic logic unit (ALU), and (3) control unit (CU). The register set differs from one computer architecture to another. It is usually a combination of general-purpose and special purpose registers. General-purpose registers are used for any purpose, hence the name general purpose. Special-purpose registers have specific functions within the CPU. For example, the program counter (PC) is a special-purpose register that is used to hold the address of the instruction to be executed next. Another example of special-purpose registers is the instruction register (IR), which is used to hold the instruction that is currently executed. The ALU provides the circuitry needed to perform the arithmetic, logical and shift operations demanded of the instruction set. The control unit is the entity responsible for fetching the instruction to be executed from the main memory and decoding and then executing it. Figure (24) shows the main components of the CPU and its interactions with the memory system and the input/output devices.

The CPU fetches instructions from memory, reads and writes data from and to memory, and transfers data from and to input/output devices.
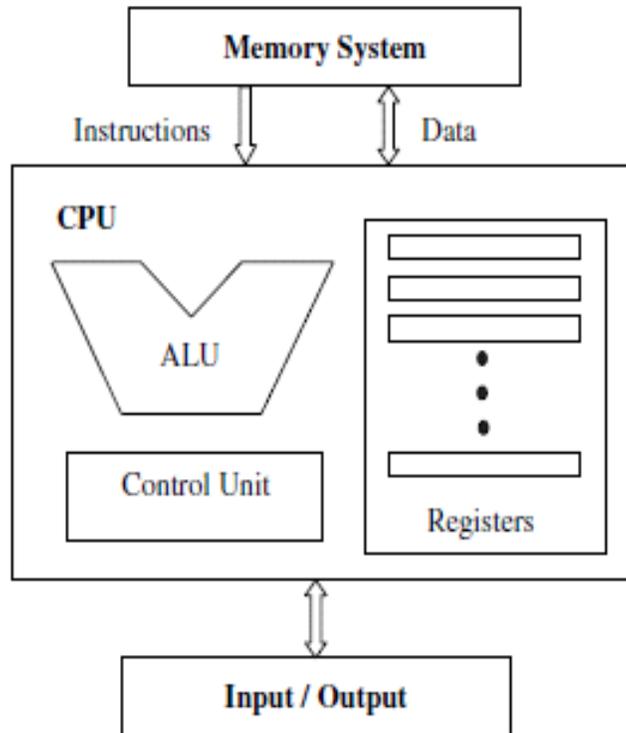
Fig (24)
**Central processing unit main components and interactions with the memory and I/O**

A typical and simple execution cycle can be summarized as follows:

1. The next instruction to be executed, whose address is obtained from the PC, is fetched from the memory and stored in the IR.
2. The instruction is decoded.
3. Operands are fetched from the memory and stored in CPU registers, if needed.
4. The instruction is executed.
5. Results are transferred from CPU registers to the memory, if needed.

The execution cycle is repeated as long as there are more instructions to execute. A check for pending interrupts is usually included in the cycle. Examples of interrupts include I/O device request, arithmetic overflow, or a page fault.

When an interrupt request is encountered, a transfer to an interrupt handling routine takes place. Interrupt handlings routines are programs that are invoked to collect the state of the currently executing program, correct the cause of the interrupt, and restore the state of the program.

The actions of the CPU during an execution cycle are defined by micro-orders issued by the control unit. These micro-orders are individual control signals sent over dedicated control lines. For example, let us assume that we want to execute an instruction that moves the contents of

register X to register Y. Let us also assume that both registers are connected to the data bus, D. The control unit will issue a control signal to tell register X to place its contents on the data bus D. After some delay, another control signal will be sent to tell register Y to read from data bus D. The activation of the control signals is determined using either hardwired control or microprogramming.

## Lecture17
### REGISTER SET
Registers are essentially extremely fast memory locations within the CPU that are used to create and store the results of CPU operations and other calculations. Different computers have different register sets. They differ in the number of registers, register types, and the length of each register. They also differ in the usage of each register. General-purpose registers can be used for multiple purposes and assigned to a variety of functions by the programmer. Special-purpose registers are restricted to only specific functions. In some cases, some registers are used only to hold data and cannot be used in the calculations of operand addresses. The length of a data register must be long enough to hold values of most data types. Some machines allow two contiguous registers to hold double-length values. Address registers may be dedicated to a particular addressing mode or may be used as address general purpose. Address registers must be long enough to hold the largest address. The number of registers in a particular architecture affects the instruction set design. A very small number of registers may result in an increase in memory references. Another type of registers is used to hold processor status bits, or flags. These bits are set by the CPU as the result of the execution of an operation. The status bits can be tested at a later time as part of another operation.

### Memory Access Registers
Two registers are essential in memory write and read operations: the memory data register (MDR) and memory address register (MAR). The MDR and MAR are used exclusively by the CPU and are not directly accessible to programmers.
In order to perform a write operation into a specified memory location, the MDR and MAR are used as follows:

1. The word to be stored into the memory location is first loaded by the CPU into MDR.
2. The address of the location into which the word is to be stored is loaded by the CPU into a MAR.
3. A write signal is issued by the CPU.

Similarly, to perform a memory read operation, the MDR and MAR are used as follows:

1. The address of the location from which the word is to be read is loaded into the MAR.
2. A read signal is issued by the CPU.
3. The required word will be loaded by the memory into the MDR ready for use by the CPU.

### Instruction Fetching Registers

Two main registers are involved in fetching an instruction for execution: the program counter (PC) and the instruction register (IR). The PC is the register that contains the address of the next instruction to be fetched. The fetched instruction is loaded in the IR for execution. After a successful instruction fetch, the PC is updated to point to the next instruction to be executed. In the case of a branch operation, the PC is updated to point to the branch target instruction after the branch is resolved, that is, the target address is known.

## Condition Registers

Condition registers, or flags, are used to maintain status information. Some architecture contains a special program status word (PSW) register. The PSW contains bits that are set by the CPU to indicate the current status of an executing program. These indicators are typically for arithmetic operations, interrupts, memory protection information, or processor status.

## Special-Purpose Address Registers

**Index Register** in index addressing, the address of the operand is obtained by adding a constant to the content of a register, called the index register. The index register holds an address displacement. Index addressing is indicated in the instruction by including the name of the index register in parentheses and using the symbol X to indicate the constant to be added.

**Segment Pointers** in order to support segmentation, the address issued by the processor should consist of a segment number (base) and a displacement (or an offset) within the segment. A segment register holds the address of the base of the segment.

**Stack Pointer**, a stack is a data organization mechanism in which the last data item stored is the first data item retrieved. Two specific operations can be performed on a stack. These are the Push and the Pop operations. A specific register, called the stack pointer (SP), is used to indicate the stack location that can be addressed. In the stack push operation, the SP value is used to indicate the location (called the top of the stack). After storing (pushing) this value, the SP is incremented

## DATAPATH

The CPU can be divided into a data section and a control section. The data section, which is also called the datapath, contains the registers and the ALU. The datapath is capable of performing certain operations on data items. The control section is basically the control unit, which issues control signals to the datapath. Internal to the CPU, data move from one register to another and between ALU and registers. Internal data movements are performed via local buses, which may carry data, instructions, and addresses. Externally, data move from registers to memory and I/O devices, often by means of a system bus. Internal data movement among registers and between the ALU and registers may be carried out using different organizations including one-bus, two-bus, or three-bus organizations. Dedicated datapaths may also be used between components that transfer data between themselves more frequently. For example, the contents of the PC are transferred to the MAR to fetch a new instruction at the beginning of each instruction cycle. Hence, a dedicated datapath from the PC to the MAR could be useful in speeding up this part of instruction execution.

**One-Bus Organization**

Using one bus, the CPU registers and the ALU use a single bus to move outgoing and incoming data. Since a bus can handle only a single data movement within one clock cycle, two-operand operations will need two cycles to fetch the operands for the ALU. Additional registers may also be needed to buffer data for the ALU. This bus organization is the simplest and least expensive, but it limits the amount of data transfer that can be done in the same clock cycle, which will slowdown the overall performance. Figure (25) shows a one-bus datapath consisting of a set of general-purpose registers, a memory address register (MAR), a memory data register (MDR), an instruction register (IR), a program counter (PC), and an ALU.
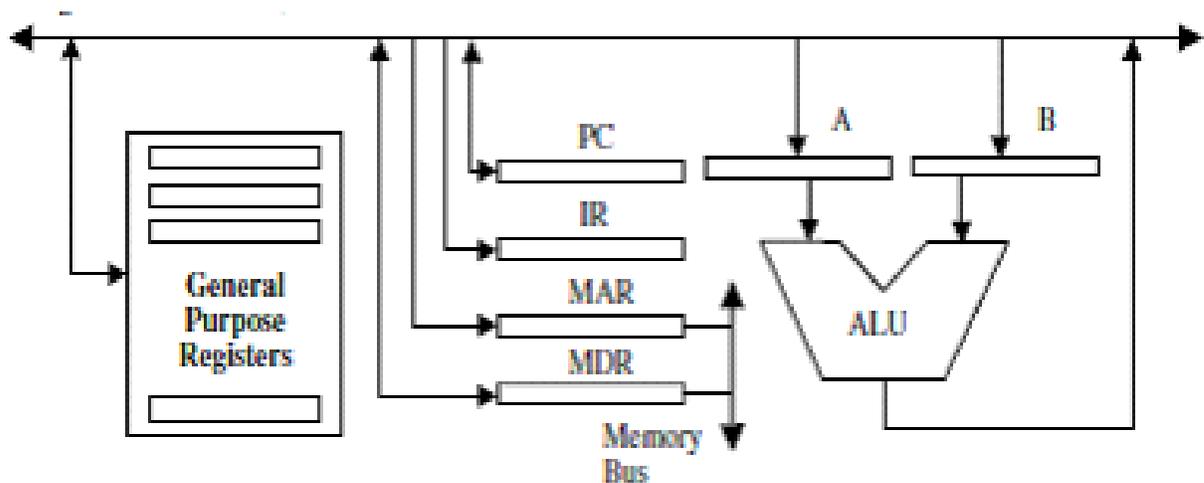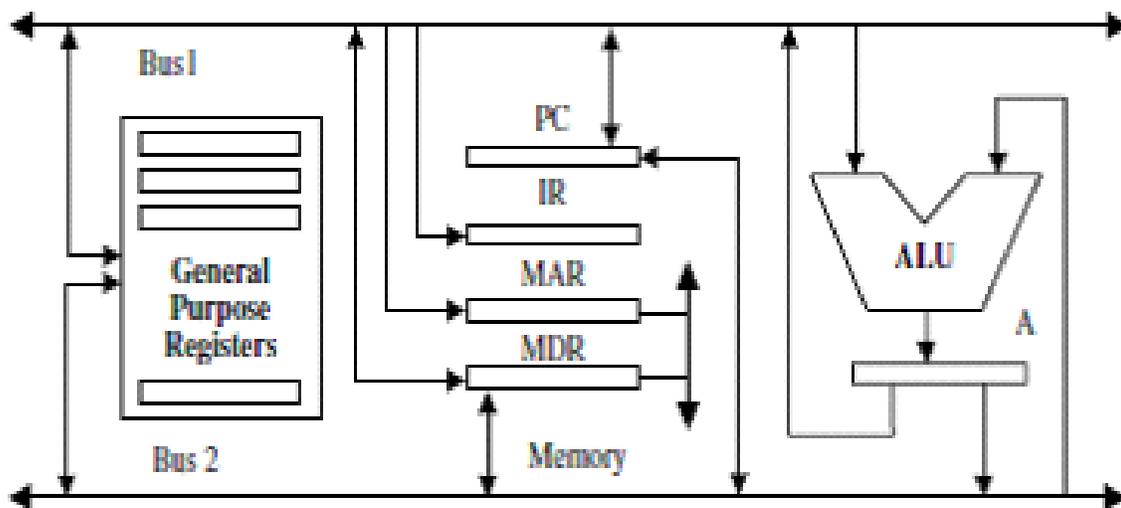


