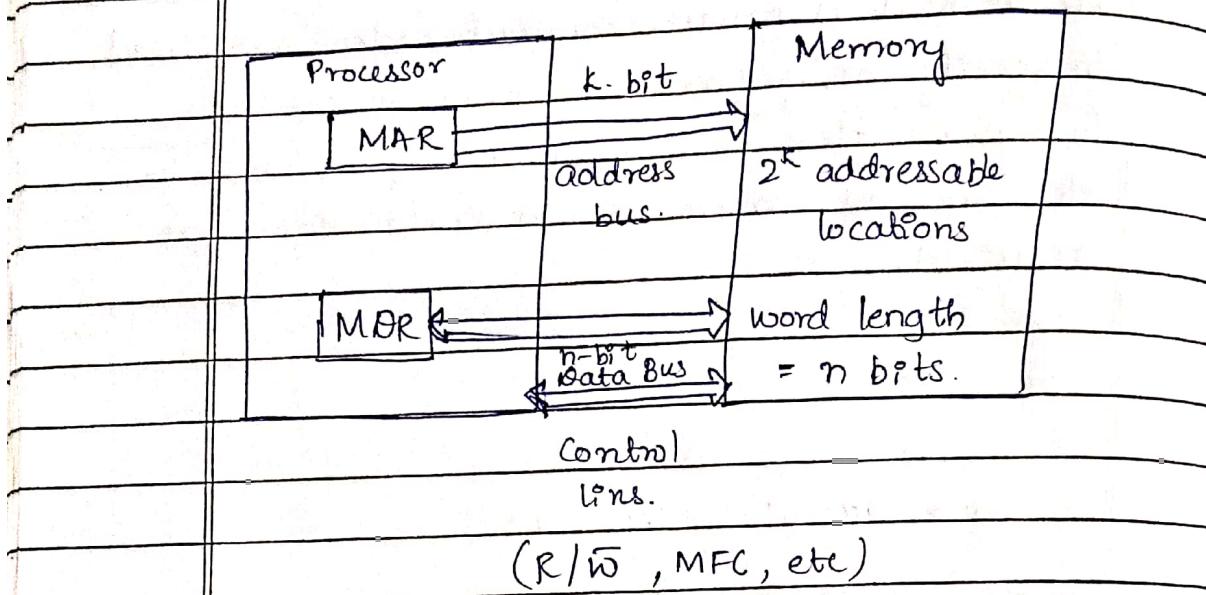


UNIT-3

MEMORY UNIT

Basic concepts



Connection between processor & Memory

- ① 32-bit = 2^{32} addressable memory locations.
- ② Memory Access Time \Rightarrow Time required to complete operation & its completion.
- ③ Memory Cycle Time \Rightarrow Time b/w two successive memory operations.
- ④ Virtual Memory.

Memory Access Time: Is the time elapsed between operation initiation and its completion.
 Ex: Time b/w Read & MFC.

Memory cycle Time:

Ex: Time b/w 2 read operations.

Memory cycle time is the minimum time delay required b/w initiation of two successive memory operation.

Virtual Memory: It translates the address specified by the program into an address that can be used to access the physical memory.

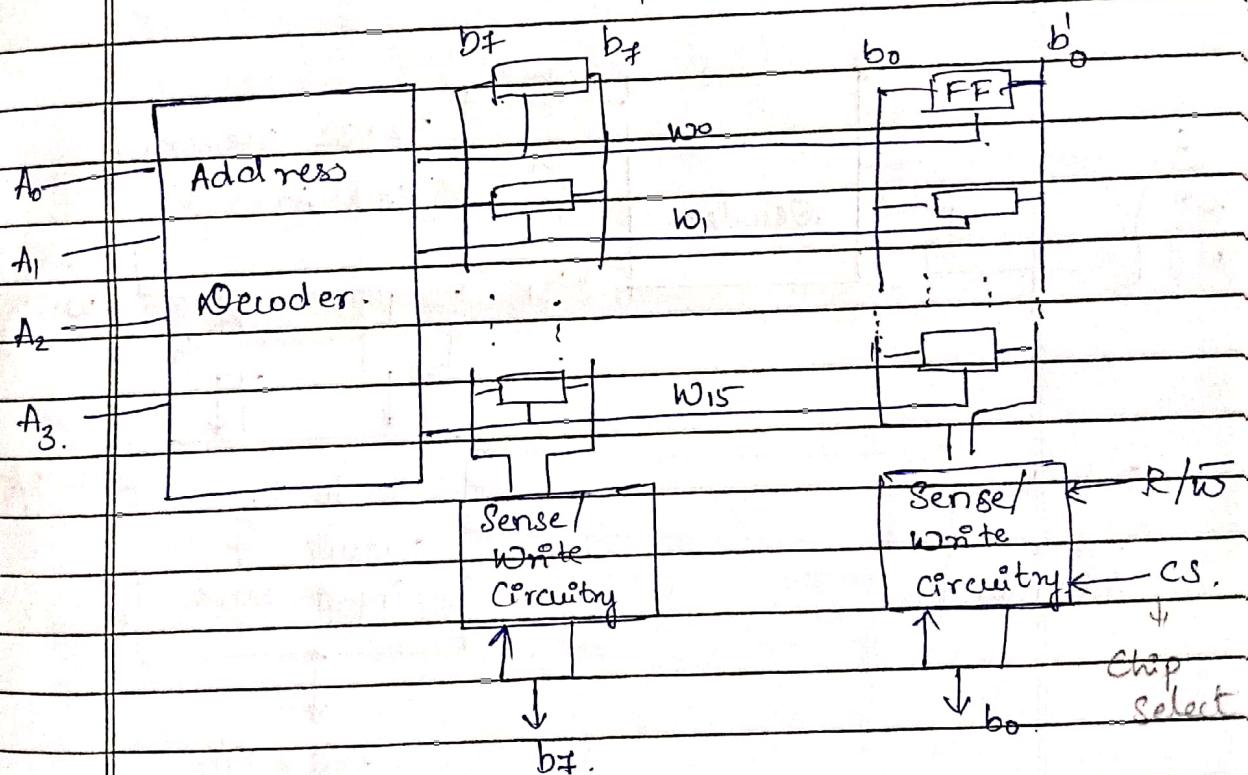
→ The address generated by the processor is called virtual address.

