## INPUT-OUTPUT ORGANIZATION

# Peripheral Devices:

The Input / output organization of computer depends upon the size of computer and the peripherals connected to it. The I/O Subsystem of the computer, provides an efficient mode of communication between the central system and the outside environment

The most common input output devices are:

      i) Monitor

      ii) Keyboard

      iii) Mouse

      iv) Printer

      v) Magnetic tapes

   The devices that are under the direct control of the computer are said to be <u>connected online.</u>

# Input - Output Interface

   Input Output Interface provides a method for transferring information between internal storage and external I/O devices.

 Peripherals connected to a computer need <u>special communication links</u> for interfacing them with the central processing unit.

  The purpose of communication link is to <u>resolve the differences that exist between the central computer and each peripheral.</u>

The Major Differences are:-

1. Peripherals are electromechnical      and electromagnetic devices and CPU and memory are electronic devices. Therefore, a conversion of signal values may be needed.

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

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

4. The operating modes of peripherals are different from each other and 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 supervises and synchronizes all input and out transfers

- These components are called Interface Units because they interface between the processor bus and the peripheral devices.

## I/O BUS and Interface Module

It defines the typical link between the processor and several peripherals.

The I/O Bus consists of data lines, address lines and control lines.

The I/O bus from the processor is attached to all peripherals interface.

To communicate with a particular device, the processor places a device address on address lines.

Each Interface decodes the address and control received from the I/O bus, interprets them for peripherals and provides signals for the peripheral controller.

It is also synchronizes the data flow and supervises the transfer between peripheral and processor.

Each peripheral has its own controller.

For example, the printer controller controls the paper motion, the print timing

The control lines are referred as I/O command. The commands are as following:
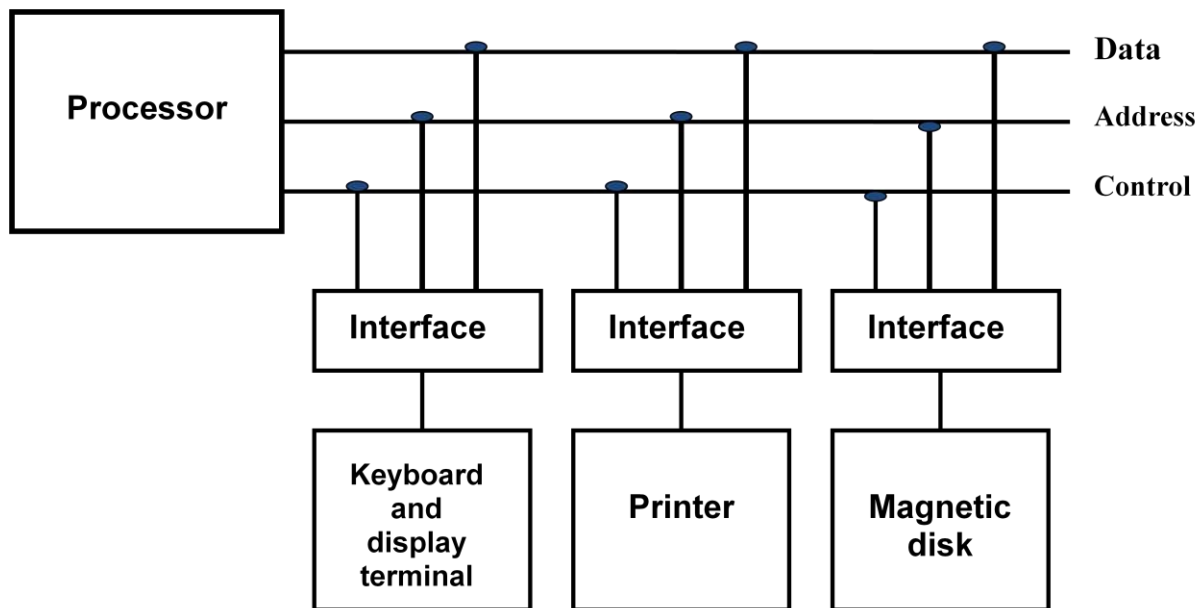
Control command- A control command is issued to activate the peripheral and to inform it what to do.

Status command- A status command is used to test various status conditions in the interface and the peripheral.

Data Output command- A data output command causes the interface to respond by transferring data from the bus into one of its registers.

Data Input command- The data input command is the opposite of the data output.

In this case the interface receives on item of data from the peripheral and places it in its buffer register. I/O Versus Memory Bus

**Connection of I/O bus to input-output devices**

To communicate with I/O, the processor must communicate with the memory unit. Like the I/O bus, the memory bus contains data, address and read/write control lines. There are 3 ways that computer buses can be used to communicate with memory and I/O:

   i.   Use two Separate buses , one for memory and other for I/O.

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

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

I/O Processor

In the first method, the computer has independent sets of data, address and control buses one for accessing memory and other for I/O. This is done in computers that provides a separate I/O processor (IOP). The purpose of IOP is to provide an independent pathway for the transfer of information between external device and internal memory.

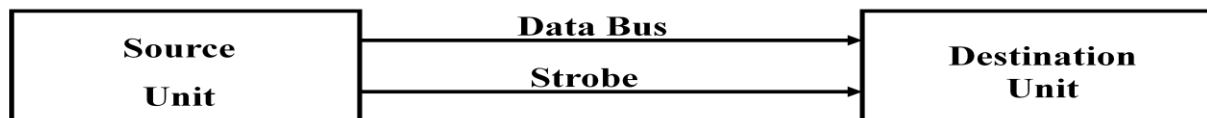## Asynchronous Data Transfer :

This Scheme is used when speed of I/O devices do not match with microprocessor, and timing characteristics of I/O devices is not predictable. In this method, process initiates the device and check its status. As a result, CPU has to wait till I/O device is ready to transfer data. When device is ready CPU issues instruction for I/O transfer. In this method two types of techniques are used based on signals before data transfer.
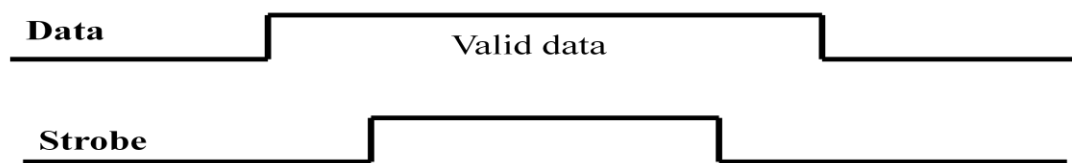
      i.  Strobe Control

      ii. Handshaking

**UNIT-V**

**Strobe Signal :**

The strobe control method of Asynchronous data transfer employs a single control line to time each transfer. The strobe may be activated by either the source or the destination unit.