**Fig (25)**
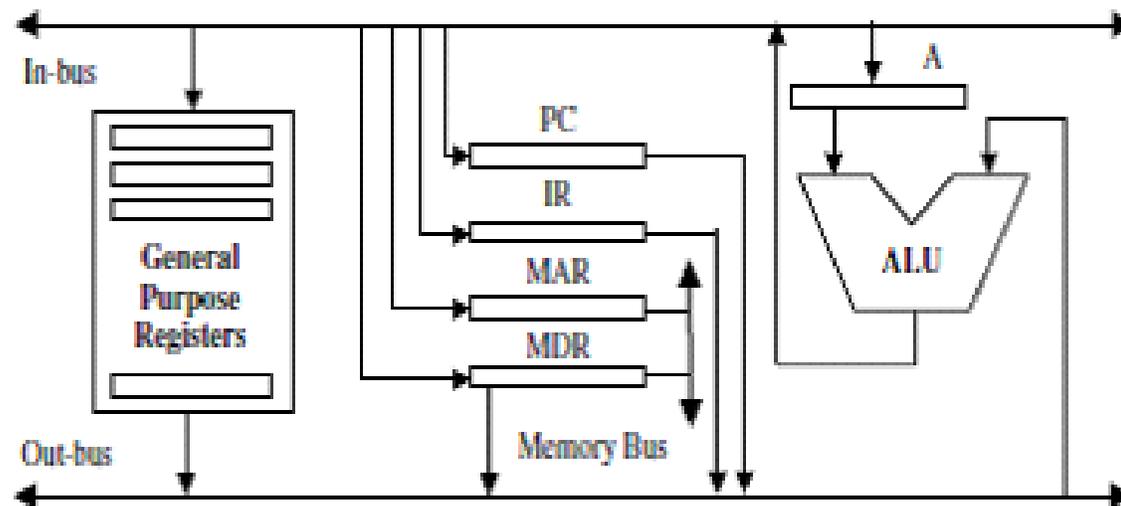**One-bus datapath**

## Lecture18

**Two-Bus Organization**

Using two buses is a faster solution than the one-bus organization. In this case, general-purpose registers are connected to both buses. Data can be transferred from two different registers to the input point of the ALU at the same time. Therefore, a two operand operation can fetch both operands in the same clock cycle. An additional buffer register may be needed to hold the output of the ALU when the two buses are busy carrying the two operands. Figure (26-a) shows a two-bus organization.

In some cases, one of the buses may be dedicated for moving data into registers (in-bus), while the other is dedicated for transferring data out of the registers (out-bus). In this case, the additional buffer register may be used, as one of the ALU inputs, to hold one of the operands. The ALU output can be connected directly to the in-bus, which will transfer the result into one of the registers. Figure (26-b) shows a two-bus organization with in-bus and out-bus.

Fig (26) Two-bus organizations.
(a) An Example of Two-Bus Datapath.
(b) Another Example of Two-Bus Datapath with in-bus and out-bus

**Three-Bus Organization**

In a three-bus organization, two buses may be used as source buses while the third is used as destination. The source buses move data out of registers (out-bus), and the destination bus may move data into a register (in-bus). Each of the two out-buses is connected to an ALU input point. The output of the ALU is connected directly to the in-bus. As can be expected, the more buses we have, the more data we can move within a single clock cycle. However, increasing the

number of buses will also increase the complexity of the hardware. Figure (27) shows an example of a three-bus datapath.



**Fig (27)**
**Three-bus datapath**

## Lecture19

**CPU INSTRUCTION CYCLE**

The sequence of operations performed by the CPU during its execution of instructions is presented in Fig. (28). As long as there are instructions to execute, the next instruction is fetched from main memory. The instruction is executed based on the operation specified in the opcode field of the instruction. At the completion of the instruction execution, a test is made to determine whether an interrupt has occurred. An interrupt handling routine needs to be invoked in case of an interrupt.

**Fig (28)**
**CPU functions**

The basic actions during fetching an instruction, executing an instruction, or handling an interrupt are defined by a sequence of micro-operations. A group of control signals must be enabled in a prescribed sequence to trigger the execution of a micro-operation.

In this section, we show the micro-operations that implement instruction fetch, execution of simple arithmetic instructions, and interrupt handling.

**Fetch Instructions**

The sequence of events in fetching an instruction can be summarized as follows:

1. The contents of the PC are loaded into the MAR.
2. The value in the PC is incremented. (This operation can be done in parallel with a memory access.)
3. As a result of a memory read operation, the instruction is loaded into the MDR.
4. The contents of the MDR are loaded into the IR.

Let us consider the one-bus datapath organization shown in Fig(26). We will see that the fetch operation can be accomplished in three steps as shown in the table below, where t0 , t1 , t2 . Note that multiple operations separated by ";" imply that they are accomplished in parallel.

| Step | Micro-operation |
|------|-----------------|
| $t_0$ | MAR ← (PC); A ← (PC) |
| $t_1$ | MDR ← Mem[MAR]; PC ← (A)+4 |
| $t_2$ | IR ← (MDR) |

Using the three-bus datapath shown in Figure 5.5, the following table shows the steps needed.

| Step | Micro-operation |
|------|-----------------|
| $t_0$ | MAR ← (PC); PC ← (PC)+4 |
| $t_1$ | MDR ← Mem[MAR] |
| $t_2$ | IR ← (MDR) |

# Lecture20

**Execute Simple Arithmetic Operation**

Add R1, R2, R0 This instruction adds the contents of source registers R1 and R2, and stores the results in destination register R0. This addition can be executed as follows:

1. The registers R0 , R1 , R2 , are extracted from the IR.
2. The contents of R1 and R2 are passed to the ALU for addition.
3. The output of the ALU is transferred to R0 .

Using the one-bus datapath shown in Figure (26), this addition will take three steps as shown in the following table, where t0 < t1 < t2 .

| Step | Micro-operation |
|------|-----------------|
| $t_0$ | A ← $(R_1)$ |
| $t_1$ | B ← $(R_2)$ |
| $t_2$ | $R_0$ ← (A)+(B) |

Using the two-bus datapath shown in Figure (27-a) this addition will take two steps as shown in the following table, where t0 < t1 .

| Step | Micro-operation |
|------|-----------------|
| $t_0$ | $A \leftarrow (R_1) + (R_2)$ |
| $t_1$ | $R_0 \leftarrow (A)$ |

Using the two-bus datapath with in-bus and out-bus shown in Figure (27-b), this addition will take two steps as shown below, where t0 > t1 .

| Step | Micro-operation |
|------|-----------------|
| $t_0$ | $A \leftarrow (R_1)$ |
| $t_1$ | $R_0 \leftarrow (A) + (R_2)$ |

Using the three-bus datapath shown in Figure (28), this addition will take only one step as shown in the following table.

| Step | Micro-operation |
|------|-----------------|
| $t_0$ | $R_0 \leftarrow (R_1) + (R_2)$ |

Add X, R0 This instruction adds the contents of memory location X to register R0 and stores the result in R0 . This addition can be executed as follows:

1. The memory location X is extracted from IR and loaded into MAR.
2. As a result of memory read operation, the contents of X are loaded into MDR.
3. The contents of MDR are added to the contents of R0 .

Using the one-bus datapath shown in Figure (26), this addition will take five steps as shown below, where t0 < t1 < t2 < t3 < t4 .

| Step | Micro-operation |
|------|-----------------|
| $t_0$ | $MAR \leftarrow X$ |
| $t_1$ | $MDR \leftarrow Mem[MAR]$ |
| $t_2$ | $A \leftarrow (R_0)$ |
| $t_3$ | $B \leftarrow (MDR)$ |
| $t_4$ | $R_0 \leftarrow (A)+(B)$ |

Using the two-bus datapath shown in Figure (27-a), this addition will take four steps as shown below, where t0 < t1 < t2 < t3 .

| Step | Micro-operation |
|------|-----------------|
| $t_0$ | $MAR \leftarrow X$ |
| $t_1$ | $MDR \leftarrow Mem[MAR]$ |
| $t_2$ | $A \leftarrow (R_0)+(MDR)$ |
| $t_3$ | $R_0 \leftarrow (A)$ |

Using the two-bus datapath with in-bus and out-bus shown in Figure (27-b), this addition will take four steps as shown below, where t0 < t1 < t2 < t3 .

| Step | Micro-operation |
|------|-----------------|
| $t_0$ | $MAR \leftarrow X$ |
| $t_1$ | $MDR \leftarrow Mem[MAR]$ |
| $t_2$ | $A \leftarrow (R_0)$ |
| $t_3$ | $R_0 \leftarrow (A)+(MDR)$ |

Using the three-bus datapath shown in Figure (28), this addition will take three steps as shown below, where t0 < t1 < t2 .

| Step | Micro-operation |
|---|---|
| $t_0$ | MAR $\leftarrow$ X |
| $t_1$ | MDR $\leftarrow$ Mem[MAR] |
| $t_2$ | $R_0 \leftarrow R_0 + (MDR)$ |

# Lecture21

**Interrupt Handling**

After the execution of an instruction, a test is performed to check for pending interrupts.