(2) SEMICONDUCTOR RAM MEMORIES.

(i) Internal organization of memory chip.

Question: Write the organization of a 16×8 memory chip.

$16 \times 8 \rightarrow 16$ words of 8 bits each



INTERNAL ORGANIZATION OF 16×8 MEMORY CHIP.

For the memory organization 32×32 identify
 - Length of data bus, length of address bus &
 number of addressable location.
 ↳ (32)

32 words of 32 bits each.

length of data bus = 32 bits.

length of address bus = 5

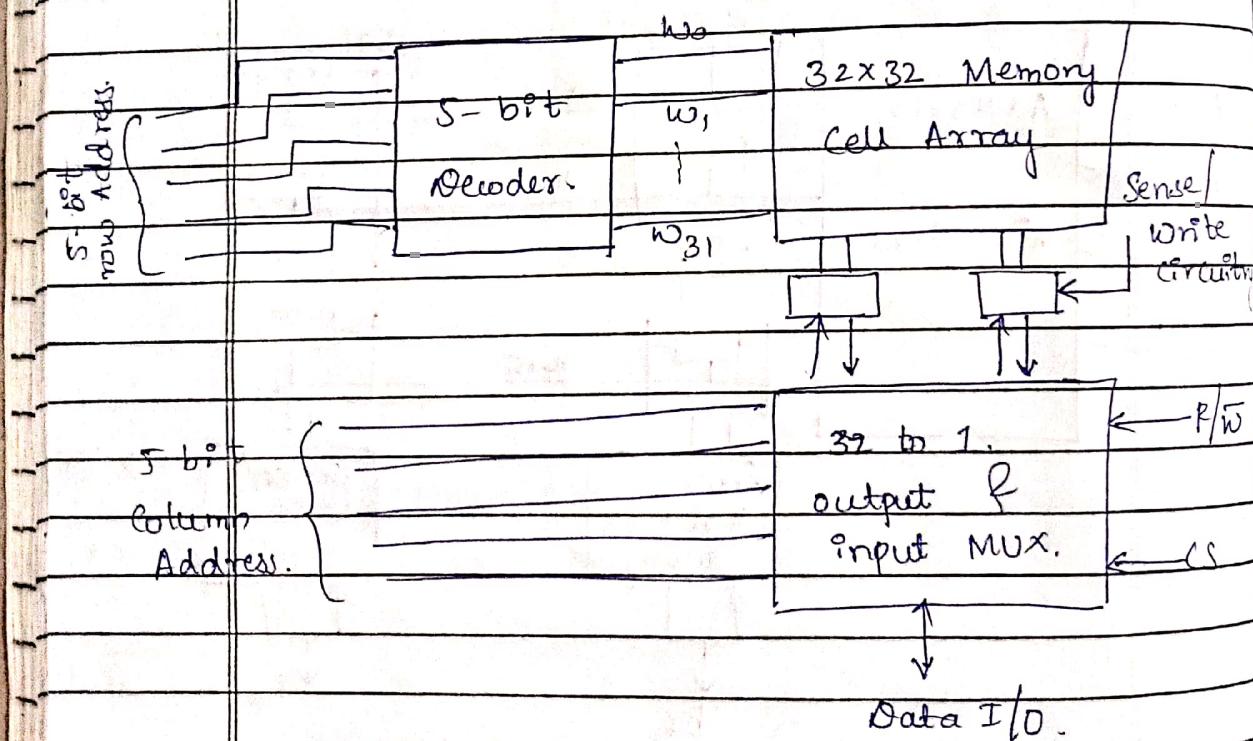
$$\{ 2^5 = 32 \}$$

Number of addressable

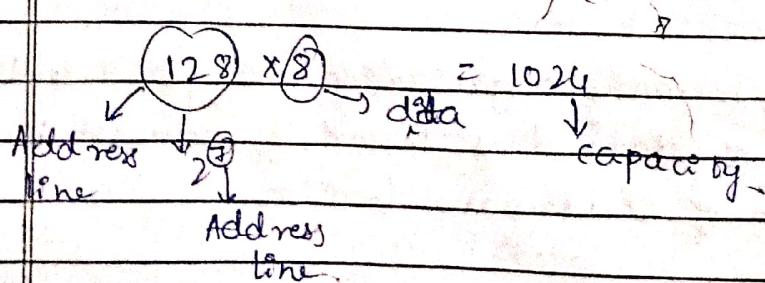
locations = 32 words.

→ During read operation the sense/read circuits, read information stored in the cells selected by the word line & transmit this information to the output data lines.

Organization of 1k x 1 Memory chip:



$$1k \times 1 \Rightarrow 1024 \times 1$$



How many external connections are required for 128×8 memory chip.

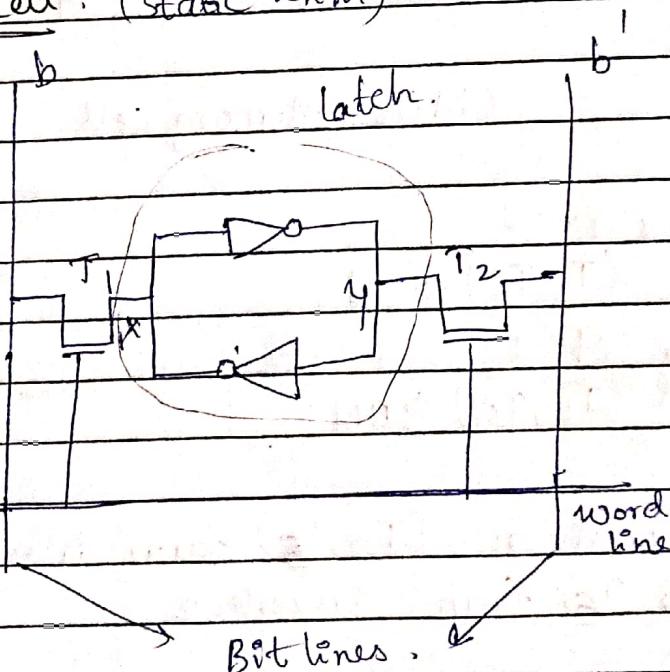
No. of address line + No. of data lines
 + Ground & power supply
 $7 + 8 + 2 = \underline{17}$

Design 128×8 Memory chip using $1k \times 1$ format. to reduce the no. of external connections.