Data Transfer Initiated by Source Unit:



(a)     **Block Diagram**

(b)     **Timing Diagram**
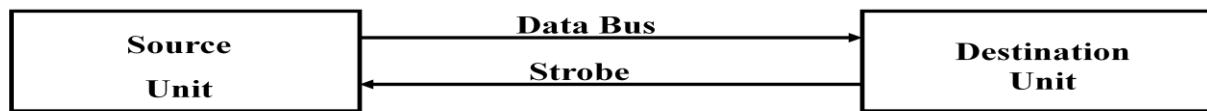
**Source-Initiated strobe for Data Transfer**

In the block diagram fig. (a), the data bus carries the binary information from source to destination unit. Typically, the bus has multiple lines to transfer an entire byte or word. The strobe is a single line that informs the destination unit when a valid data word is available.

The timing diagram fig. (b) the source unit first places the data on the data bus. The information on the data bus and strobe signal remain in the active state to allow the destination unit to receive the data.
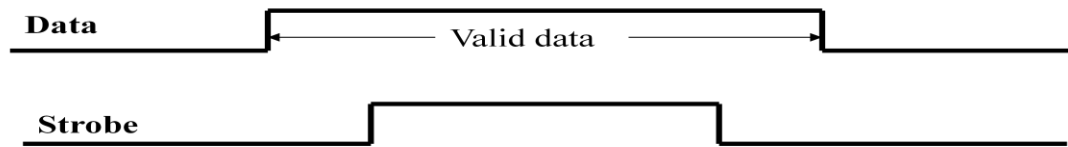
Data Transfer Initiated by Destination Unit:

In this method, the destination unit activates the strobe pulse, to informing the source to provide the data. The source will respond by placing the requested binary information on the data bus.

The data must be valid and remain in the bus long enough for the destination unit to accept it. When accepted the destination unit then disables the strobe and the source unit removes the data from the bus.

**UNIT-V**

(a)    Block Diagram

(b)    Timing Diagram

**Destination–Initiated strobe for Data Transfer**

**Disadvantage of Strobe Signal :**

The disadvantage of the strobe method is that, the source unit initiates the transfer has no way of knowing whether the destination unit has actually received the data item that was places in the bus. Similarly, a destination unit that initiates the transfer has no way of knowing whether the source unit has actually placed the data on bus. The Handshaking method solves this problem.

**Handshaking:**

The handshaking method solves the problem of strobe method by introducing a second control signal that provides a reply to the unit that initiates the transfer.
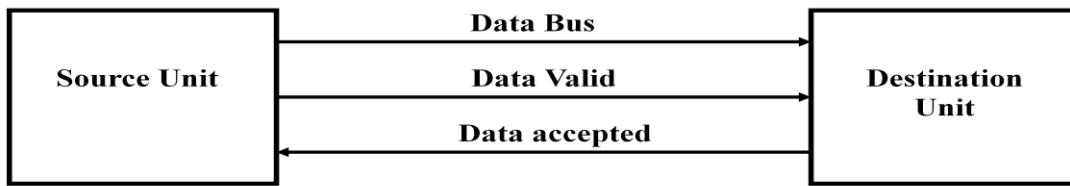
Principle of Handshaking:

The basic principle of the two-wire handshaking method of data transfer is as follow:
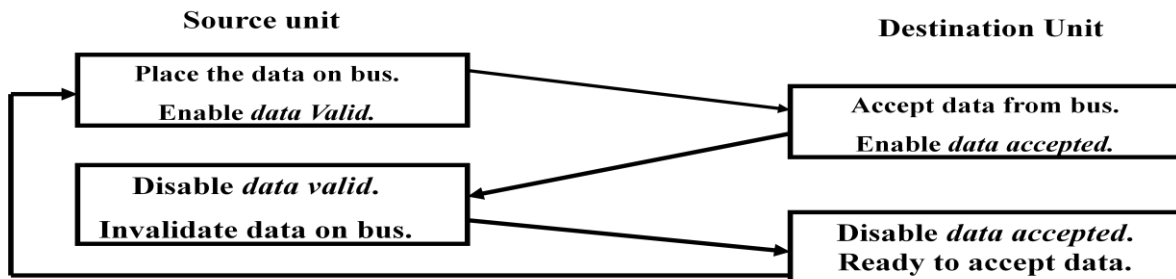
One control line is in the same direction as the data flows in the bus from the source to destination. It is used by source unit to inform the destination unit whether there a valid data in the bus. The other control line is in the other direction from the destination to the source. It is used by the destination unit to inform the source whether it can accept the data. The sequence of control during the transfer depends on the unit that initiates the transfer.

Source Initiated Transfer using Handshaking:

The sequence of events shows four possible states that the system can be at any given time. The source unit initiates the transfer by placing the data on the bus and enabling its *data valid* signal. The *data accepted* signal is activated by the destination unit after it accepts the data from the bus. The source unit then disables its *data accepted* signal and the system goes into its initial state.

**UNIT-V**

(a)     **Block Diagram**



**(b) Sequence of events**

Destination Initiated Transfer Using Handshaking:

The name of the signal generated by the destination unit has been changed to *ready for data* to reflects its new meaning. The source unit in this case does not place data on the bus until after it receives the *ready for data* signal from the destination unit. From there on, the handshaking procedure follows the same pattern as in the source initiated case.

   The only difference between the Source Initiated and the Destination Initiated transfer is in their choice of Initial sate.



(a)     **Block Diagram**



**(b) Sequence of events**

**Destination-Initiated transfer using Handshaking**

Advantage of the Handshaking method:

- ➢ The Handshaking scheme provides degree of flexibility and reliability because the successful completion of data transfer relies on active participation by both units.

- ➢ If any of one unit is faulty, the data transfer will not be completed. Such an error can be detected by means of a *Timeout mechanism* which provides an alarm if the data is not completed within time.

