

→ The printing character consists of 26 - upper cases (A - Z)

26 - lower case letters

10 - numericals (0 - 9)

32 - special printable characters like %, *, \$... etc.

(i) Control functions:-

NUL - NULL

SOH - Start of heading

STX - Start of text

ETX - End of text

EOT - End of transmission

ENQ - Enquiry

ACK - Acknowledge

BEL - Bell

BS - Back space

HT - Horizontal Tab

LF - Line feed

VT - Vertical tab

FF - Form feed

CR - Carriage return

SO - Shift out

SI - Shift in

SP - Space

format
control
(Layout
printing)

DLE - Data Link Escape

DC1 - Device Control 1

DC2 - " " 2

DC3 - " " 3

DC4 - " " 4

NAK - Negative Acknowledge

SYN - Synchronization

ETB - End of transmission block

CAN - Cancel

EM - End of medium

SUB - Substitute

ESC - Escape

Information separator { FS - file separator

{ RS - Record separator

{ GS - Group separator

{ US - Unit separator

DEL - Delete.

There are 3-types of Control characters.

(a) Format control :- are characters that control the layout printing

which include initial type write controls like

BS (back space), HT ; CR ;

(b) Information separator :- are used to separate the data into division

like paragraphs and pages which include

RS (Record separator) ; file separator (FS)

(c) Communication Control characters :- are useful during tx of text. They denote

terminals

e.g. - STX ; ETX which are used to frame a text message.

3

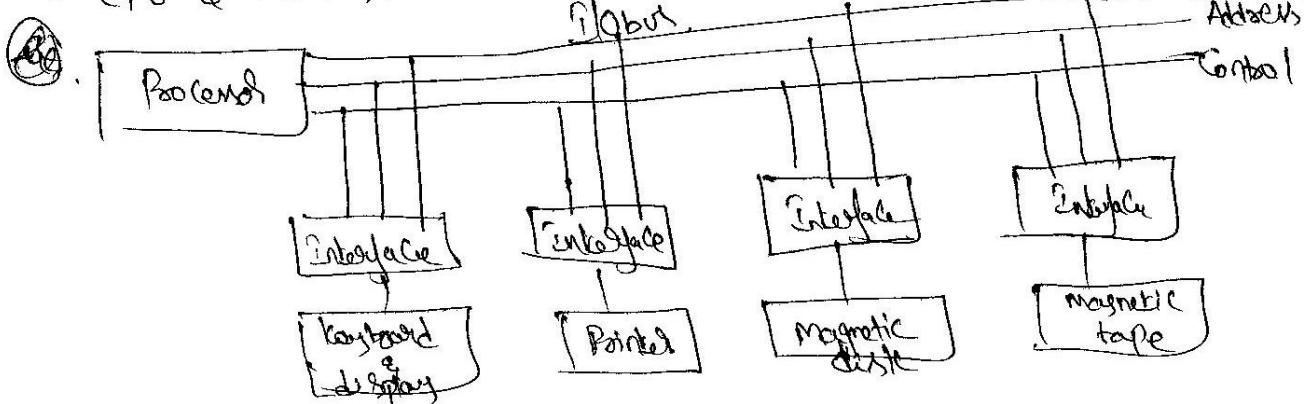
⑥ Input-Output Interface:- is a method of transferring information.
How Internal storage & external I/O devices.

→ Major Difference b/w Peripheral devices & CPU is.

① Peripheral are electro-mechanical & electro-magnetic devices and
CPU & memory are electronic devices which have major
difference from operation.

② The Data transfer rate of peripheral is usually slower than transfer rate of CPU.

which needs Synchronization mechanism.
③ Data codes and formats in peripherals differ from word formats in
CPU & memory.



There are 4-types of Commands

① Control ② Status

is issued to activate
Peripheral.

is used to test
various status
conditions

③ Data output

↓
respond by
transferring data from
bus into one of
its registers

④ Data input



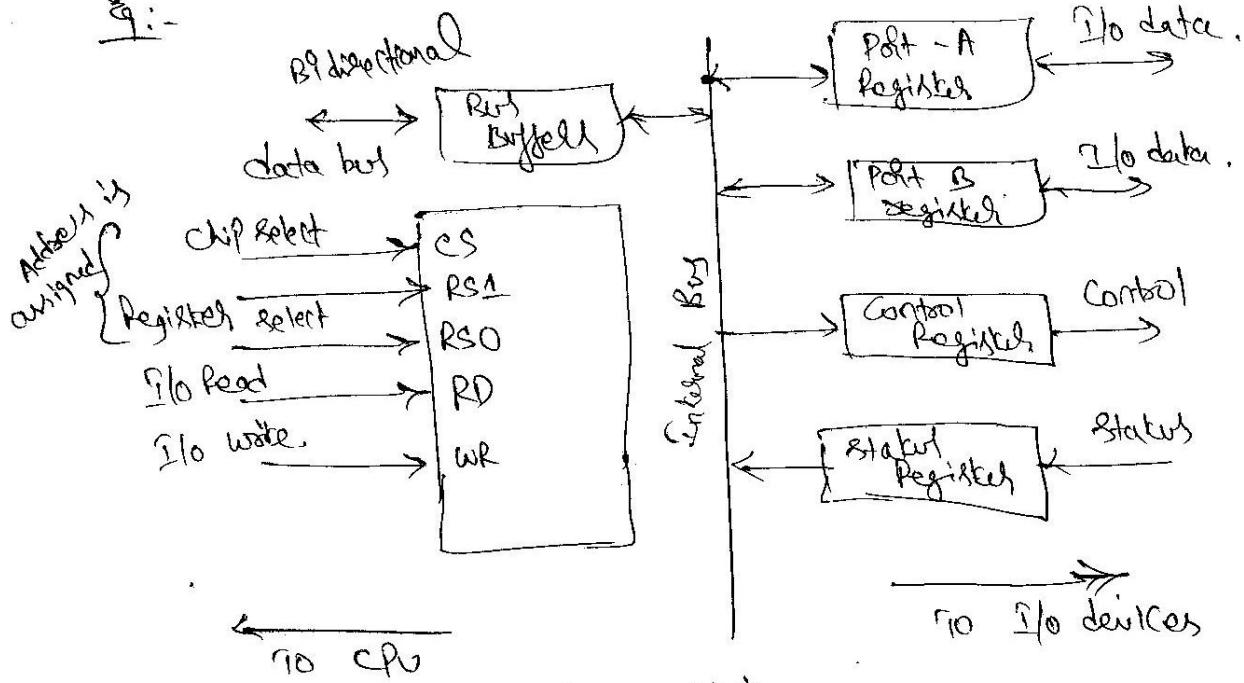
→ 3-Ways that Computer buses used to communicate memory I/O.

- ① use two separate buses one for memory & other for I/O
- ② use ~~one~~ common bus for both I/O but have separate control lines for each.
- ③ use one common bus for Memory & I/O with Common Control lines.

③ Isolated I/O memory mapping I/o:-

- Isolated I/o is method for assigning address in common bus
- the memory mapped I/o is to use the same address space for both memory and I/O.
 - The I/O read & write & enables control lines during I/o transfer.
 - memory read & memory write enables " " " I/O "

E:-



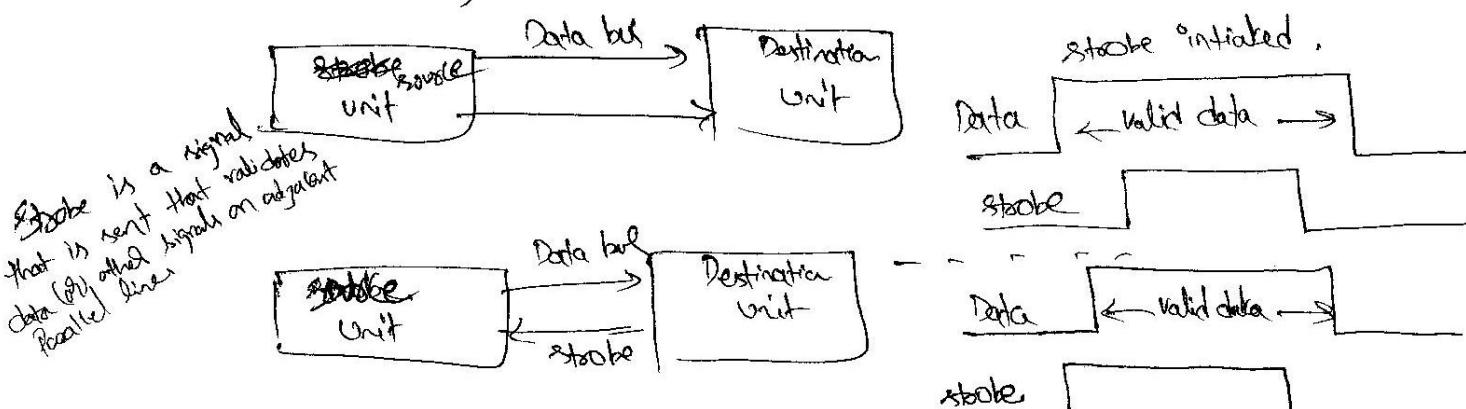
chip select CS	Register Select		Port Selected more data bus in high preference
	RS1	RS0	
0	X	X	Port - A Registered
1	0	0	Port - B "
1	0	1	Control Register
1	1	0	
1	1	1	Status Register.

a) Asynchronous Data transfer :-

" " b/w two independent units require control signals to be transmitted via the communicating units to indicate the time at which being transmitted.

→ The strobe control method of Asynchronous data transfer employs a single control line to time each transfer.

Strobe may be activated either source (or) destination unit

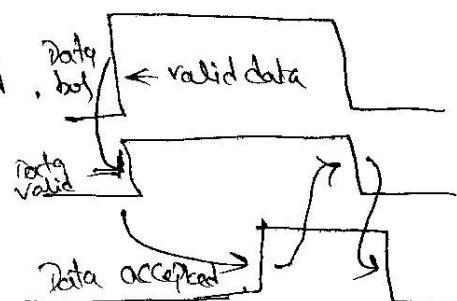
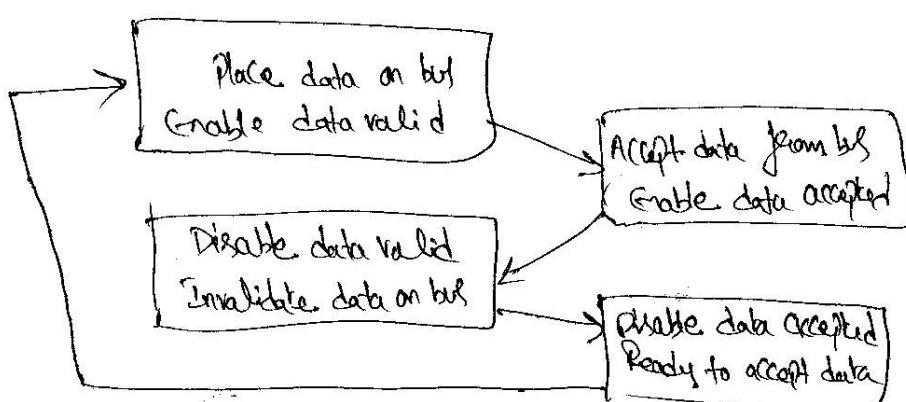


b) Handshaking :-

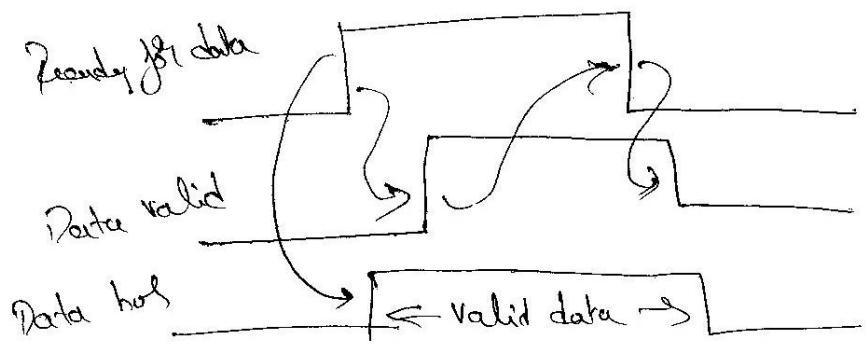
The dis-advantage of strobe method is source / destination unit initiates transfer has no way of knowing whether data is placed on bus or not.