\Rightarrow 10-bit address is divided into two groups of 5-bits each to form the row & column addresses for the cell array.

STATIC MEMORY

SRAM Cell: (Static RAM)



- \Rightarrow T₁ & T₂ act as switches that can be opened or closed by controlling the word line.
- \Rightarrow When word line is grounded, transistors are

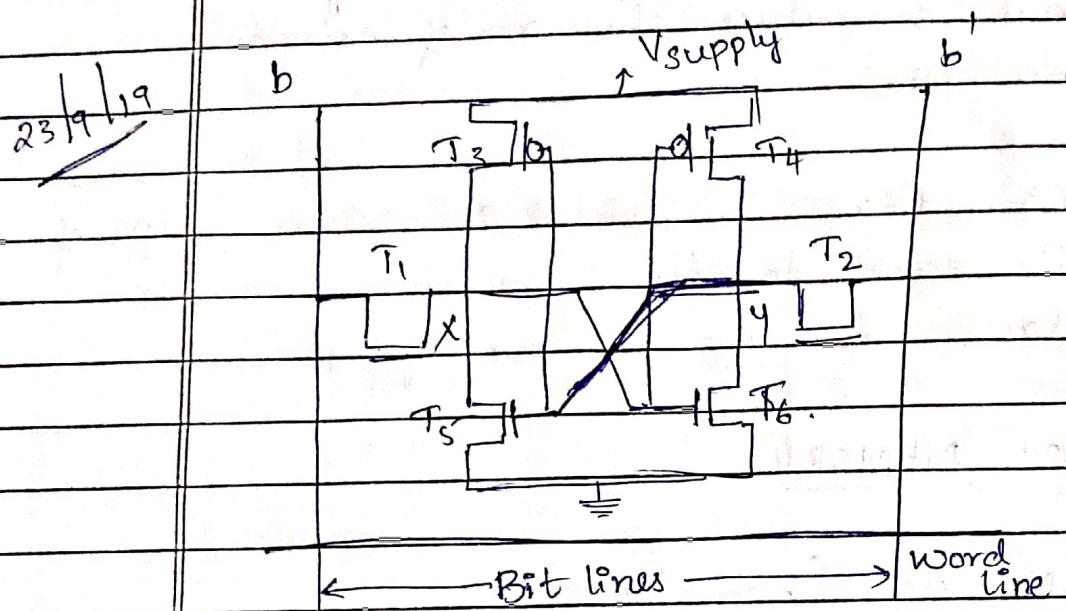
Asha

closed and the latch retains its state.

① Read - $\{ b=1, b'=0 \}$ or $\{ b=0, b'=1 \}$

$\hookrightarrow b \& b'$ monitored by sense / write circuits.

② Write - Appropriate value set on $b \& b'$ & write word line is activated



CMOS Memory cell.

For $x=1$

$(T_3, T_6) \Rightarrow ON$.

For $y=1$

$(T_4, T_5) \Rightarrow ON$.

→ Transistors (T_3, T_6) form inverter 1 and (T_4, T_5) forms inverter 2.

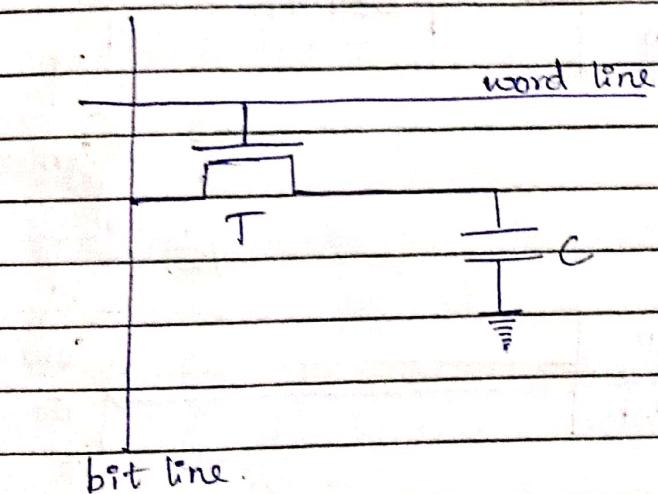
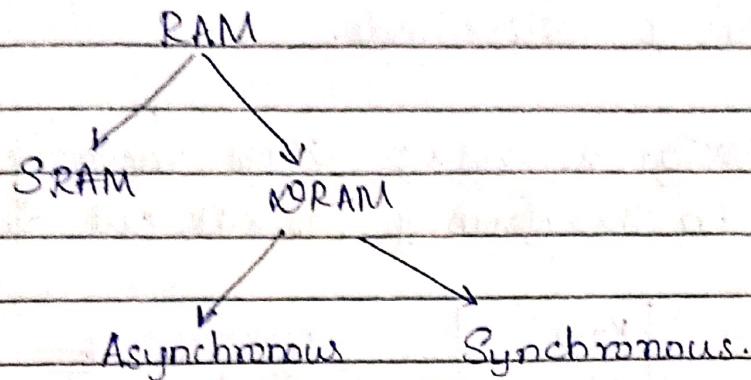
→ Voltage at X is maintained high by having $T_3 \& T_6$ in ON state & T_4, T_5 in OFF state

If T_1 & T_2 are turned on (closed) it will have high signal.

Advantages of CMOS Memory cell

- * low power consumption
- + can be accessed quickly (in nanoseconds)

ASYNCHRONOUS DYNAMIC RAM (DRAM).