**Asynchronous Serial Transmission:**

The transfer of data between two units is serial or parallel. In parallel data transmission, n bit in the message must be transmitted through n separate conductor path. In serial transmission, each bit in the message is sent in sequence one at a time.
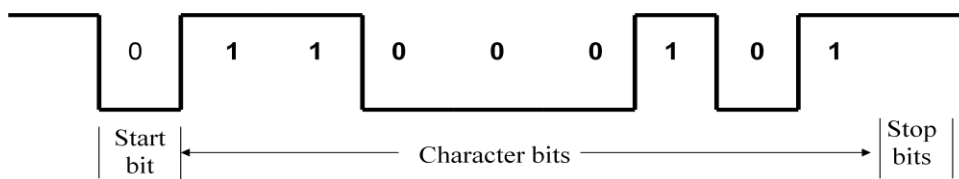
Parallel transmission is faster but it requires many wires. It is used for short distances and where speed is important. Serial transmission is slower but is less expensive.

In Asynchronous serial transfer, each bit of message is sent a sequence at a time, and binary information is transferred only when it is available. When there is no information to be transferred, line remains idle.

In this technique each character consists of three points :

<div style="text-align:center">

i. Start bit

ii. Character bit

iii. Stop bit

</div>

i.    Start Bit-  First bit, called start bit is always zero and used to indicate the beginning character.

ii.   Stop Bit-   Last bit, called stop bit is always one and used to indicate end of characters. Stop bit is always in the 1- state and frame the end of the characters to signify the idle or wait state.

iii.  Character Bit-  Bits in between the start bit and the stop bit are known as character bits. The character bits always follow the start bit.



**Asynchronous Serial Transmission**

Serial Transmission of Asynchronous is done by two ways:

**UNIT-V**

a) Asynchronous Communication Interface

b) First In First out Buffer

**Asynchronous Communication Interface:**

It works as both a receiver and a transmitter. Its operation is initialized by CPU by sending a byte to the control register.

The transmitter register accepts a data byte from CPU through the data bus and transferred to a shift register for serial transmission.

The receive portion receives information into another shift register, and when a complete data byte is received it is transferred to receiver register.

CPU can select the receiver register to read the byte through the data bus. Data in the status register is used for input and output flags.

**First In First Out Buffer (FIFO):**

A First In First Out (FIFO) Buffer is a memory unit that stores information in such a manner that the first item is in the item first out. A FIFO buffer comes with separate input and output terminals. The important feature of this buffer is that it can input data and output data at two different rates.

When placed between two units, the FIFO can accept data from the source unit at one rate, rate of transfer and deliver the data to the destination unit at another rate.

If the source is faster than the destination, the FIFO is useful for source data arrive in bursts that fills out the buffer. FIFO is useful in some applications when data are transferred asynchronously.

## Modes of Data Transfer :

Transfer of data is required between CPU and peripherals or memory or sometimes between any two devices or units of your computer system. To transfer a data from one unit to another one should be sure that both units have proper connection and at the time of data transfer the receiving unit is not busy. This data transfer with the computer is Internal Operation.

All the internal operations in a digital system are synchronized by means of clock pulses supplied by a common clock pulse Generator. The data transfer can be

i. Synchronous or

ii. Asynchronous

When both the transmitting and receiving units use same clock pulse then such a data transfer is called Synchronous process. On the other hand, if the there is not concept of clock pulses

and the sender operates at different moment than the receiver then such a data transfer is called Asynchronous data transfer.
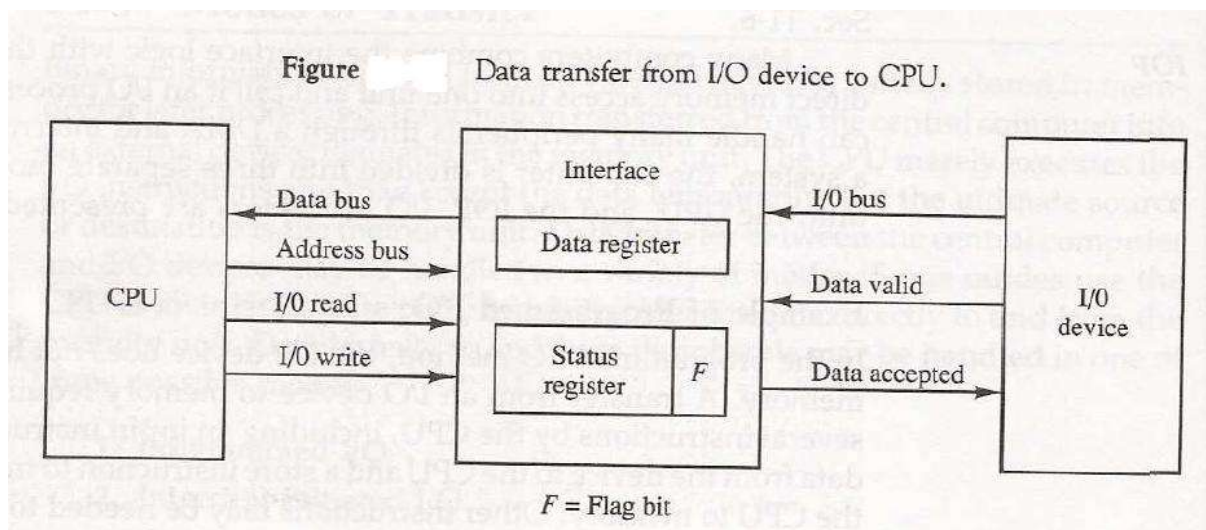
The data transfer can be handled by various modes. some of the modes use CPU as an intermediate path, others transfer the data directly to and from the memory unit and this can be handled by 3 following ways:

<div style="text-align:center">

i. Programmed I/O

ii. Interrupt-Initiated I/O

iii. Direct Memory Access (DMA)

</div>

### Programmed I/O Mode:

In this mode of data transfer the operations are the results in I/O instructions which is a part of computer program. Each data transfer is initiated by a instruction in the program. Normally the transfer is from a CPU register to peripheral device or vice-versa.

Once the data is initiated the CPU starts monitoring the interface to see when next transfer can made. The instructions of the program keep close tabs on everything that takes place in the interface unit and the I/O devices.