The handshaking method solve by using

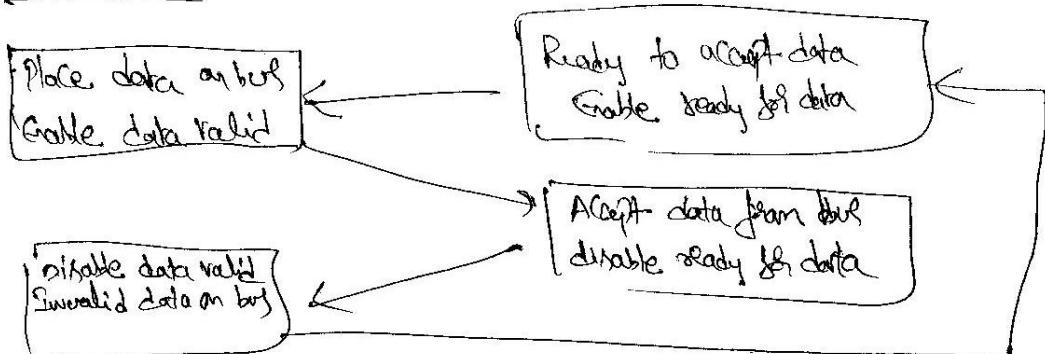
(a) Data valid ; (b) Data accepted



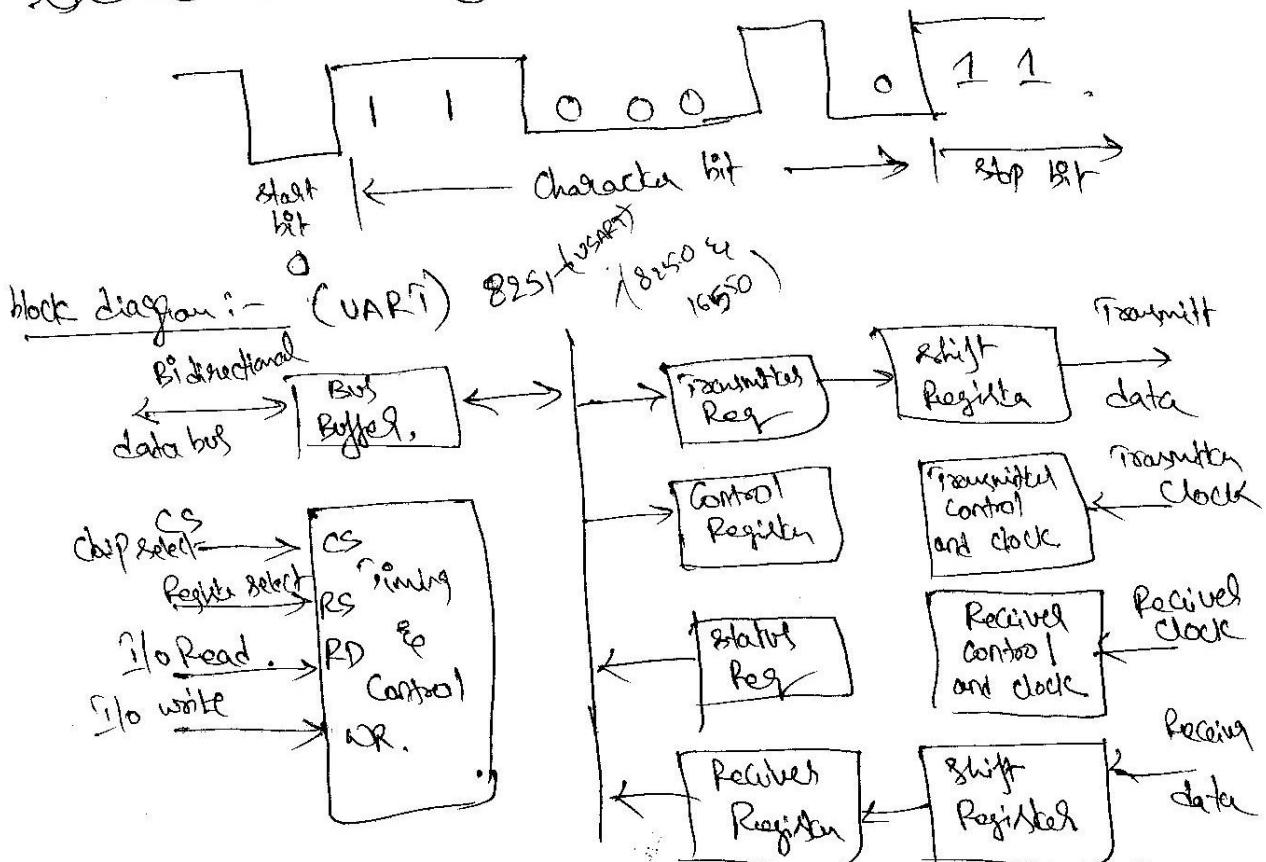
Source initiated transfer using handshaking



Source unit Destination unit



(c) Asynchronous Serial Transfer:-



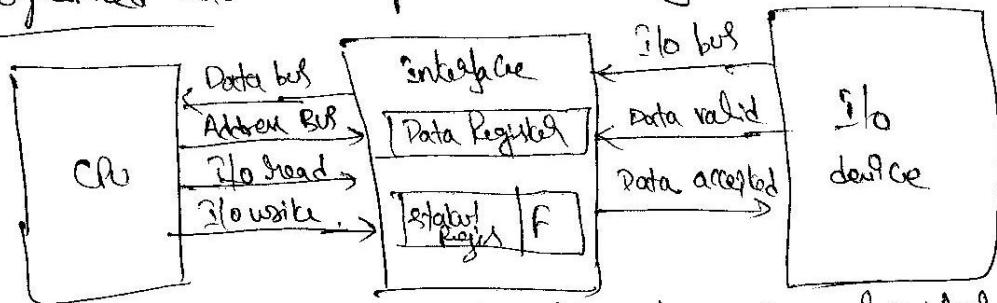
<u>CS</u>	<u>RS</u>	<u>operation</u>	<u>Register Selected</u>
0	X	X	none data bus in high impedance
1	0	WR	Transmitter Register
1	1	WR	Control Register
1	0	RD	Received Register
1	1	RD	Status Register

3) Model of I/O:-

Data transfer to & from peripherals may be handled in 3 ways.

- ① Programmed I/O
- ② Interrupt-initiated I/O
- ③ Direct Memory Access (DMA)

① Programmed I/O :- Requires execution of several instructions by the CPU.

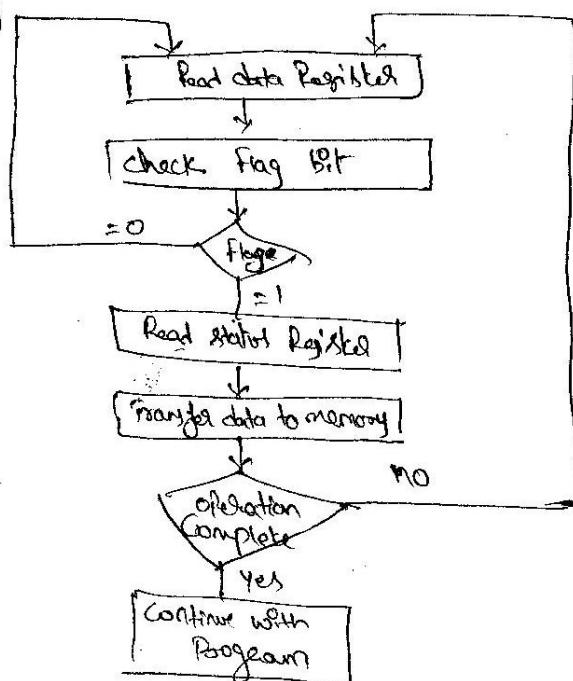


→ If flag = 1 the CPU reads data from data register otherwise it cleared

→ The transfer of each byte requires following steps

- a) Read the status register
- b) Check status of flag bit and branch to step 1 if not set (if) to Step -3
- c) Read the data register

Flowchart



② Interrupt Initiated I/O:- CPU responds to interrupt signal. by setting the return address from the program counter in to memory ~~data~~, stack.