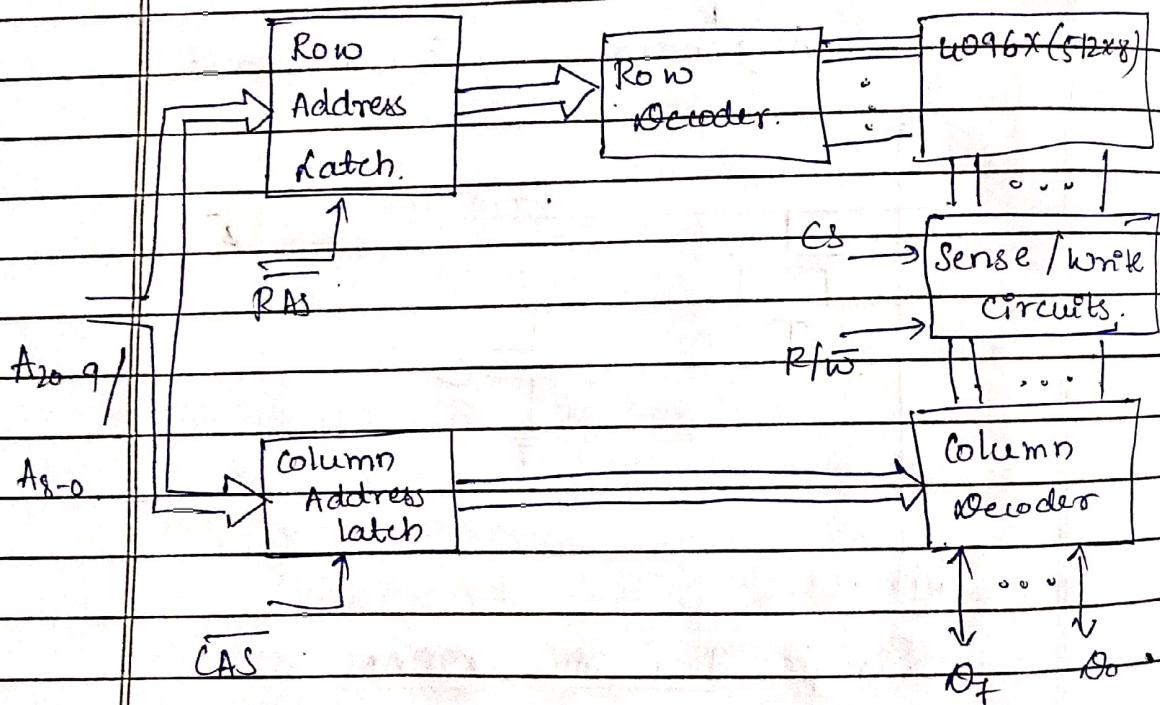
Single Transistor DRAM Cell

If capacitor \leq Threshold $\Rightarrow 0$.
 If capacitor charge

If capacitor \geq Threshold $\Rightarrow 1$.
 If capacitor charge

- DRAM is cheaper.
- Main memory is made up of DRAM.
- Cache memory is made up of SRAM.
- DRAM doesn't retain its state indefinitely.
- Information is stored in the dynamic memory cell in the form of capacitor charge & this charge will maintained only for 10's of milliseconds.

Design a $2M \times 8$ DRAM memory chip.
in the form of $4K \times 4K$ cell array.



$$2M \times 8 \rightarrow 2^{21} \times 8$$

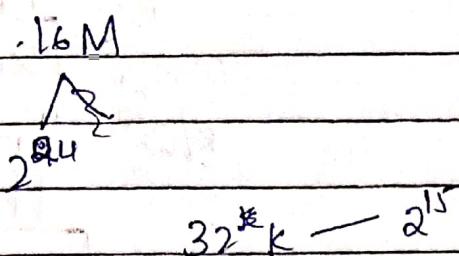
$$4K \times 4K \Rightarrow 4096 \times 4096$$

$$4096 \xrightarrow{6 bits} 512 \times 8$$

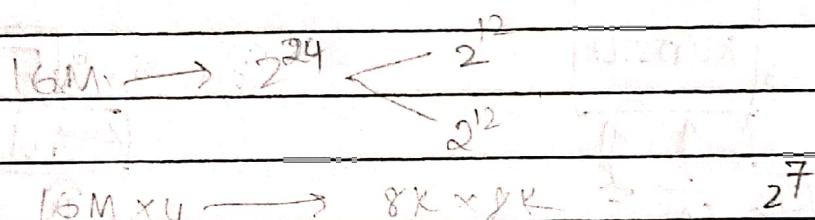
2¹²

During R/W operation row address is applied first to select the word & column address is applied to select the byte.

Design a $16M \times 8$ using $32K \times 4K$ Cell Array.



Design $64M$ Memory using $16M \times 4$, $8M \times 8$, $4M \times 16$



Design a 16×8 Memory chip

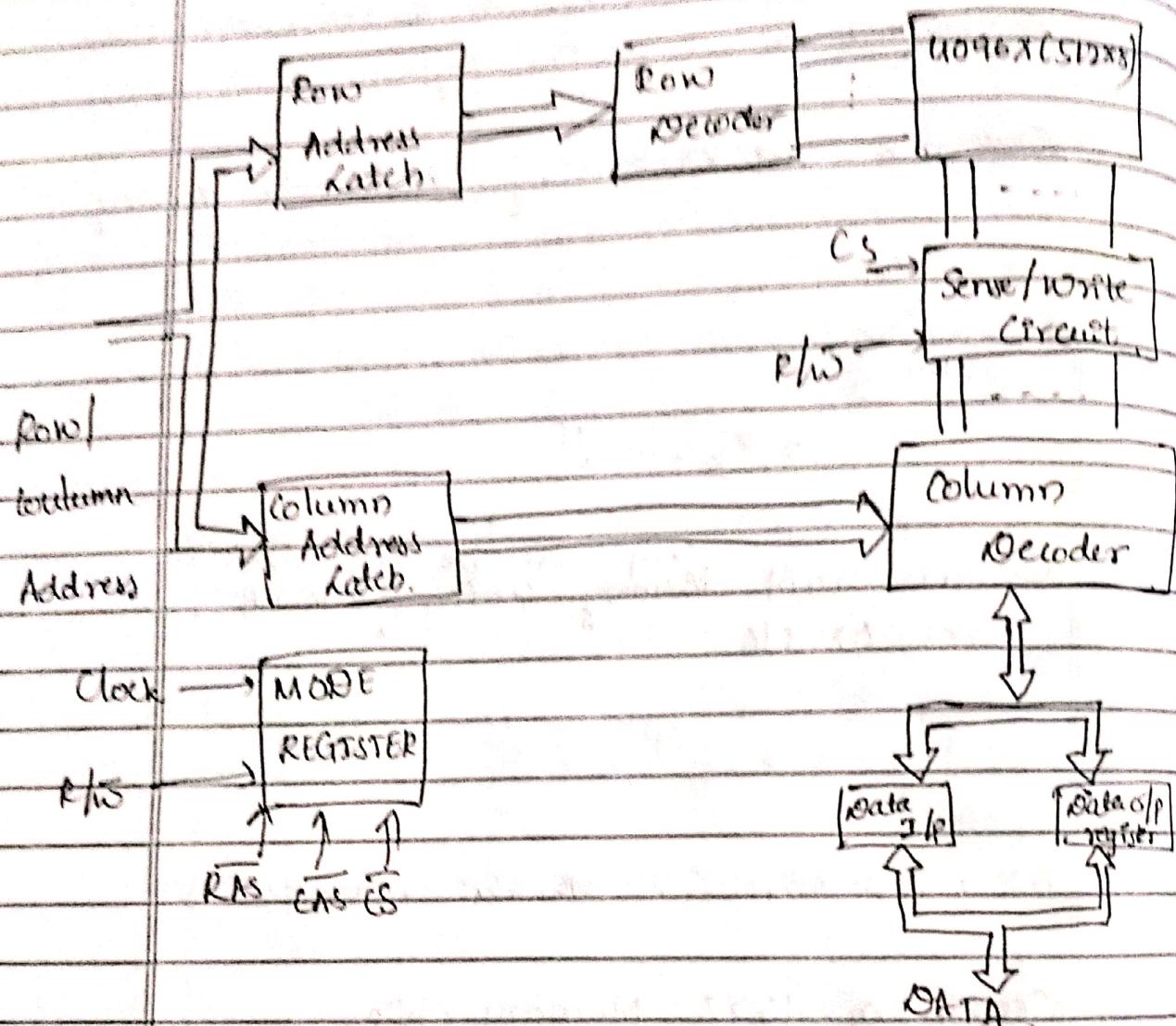
→ DRAM chip is organised to read or write no. of data bits in parallel to reduce no. of memory chip.