**Figure**      Data transfer from I/O device to CPU.



F = Flag bit

- The transfer of data requires three instructions:

1. Read the status register.

2. Check the status of the flag bit and branch to step 1 if not set or to step 3 if set.

3. Read the data register.

In this technique CPU is responsible for executing data from the memory for output and storing data in memory for executing of Programmed I/O as shown in Flowchart-:

**Programmed I/O**

CPU issues the read or write command to I/O module

I/O module informs about its status to CPU

Status → Error

Busy

Ready

CPU reads word from I/O module & writes it to memory or CPU reads word from memory & writes it to I/O module

Is transfer complete

NO

Execute next instruction

Drawback of the Programmed I/O :

The main drawback of the Program Initiated I/O was that the CPU has to monitor the units all the times when the program is executing. Thus 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 and the CPU time is wasted a lot in keeping an eye to the executing of program.

To remove this problem an Interrupt facility and special commands are used.

**Interrupt-Initiated I/O :**

In this method an interrupt facility an interrupt command is used to inform the device about the start and end of transfer. In the meantime the CPU executes other program. When the interface determines that the device is ready for data transfer it generates an Interrupt Request and sends it to the computer.

When the CPU receives such an signal, it temporarily stops the execution of the program and branches to a service program to process the I/O transfer and after completing it returns back to task, what it was originally performing.

- In this type of IO, computer does not check the flag. It continue to perform its task.

- Whenever any device wants the attention, it sends the interrupt signal to the CPU.

- CPU then deviates from what it was doing, store the return address from PC and branch to the address of the subroutine.

- There are two ways of choosing the branch address:

  - Vectored Interrupt

  - Non-vectored Interrupt

- In vectored interrupt the source that interrupt the CPU provides the branch information. This information is called interrupt vectored.

- In non-vectored interrupt, the branch address is assigned to the fixed address in the memory.

**Priority Interrupt:**

- There are number of IO devices attached to the computer.

- They are all capable of generating the interrupt.

- When the interrupt is generated from more than one device, priority interrupt system is used to determine which device is to be serviced first.

- Devices with high speed transfer are given higher priority and slow devices are given lower priority.

- Establishing the priority can be done in two ways:

  - Using Software

  - Using Hardware

- A pooling procedure is used to identify highest priority in software means.

**Polling Procedure :**

- There is one common branch address for all interrupts.

- Branch address contain the code that polls the interrupt sources in sequence. The highest priority is tested first.

- The particular service routine of the highest priority device is served.

- The disadvantage is that time required to poll them can exceed the time to serve them in large number of IO devices.
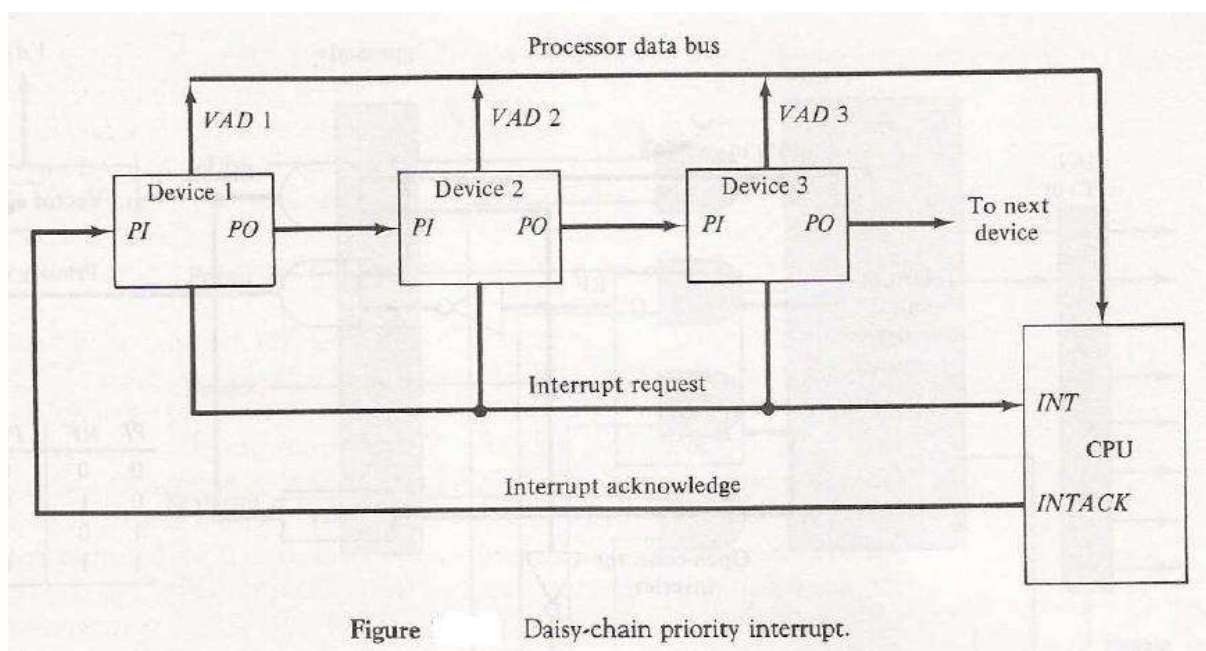
Using Hardware:

- Hardware priority system function as an overall manager.

**UNIT-V**

- It accepts interrupt request and determine the priorities.

- To speed up the operation each interrupting devices has its own interrupt vector.

- No polling is required, all decision are established by hardware priority interrupt unit.

- It can be established by serial or parallel connection of interrupt lines.

**Serial or Daisy Chaining Priority:**

- Device with highest priority is placed first.

- Device that wants the attention send the interrupt request to the CPU.

- CPU then sends the INTACK signal which is applied to PI(priority in) of the first device.

- If it had requested the attention, it place its VAD(vector address) on the bus. And it block the signal by placing 0 in PO(priority out)

- If not it pass the signal to next device through PO(priority out) by placing 1.

- This process is continued until appropriate device is found.

- The device whose PI is 1 and PO is 0 is the device that send the interrupt request.



Figure  Daisy-chain priority interrupt.

**Parallel Priority Interrupt :**

- It consist of interrupt register whose bits are set separately by the interrupting devices.

- Priority is established according to the position of the bits in the register.

**UNIT-V**

- Mask register is used to provide facility for the higher priority devices to interrupt when lower priority device is being serviced or disable all lower priority devices when higher is being serviced.

- Corresponding interrupt bit and mask bit are ANDed and applied to priority encoder.

- Priority encoder generates two bits of vector address.

- Another output from it sets IST(interrupt status flip flop).