If there is an interrupt request waiting, the following steps take place:

1. The contents of PC are loaded into MDR (to be saved).
2. The MAR is loaded with the address at which the PC contents are to be saved.
3. The PC is loaded with the address of the first instruction of the interrupt handling routine.
4. The contents of MDR (old value of the PC) are stored in memory.

The following table shows the sequence of events, where t1 < t2 < t3 .

| Step | Micro-operation | |
|---|---|---|
| $t_1$ | MDR $\leftarrow$ | (PC) |
| $t_2$ | MAR $\leftarrow$ | address1 (where to save old PC); |
| | PC $\leftarrow$ | address2 (interrupt handling routine) |
| $t_3$ | Mem[MAR] $\leftarrow$ | (MDR) |

**CONTROL UNIT**

The control unit is the main component that directs the system operations by sending control signals to the datapath. These signals control the flow of data within the CPU and between the CPU and external units such as memory and I/O. Control buses generally carry signals between the control unit and other computer components in a clock-driven manner. The system clock produces a continuous sequence of pulses in a specified duration and frequency. A sequence of steps t0 , t1 , t2 , . . . ,
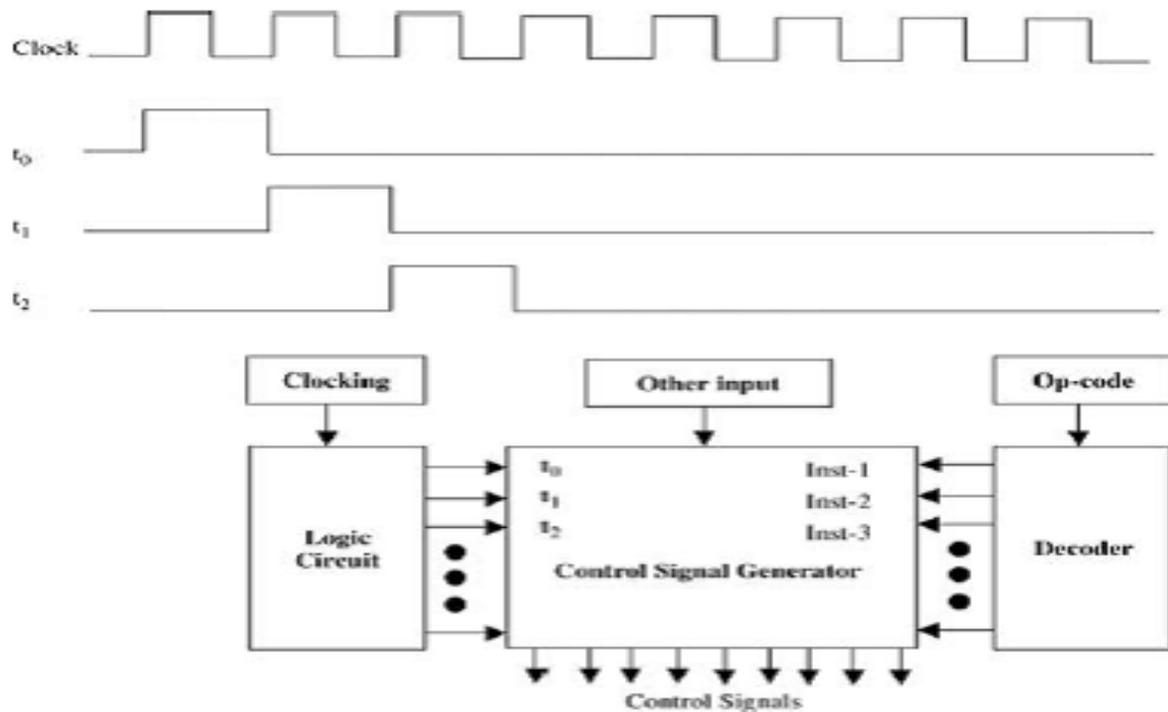
**Fig (29)**

**Timing of control signals**

(t0 , t1 , t2 , . . .) are used to execute a certain instruction. The op-code field of a fetched instruction is decoded to provide the control signal generator with information about the instruction to be executed. Step information generated by a logic circuit module is used with other inputs to generate control signals. The signal generator can be specified simply by a set of Boolean equations for its output in terms of its inputs. Figure (29) shows a block diagram that describes how timing is used in generating control signals.

There are mainly two different types of control units: **microprogrammed** and **hardwired**. In microprogrammed control, the control signals associated with operations are stored in special memory units inaccessible by the programmer as control words. A control word is a microinstruction that specifies one or more microoperations.

A sequence of microinstructions is called a microprogram, which is stored in a ROM or RAM called a control memory CM.

In hardwired control, fixed logic circuits that correspond directly to the Boolean expressions are used to generate the control signals. Clearly hardwired control is faster than microprogrammed control. However, hardwired control could be very expensive and complicated for complex systems. Hardwired control is more economical for small control units. It should also be noted that microprogrammed control could adapt easily to changes in the system design. We can easily add new instructions without changing hardware. Hardwired control will require a redesign of the entire systems in the case of any change.

**Example 1** Let us revisit the add operation in which we add the contents of source registers R1 , R2 , and store the results in destination register R0 . We have shown earlier that this operation can be done in one step using the three-bus datapath shown in Figure (26).

Let us try to examine the control sequence needed to accomplish this addition at step t0 . Suppose that the op-code field of the current instruction was decoded to Inst-x type. First, we need to select the source registers and the destination register, then we select Add as the ALU function to be performed. The following table shows the needed step and the control sequence.

| Step | Instruction type | Micro-operation | Control |
|---|---|---|---|
| $t_0$ | Inst-x | $R_0 \leftarrow (R_1) + (R_2)$ | Select $R_1$ as source 1 on out-bus1 ($R_1$ out-bus1) Select $R_2$ as source 2 on out-bus2 ($R_2$ out-bus2) Select $R_0$ as destination on in-bus ($R_0$ in-bus) Select the ALU function Add (Add) |

Figure (30)  shows the signals generated to execute Inst-x during time period t0 . The AND gate ensures that these signals will be issued when the op-code is decoded into Inst-x and during time period t0 . The signals (R1 out-bus 1), (R2 out-bus2), (R0 in-bus), and (Add) will select R1 as a source on out-bus1, R2 as a source on outbus2, R0 as destination on in-bus, and select the ALUs add function, respectively.
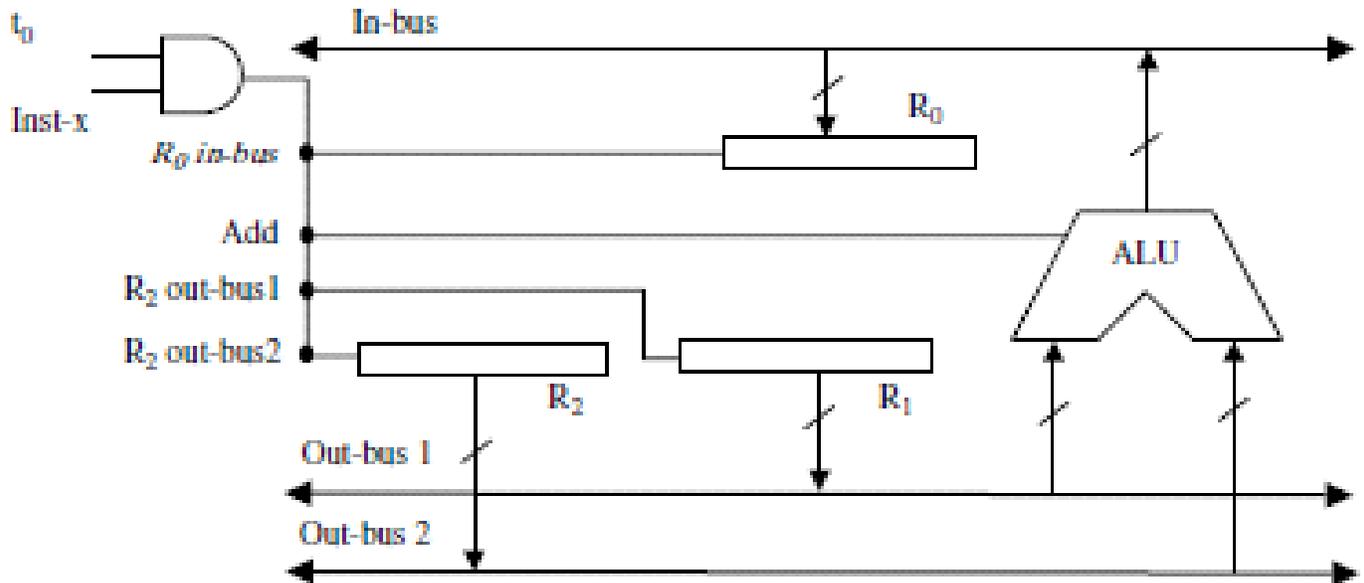
**Fig (30)**

**Signals generated to execute Inst-x on three-bus datapath during time period t0**

## Lecture22

**Example 2** Let us repeat the operation in the previous example using the one-bus datapath shown in Fig. (26). We have shown earlier that this operation can be carried out in three steps using the one-bus datapath. Suppose that the op-code field of the current instruction was decoded to Inst-x type. The following table shows the needed steps and the control sequence.

| Step | Instruction type | Micro-operation | |
|------|------------------|-----------------|---|
| $t_0$ | Inst-x | $A \leftarrow (R_1)$ | Select $R_1$ as source ($R_1$ out) |
| | | | Select $A$ as destination ($A$ in) |
| $t_1$ | Inst-x | $B \leftarrow (R_2)$ | Select $R_2$ as source ($R_2$ out) |
| | | | Select $B$ as destination ($B$ in) |
| $t_2$ | Inst-x | $R_0 \leftarrow (A) + (B)$ | Select the ALU function Add (Ac |
| | | | Select $R_0$ as destination ($R_0$ in) |

Figure (31) shows the signals generated to execute Inst-x during time periods t0 , t1 , and t2 . The AND gates ensure that the appropriate signals will be issued when the op-code is decoded into Inst-x and during the appropriate time period. During t0 , the signals (R1 out) and (A in) will

be issued to move the contents of R1 into A. Similarly during t1 , the signals (R2 out) and (B in) will be issued to move the contents of R2 into B. Finally, the signals (R0 in) and (Add) will be issued during t2 to add the contents of A and B and move the results into R0 .