FAST PAGE MODE.

↳ Burst (words).

In fast page mode bytes are transferred sequentially by applying consecutive column addresses under the control of successive CAS (Column Address Store)

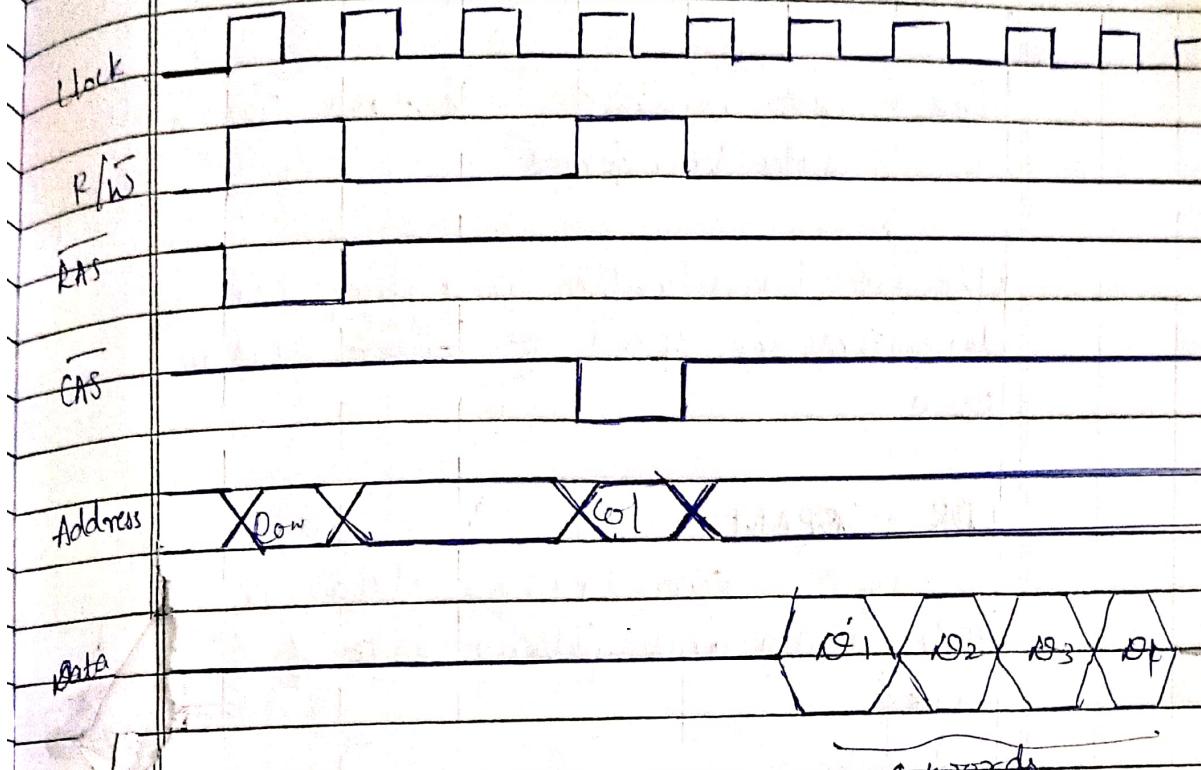
SYNCHRONOUS DRAM



\Rightarrow DRAM's are directly synchronised by clock signal mode register is used to store the control information.

Here, counters are used to store the no. of clock pulses required for data transfer operation.

Timing Diagrams for read operation for a burst read of length 4 in SRAM.



Write the timing diagrams for a burst read of length 4 in SRAM.

- Major Manufacturers like intel produces SRAM chip in which the memory takes two or three clock cycles to activate the selected row.

Latency: Memory latency is a amount of time taken to transfer one word of data to /from memory. (must be low)

Latency = 5 clock cycle to transfer 1st word.

In block transfers latency is the time taken to transfer first word of data.

Memory Bandwidth: (High)

Number of bits / bytes that are transferred to / from memory per second.

$B/W = \text{rate at which data transferred} \times \text{data Bus width}$

Effective band width in a computer is dependent on speed of ~~burst~~ and memory speed.