## Priority Encoder Truth Table

| Inputs | | | | Outputs | | | Boolean functions |
|---|---|---|---|---|---|---|---|
| $I_0$ | $I_1$ | $I_2$ | $I_3$ | $x$ | $y$ | IST | |
| 1 | × | × | × | 0 | 0 | 1 | |
| 0 | 1 | × | × | 0 | 1 | 1 | $x = I_0'I_1'$ |
| 0 | 0 | 1 | × | 1 | 0 | 1 | $y = I_0'I_1 + I_0'I_2'$ |
| 0 | 0 | 0 | 1 | 1 | 1 | 1 | $(IST) = I_0 + I_1 + I_2 + I_3$ |
| 0 | 0 | 0 | 0 | × | × | 0 | |

The Execution process of Interrupt–Initiated I/O is represented in the flowchart:

**Interrupt –Initiated I/O**

```
CPU issues the read or write
command to I/O module
        ↓
I/O module informs about its
status to CPU
        ↓
    Status ──────────────→ Error
        ↓
CPU reads word from I/O module & writes it to memory or
CPU reads word from memory & writes it to I/O module
        ↓
      Is
    transfer
    complete
     ↓ yes
Execute next instruction
```

NO

**Direct Memory Access (DMA):**

In the Direct Memory Access (DMA) the interface transfer the data into and out of the memory unit through the memory bus. The transfer of data between a fast storage device such as magnetic disk and memory is often limited by the speed of the CPU. Removing the CPU from the path and letting the peripheral device manage the memory buses directly would improve the speed of transfer. This transfer technique is called Direct Memory Access (DMA).

During the DMA transfer, the CPU is idle and has no control of the memory buses. A DMA Controller takes over the buses to manage the transfer directly between the I/O device and memory.

The CPU may be placed in an idle state in a variety of ways. One common method extensively used in microprocessor is to disable the buses through special control signals such as:

- Bus Request (BR)

- Bus Grant (BG)

These two control signals in the CPU that facilitates the DMA transfer. The *Bus Request (BR)* input is used by the *DMA controller* to request the CPU. When this input is active, the CPU terminates the execution of the current instruction and places the address bus, data bus

and read write lines into a *high Impedance state.* High Impedance state means that the output is disconnected.



**CPU bus Signals for DMA Transfer**

The CPU activates the *Bus Grant (BG)* output to inform the external DMA that the Bus Request (BR) can now take control of the buses to conduct memory transfer without processor.

When the DMA terminates the transfer, it disables the *Bus Request (BR)* line. The CPU disables the *Bus Grant (BG)*, takes control of the buses and return to its normal operation.

The transfer can be made in several ways that are:

<div align="center">

i. DMA Burst

ii. Cycle Stealing

</div>

   i)    DMA Burst :- In DMA Burst transfer, a block sequence  consisting of a number of memory words is transferred in continuous burst while the DMA controller is master of the memory buses.

  ii)    Cycle Stealing :- Cycle stealing allows the DMA controller to transfer one data word at a time, after which it must returns control of the buses to the CPU.

DMA Controller:

The DMA controller needs the usual circuits of an interface to communicate with the CPU and I/O device. The DMA controller has three registers:

<div align="center">

i.  Address Register

ii. Word Count Register
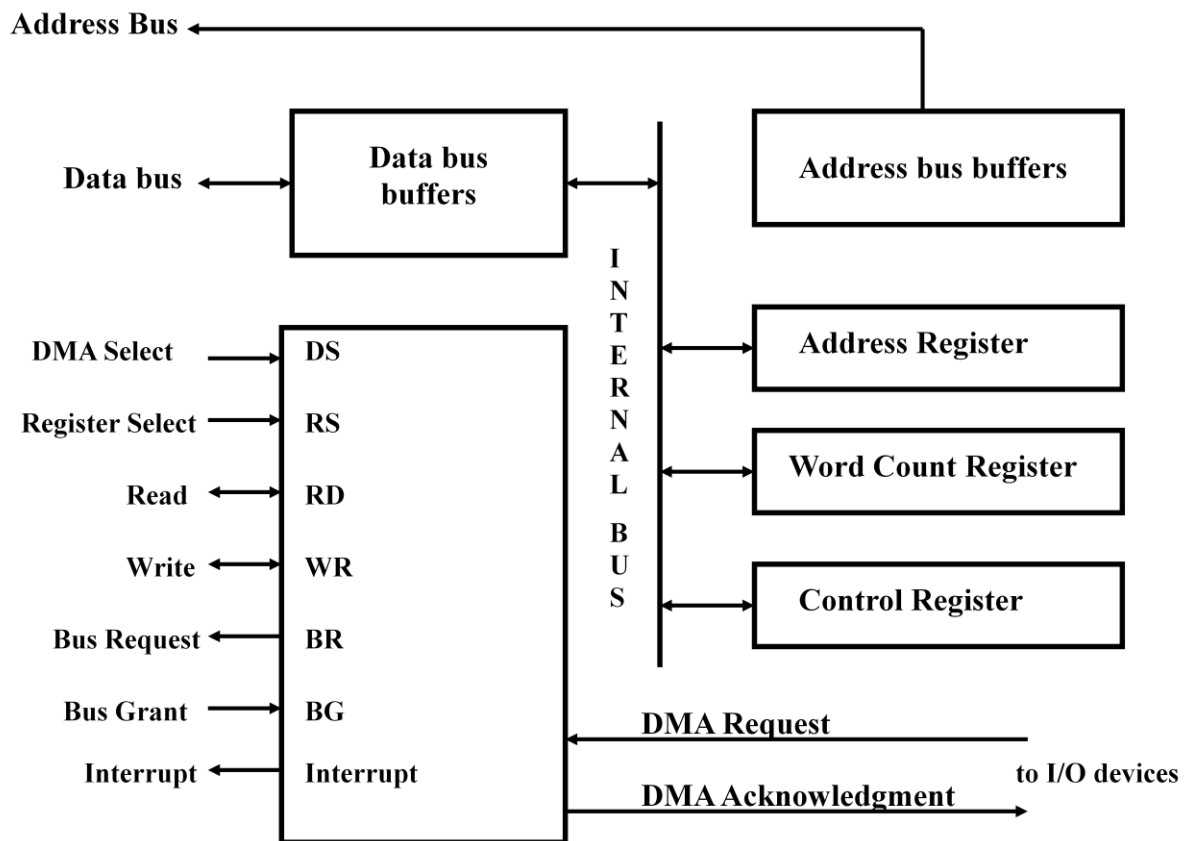
iii. Control Register

</div>

i. Address Register :- Address Register contains an address to specify the desired location in memory.

ii. Word Count Register :- WC holds the number of words to be transferred. The register is incre/decre by one after each word transfer and internally tested for zero.

i. Control Register :- Control Register specifies the mode of transfer

The unit communicates with the CPU via the data bus and control lines. The registers in the DMA are selected by the CPU through the address bus by enabling the DS (DMA select) and RS (Register select) inputs. The RD (read) and WR (write) inputs are bidirectional.