## Hardwired Implementation
In hardwired control, a direct implementation is accomplished using logic circuits. For each control line, one must find the Boolean expression in terms of the input to the control signal generator.



**Fig (31)**

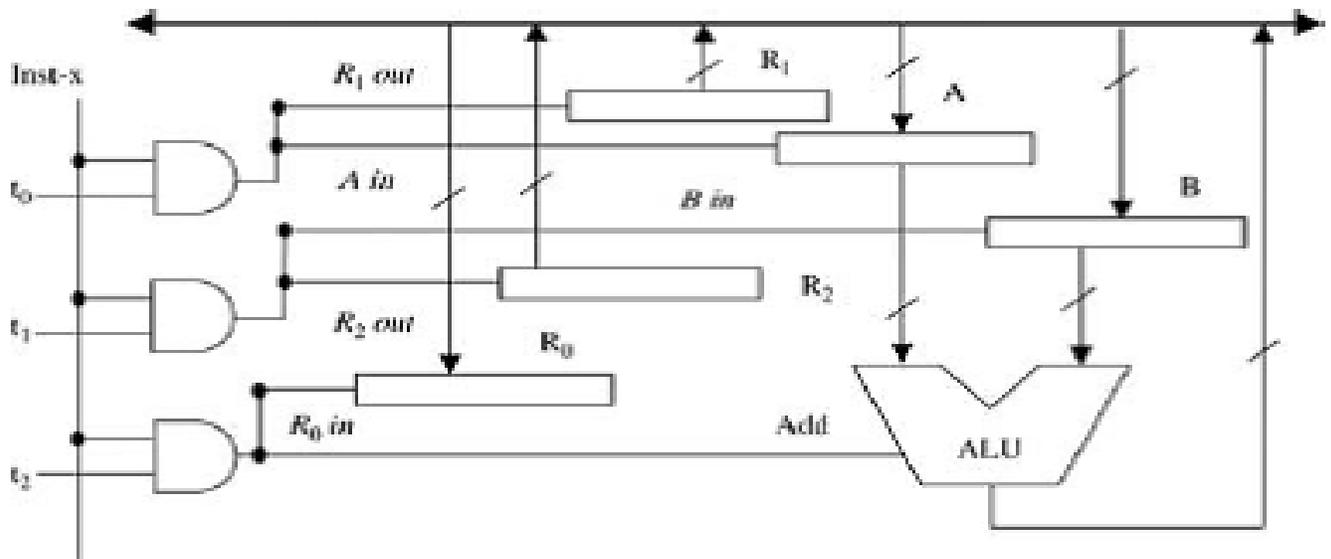**Signals generated to execute Inst-x on one-bus datapath during time period t0, t1, t2**

# Lecture23

## Microprogrammed Control Unit
The idea of microprogrammed control units was introduced by M. V. Wilkes in the early 1950s. Microprogramming was motivated by the desire to reduce the complexities involved with hardwired control. As we studied earlier, an instruction is implemented using a set of micro-operations.
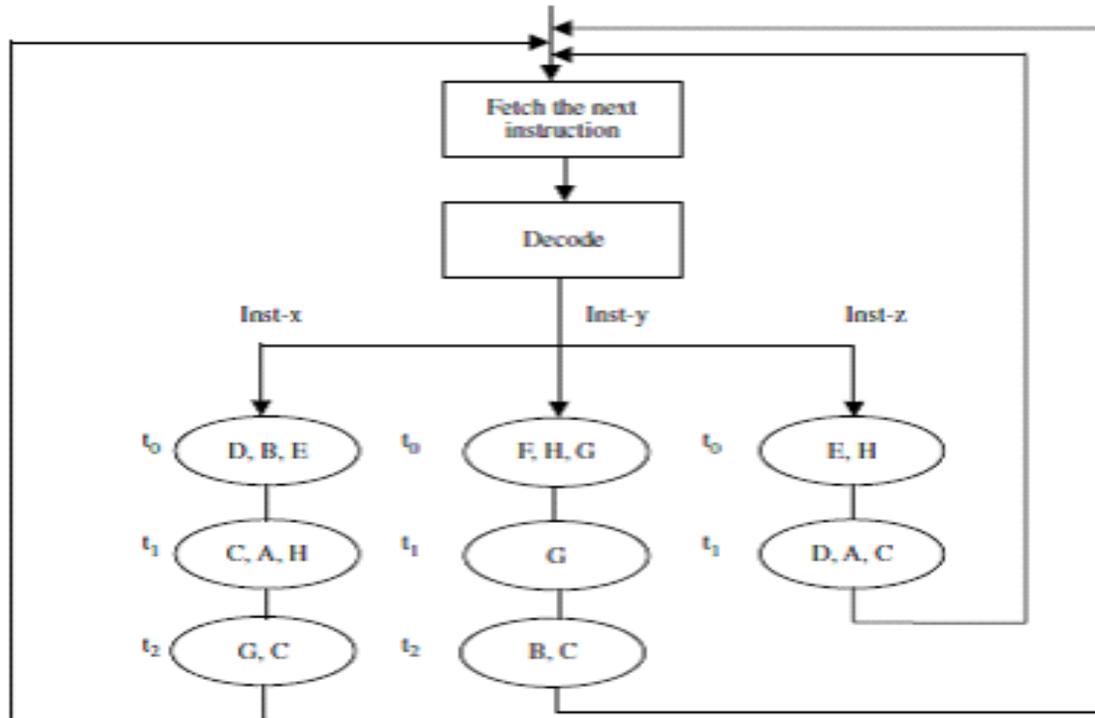
**Fig (32)**
**Instruction execution state diagram**

Associated with each micro-operation is a set of control lines that must be activated to carry out the corresponding micro-operation. The idea of microprogrammed control is to store the control signals associated with the implementation of a certain instruction as a microprogram in a special memory called a control memory (CM). A micro program consists of a sequence of microinstructions. A microinstruction is a vector of bits, where each bit is a control signal, condition code, or the address of the next microinstruction.

Microinstructions are fetched from CM the same way program instructions are fetched from main memory Fig. (33). When an instruction is fetched from memory, the op-code field of the instruction will determine which microprogram is to be executed. In other words, the op-code is mapped to a microinstruction address in the control memory. The microinstruction processor uses that address to fetch the first microinstruction in the microprogram. After fetching each microinstruction, the appropriate control lines will be enabled. Every control line that corresponds to a "1" bit should be turned on. Every control line that corresponds to a "0" bit should be left off.

After completing the execution of one microinstruction, a new microinstruction will be fetched and executed. If the condition code bits indicate that a branch must be taken, the next microinstruction is specified in the address bits of the current microinstruction. Otherwise, the

next microinstruction in the sequence will be fetched and executed. The length of a microinstruction is determined based on the number of micro-operations specified in the microinstructions, the way the control bits will be interpreted, and the way the address of the next microinstruction is obtained. A microinstruction may specify one or more micro-operations that will be activated simultaneously. The length of the microinstruction will increase as the number of parallel micro-operations per microinstruction increases. Furthermore, when each control bit in the microinstruction corresponds to exactly one control line, the length of microinstruction could get bigger. The length of a microinstruction could be reduced if control lines are coded in specific fields in the microinstruction. Decoders will be needed to map each field into the individual control lines. Clearly, using the decoders will reduce the number of control lines that can be activated simultaneously. There is a tradeoff between the length of the microinstructions and the amount of parallelism. It is important that we reduce the length of microinstructions to reduce the cost and access time of the control memory. It may also be desirable that more micro-operations be performed in parallel and more control lines can be activated simultaneously.



**Fig (33)**

**Fetching microinstructions (control words)4**

**Horizontal Versus Vertical Microinstructions**: Micro-instructions can be classified as horizontal or vertical. Individual bits in horizontal microinstructions correspond to individual control lines. Horizontal microinstructions are long and allow maximum parallelism since each bit controls a single control line. In vertical microinstructions, control lines are coded into specific fields within a microinstruction. Decoders are needed to map a field of k bits to 2k possible combinations of control lines. For example, a 3-bit field in a microinstruction could be used to specify any one of eight possible lines. Because of the encoding, vertical microinstructions are much shorter than horizontal ones. Control lines encoded in the same field cannot be activated simultaneously. Therefore, vertical microinstructions allow only limited parallelism. It should be noted that no decoding is needed in horizontal microinstructions while decoding is necessary in the vertical case.