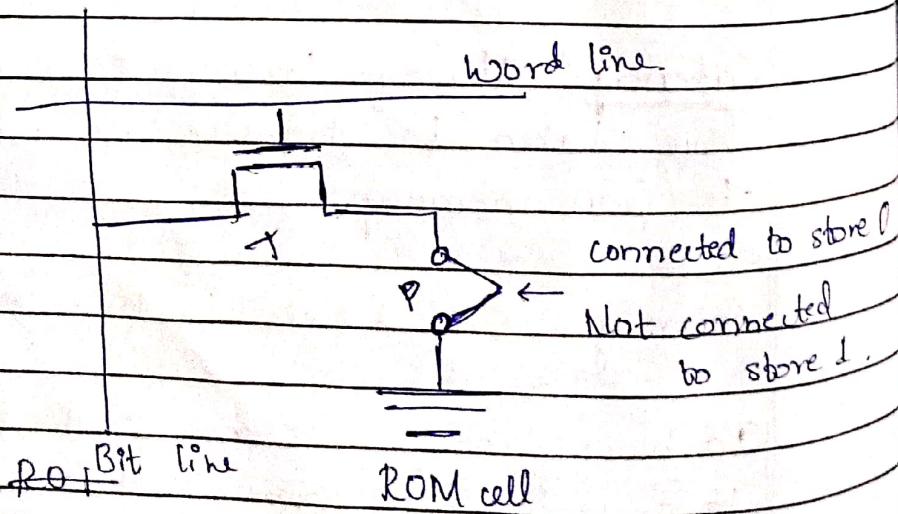
DDR - SDRAM

can transfer data on both rising and falling edge of clock.

$B/W \text{ of DDR-SDRAM} = 2 \times B/W \text{ of SDRAM}$

→ It is used in graphic terminals & machine learning systems.

ROM: (Read only- Memory).





- ROM is non-volatile
- Small amount of non-volatile memory called ROM is provided that holds the instructions whose execution results in loading the boot program from the list.
- Bit line tops connected to the power supply through a register.
- Cells circuit at the end of bit line generates proper output.

27/9/19

Types of ROM:

- PROM
- EEPROM
- EEPROM.
- Flash Memory (Flash Card & Drives)

PROM: Programmable Read only Memory.

1. Data loaded by user.
2. Programmability achieved by inserting Fuse at "P".
3. Memory contains 0's before programming.
4. User inserts 1's by burning out Fuse using high current pulse.
5. Irreversible process.

Advantages:

1. Flexibility, faster.
2. Less expensive.

- Programmable ROM is used in cell phones where small no. of ROM's are required

EPROM: Erasable PROM

- 1) Allows stored data to be erased by exposing chip to uv light.
- 2) Charge injected into special transistor.

Advantages:

- 1) Flexibility, Retains information for long time.

Disadvantage:

- 1) Chip must be physically removed from circuit for re-programming.
 - 2) Entire contents erased.
- The EPROM chip allows stored data to be erased & new data to be loaded.

EEPROM: Electrically Erasable PROM

- 1) Can be programmed & erased electrically.
- 2) Erasing can be done selectively.

Disadvantage:

- 1) Requires different voltages for reading, writing & erasing stored data.
- EEPROM doesn't have to be physically removed & cell contents can be electrically erased.

Flash Memory:

Flash Card:

- 1) Larger module to mount flash chips on small card.
- 2) Card is plugged into accessible slot.
- 3) 1 minute of music → 1 MB
1 hour of music → 64 MB



Flash Drive

- 1) Emulates harddisk
- 2) Shorter seek & access times
- 3) Low power consumption
- 4) Weakens after it has been written 1 million times

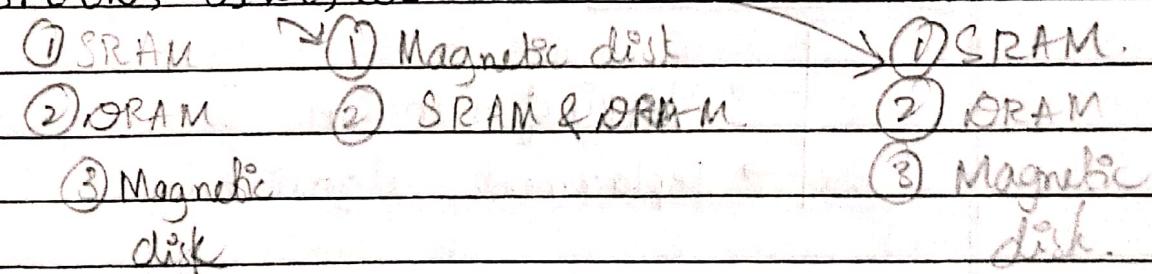
Examples: Pendrive, External hard drive

- ⇒ hand-held computers, cell-phones, digital cameras
- MP3 players ⇒ Stores Music.
- Digital cameras ⇒ Image data.

Disadvantage of flash drive over HHD (Hard Disk Drive)

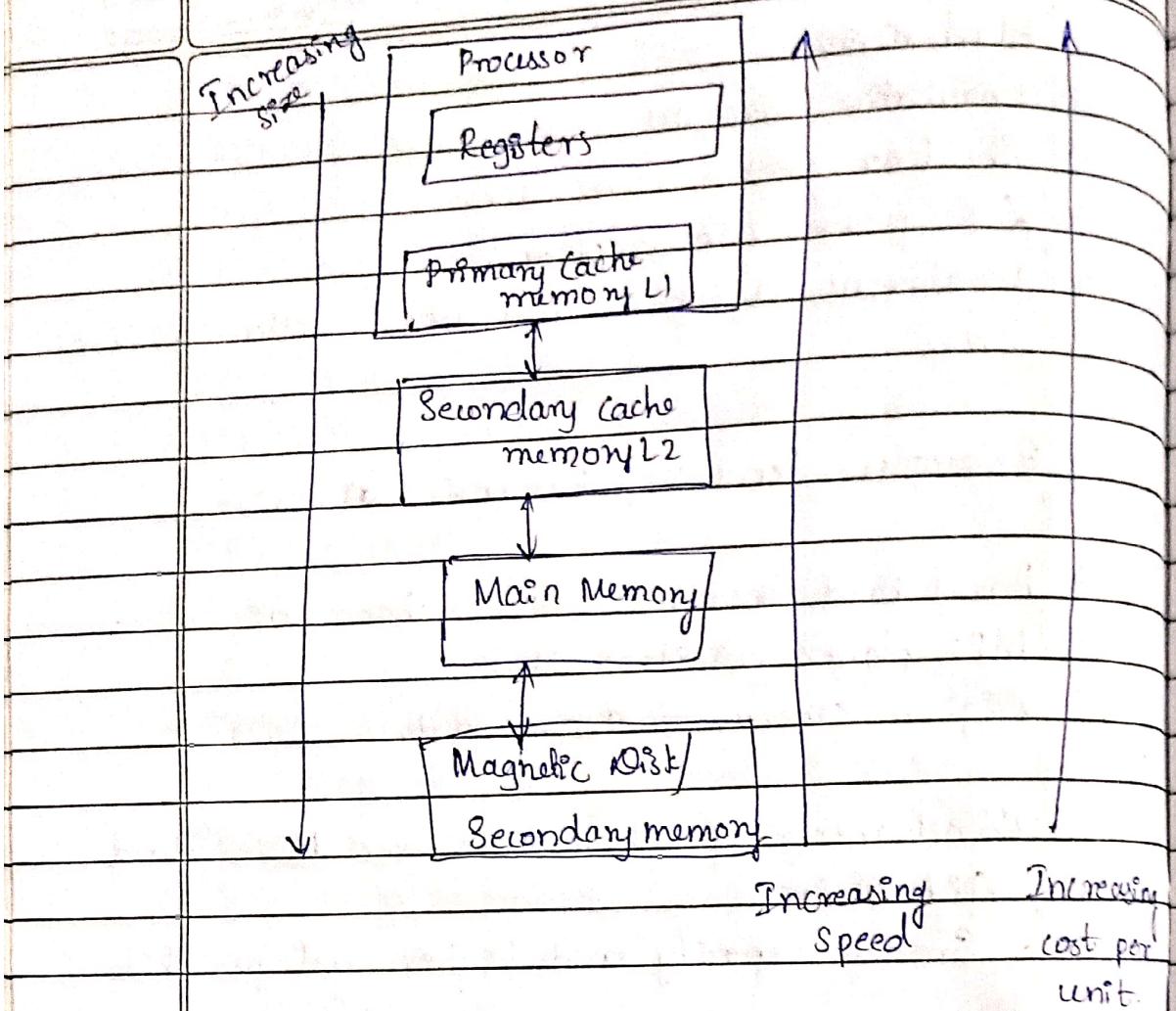
- ⇒ Small capacity and higher cost per bit.