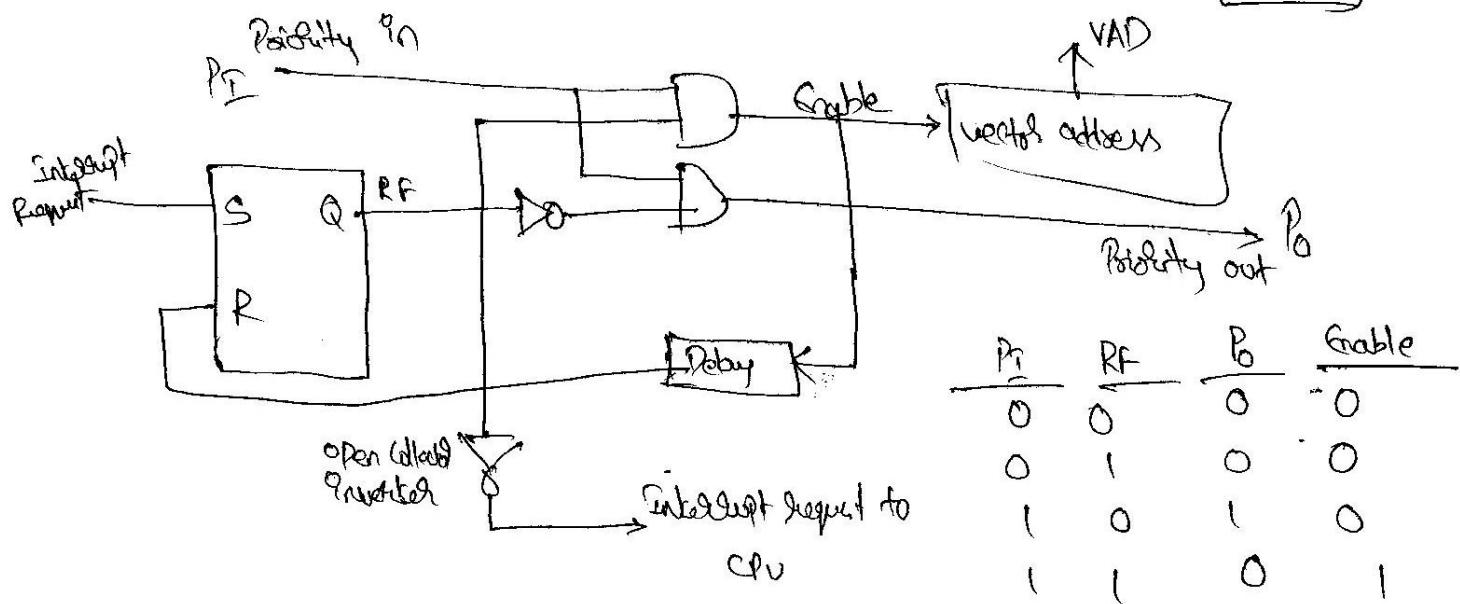
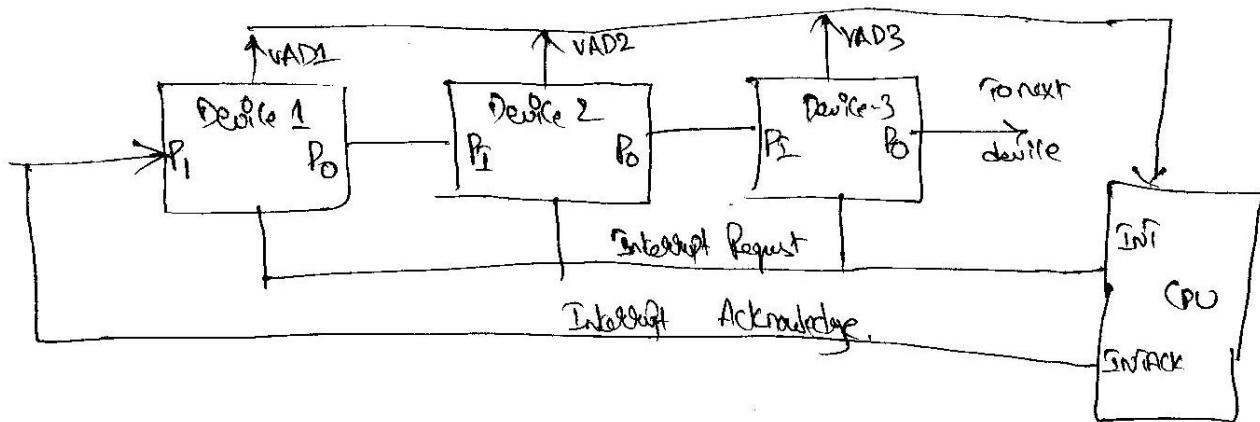
4) Priority Interrupt:- is a system that establishes a Priority level which sources to determine which condition is to be serviced first. when two request arrives simultaneously.

(a) more request arrives simultaneously.
→ The High Priority interrupt levels are assigned to request which if delayed will be interrupted.

→ Devices with high speed transfer such as magnetic disks are given high Priority
→ " " " slow " " " keyboards receive low "

(a) Daisy Chaining Priority:-

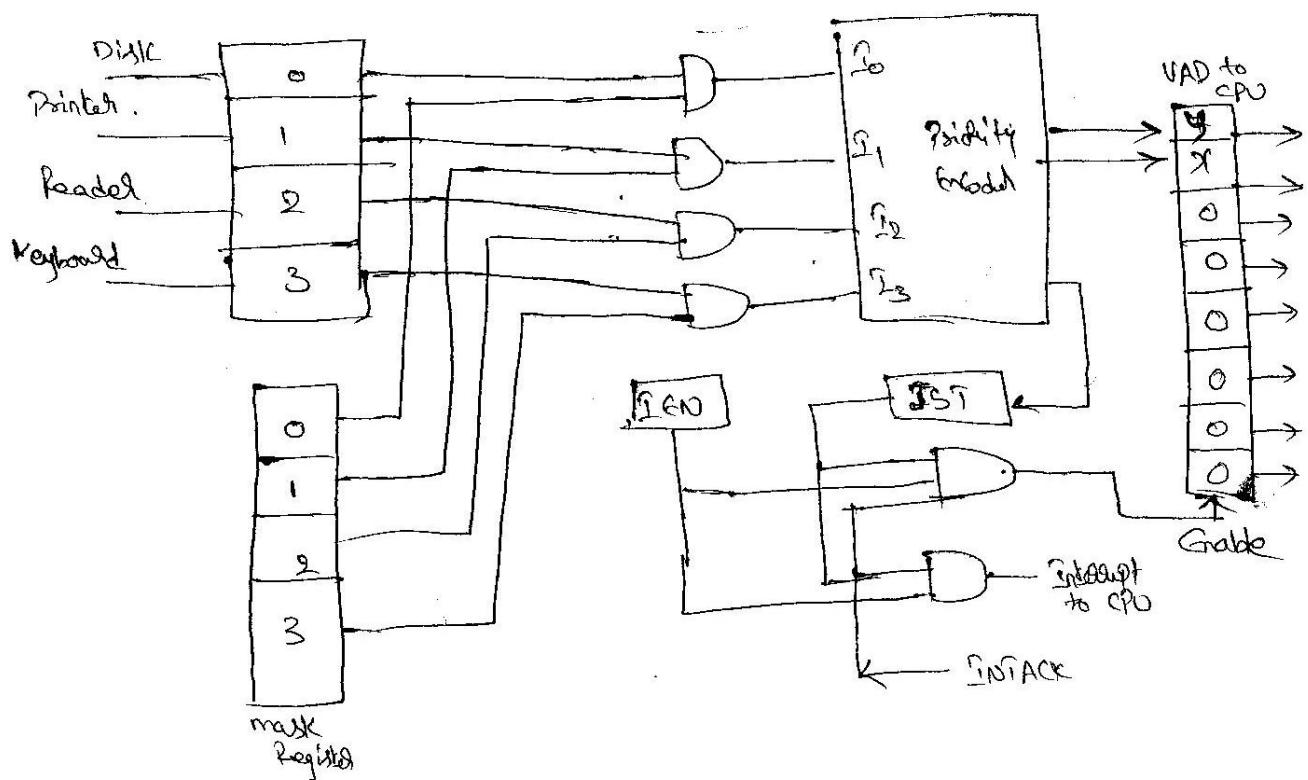
→ Method consists of serial connection of all devices that support interrupt
→ The device with high priority is placed in first position followed by lower priority.
→ Each device has its own interrupt vector (VAD) into data bus



9

~~Parallel Priority interrupt:-~~

The parallel priority interrupt method uses a register whose bits are set separately by the interrupt signal from each device. Priority is established according to the position of the bits in register. In addition to interrupt register, the circuit may include a mask register whose purpose is to control the status of each interrupt request. The mask register can be programmed to disable lower priority interrupts while a higher priority device is being serviced.

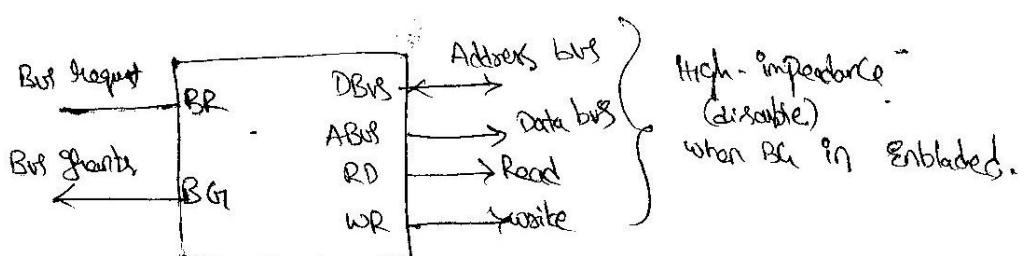


~~Direct Memory Access (DMA):-~~

The transfer of data b/w a fast storage device such as magnetic disk and memory often limited by the speed of the CPU, removing the CPU from the path and letting the peripheral device manage the memory busses directly. This transfer technique is called direct memory access (DMA).

Data formats of peripheral devices differ from memory and CPU data formats.

This transfer technique is called direct memory access (DMA).



- The direct memory Access (DMA) I/O technique provides direct access to memory while the CPU is temporarily disable.
- I/O device is connected to system bus via a special interface circuit known as DMA Controller.
- In DMA both CPU and DMA Controller have access to main memory via a shared system bus having data, address and control lines.
- A DMA Controller temporally borrows the address bus, data bus, and control bus from microprocessor and transfers the data bytes directly between an I/O port and series of memory locations.
- DMA transfer is also used to do high speed memory to memory transfers.

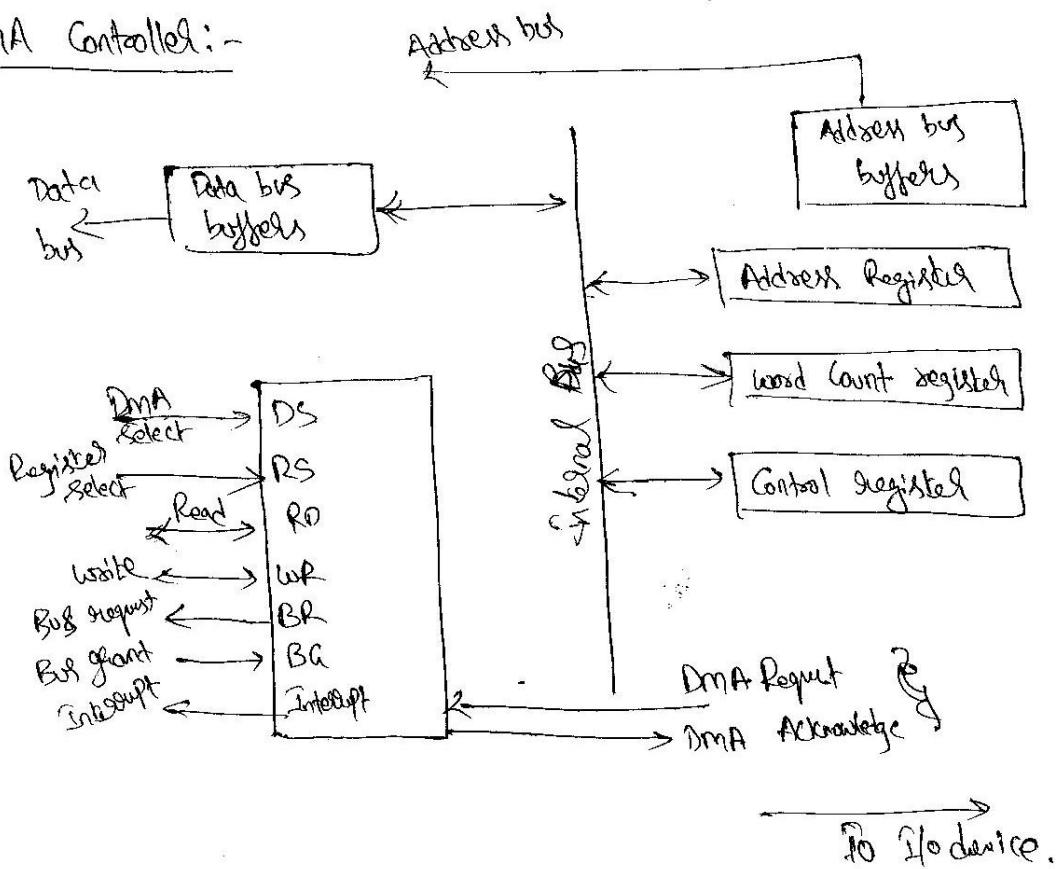
With out DMA:

→ using Programmed I/O it is typically fully occupied for entire duration of read (or) write operations and thus unavailable to perform other work.