**Example 3** Consider the three-bus datapath shown in Figure (28). In addition to the PC, IR, MAR, and MDR, assume that there are 16 general-purpose registers numbered R0–R15 . Also, assume that the ALU supports eight functions (add, subtract, multiply, divide, AND, OR, shift left, and

shift right). Consider the add operation Add R1 , R2 , R0 , which adds the contents of source registers R1 , R2 , and store the results in destination register R0 . In this example, we will study the format of the microinstruction under horizontal organization. We will use horizontal microinstructions, in which there is a control bit for each control line. The format of the microinstruction should have control bits for the following:

. ALU operations
. Registers that output to out-bus1 (source 1)
. Registers that output to out-bus2 (source 2)
. Registers that input from in-bus (destination)
. Other operations that are not shown here

The following table shows the number of bits needed for ALU, Source 1, Source 2, and destination:

| Purpose | Number of bits | Explanations |
|---|---|---|
| ALU | 8 bits | 8 functions |
| Source 1 | 20 bits | 16 general-purpose registers + 4 special-purpose registers |
| Source 2 | 16 bits | 16 general-purpose registers |
| Destination | 20 bits | 16 general-purpose registers + 4 special-purpose registers |



**Microinstruction for Add R1, R2, R0**

## Lecture24

**Example 4** In this example, we will use vertical microinstructions, in which decoders will be needed. We will use a three-bus datapath as shown in Figure (28). Assume that there are 16 general-purpose registers and that the ALU supports eight functions. The following tables show the encoding for ALU functions, registers connected to out-bus 1 (Source 1), registers connected to out-bus 2 (Source 2), and registers connected to in-bus (Destination).

| Purpose | Number of bits | Explanations |
| --- | --- | --- |
| ALU | 4 bits | 8 functions + none |
| Source 1 | 5 bits | 16 general-purpose registers + 4 special-purpose registers + none |
| Source 2 | 5 bits | 16 general-purpose registers + none |
| Destination | 5 bits | 16 general-purpose registers + 4 special-purpose registers + none |

| Encoding | ALU function |
| --- | --- |
| 0000 | None (ALU will connect out-bus1 to in-bus) |
| 0001 | Add |
| 0010 | Subtract |
| 0011 | Multiple |
| 0100 | Divide |
| 0101 | AND |
| 0110 | OR |
| 0111 | Shift left |
| 1000 | Shift right |

| Encoding | Source 1 | Destination | Encoding | Source 2 |
| --- | --- | --- | --- | --- |
| 00000 | $R_0$ | $R_0$ | 00000 | $R_0$ |
| 00001 | $R_1$ | $R_1$ | 00001 | $R_1$ |
| 00010 | $R_2$ | $R_2$ | 00010 | $R_2$ |
| 00011 | $R_3$ | $R_3$ | 00011 | $R_3$ |
| 00100 | $R_4$ | $R_4$ | 00100 | $R_4$ |
| 00101 | $R_5$ | $R_5$ | 00101 | $R_5$ |
| 00110 | $R_6$ | $R_6$ | 00110 | $R_6$ |
| 00111 | $R_7$ | $R_7$ | 00111 | $R_7$ |
| 01000 | $R_8$ | $R_8$ | 01000 | $R_8$ |
| 01001 | $R_9$ | $R_9$ | 01001 | $R_9$ |
| 01010 | $R_{10}$ | $R_{10}$ | 01010 | $R_{10}$ |

| Encoding | Source 1 | Destination | Encoding | Source 2 |
|---|---|---|---|---|
| 01011 | $R_{11}$ | $R_{11}$ | 01011 | $R_{11}$ |
| 01100 | $R_{12}$ | $R_{12}$ | 01100 | $R_{12}$ |
| 01101 | $R_{13}$ | $R_{13}$ | 01101 | $R_{13}$ |
| 01110 | $R_{14}$ | $R_{14}$ | 01110 | $R_{14}$ |
| 01111 | $R_{15}$ | $R_{15}$ | 01111 | $R_{15}$ |
| 10000 | PC | PC | 10000 | None |
| 10001 | IR | IR | | |
| 10010 | MAR | MAR | | |
| 10011 | MDR | MDR | | |
| 10100 | NONE | NONE | | |

**Example 5** Using the same encoding as Example 4, let us find vertical microinstructions used in fetching an instruction.

MAR PC First, we need to select PC as source 1 by using "10000" for source 1 field. Similarly, we select MAR as our destination by using "10010" in the destination field. We also need to use "0000" for the ALU field, which will be decoded to "NONE". As shown in the ALU encoding table (Example 4), "NONE" means that out-bus1 will be connected to in-bus. The field source 2 will be set to "10000", which means none of the registers will be selected.

Memory Read and Write Memory operations can easily be accommodated by adding 1 bit for read and another for write. The two microinstructions perform memory read and write, respectively.

Fetching an instruction can be done using the three microinstructions. The first and second microinstructions have been shown above. The third microinstruction moves the contents of the MDR to IR (IR MDR). MDR is selected as source 1 by using "10011" for source 1 field. Similarly, IR is selected as the destination by using "10001" in the destination field. We also need to use "0000" ("NONE")

| ALU | | | | Source 1 | | | | | Source 2 | | | | Destination | | | | | | Others |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | Others |

**Figure 5.14    Microinstruction for Add $R_1$, $R_2$, $R_0$**

| ALU | | | | Source 1 | | | | | Source 2 | | | | | Destination | | | | | Others |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | Others |

| | | | | R | W | |
|---|---|---|---|---|---|---|
| ALU | Source 1 | Source 2 | Destination | 1 | 0 | Others |

MDR ← Mem[MAR]

| | | | | R | W | |
|---|---|---|---|---|---|---|
| ALU | Source 1 | Source 2 | Destination | 0 | 1 | Others |

Mem[MAR] ← MDR

**Figure 5.16    Microinstructions for memory read and write**

| | ALU | | | | Source 1 | | | | | Source 2 | | | | | Destination | | | | R | W | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| MAR ← (PC) | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | Others |
| MDR ← Mem[MAR] | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | Others |
| IR ← MDR | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | Others |

**Microinstructions for memory read and write**

| | ALU | | | | | Source 1 | | | | | Source 2 | | | | | Destination | | | | R | W | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| MAR ← (PC) | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | Others |
| MDR ← Mem[MAR] | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | Others |
| IR ← MDR | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | Others |

**Microinstructions for fetching an instruction**

in the ALU field, which means that out-bus1 will be connected to in-bus. The field source 2 will be set to "10000", which means none of the registers will be selected.

## Lecture25
## Parallel processing and Pipelining

Pipelining is a technique of decomposing a sequential process into sub-operations, with each sub-process being executed in a special dedicated segment that operates concurrently with all other segments. A pipeline can be visualized as a collection of processing segments through which binary information flows.

Each segment performs partial processing dictated by the way the task is partitioned. The result obtained from the computation in each segment is transferred to the next segment in the pipeline. The final result is obtained after the data have passed through all segments. The name "pipeline" implies a flow of information analogous to an industrial assembly line. It is characteristic of pipelines that several computations can be in progress in distinct segments at the same time. The overlapping of computation is made possible by associating a register with each segment in the pipeline. The registers provide isolation between each segment so that each can operate on distinct data simultaneously.

Perhaps the simplest way of viewing the pipeline structure is to imagine that each segment consists of an input register followed by a combinational circuit. The register holds the data and the combinational circuit performs the sub-operation in the particular segment. The output of

the combinational circuit in a given segment is applied to the input register of the next segment. A clock is applied to all registers after enough time has elapsed to perform all segment activity. In this way the information flows through the pipeline one step at a time.

The pipeline organization will be demonstrated by means of a simple example. Suppose that we want to perform the combined multiply and add operations with a stream of numbers.

Ai*Bi + C for i = 1,2,3,... ,7

Each sub-operation is to be implemented in a segment within a pipeline. Each segment has one or two registers and a combinational circuit as shown in Fig. (33). Rl through R5 are registers that receive new data with every clock pulse. The multiplier and adder are combinational circuits. The sub-operations performed in each segment of the pipeline are as follows:

$$R1 \leftarrow A_i, \quad R2 \leftarrow B_i \qquad \text{Input } A_i \text{ and } B_i$$

$$R3 \leftarrow R1 * R2, \quad R4 \leftarrow C_i \qquad \text{Multiply and input } C_i$$

$$R5 \leftarrow R3 + R4 \qquad \text{Add } C_i \text{ to product}$$

The five registers are loaded with new data every clock pulse. The effect of each clock is shown in Table 9-1. The first clock pulse transfers Ai and Bi into R1 and R2.

The second clock pulse transfers the product of Rl and R2 into R3 and Q into R4. The same clock pulse transfers A2 and B2 into Rl and R2. The third clock pulse operates on all three segments simultaneously. It places A3 and B3 into Rl and R2, transfers the product of Rl and R2 into R3, transfers C2 into R4, and places the sum of R3 and R4 into R5. It takes three clock pulses to fill up the pipe and retrieve the first output from R5. From there on, each clock produces a new output and moves the data one step down the pipeline. This happens as long as new input data flow into the system. When no more input data are available, the clock must continue until the last output emerges out of the pipeline.
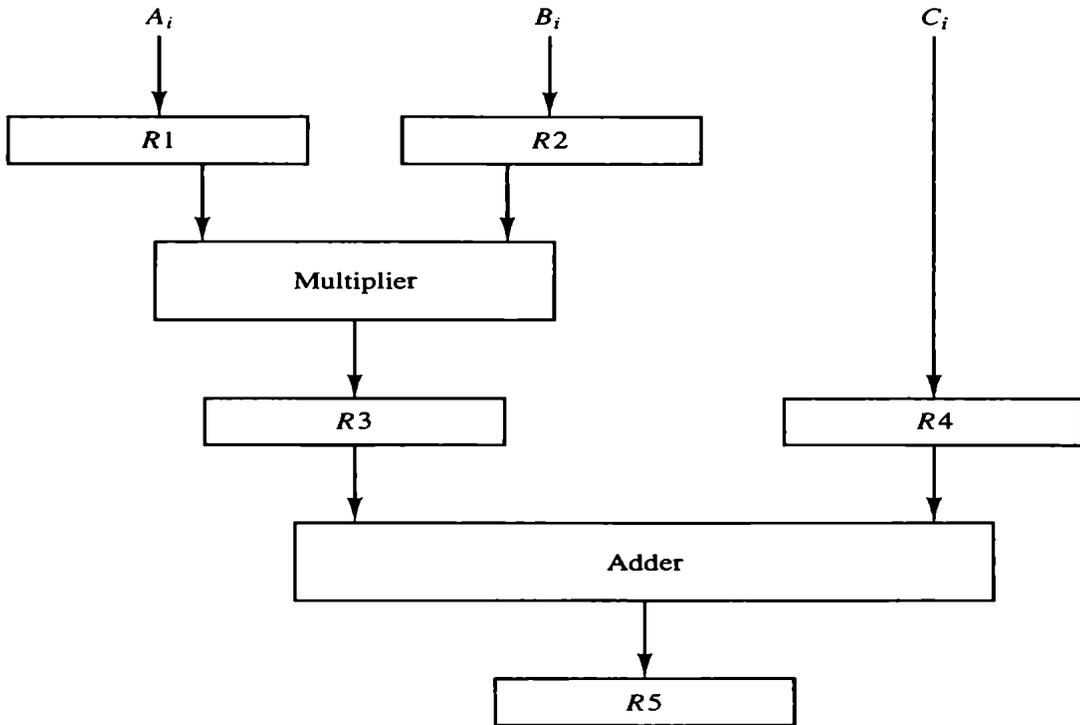
$A_i$      $B_i$      $C_i$