SPEED, SIZE, COST



MEMORY HIERARCHY

- ⇒ To speed up the program execution, instructions and data will be stored in cache.

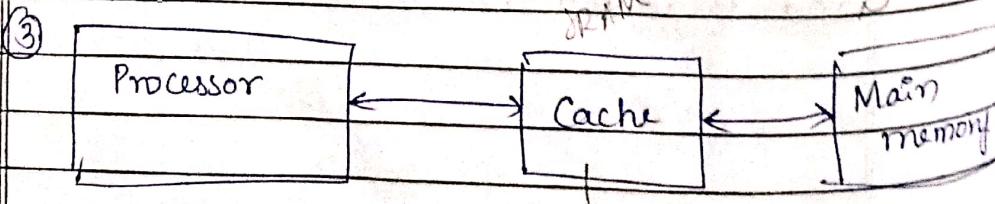


CACHE - MEMORY

- Mapping function.
 - Replacement algorithms

① Locality of Reference → temporal
→ Spatial

② Block - Set of instructions



4 Mapping function

It stores instruction & data

(5) Replacement Algorithm.

(6) Read Hit.

(7) Read Miss

(8) Write Hit.

(9) Write Miss.

(10) Write through protocol.

(11) write - Back protocol / copy - back

(12) Load through / Early restart.

(1) Locality of reference is a concept that many instructions in the localised area of the program are executed repeatedly during same time period & the remainder of the program is accessed relatively infrequently.

Temporal: Recently executed instruction is likely to be executed very soon.

Spatial: Instructions closed to the recently executed instruction are likely to be executed soon.

(2) Block: Block refers to set of contiguous address location of same size.

(3) Mapping function: It specifies the correspondence b/w main memory blocks & cache blocks.

(4) Replacement Algorithm: are used to replace cache block when a word not in cache is referenced.

(5) Write through protocol: Cache location & main memory access are updated simultaneously.

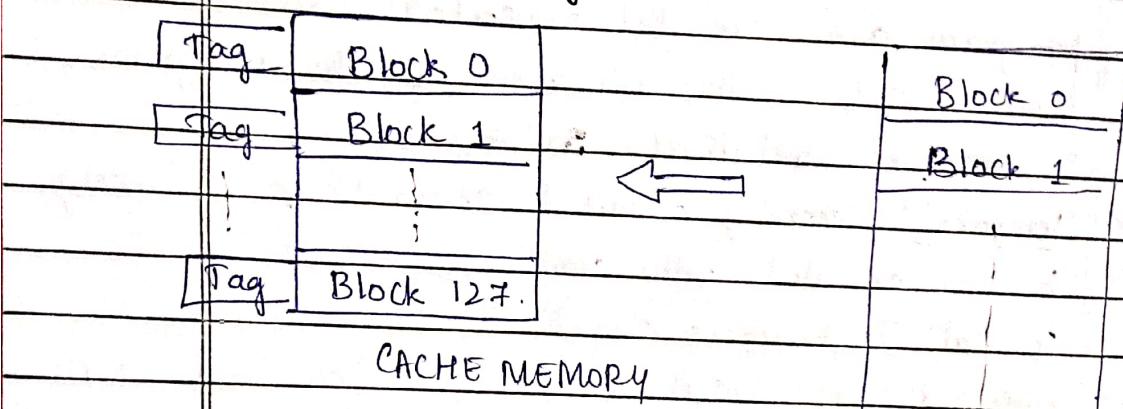
(ii)

Write back protocol: Only cache location is updated & it is ~~marked~~ modified with dirty or modified bit, main memory location is updated later when the cache block is being replaced.

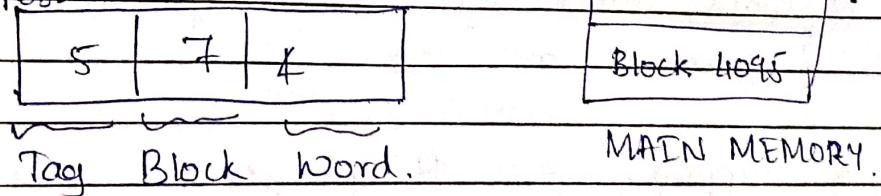
MAPPING FUNCTIONS

- (1) DIRECT MAPPING
- (2) ASSOCIATIVE MAPPING.
- (3) SET ASSOCIATIVE MAPPING.

Direct Mapping:



Memory Address.



Tag Block Word.

Main memory \Rightarrow 4096 blocks.