With DMA

→ The CPU initiates the transfer, does other operations while the transfer is in progress and receives an interrupt from DMA controller when operation is done.

DMA Controller:-



4-channels
Specify
16-bit Address
and
(Data Reg.) 8 bit

The (BR) Bus Request bit is used by DMA Controller to request the CPU to relinquish control of the bus.

The transfer can be made in several ways

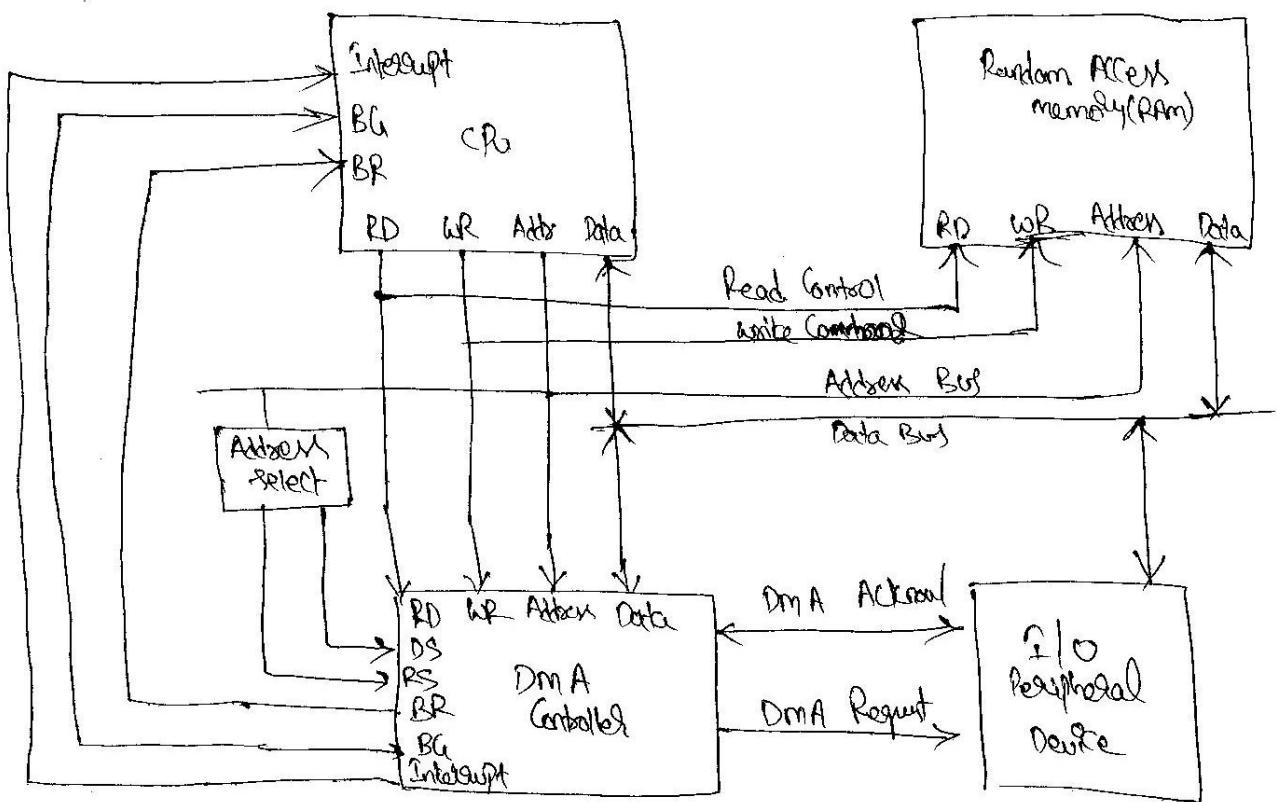
1. Burst transfer : allows DMA Controller to transfer entire block at a time.
2. Cycle Stealing:- allows " " " " one data at a time

DMA Controlled:-

- * When (BG) Bus Grant is '0' the CPU can communicate with DMA Registers through the Data bus to read from (or) write to DMA Registers.
- * When BG is '1' the CPU has relinquished the bus and DMA can communicate directly with memory by specifying an address in the address bus and activating the RD (or) WR Control.
- * DMA Controller has 3 Registers (Address Register), (Word Count Register), (Control Register)
 - Address Register → Contains an address to specify the desired location in memory
 - " " is incremented after each register holds the no. of words and internally to be transferred
 - It is decremented by one after each word transferred and internally tested for zero.
- * CPU initialize the DMA by sending \$ information through data bus:-
 - (1) The starting address of memory block where data are available (or) where data are to be stored.
 - (2) The word Count, which is the no. of words in the memory.
 - (3) Control to specify the mode of transfer such as Read (or) write.
 - (4) A Control to start the DMA transfer.

DMA Transfer:-

(12)



The CPU communicates with the DMA through address and data buses as with interface unit.

- The DMA has its own address, which activates DS and RS lines
- When Peripheral device sends the DMA Request the DMA Controller Activates the BR line, informing the CPU relinquish the bus.
- When Peripheral device receives DMA Acknowledgment it puts a word in the Data bus (for write) or recall a word from the data bus (for Read).



Multiprocessors:-

Is an interconnection of 2-(d) more CPU with memory & I/O devices
means CPU (d) Input-output processor (IOP)

- A multiprocessor system with common shared memory is classified as shared memory (d) tightly coupled multiprocessor.

Multiprocessors:- (CPU & I/O Processor)

(13)

Inter Connection Structure:-

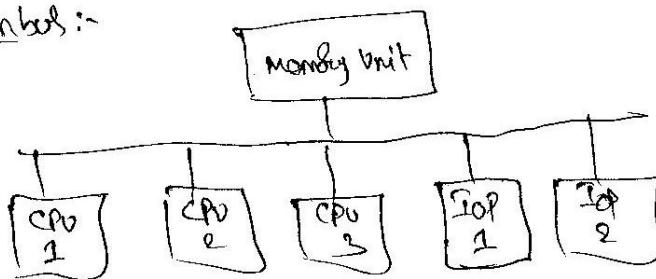
How the components can have different physical configuration depending on no. of transfer path that are available b/w processor and memory in shared memory system.

→ Several physical form available for establish an inter connection like

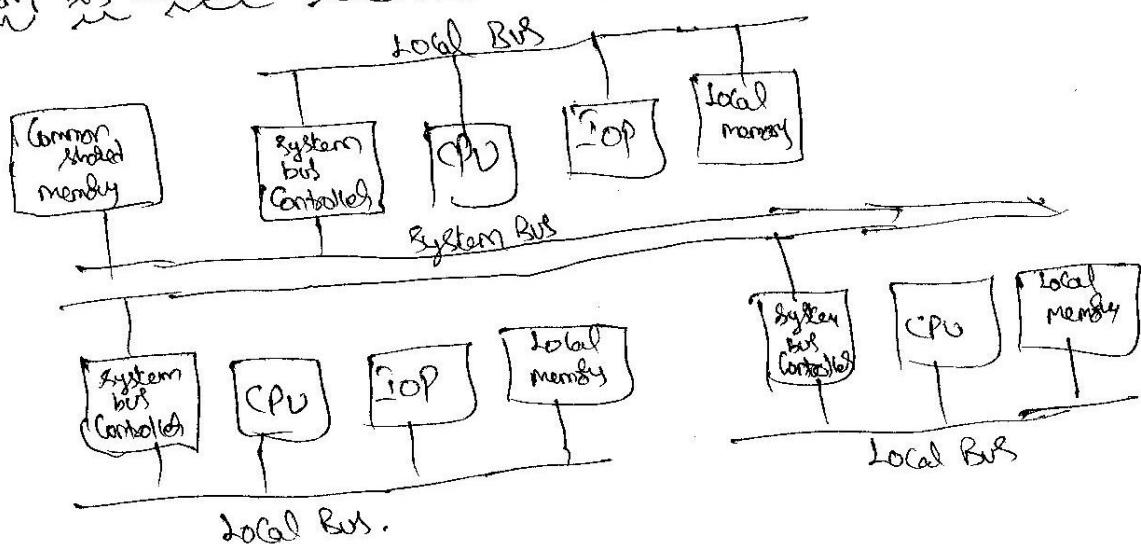
- ① Time Shared Common bus
- ② Multi Port memory
- ③ Cross bar Switch
- ④ Multi Stage switching
- ⑤ Hyper Cube System

① Time Shared Common bus:-

only one processor can communicate with memory at any time.



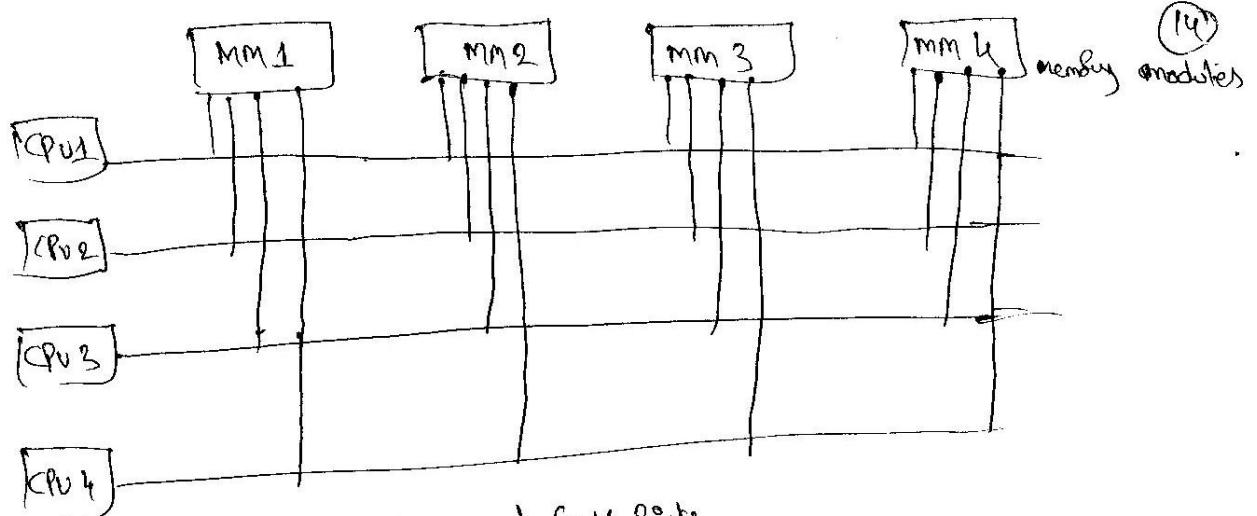
System bus Structure for multiprocessor:-