| $R1$ | | $R2$ |
|---|---|---|

Multiplier

| $R3$ | | $R4$ |
|---|---|---|

Adder

$R5$

**Figure (33)**

**Example of pipeline processing.**

| Clock Pulse Number | Segment 1 | | Segment 2 | | Segment 3 |
|---|---|---|---|---|---|
| | $R1$ | $R2$ | $R3$ | $R4$ | $R5$ |
| 1 | $A_1$ | $B_1$ | — | — | — |
| 2 | $A_2$ | $B_2$ | $A_1 * B_1$ | $C_1$ | — |
| 3 | $A_3$ | $B_3$ | $A_2 * B_2$ | $C_2$ | $A_1 * B_1 + C_1$ |
| 4 | $A_4$ | $B_4$ | $A_3 * B_3$ | $C_3$ | $A_2 * B_2 + C_2$ |
| 5 | $A_5$ | $B_5$ | $A_4 * B_4$ | $C_4$ | $A_3 * B_3 + C_3$ |
| 6 | $A_6$ | $B_6$ | $A_5 * B_5$ | $C_5$ | $A_4 * B_4 + C_4$ |
| 7 | $A_7$ | $B_7$ | $A_6 * B_6$ | $C_6$ | $A_5 * B_5 + C_5$ |
| 8 | — | — | $A_7 * B_7$ | $C_7$ | $A_6 * B_6 + C_6$ |
| 9 | — | — | — | — | $A_7 * B_7 + C_7$ |

# Lecture26

**General Considerations and pipeline processor**

Any operation that can be decomposed into a sequence of sub-operations of about the same complexity can be implemented by a <mark>pipeline processor</mark>. The technique is efficient for those applications that need to repeat the same task many times with different sets of data. The general structure of a four-segment pipeline is illustrated in Fig (34). The operands pass through all four segments in a fixed sequence. Each segment consists of a combinational circuit S, which performs a sub-operation over the data stream flowing through the pipe. The segments are separated by registers K, which hold the intermediate results between the stages. Information flows between adjacent stages under the control of a common clock applied to all the registers simultaneously. We define a task as the total operation performed going through all the segments in the pipeline. The behavior of a pipeline can be illustrated with a space-time diagram. This is a diagram that shows the segment utilization as a function of time. The space-time diagram of a four-segment pipeline is demonstrated in Fig. (34).
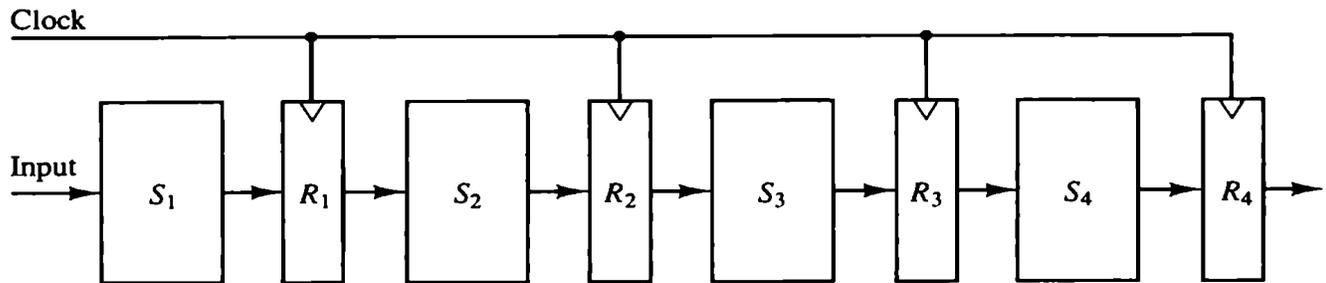


**Figure (34)**
**Four-segment pipeline.**

The horizontal axis displays the time in clock cycles and the vertical axis gives the segment number. The diagram shows six tasks T1 through T6 executed in four segments. Initially, task T1 is handled by segment 1. After the first clock, segment 2 is busy with T1, while segment 1 is busy with task T2. Continuing in this manner, the first task T1 is completed after the fourth clock cycle. From then on, the pipe completes a task every clock cycle. No matter how many segments there are in the system, once the pipeline is full, it takes only one clock period to obtain an output. Now consider the case where a k-segment pipeline with a clock cycle time tp is used to execute n tasks. The first task T1 requires a time equal to ktp to complete its operation since there are k segments in the pipe. The remaining n - 1 task emerge from the pipe

at the rate of one task per clock cycle and they will be completed after a time equal to (n - 1)tp. Therefore, to complete n tasks using a k-segment pipeline requires k + (n - 1) clock cycles.

**For example**, the diagram of Fig. (35) Shows four segments and six tasks. The time required to complete all the operations is 4 + (6 - 1) = 9 clock cycles, as indicated in the diagram.

Next consider a non-pipeline unit that performs the same operation and takes a time equal to $t_n$ to complete each task. The total time required for n tasks is $nt_n$. The speedup of a pipeline processing over an equivalent non-pipeline processing is defined by the ratio :
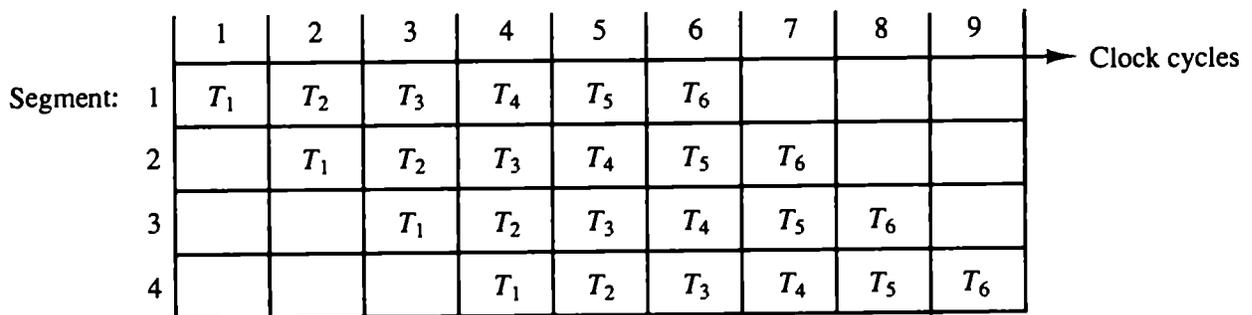
$$S = \frac{nt_n}{(k + n - 1)t_p}$$



**Figure (35)**
**Space-time diagram for pipeline.**

As the number of tasks increases, n becomes much larger than k - 1, and k + n - 1 approaches the value of n. Under this condition, the speedup becomes

$$S = \frac{t_n}{t_p}$$

If we assume that the time it takes to process a task is the same in the pipeline and non-pipeline circuits, we will have $t_n$ = ktp. Including this assumption, the speedup reduces to :

$$S = \frac{kt_p}{t_p} = k$$

This shows that the theoretical maximum speed up that a pipeline can provide is k, where k is the number of segments in the pipeline. To clarify the meaning of the speedup ratio, consider the following numerical example:

Let the time it takes to process a sub-operation in each segment be equal to tp = 20 ns. Assume that the pipeline has k = 4 segments and executes n = 100 tasks in sequence.

The pipeline system will take (k + n - 1)tp = (4 + 99) x 20 = 2060 ns to complete. Assuming that t,, = ktp = 4 x 20 = 80 ns, a non-pipeline system requires nktp = 100 x 80 = 8000 ns to complete the 100 tasks. The speedup ratio is equal to 8000/2060 = 3.88. As the number of tasks increases, the speedup will approach 4, which is equal to the number of segments in the pipeline. If we assume that t,, = 60 ns, the speedup becomes 60/20 = 3.

To duplicate the theoretical speed advantage of a pipeline process by means of multiple functional units, it is necessary to construct k identical units that will be operating in parallel. The implication is that a /c-segment pipeline processor can be expected to equal the performance of k copies of an equivalent non-pipeline circuit under equal operating conditions. This is illustrated in Fig. (36), where four identical circuits are connected in parallel. Each P circuit performs the same task of an equivalent pipeline circuit. Instead of operating with the input data in sequence as in a pipeline, the parallel circuits accept four input data items simultaneously and perform four tasks at the same time. As far as the speed of operation is concerned, this is equivalent to a four segment pipeline. Note that the four-unit circuit of Fig. 9-5 constitutes a single-instruction multiple-data (SIMD) organization since the same instruction is used to operate on multiple data in parallel. There are various reasons why the pipeline cannot operate at its maximum theoretical rate. Different segments may take different times to complete their sub-operation. The clock cycle must be chosen to equal the time delay of the segment with the maximum propagation time.
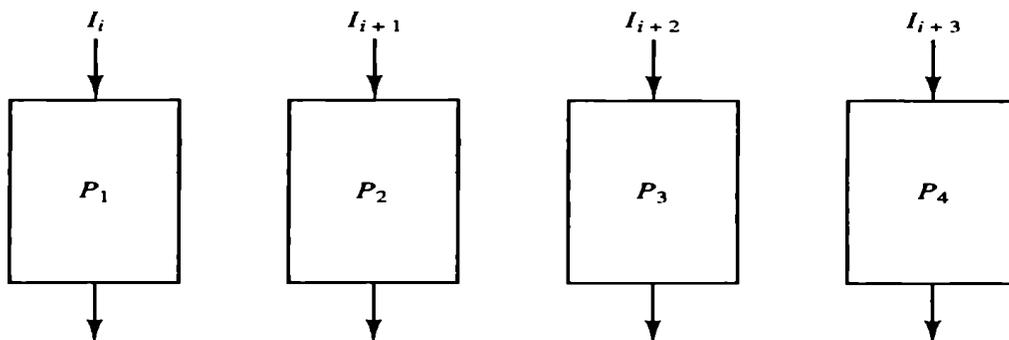
**Figure (36)**
**Multiple functional units in parallel.**

This causes all other segments to waste time while waiting for the next clock. Moreover, it is not always correct to assume that a nonpipe circuit has the same time delay as that of an equivalent pipeline circuit. Many of the intermediate registers will not be needed in a single-unit circuit, which can usually be constructed entirely as a combinational circuit. Nevertheless, the pipeline technique provides a faster operation over a purely serial sequence even though the maximum theoretical speed is never fully achieved. There are two areas of computer design where the pipeline organization is applicable. An arithmetic pipeline divides an arithmetic operation into sub- operations for execution in the pipeline segments. An instruction pipeline operates on a stream of instructions by overlapping the fetch, decode, and execute phases of the instruction cycle.