When the BG (Bus Grant) input is 0, the CPU can communicate with the DMA registers through the data bus to read from or write to the DMA registers. When BG =1, the DMA can communicate directly with the memory by specifying an address in the address bus and activating the RD or WR control.



**Block Diagram of DMA Controller**

DMA Transfer:

The CPU communicates with the DMA through the address and data buses as with any interface unit. The DMA has its own address, which activates the DS and RS lines. The CPU initializes the DMA through the data bus. Once the DMA receives the start control command, it can transfer between the peripheral and the memory.
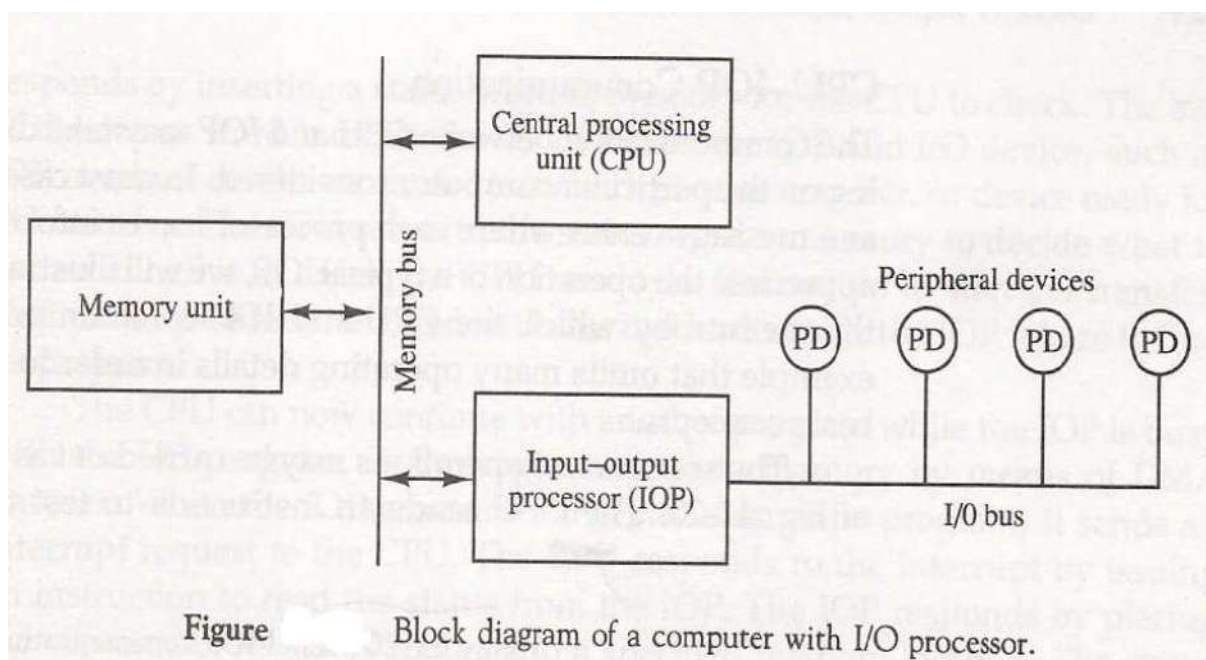
When BG = 0 the RD and WR are input lines allowing the CPU to communicate with the internal DMA registers. When BG=1, the RD and WR are output lines from the DMA controller to the random access memory to specify the read or write operation of data.

Summary :
- Interface is the point where a connection is made between two different parts of a system.
- The strobe control method of Asynchronous data transfer employs a single control line to time each transfer.
- The handshaking method solves the problem of strobe method by introducing a second control signal that provides a reply to the unit that initiates the transfer.
- Programmed I/O mode of data transfer the operations are the results in I/O instructions which is a part of computer program.
- In the Interrupt Initiated I/O method an interrupt facility an interrupt command is used to inform the device about the start and end of transfer.
- In the Direct Memory Access (DMA) the interface transfer the data into and out of the memory unit through the memory bus.