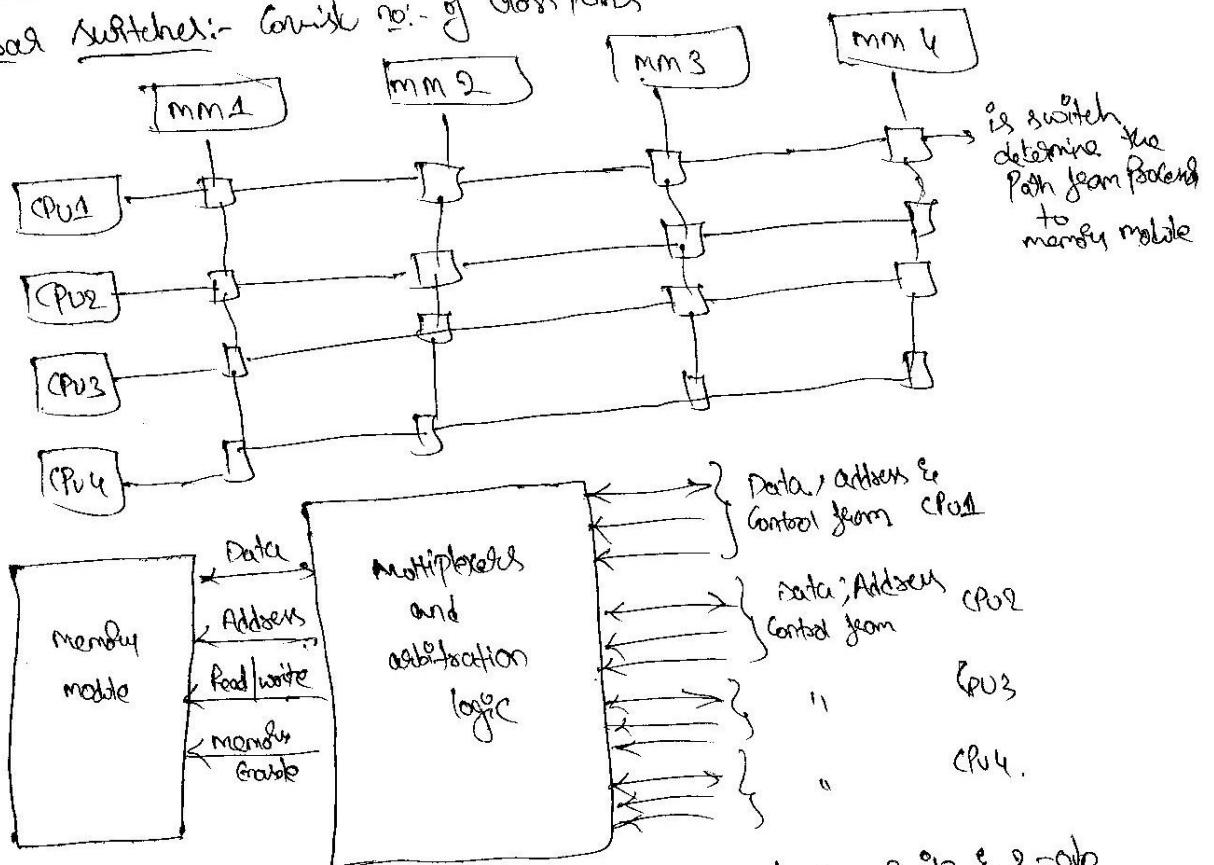
② Multi Port memory:-

Employed separate bus b/w each memory module and CPU.

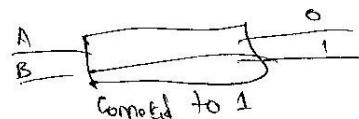
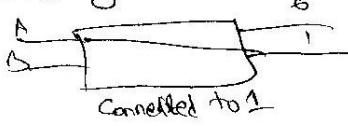
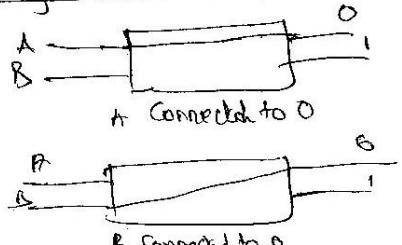
- A protocol bus consists of address, data and control lines required to communicate with memory.
- Advantage of multiport memory organization is high transfer rate is achieved.
- Disadvantage :- requires expensive memory control logic and large no. of cables.



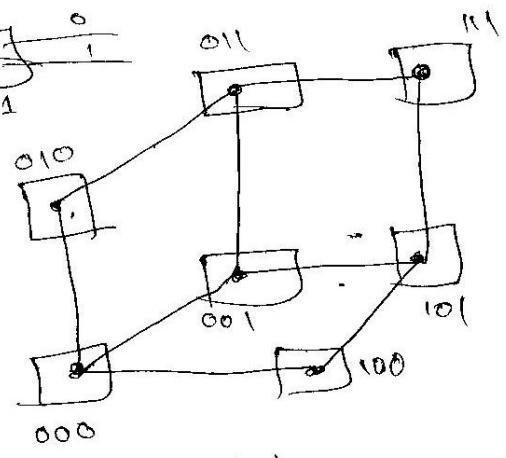
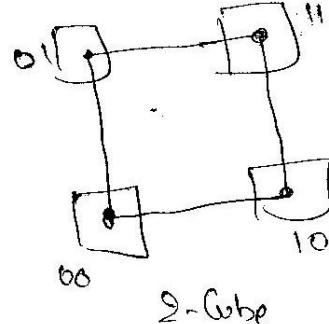
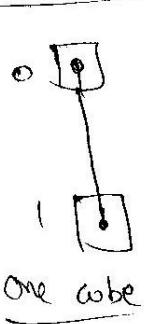
(2) Cross Bar Switched:- Count no. of Cross Points



(3) Multistage Switching:- Basic component of multistage nw is 2-2P & 2-0P

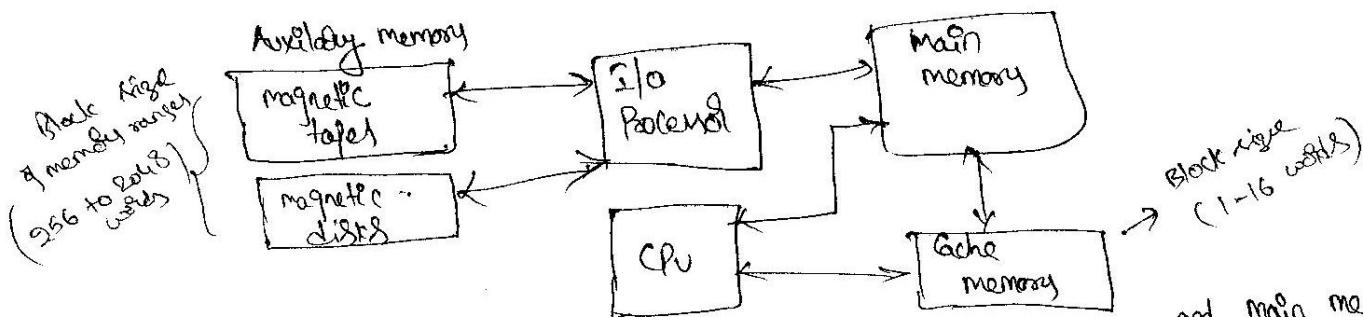


Hyper Cube Interconnection:-



Q) Memory Hierarchy:-

- The memory unit communicate directly to CPU is called Main memory.
 - Devices that provide Backup Storage are called Auxiliary Memory.
E.g.: magnetic disks & magnetic tapes. (used for storing system programs, large data files, and other backup information).
 - A special very high speed memory called Cache is sometimes used to increase the speed of processing by making current programs and data.



→ While I/O processor manages data transfer b/w Auxiliary memory and main memory, the Cache organization is concerned with the transfer of information b/w main memory and CPU.

→ The Auxiliary memory has large storage capacity but has low access speed. Compiled

to main " " small and has very high Access speed.

→ The Cache memory is very small and has 1 MB.

- No cache memory
- Typical Access time ratio b/w Cache & main memory is 1:100

E.g.: A typical Cache memory may have Access time 100ns
 main ~~Auxiliary~~ " " " " 1000ns
 Auxiliary " " " " 1000 times than main memory.

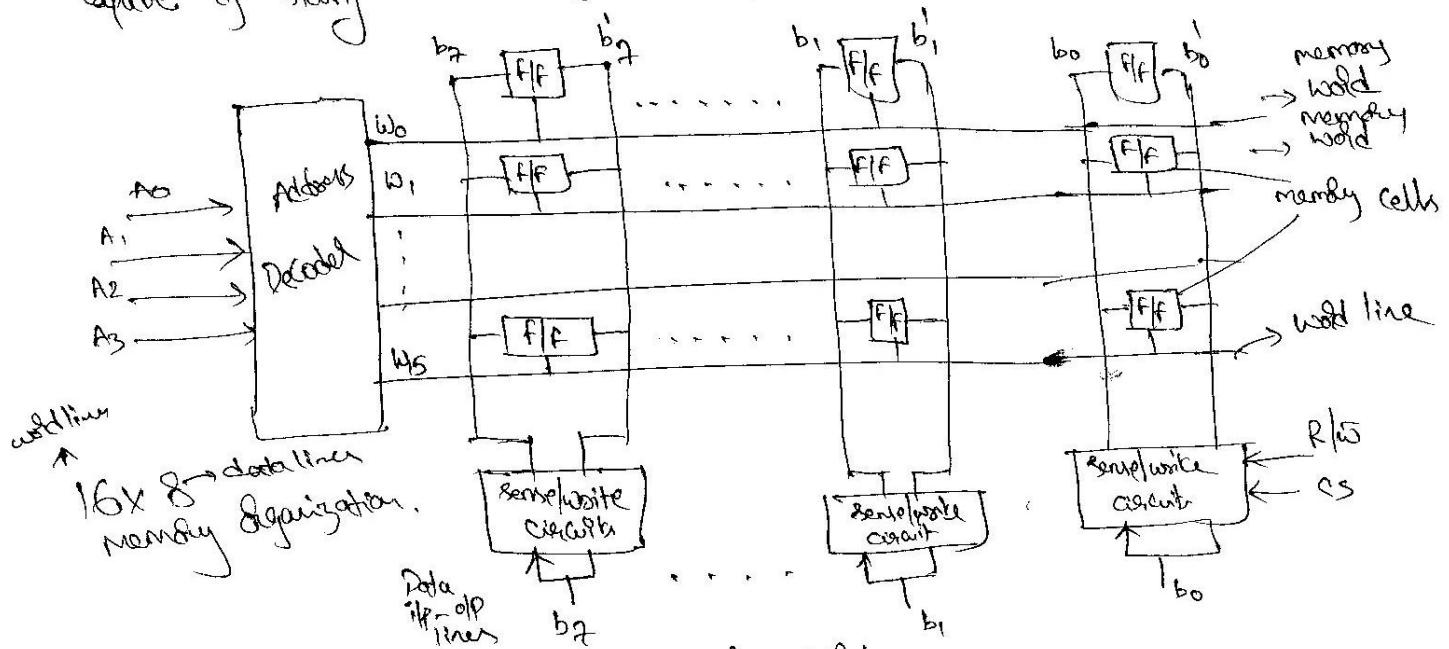
Ranges:- Block size in Auxiliary memory ranges 256 to 2048 words.
" Cache " 1 to 16 words.

→ The part of the computer system that supervised the flow of information between Auxiliary memory and Main memory is called Memory management system.

Main memory :- is central storage unit in computer system, which is relatively large and fast memory used to store programs and data during computer operation.

→ Internal organization of memory chip :-

memory cells are usually organized in form of array in which each cell is capable of storing 1-bit information.



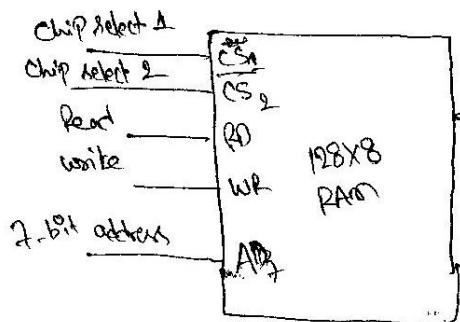
Each row of cells constitutes memory word.