## Lecture27

### Arithmetic Pipeline

Pipeline arithmetic units are usually found in very high speed computers. They are used to implement floating-point operations, multiplication of fixed-point numbers, and similar computations encountered in scientific problems. A pipeline multiplier is essentially an array multiplier with special adders designed to minimize the carry propagation time through the partial products. Floating-point operations are easily decomposed into sub-operations as demonstrated in Sec. (37). We will now show an example of a pipeline unit for floating-point addition and subtraction. The inputs to the floating-point adder pipeline are two normalized floating-point binary numbers.

$$X = A \times 2^a$$
$$Y = B \times 2^b$$

A and B are two fractions that represent the mantissas and a and b are the exponents. The floating-point addition and subtraction can be performed in four segments. The registers labeled R are placed between the segments to store intermediate results. The sub-operations that are performed in the four segments are:

1. Compare the exponents.
2. Align the mantissas.
3. Add or subtract the mantissas.

4. Normalize the result.

The exponents are compared by subtracting them to determine their difference. The larger exponent is chosen as the exponent of the result. The exponent difference determines how many times the mantissa associated with the smaller exponent must be shifted to the right. This produces an alignment of the two mantissas. It should be noted that the shift must be designed as a combinational circuit to reduce the shift time. The two mantissas are added or subtracted in segment 3. The result is normalized in segment 4. When an overflow occurs, the mantissa of the sum or difference is shifted right and the exponent incremented by one. If an underflow occurs, the number of leading zeros in the mantissa determines the number of left shifts in the mantissa and the number that must be subtracted from the exponent.

**Example:** Consider the two normalized floating-point numbers:

$$X = 0.9504 \times 10^3$$
$$Y = 0.8200 \times 10^2$$

The two exponents are subtracted in the first segment to obtain 3-2 = 1. The larger exponent 3 is chosen as the exponent of the result. The next segment shifts the mantissa of Y to the right to obtain

$$X = 0.9504 \times 10^3$$
$$Y = 0.0820 \times 10^3$$

This aligns the two mantissas under the same exponent. The addition of the two mantissas in segment 3 produces the sum
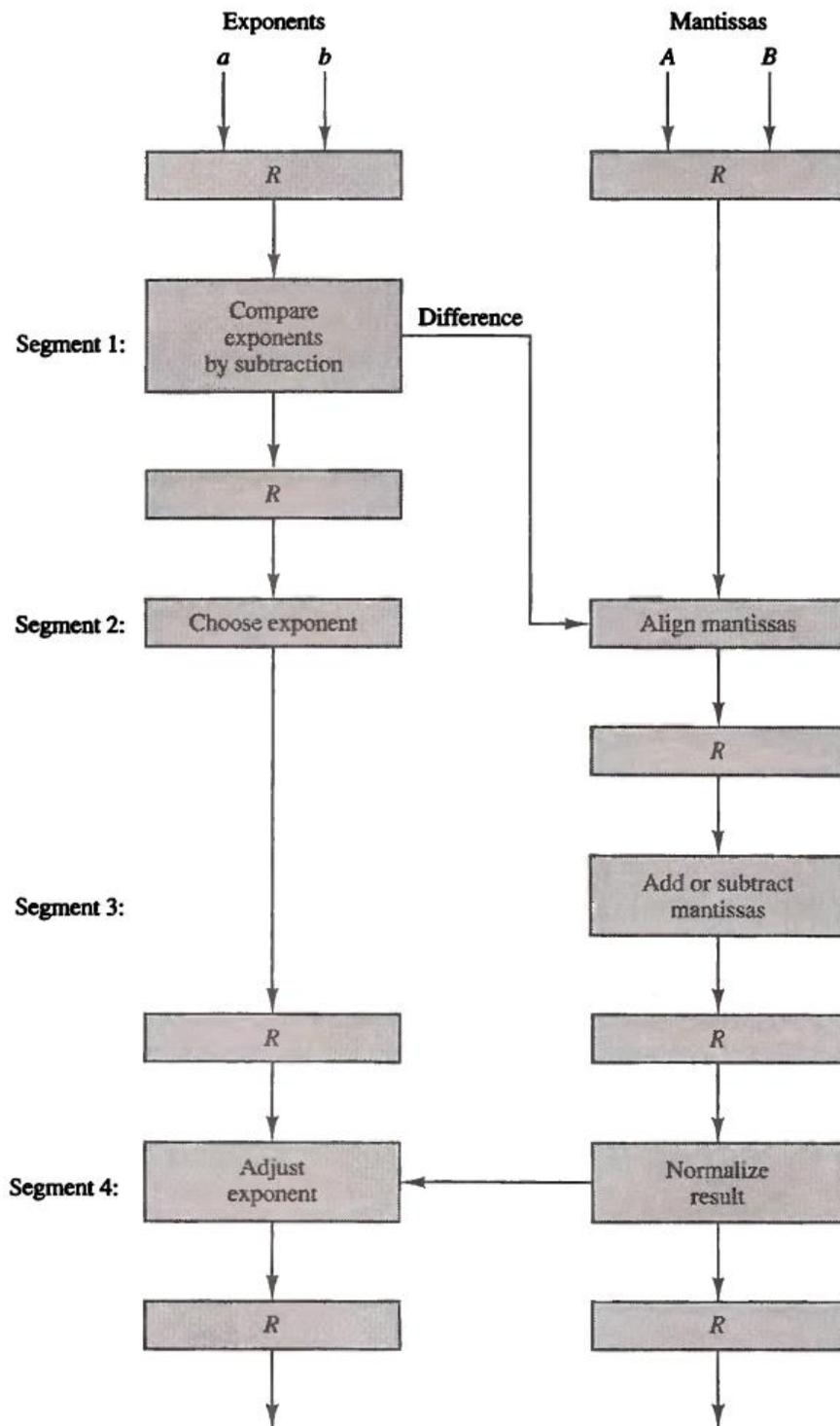
$$Z = 1.0324 \times 10^3$$

Exponents

a    b

Mantissas

A    B

| R | R |

**Segment 1:**

Compare exponents by subtraction → Difference

| R |

**Segment 2:** Choose exponent → Align mantissas

| R |

**Segment 3:** Add or subtract mantissas

| R | R |

**Segment 4:** Adjust exponent ← Normalize result

| R | R |

**Figure (37)**
**Pipeline for floating-point addition and subtraction.**

The sum is adjusted by normalizing the result so that it has a fraction with a nonzero first digit. This is done by shifting the mantissa once to the right and incrementing the exponent by one to obtain the normalized sum.

$$Z = 0.10324 \times 10^4$$

The comparator, shifter, adder-subtractor, incrementer, and decrementer in the floating-point pipeline are implemented with combinational circuits. Suppose that the time delays of the four segments are $f$, = 60 ns, $t2$ = 70 ns, $t3$ = 100 ns, $t4$ = 80 ns, and the interface registers have a delay of $tr$ = 10 ns. The clock cycle is chosen to be $tp = t3 + tT$ = 110 ns. An equivalent nonpipeline floatingpoint adder-subtractor will have a delay time $t,, = f, + t2 + t3 + t4 + tr$ = 320 ns. In this case the pipelined adder has a speedup of 320/110 = 2.9 over the nonpipelined adder.

## Lecture28

### Instruction Pipeline

Pipeline processing can occur not only in the data stream but in the instruction stream as well. An instruction pipeline reads consecutive instructions from memory while previous instructions are being executed in other segments. This causes the instruction fetch and executes phases to overlap and perform simultaneous operations. One possible digression associated with such a scheme is that an instruction may cause a branch out of sequence. In that case the pipeline must be emptied and all the instructions that have been read from memory after the branch instruction must be discarded.

Consider a computer with an instruction fetch unit and an instruction execution unit designed to provide a two-segment pipeline. The instruction fetch segment can be implemented by means of a first-in, first-out (FIFO) buffer. This is a type of unit that forms a queue rather than a stack. Whenever the execution unit is not using memory, the control increments the program counter and uses its address value to read consecutive instructions from memory. The instructions are inserted into the FIFO buffer so that they can be executed on a first-in, first-out basis. Thus an instruction stream can be placed in a queue, waiting for decoding and processing by the execution segment. The instruction stream queuing mechanism provides an efficient way for reducing the average access time to memory for reading instructions. Whenever there is space in the FIFO buffer, the control unit initiates the next instruction fetch phase. The buffer acts as a queue from which control then extracts the instructions for the execution unit.

Computers with complex instructions require other phases in addition to the fetch and execute to process an instruction completely. In the most general case, the computer needs to process each instruction with the following sequence of steps.

1. Fetch the instruction from memory.
2. Decode the instruction.
3. Calculate the effective address.
4. Fetch the operands from memory.
5. Execute the instruction.
6. Store the result in the proper place.

There are certain difficulties that will prevent the instruction pipeline from operating at its maximum rate. Different segments may take different times to operate on the incoming information. Some segments are skipped for certain operations. For example, a register mode instruction does not need an effective address calculation. Two or more segments may require memory access at the same time, causing one segment to wait until another is finished with the memory. Memory access conflicts are sometimes resolved by using two memory buses for accessing instructions and data in separate modules. In this way, an instruction word and a data word can be read simultaneously from two different modules.

The design of an instruction pipeline will be most efficient if the instruction cycle is divided into segments of equal duration. The time that each step takes to fulfill its function depends on the instruction and the way it is executed.

**Example: Four-Segment Instruction Pipeline**

Assume that the decoding of the instruction can be combined with the calculation of the effective address into one segment. Assume further that most of the instructions place the result into a processor registers so that the instruction execution and storing of the result can be combined into one segment. This reduces the instruction pipeline into four segments.

Figure 9-7 shows how the instruction cycle in the CPU can be processed with a four-segment pipeline. While an instruction is being executed in segment 4, the next instruction in sequence is busy fetching an operand from memory in segment 3. The effective address may be calculated in a separate arithmetic circuit for the third instruction, and whenever the memory is available, the fourth and all subsequent instructions can be fetched and placed in an instruction FIFO. Thus up to four sub-operations in the instruction cycle can overlap and up to four different instructions can be in progress of being processed at the same time.