## Input-Output Processor:

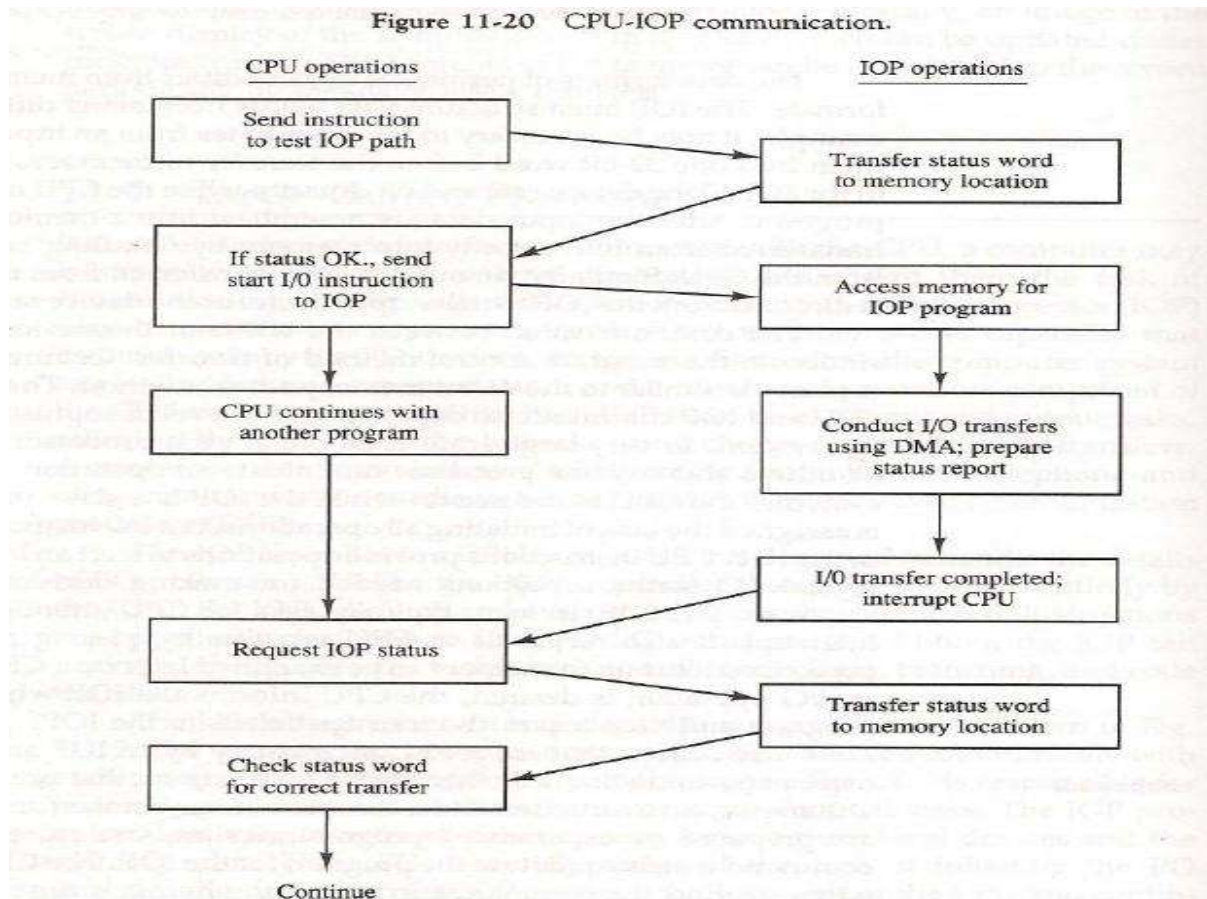- It is a processor with direct memory access capability that communicates with IO devices.

- IOP is similar to CPU except that it is designed to handle the details of IO operation.

- Unlike DMA which is initialized by CPU, IOP can fetch and execute its own instructions.

- IOP instruction are specially designed to handle IO operation.



Figure    Block diagram of a computer with I/O processor.

- Memory occupies the central position and can communicate with each processor by DMA.

- CPU is responsible for processing data.

- IOP provides the path for transfer of data between various peripheral devices and memory.

- Data formats of peripherals differ from CPU and memory. IOP maintain such problems.

- Data are transfer from IOP to memory by stealing one memory cycle.

- Instructions that are read from memory by IOP are called commands to distinguish them from instructions that are read by the CPU.

**Figure 11-20    CPU-IOP communication.**

CPU operations                                          IOP operations

```
  ┌──────────────────────┐
  │  Send instruction    │────────┐      ┌──────────────────────┐
  │  to test IOP path    │        │      │  Transfer status word │
  └──────────────────────┘        └─────▶│  to memory location   │
                                          └──────────────────────┘
  ┌──────────────────────┐
  │  If status OK., send │        ┌──────────────────────┐
  │  start I/0 instruction│──────▶│  Access memory for    │
  │  to IOP              │        │  IOP program          │
  └──────────────────────┘        └──────────────────────┘
  ┌──────────────────────┐        ┌──────────────────────┐
  │  CPU continues with  │        │  Conduct I/O transfers│
  │  another program     │        │  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
```

Instruction *that are read from memory by an IOP*

&raquo; Distinguish from instructions that are read by the CPU

&raquo; Commands are prepared by experienced programmers and are stored in memory
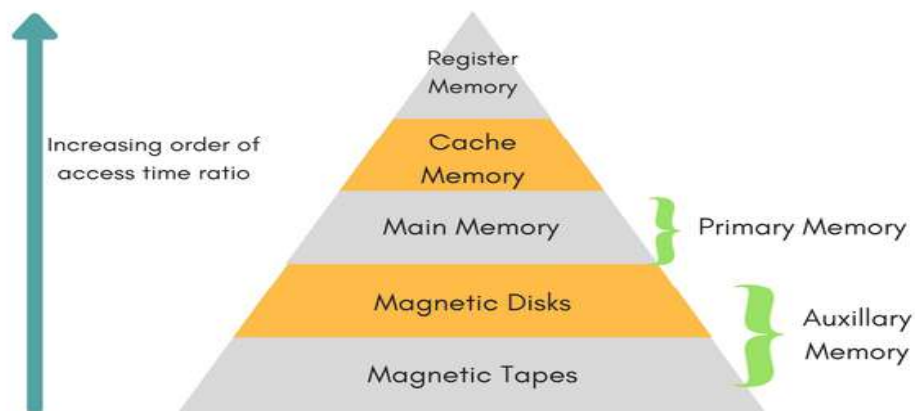
&raquo; Command word = IOP program

# Memory Organization

A memory unit is the collection of storage units or devices together. The memory unit stores the binary information in the form of bits. Generally, memory/storage is classified into 2 categories:

**Volatile Memory**: This loses its data, when power is switched off.

- **Non-Volatile Memory**: This is a permanent storage and does not lose any data when power is switched off.
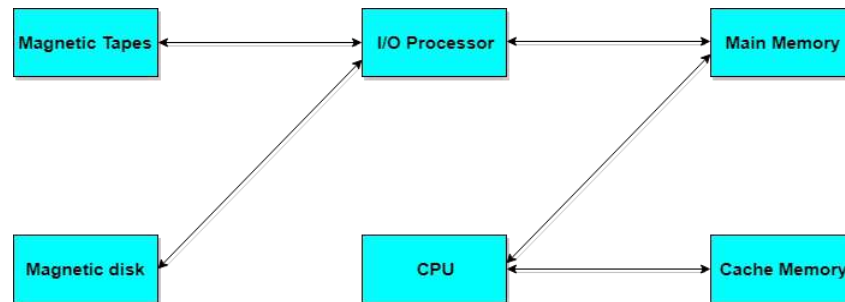
**Memory Hierarchy**



The total memory capacity of a computer can be visualized by hierarchy of components. The memory hierarchy system consists of all storage devices contained in a computer system from the slow Auxiliary Memory to fast Main Memory and to smaller Cache memory.

**Auxillary memory** access time is generally **1000 times** that of the main memory, hence it is at the bottom of the hierarchy.