- Read operation :- the information stored in the cells selected by a word line and transmits the information to op data line.
- Write operation :- the circuit recv/write, receives op information and store it in Cells of selected word.

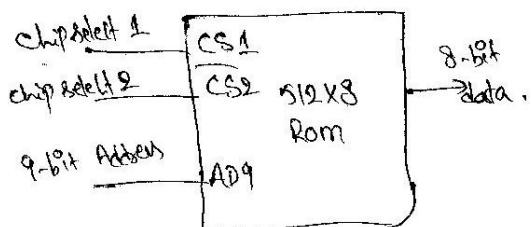
- Integrated circuit RAM chips are available in 2 operating models → static and dynamic.
- The static RAM consists essentially of internal flip-flops to store binary information. The stored information remains valid as long as power is applied to unit.
- The dynamic RAM stores binary information in form of electric charges that are applied to capacitors.
- The ROM portion of main memory is needed for storing an initial program called boot strap loaded.

RAM and ROM Switches:-

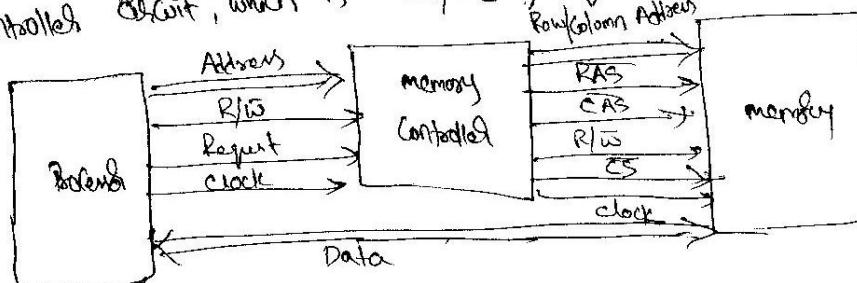
(17)



CS1	CS2	RD	WR	function	State of data bus
0	0	X	X	Inhibit	High impedance
0	1	X	X	Inhibit	High "
1	0	0	0	Inhibit	High "
1	0	0	1	write	Input data to RAM
1	0	1	X	Read	Output data from RAM
1	1	X	X	Inhibit	High impedance.



Memory Controller:- The required multiplexing of address bits is usually performed by a memory controller circuit, which is integrated b/w processor and the dynamic memory.



→ Virtual Memory:- is used to increase the apparent size of physical memory. → same time data may be stored in physical memory location that have address different from those specified by the program.

In such cases address generated by processor is referred as virtual (or) logical address.

→ the virtual address is mapped into physical memory where data are actually stored. mapping function is implemented by a special memory unit called memory management unit.

UNIT - 4

THE MEMORY SYSTEM

))

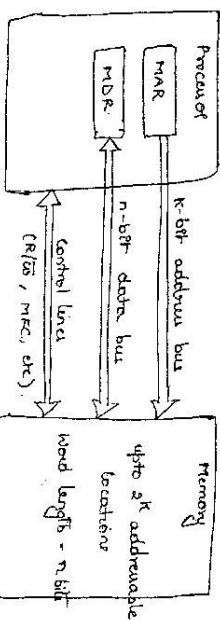
* Some basic concepts

- Maximum size of memory that can be used in any computer is defined by addressing scheme.

Ex- Computer that generates 16-bit address is capable of addressing upto 2^{16} = 64K (Kilo) memory locations.

16 bit address $\Rightarrow 2^16 = 64\text{K}$ (Kilo) memory locations

- connections between processor & memory is shown below:



Data transfer between processor and memory takes place through 2 processor registers, MAR (Memory Address Register) and MDR (Memory Data Register).

- if $\text{MAR} = \text{n-bit register}$ \Rightarrow K-bit address bus
- \Rightarrow 2^K addressable locations

- if $\text{MDR} = \text{n-bit register}$ \Rightarrow n-bit data bus

\Rightarrow word length = n-bit

\Rightarrow during a memory cycle, n-bits of data are transferred between memory and processor.

- Word length defines the no. of bits stored or retrieved in a memory access. Control lines include R/W (Read / write) and MFC (Memory Function control) for co-ordinating data transfer.
- Other control lines can be included to indicate indicate the no. of bytes to be transferred.

- Processor reads data from memory by loading the address of the required memory location into MAR register and setting R/W line to 1.
- Memory responds by placing data from the addressed location onto data lines, and confirms this action by outputting MFC signal.
- Upon receipt of MFC signal, the processor loads the data on the data lines into MDR register.

b) Write Operations

- Processor writes data into memory location by loading the address of the location into MAR and loading data into MDR.
- Set R/W = 0 to indicate a write operation.
- Word in MDR is written into memory location indicated by MAR.

c) Read or Write operations involve consecutive address locations in the memory. Thus block transfer operation can be performed, in which the 1st location address is only sent.

→ measure of speed of memory unit is defined by

- Memory Access time
- Memory cycle time

- Memory Access time
 - It is the time that elapses between the initiation of an operation and the completion of their operation.
 - Ex- time between Read and R/W signals.
- Memory cycle time
 - It is the minimum time delay required between 2 successive Read operations.

- a) How many chips are needed, and how should their addresses lines be connected to provide a memory capacity of 1024 bytes?

Sol: Memory capacity = $1024 \text{ bytes} = 1024 \times 8$.

Size of each RAM chip = 1024×1 .

\Rightarrow We requires 8 chips with their address lines connected in parallel.

- b) How many chips are needed to provide a memory capacity of 16K bytes? Explain in words how the chips are to be connected to the address bus?

Sol: memory capacity = $16K \text{ bytes} = 16 \times 1024 \times 8$.

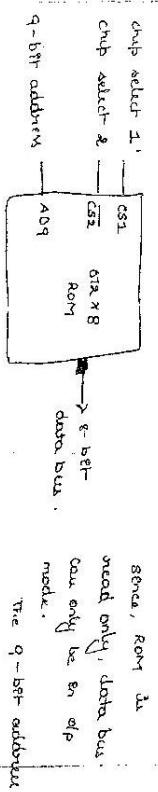
Given size of RAM = 1024×1 .

$$\text{No. of chips required} = \frac{16 \times 1024 \times 8}{1024 \times 1} = 16 \times 8 = 128 \text{ chips}$$

No. of address lines required = 14 ($2^4 \times 2^{10} = 16 \times 1024$).

14 address lines assigned to each chip
and 4 lines are decoded into 8 chip-select pins.

* ROM chips



Selecting only one of the 612 bytes stored in ROM

For ROM chip to operate, $CS1 = 1$ and $CS2 = 0$,

otherwise, data bus is in high-impedance state.

There is no need for word control because the can only read. \Rightarrow When chip is enabled, the byte selected by the address bus appears on the bus. (data bus).

- (P) A ROM chip of 1024×8 bit has 4 select pins and operates from a 5V power supply. How many pins are needed for IC package?

Sol: Draw a block diagram and label all 4 pin and chip manually on the ROM 1024×8 ROM as constructed from 8 ~~stacked~~ chips with 4 select pins and 5V power supply.

\Rightarrow No. of select pins = 4.

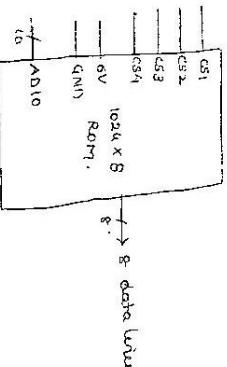
No. of power supply pins = 2 (Vcc, GND)

Ques

No. of address lines > 10.

No. of chip's & No. of data lines > 8.

Total no. of pins in 1024×8 ROM = $4 + 2 + 10 + 8 = 24$ pins.



* memory address map

The design 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 the knowledge of the size of the memory needed and type of RAM or 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 is a pictorial representation of assigned address space for each chip in the system and is called memory address map.

- (P) Design a computer system which needs 512 bytes of RAM and 512 bytes of ROM.

Sol: Size of RAM chip = 128×8
Size of ROM chip = 64×8

\Rightarrow To design a memory of 512 bytes of RAM and 512 bytes of ROM, we require 4 ($\frac{512}{128} = 4$), RAM chips and 1 ROM chip.

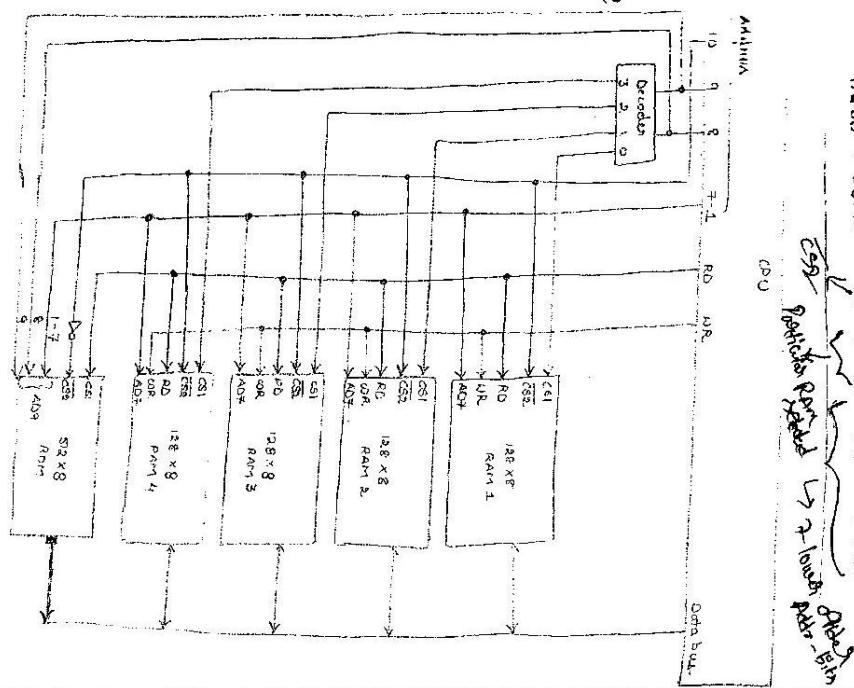
Total no. of address locations = 512 in RAM + 512 in ROM = $1024 = 2^{10}$.

\Rightarrow No. of address lines required = 10.

CPU has 16 address lines \Rightarrow higher order 6 address lines 16 to 21 are always 0's.

Address lines 8 and 9 act as chip-select for RAM₁, RAM₂, RAM₃ and RAM₄. Address line 10 acts as chip-select for ROM.

Component	Maskable	Address Bus
RAM ₁	000 - 011	0 0 0 x x x x x x x x
RAM ₂	000 - 011	0 0 1 x x x x x x x x
RAM ₃	0111	0 1 0 x x x x x x x x
RAM ₄	0110 - 0111	0 1 1 x x x x x x x x
ROM	0000 - 03FF	1 0 1 x x x x x x x x



- (a) The bus which carries addresses to and from memory is called address bus. It has two parts - address bus and data bus.
- The lower-order lines in the address bus select the byte within a particular chip and other lines in the address bus select a particular chip through its chip select pins.
 - Each RAM receives 4 low-order bits of the address bus to select one of 128 possible bytes.
 - A particular RAM is selected from lines 8 and 9 in the address bus. Thus 8 lines through a 2-to-1 decoder whose output go to CS inputs in each RAM chip.

q8
00 - RAM₁
01 - RAM₂
10 - RAM₃
11 - RAM₄

→ RD and WR pins from the microprocessor are applied to pins of each RAM chip.