Once in a while, an instruction in the sequence may be a program control type that causes a branch out of normal sequence. In that case the pending operations in the last two segments are completed and all information stored in the instruction buffer is deleted. The pipeline then restarts from the new address stored in the program counter. Similarly, an interrupt request, when

Acknowledged, will cause the pipeline to empty and start again from a new address value.
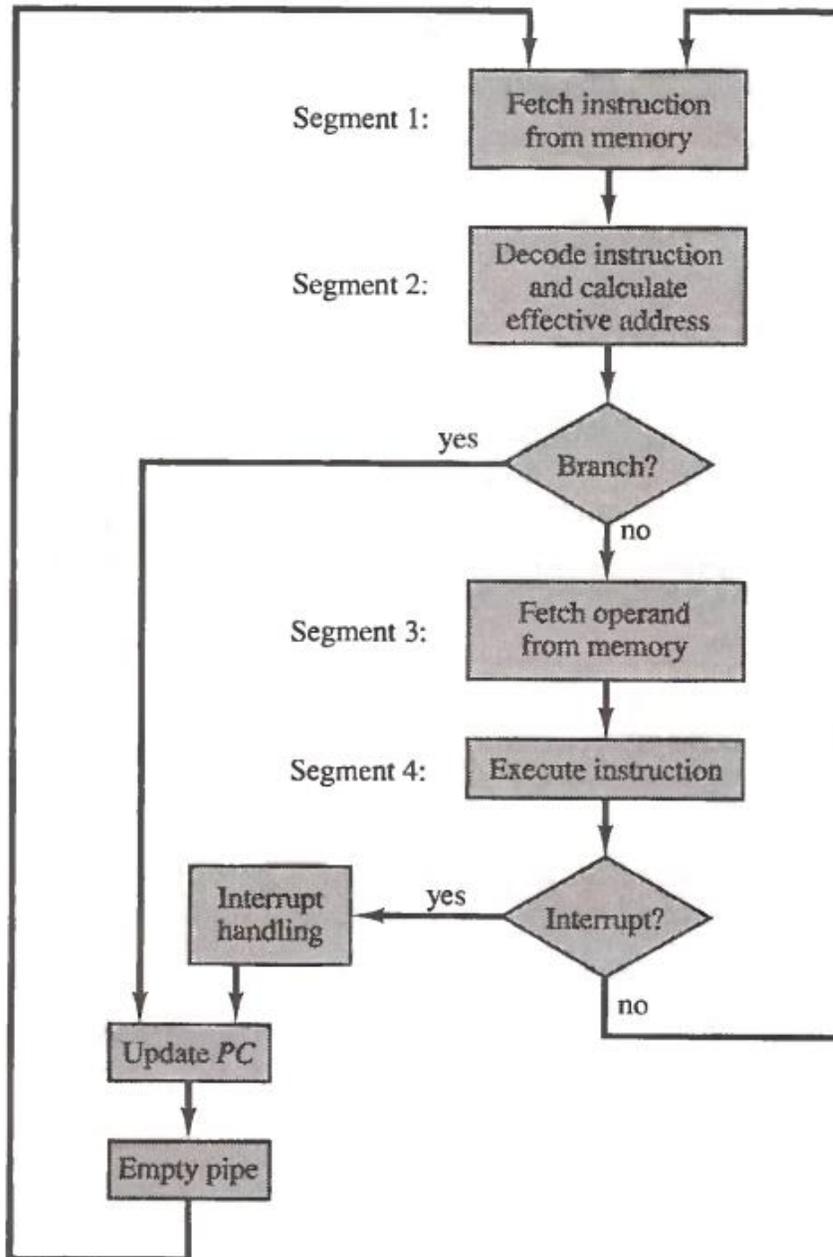
**Figure (38)**
**Four-segment CPU pipeline.**

Figure (38) shows the operation of the instruction pipeline. The time in the horizontal axis is divided into steps of equal duration. The four segments are represented in the diagram with an abbreviated symbol.

1. FI is the segment that fetches an instruction.
2. DA is the segment that decodes the instruction and calculates the effective address.

3. FO is the segment that fetches the operand.

4. EX is the segment that executes the instruction.

It is assumed that the processor has separate instruction and data memories so that the operation in FI and FO can proceed at the same time. In the absence of a branch instruction, each segment operates on different instructions.

| Step: | | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Instruction: | 1 | FI | DA | FO | EX | | | | | | | | | |
| | 2 | | FI | DA | FO | EX | | | | | | | | |
| (Branch) | 3 | | | FI | DA | FO | EX | | | | | | | |
| | 4 | | | | FI | – | – | FI | DA | FO | EX | | | |
| | 5 | | | | | – | – | – | FI | DA | FO | EX | | |
| | 6 | | | | | | | | | FI | DA | FO | EX | |
| | 7 | | | | | | | | | | FI | DA | FO | EX |

**Timing of instruction pipeline.**

Thus, in step 4, instruction 1 is being executed in segment EX; the operand for instruction 2 is being fetched in segment FO; instruction 3 is being decoded in segment DA; and instruction 4 is being fetched from memory in segment FI. Assume now that instruction 3 is a branch instruction. As soon as this instruction is decoded in segment DA in step 4, the transfer from FI to DA of the other instructions is halted until the branch instruction is executed in step 6. If the branch is taken, a new instruction is fetched in step 7. If the branch is not taken, the instruction fetched previously in step 4 can be used. The pipeline then continues until a new branch instruction is encountered. Another delay may occur in the pipeline if the EX segment needs to store the result of the operation in the data memory while the FO segment needs to fetch an operand. In that case, segment FO must wait until segment EX has finished its operation. In general, there are three major difficulties that cause the instruction pipeline to deviate from its normal operation.

1. Resource conflicts caused by access to memory by two segments at the same time. Most of these conflicts can be resolved by using separate instruction and data memories.

2. Data dependency conflicts arise when an instruction depends on the result of a previous instruction, but this result is not yet available.

3. Branch difficulties arise from branch and other instructions that change the value of PC.

## Lecture29

**Data Dependency**

A difficulty that may cause a degradation of performance in an instruction pipeline is due to possible collision of data or address. A collision occurs when an instruction cannot proceed because previous instructions did not complete certain operations. A data dependency occurs when an instruction needs data that are not yet available. For example, an instruction in the FO segment may need to fetch an operand that is being generated at the same time by the previous instruction in segment EX. Therefore, the second instruction must wait for data to become available by the first instruction. Similarly, an address dependency may occur when an operand address cannot be calculated because the information needed by the addressing mode is not available. For example, an instruction with register indirect mode cannot proceed to fetch the operand if the previous instruction is loading the address into the register. Therefore, the operand access to memory must be delayed until the required address is available. Pipelined computers deal with such conflicts between data dependencies in a variety of ways.

The most straightforward method is to insert hardware interlocks. An interlock is a circuit that detects instructions whose source operands are destinations of instructions farther up in the pipeline. Detection of this situation causes the instruction whose source is not available to be delayed by enough clock cycles to resolve the conflict. This approach maintains the program sequence by using hardware to insert the required delays.

Another technique called operand forwarding uses special hardware to detect a conflict and then avoid it by routing the data through special paths between pipeline segments. For example, instead of transferring an ALU result into a destination register, the hardware checks the destination operand, and if it is needed as a source in the next instruction, it passes the result directly into the ALU input, bypassing the register file. This method requires additional hardware paths through multiplexers as well as the circuit that detects the conflict.

A procedure employed in some computers is to give the responsibility for solving data conflicts problems to the compiler that translates the high-level programming language into a machine language program. The compiler for such computers is designed to detect a data conflict and reorder the instructions as necessary to delay the loading of the conflicting data by inserting no-operation instructions. This method is referred to as delayed load.

**Handling of Branch Instructions**

One of the major problems in operating an instruction pipeline is the occurrence of branch instructions. A branch instruction can be **conditional** or **unconditional**. An unconditional branch always alters the sequential program flow by loading the program counter with the target address. In a conditional branch, the control selects the target instruction if the condition is satisfied or the next sequential instruction if the condition is not satisfied. As mentioned

previously, the branch instruction breaks the normal sequence of the instruction stream, causing difficulties in the operation of the instruction pipeline.

Pipelined computers employ various hardware techniques to minimize the performance degradation caused by instruction branching. One way of handling a conditional branch is to pre-fetch the target instruction in addition to the instruction following the branch. Both are saved until the branch is executed. If the branch condition is successful, the pipeline continues from the branch target instruction. An extension of this procedure is to continue fetching instructions from both places until the branch decision is made. At that time control chooses the instruction stream of the correct program flow.

Another possibility is the use of a branch target buffer or BTB. The BTB is an associative memory included in the fetch segment of the pipeline. Each entry in the BTB consists of the address of a previously executed branch instruction and the target instruction for that branch. It also stores the next few instructions after the branch target instruction. When the pipeline decodes a branch instruction, it searches the associative memory BTB for the address of the instruction. If it is in the BTB, the instruction is available directly and prefetch continues from the new path. If the instruction is not in the BTB, the pipeline shifts to a new instruction stream and stores the target instruction in the BTB. The advantage of this scheme is that branch instructions that have occurred previously are readily available in the pipeline without interruption. A variation of the BTB is the loop buffer. This is a small very high speed register file maintained by the instruction fetch segment of the pipeline. When a program loop is detected in the program, it is stored in the loop buffer in its entirety, including all branches. The program loop can be executed directly without having to access memory until the loop mode is removed by the final branching out.

Another procedure that some computers use is branch prediction. A pipeline with branch prediction uses some additional logic to guess the outcome of a conditional branch instruction before it is executed. The pipeline then begins prefetching the instruction stream from the predicted path. A correct prediction eliminates the wasted time caused by branch penalties.

A procedure employed in most RISC processors is the delayed branch. In this procedure, the compiler detects the branch instructions and rearranges the machine language code sequence by inserting useful instructions that keep the pipeline operating without interruptions. An example of delayed branch is the insertion of a no-operation instruction after a branch instruction. This causes the computer to fetch the target instruction during the execution of the no- operation instruction, allowing a continuous flow of the pipeline. An example of delayed branch is presented in the next section.