The **main memory** occupies the central position because it is equipped to communicate directly with the CPU and with auxiliary memory devices through Input/output processor (I/O).

When the program not residing in main memory is 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 in main memory for other programs that are currently in use.

The **cache memory** is used to store program data which is currently being executed in the CPU. Approximate access time ratio between cache memory and main memory is about **1 to 7~10**



**Memory Access Methods**

Each memory type, is a collection of numerous memory locations. To access data from any memory, first it must be located and then the data is read from the memory location. Following are the methods to access information from memory locations:

1. **Random Access**: Main memories are random access memories, in which each memory location has a unique address. Using this unique address any memory location can be reached in the same amount of time in any order.

2. **Sequential Access**: This methods allows memory access in a sequence or in order.

3. **Direct Access**: In this mode, information is stored in tracks, with each track having a separate read/write head.

**Main Memory**

The memory unit that communicates directly within the CPU, Auxillary memory and Cache memory, is called main memory. It is the central storage unit of the computer system. It is a large and fast memory used to store data during computer operations. Main memory is made up of **RAM** and **ROM**, with RAM integrated circuit chips holing the major share.

- RAM: Random Access Memory

- o **DRAM**: Dynamic RAM, is made of capacitors and transistors, and must be refreshed every 10~100 ms. It is slower and cheaper than SRAM.

- o **SRAM**: Static RAM, has a six transistor circuit in each cell and retains data, until powered off.

- o **NVRAM**: Non-Volatile RAM, retains its data, even when turned off. Example: Flash memory.

- ROM: Read Only Memory, is non-volatile and is more like a permanent storage for information. It also stores the **bootstrap loader** program, to load and start the operating system when computer is turned on. **PROM**(Programmable ROM), **EPROM**(Erasable PROM) and **EEPROM**(Electrically Erasable PROM) are some commonly used ROMs.

**Auxiliary Memory**

Devices that provide backup storage are called auxiliary memory. **For example:** Magnetic disks and tapes are commonly used auxiliary devices. Other devices used as auxiliary memory are magnetic drums, magnetic bubble memory and optical disks.

It is not directly accessible to the CPU, and is accessed using the Input/Output channels.

**Cache Memory**

The data or contents of the main memory that are used again and again by CPU, are stored in the cache memory so that we can easily access that data in shorter time.

Whenever the CPU needs to access memory, it first checks the cache memory. If the data is not found in cache memory then the CPU moves onto the main memory. It also transfers block of recent data into the cache and keeps on deleting the old data in cache to accomodate the new one.

**Hit Ratio**

The performance of cache memory is 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 then it counts as a **miss**.

The ratio of the number of hits to the total CPU references to memory is called hit ratio.

Hit Ratio = Hit/(Hit + Miss)
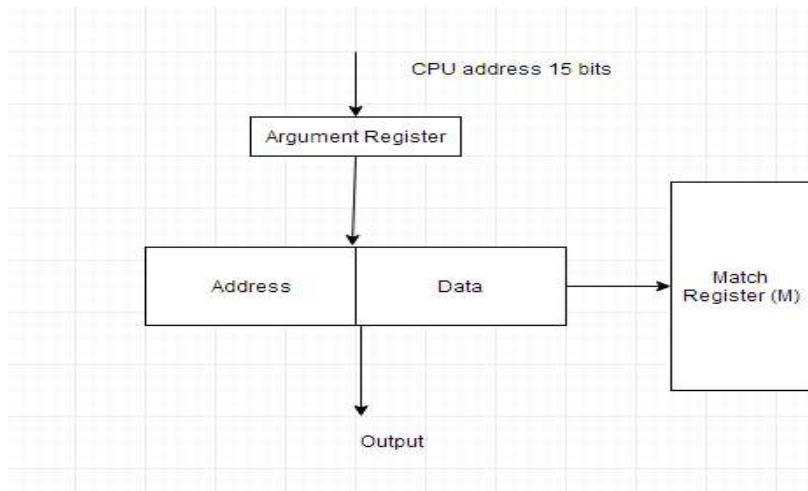
**Associative Memory**

It is also known as **content addressable memory (CAM)**. It is a memory chip in which each bit position can be compared. In this the content is compared in each bit cell which allows very fast table lookup. Since the entire chip can be compared, contents are randomly stored without considering addressing scheme. These chips have less storage capacity than regular memory chips.

The transformation of data from main memory to cache memory is called mapping. There are 3 main types of mapping:

- Associative Mapping
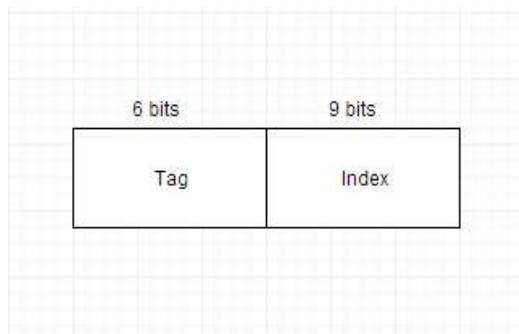
- Direct Mapping

- Set Associative Mapping

**Associative Mapping**

The associative memory stores both address and data. The address value of 15 bits is 5 digit octal numbers and data is of 12 bits word in 4 digit octal number. A CPU address of 15 bits is placed in **argument register** and the associative memory is searched for matching address.

**Direct Mapping**

The CPU address of 15 bits is divided into 2 fields. In this the 9 least significant bits constitute the **index** field and the remaining 6 bits constitute the **tag** field. The number of bits in index field is equal to the number of address bits required to access cache memory.



**Set Associative Mapping**

The disadvantage of direct mapping is that two words with same index address can't reside in cache memory at the same time. This problem can be overcome by set associative mapping.

In this we can store two or more words of memory under the same index address. Each data word is stored together with its tag and this forms a set.

**Replacement Algorithms**

Data is continuously replaced with new data in the cache memory using replacement algorithms. Following are the 2 replacement algorithms used:

- FIFO - First in First out. Oldest item is replaced with the latest item.

- LRU - Least Recently Used. Item which is least recently used by CPU is removed.

**Virtual Memory**

Virtual memory is the separation of logical memory from physical memory. This separation provides large virtual memory for programmers when only small physical memory is available.

Virtual memory is used to give programmers the illusion that they have a very large memory even though the computer has a small main memory. It makes the task of programming easier because the programmer no longer needs to worry about the amount of physical memory available.