- Selection between RAM and ROM is achieved through bus line 10. RAMs are selected when bit line is 0 and ROM is selected when 1.
- The ROM chip to be enabled only during a read control operation.
- Address bus lines 1 to 4 are applied to the 8-to-1 address of ROM without going through the decoder. Thus unique addresses 0 to 51 to RAM and 612 to 1023 to ROM.
- The data bus of the ROM has only an output capability, whereas the data bus connected to the RAMs can transfer information in both directions.

- (b) A computer employs RAM chips of 256x8 and ROM chips of 1024x8. The computer memory needs 4 bytes of RAM, 4K bytes of ROM and 4 bytes of ROM interface units, each of 4 registers. A memory-mapped I/O configuration is used. The two highest-order bits of the address bus are assigned to RAM, 01 to ROM and 10 to ROM interface units.
- a) How many RAM and ROM chips are needed?
 - b) Draw a memory-address map for the system.
 - c) Give the address ranges in hexadecimal for RAM, ROM and ROM interface.

Sol: (a) Size of RAM chips = 256x8.

Required RAM capacity = 2Kx8

\Rightarrow No. of RAM chips = $\frac{2K \times 8}{256 \times 8} = 2 \times 1024$

g = 512 Kbytes
g = 128 Kbytes

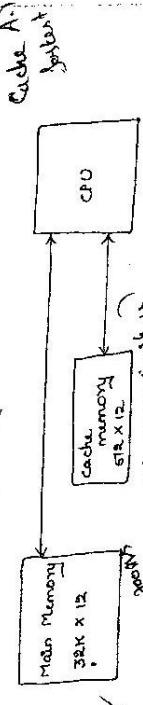
Memory Connection to CPU

Cache memory:

When a program loop is executed, the CPU \rightarrow repeatedly refers to the set of instructions in memory that constitute the loop. Eventually a given subroutine is called, its set of instructions are fetched from memory, thus loops and subroutines tend to localize the references to memory for fetching instructions.

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 **Locality of Reference**.

If the active portions of the program and data are placed in a fast small memory, the average access time can be reduced, thus reducing the total execution time of the program. Such a fast and small memory is called **Cache memory**. It is placed between CPU and main memory as shown below.



The cache access time is less than the access time of main memory by a factor of 6 to 10. The cache is the fastest component to the memory hierarchy and approaches the access time of the CPU.

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 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.

Basic operation of Cache: 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 obtain the word.

transformed from main memory to cache memory. The block size may vary from 2 word to 16 words.

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 called **hit**. If word is not found in cache, it is in main memory and it counts a **miss**.

$$\text{Hit Ratio} = \frac{\text{No. of hits}}{\text{Total no. references to main memory}}$$

$$= \frac{\text{No. of hits}}{\text{No. of hits} + \text{No. of misses}}$$

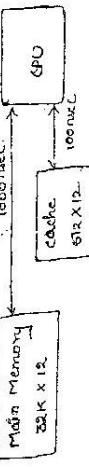
Hit ratio of 0.9 is experimentally seen. This hit ratio verifies the validity of the locality of reference property.

(Pb) The access time of a cache memory is 10 nsec and that of main memory is 100 nsec. It is estimated that 80% of memory requests are for read and the remaining 20% for write. The hit ratio for read access only is 0.9. A suitable though procedure is used.

a) what is the average access time of the system considering only memory read cycles?

b) what is the average access time of the system considering both read & write requests?

c) what is the hit ratio taking into consideration the write cycles?



Hit ratio of read access = 0.7.

\Rightarrow probability of hit = 0.9, probability of miss = 0.1.

If hit, data is accessed from cache \Rightarrow 100 nsec.
If miss, data is accessed from main memory after checking cache \Rightarrow 100 nsec + 100 nsec.

Average access time for read access = $0.9 \times 100\text{nsec} + 0.1 \times 200\text{nsec}$
 $= 110\text{nsec}$.

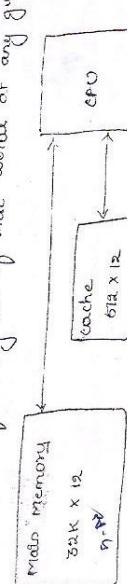
\Rightarrow 110 nsec

\downarrow = % of read cycles : Avg. access time for read cycles
 \downarrow + 1. 0% write cycles x Avg access time of cache cycles
 $= 0.8 \times 200 \text{ nsec} + 0.2 \times 100 \text{ nsec}$
 Cache
 write cycle - Through procedure.
 updates both main memory & cache.
 simultaneously \Rightarrow Avg access time = 100 nsec
 c) After conducting write cycles = $0.8 \times 0.9 = 0.72$

The transformation of data from main memory to cache memory is referred to as a mapping procedure. There 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.

Consider, The main memory can store 32K words of 12 bit each. The cache is capable of storing 64 of these words at any given time.

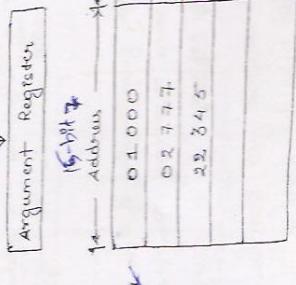


CPU generates cache address. It first sends a 16-bit address to cache. If there is a hit, the CPU receives 12-bit data from main memory and the word is then transferred to cache.

1. Associative mapping: The fastest and most flexible cache organization uses an associative memory. An associative memory stores both address and data of the memory word. This permits any location in cache to store any word from main memory.



cpu generates 16-bit address



12-bit T

The 16-bit address is shown as 6 digit octal number and the 12-bit data is shown as 4 digit octal numbers. A CPU address of 16-bit is placed in the argument register and the associative memory is searched for a matching address.

1) If address is found, the corresponding 12-bit data is read and sent to CPU.

2) If no match occurs, The main memory is accessed for the word. The address-data pair is then transferred to the cache.

3) If cache is full, an address-data pair must be displaced to make room for a pair that is read and not present in cache. The decision as to what pair is replaced depends on the replacement algorithm.

In general, FIFO or LRU displacement algorithms are used.

Disadv.: Associative memories are expensive, compound to random-access memories because of the added logic associated with each cell.

Random-access memories are used in

- ✓ 1) Direct mapping
- ✓ 2) Set-Associative mapping

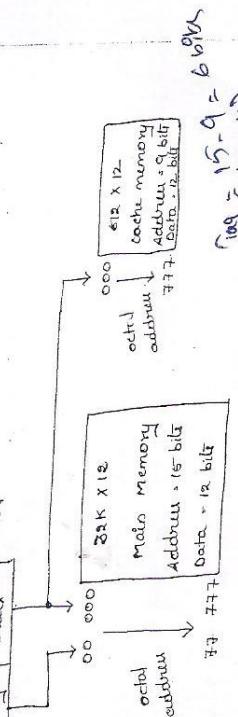
5×2^10
 Address 15-bit

X

✓

Direct mapping

The addressing relationships between main memory and cache in direct mapping are shown below:



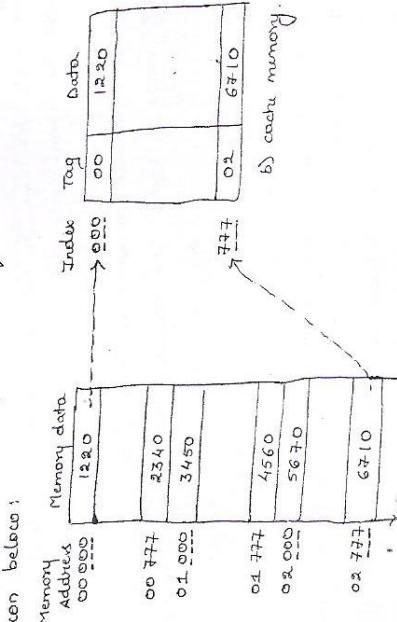
The CPU address of 16 bits is divided into 3 fields. The least significant 9 bits constitute the index field and the remaining 6 bits form the tag field.

No. of bits in the index field = No. of address bits required to access cache memory.

In general, there are 2^k words in cache and 2^n words in main memory. The n -bit memory address is divided into 2 fields: k bits for index field and $(n-k)$ bits for the tag field.

The direct mapping uses n -bit address to access main memory and k -bit address to access cache memory.

The internal organization of words in the cache is shown below:



(q. 2e)

Each cell in the cache consists of the data word and its associated tag. When a new word is brought into the cache, the tag bits stored alongside the data bits. When the CPU generates a memory request, the index field is used for the address to access the cache.

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 sourced from main memory. It is then stored in the cache together with the new tag, replacing the previous value.

Disadvantages: We cannot store two words with same index and different tag values. \Rightarrow Hit ratio will fall if two words with same index and different tag values are accessed repeatedly.

However, this possibility is minimized by the fact that such words are unlikely to appear in the address range multiples of 64 (in our example) based on "locality of reference".

Direct mapping: Direct mapping stores above uses a block size of 1 word. The same organization but block size of 8 words is shown below.

Index	Tag	Data
0000	01	3456
Block 0	004	01 6578
	040	
	014	
Block 1	014	
Block 2	014	
	040	
	076	
Block 3	02	
	074	
	047	
Block 4	02	
	074	
	047	
Block 5	02	
	074	
	047	
Block 6	02	
	074	
	047	
Block 7	02	
	074	
	047	

The index field is divided into 2 parts:

- block field
- word field

In a 64-word cache, there are 64 blocks ($\frac{512}{8} = 64$) of 8 words each.

\Rightarrow block field \times 6 bits
 \Rightarrow word field \times 3 bits

Index Block word
9 bits

Tag field stored in the cache is common for all 8 words in the same block.

memory of 1K words. The cache uses direct mapping with block size of 4 words.

a) How many bits are there in the tag, index, block and word fields of the address format?

b) How many bits are there in each word of cache, and how are they divided into functions? Include a valid bit.

c) How many blocks can the cache accommodate.

Sol: Size of main memory = $64 \text{ K} \times 16 = 2^{16} \times 16$.

\Rightarrow No. of address bits = 16 bits

Size of cache = 1K words = 1024 words .

\Rightarrow No. of address bits = 10 bits

Block size = 4 words.

Index : no. of bits in cache address.
Tag : 16 - index.

Block field : Index - Word.
Word field : 4 words = 2^2 .

Block field = Index - Word.
 \Rightarrow word field = 2 bits

Block field = Index - Word.
 \Rightarrow 10 - 2 + 8 bits

c) Cache : 16 bits
word of cache = $1+6+16$
 $\Rightarrow 23$ bits

Valid bit is used during working with write-back policy.

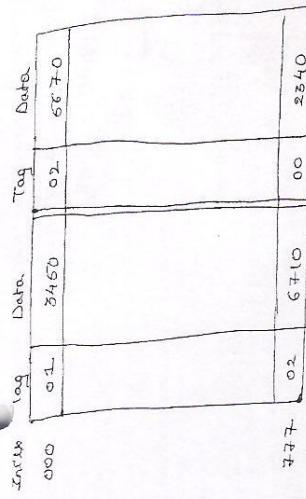
c) Cache can accommodate 2^8 blocks as block field = 8 bits

= 256 blocks of 4 words each.

c) Set-Associative Mapping

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. 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 and

(9.28)



Two-way set-associative mapping cache.

A two-way set-associative mapping accommodates 2 words of main memory in each word of cache.

In general, a set-associative cache of "set size k" will accommodate "k" words of main memory in each word of cache.

In above example, the words stored at addresses 01000 and 02000 of main memory are stored in cache memory at index address 000.illy, words at addresses 02747 and 00747 are stored in cache at index address 747.

When 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 thus 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 :

\Rightarrow The name "Set-associative" hit ratio increases as set-size increases because more words with the same index but different tags can reside in cache.

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

Replacement Algorithms

when a miss occurs in a set-associative cache and the set is full, it is necessary to replace 1 of the tag-data items with a new value.

The most common replacement algorithm is LRU (Least Recently Used).

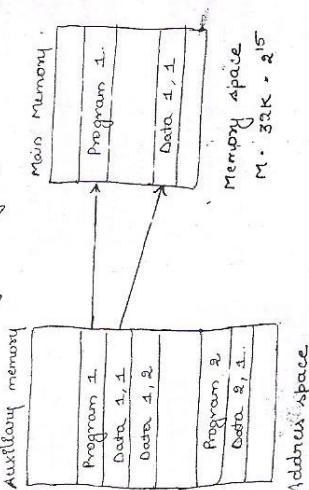
Cache is not replaced by another word unless the valid bit is set to 1 and a mismatch of tags occurs. If valid bit happens to be 0, the new word automatically replaces the previous data.

* Virtual memory

Virtual memory is a concept used in some large computer systems. It permits the user to construct programs as though the large memory space is available, equal to the totality of auxiliary 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.

Consider a system with a main memory capacity of 64K words (2¹⁶ words) and auxiliary memory of 1024K words (2¹⁰).



$$N = 1024 K = 2^{10}$$

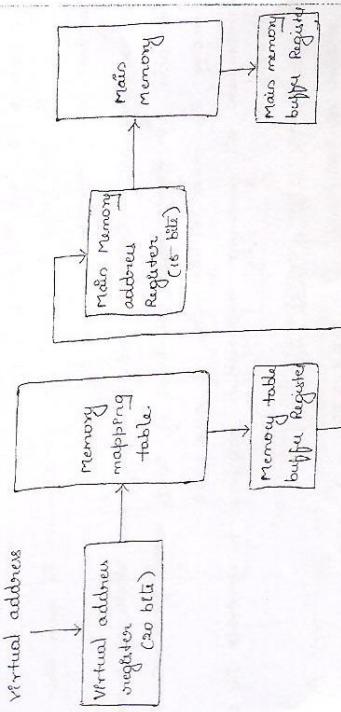
Address in auxiliary memory is called virtual address. An address used by the programmer is virtual address and set of all virtual addresses is called Address space.

An address in main memory is called physical address. The set of all physical addresses is called Memory space.

The address space is allowed to be larger than memory space so computer can use virtual memory.

Each address referenced by the CPU (virtual address) goes through an address mapping from the so-called virtual address to a physical address in main memory. A table is thus needed to map a virtual address of 20 bits to a physical address of 16 bits as shown below:

Virtual address

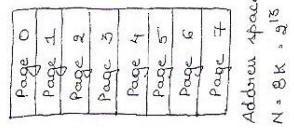


The mapping table may be stored in a separate memory or in main memory. If stored in a separate memory, an additional memory unit is required as well as one extra memory access time. If stored in main memory, the table takes space from main memory and two access to memory are required with the program running at half speed. A 3rd alternative is to use associative memory.

The address space is broken down into groups of equal size called Pages, and the physical memory or memory space is broken into groups of equal sizes called Blocks. The size of page & block are same. Portions of programs are moved from auxiliary memory to main memory in blocks equal to size of a page.

Consider a system with memory space of 4K and address space of 8K.

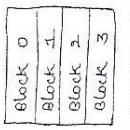
If we split each into groups of 1K words, we get 8 pages and 4 blocks as shown below:



Memory space
m. 4K = 2¹²

Address space
N = 8K = 2¹³

At any point of time, 4 pages of address space can override in main memory in any 1 of the 4 blocks.



Facilitated by each virtual address as considered can be represented by 2 numbers:

- a page number address

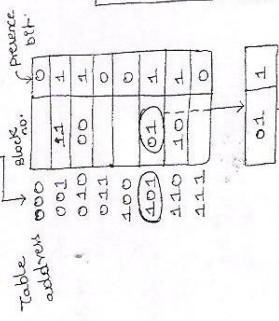
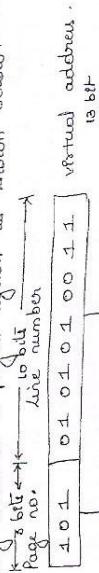
- a line number within a page.

In a computer with 2¹⁰ words per page, 10 bits are used to specify a line number and remaining higher-order bits of the virtual address specifying the page number.

$$\text{Ex:- } N = \log_2 8 \text{ k.} = 2^{10}$$

page size = 1024 words $\Rightarrow 2^{10}$
 \Rightarrow Line no = 10 bits
 page no = 13 - 10 = 3 bits to specify 4 of the 8 pages.

memory table in a paged system is shown below.



Memory page table.

The memory page table consists of 8 words, one per each page. The address to the page table denotes the page number and the content of the word gives the block number. The page is stored in main memory.

The table shows that pages 4, 2, 5 and 6 are now available in main memory. 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.

Virtual address of 15 bits. The 5 higher-order bits of the virtual memory address specify the page no. and also an address for memory-page table.

The content of the word in the memory-page table at the page number address is read out into memory table buffer register if presence bit = "1". The block no. thus read is transformed to the two higher-order bits of the main memory address register. The line no. from the virtual address is transformed into the lower-order bits of the memory address register. A read signal to main memory thus transmits the content of the word in main memory to main memory buffer register which is then used by CPU.

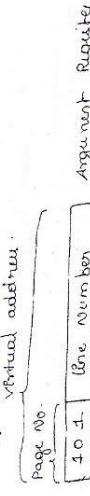
If presence bit = "0", it signifies that the content of the word referenced by the virtual address does not reside in main memory. A call to OS is then generated to fetch word from auxiliary memory and place it into main memory.

Disadv: A random access memory page table is inefficient way of storage utilization. Here, we need 8 words of memory, i.e. for each page, but at least 4 words will always be marked empty because main memory cannot accommodate more than 4 blocks.

In general, a system with "n" pages and "m" blocks require a memory-page-table of "n" locations of which up to "m" locations will be marked with block numbers and all others will be empty.

So:- A more efficient way is to organize the page table. It is to construct words equal to no. of blocks in main memory. In this case of memory is divided & it is fully utilized.

This is implemented using consecutive memory cells each word in the memory containing a page no. together with its corresponding block number. Virtual address.



Page Number Register

Page Number Register

Page No.	Page Number Register
40 4	Page Number Register
40 0	Page Number Register
40 1	Page Number Register
40 2	Page Number Register
40 3	Page Number Register
40 4	Page Number Register
40 5	Page Number Register
40 6	Page Number Register
40 7	Page Number Register
40 8	Page Number Register
40 9	Page Number Register
40 10	Page Number Register
40 11	Page Number Register
40 12	Page Number Register
40 13	Page Number Register
40 14	Page Number Register
40 15	Page Number Register
40 16	Page Number Register
40 17	Page Number Register
40 18	Page Number Register
40 19	Page Number Register
40 20	Page Number Register
40 21	Page Number Register
40 22	Page Number Register
40 23	Page Number Register
40 24	Page Number Register
40 25	Page Number Register
40 26	Page Number Register
40 27	Page Number Register
40 28	Page Number Register
40 29	Page Number Register
40 30	Page Number Register
40 31	Page Number Register

Page No.	Page Number Register
1 1 1 Key Register	Key Register
1 1 0	Memory
1 0 1	Memory
1 0 0	Memory
0 0 1	Memory

range memory

A virtual memory system is a combination of hardwired and software techniques. The memory management hardware 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 has to be transferred from auxiliary memory to main memory.
3. where the page is to be placed in main memory.

when a program starts execution, 1 or more pages are transferred into main memory and the page table is set to indicate their positions. The program is executed from main memory until it attempts to reference a page that is still in auxiliary memory. Thus condition is called page fault.

when a page fault occurs in a virtual memory, it signifies that the page referenced by CPU is not in main memory. A new page is then transferred from auxiliary memory to main memory. If main memory is full, a page is removed from main memory to make room for the new page.

The policy for choosing pages to remove is determined from the replacement algorithm. The goal of a replacement algorithm is to try to remove the page least likely to be referenced in the immediate future.

1. FIFO (first in first out)
2. LRU (Least Recently used).

(b)

A virtual memory system has an address space of 8K words, a memory space of 4K words, and page + block size of 2K words. The following page references occur during a given interval of time.

4 2 0 1 2 6 1 4 0 2 0 2 3 5 7
Determine the page that resides in main memory after each page reference if the replacement algorithm used is a) FIFO b) LRU

(c)

A virtual memory has a page size of 4K words. There are 8 pages and 4 blocks. The associative memory page table contains the following entries.

Make a list of all virtual addresses in decimal that will cause a page fault to be issued by CPU.

Page	Block
0	3
1	2
4	4
6	0

available = 0 to 1023 (2K-1)
to main = 1024 (1K) to 2047 (2K-1)
memory = 2K to (3K-1)

(3) — 3K to (4K-1)

4 — 4K to (5K-1)

(5) — 5K to (6K-1)

6 — 6K to (7K-1)

(7) — 7K to (8K-1)

Page	Block
0	4
1	2
2	0
3	1
4	3
5	5
6	7
7	6

Page	Block
0	3
1	2
2	0
3	5
4	4
5	6
6	1
7	7

a) LRU b) FIFO

1 2 0 1 2 0 1 4 0 1 0 2 3 5 7

4 1 4 2 0 2 0 1 4 2 0 1 4 2 0 1

4 1 4 2 0 2 0 1 4 2 0 1 4 2 0 1

4 1 4 2 0 2 0 1 4 2 0 1 4 2 0 1

4 1 4 2 0 2 0 1 4 2 0 1 4 2 0 1

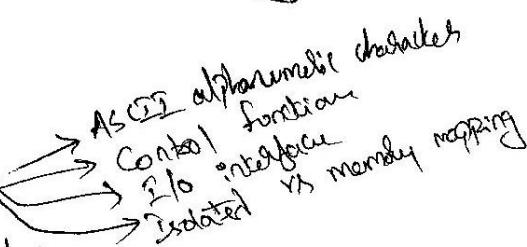
UNIT - 4

(1)

Input - output Organization & Memory System

1) Input - output organization:-

↳ Peripheral Devices



↳ Asynchronous Data transfer.

↳ Modes of transfer

↳ Priority interrupt

↳ Direct memory Access (DMA)

↳ serial Communication.

↳ Input - output Processor.

2) Memory System:-

→ Memory hierarchy

→ main memory

→ Auxiliary memory

→ Associative memory

→ Cache memory

→ Virtual memory

→ Memory management hardware.

3) Input - output organization:-

(a) Peripheral Devices:- I/O devices attached to Computer are called
Peripherals like keyboards ; display units ; ~~printers~~ printers ; magnetic disks ;
^(flo) Hard disks ; magnetic tapes.

(i) ASCII alphanumeric characters:-

(American Standard Code for Information Interchange) 91 res

7-bit code 128 characters

(b, - b7) (as shown in table),

→ out of 128 characters The ASCII Code contain 94 characters
Can be printed and 34 non printing characters used for various
control functions.