1 blocks = 16 words \Rightarrow Total $= \frac{4096 \times 16}{\text{Memory}}$

Cache memory \rightarrow 128 Blocks. $= 128 \times 16 = 2048$ words.
 Memory address size $\Rightarrow 16 = 2^{16} = 65536$

$j \bmod 128$

$j \rightarrow \text{Block}$

19 54 6
10011 0110110 0110

Processor wants to read 6th word from
54th block present in 19th tag.

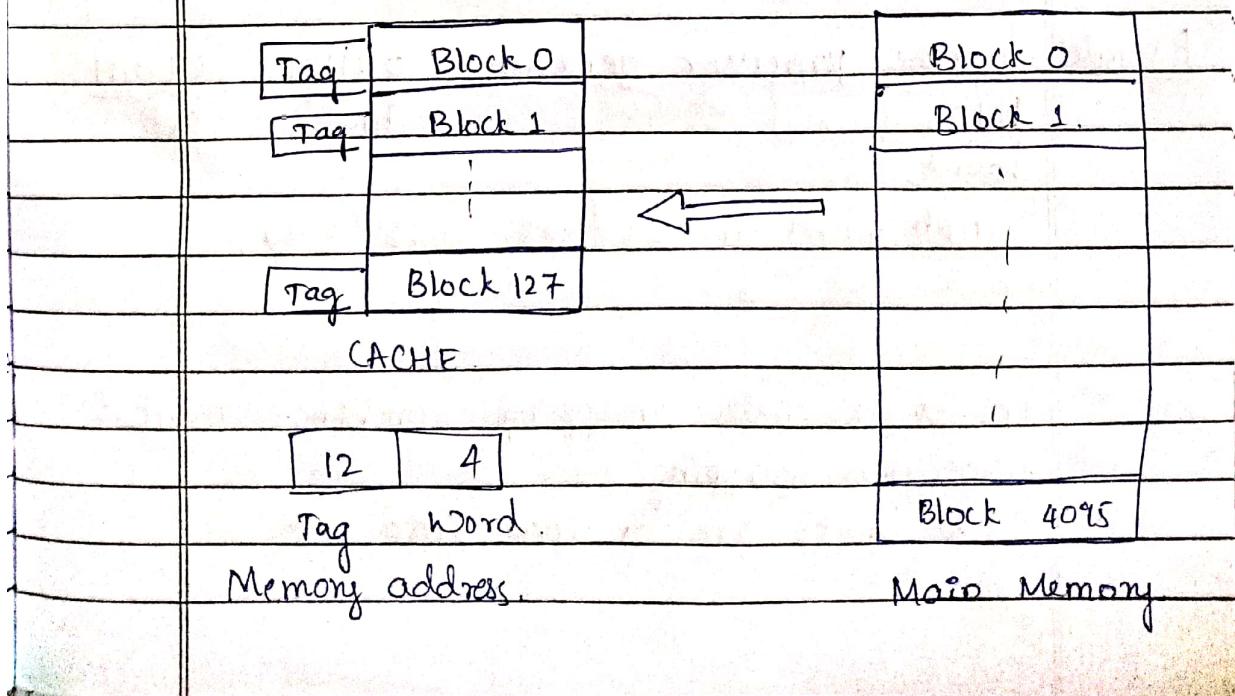
$j \bmod \text{No. of cache blocks}$

- Tag bits are used to identify whether the addressed word is present in cache.
- Direct mapping is easy to implement but not very flexible.
- Memory wastage happens because 128 blocks has to be stored in block 0 even if other cache blocks are empty.

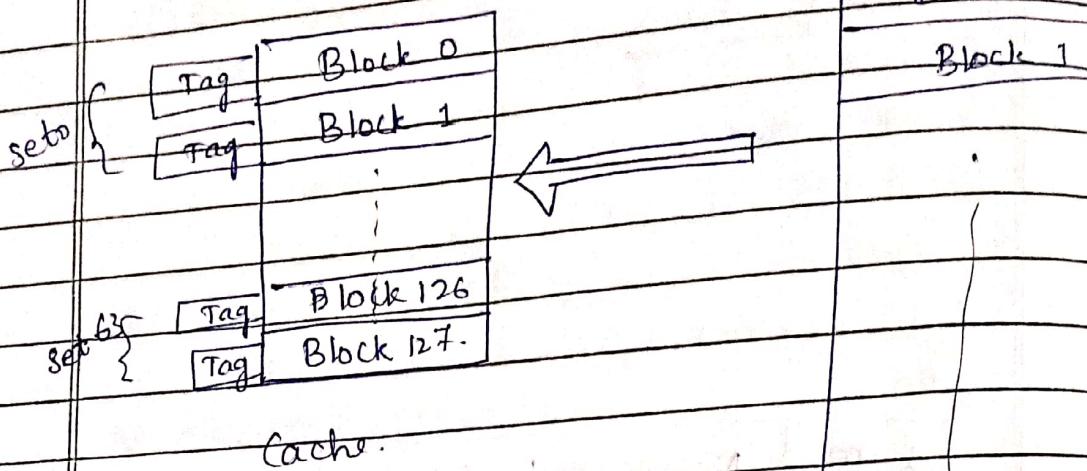
1) 19

MAPPING FUNCTION

2) ASSOCIATIVE MAPPING.



3) Set-Associative Mapping.



- In Associative Mapping, here main memory can be placed into any cache block position.
- 12 tag bits are used to identify memory block when it is resident in the main cache memory.

Example: Suppose processor generates 011101 010111
 1010 : tag Set
word

10th word in 23rd set, 13th tag.

- In Set associative mapping, combines direct & associative mapping
- Cache blocks are grouped into sets.

- In the example given two tags in set 23 searched to see if either one matches tag 13
 - Valid Bit is provided for each word indicating whether the block contains valid data.
- Valid Bit = 1 \Rightarrow Data in cache & main-memory are consistent.
- Cache coherence problem need to ensure that two different entries use the same copies of data.

DMA \rightarrow Data Memory Access.

REPLACEMENT ALGORITHM:

1) LRU (Least Recently Used).

	Block 0
	Block 1
	Block 2
	Block 3.

Replacement algorithms are not required for direct mapped cache.

The aim is to keep track of the no. of blocks that are referenced. LRU algorithm overwrites the cache block that are not been referenced for longest time.

So a counter is kept with each block to

track the no. of references to the cache block.

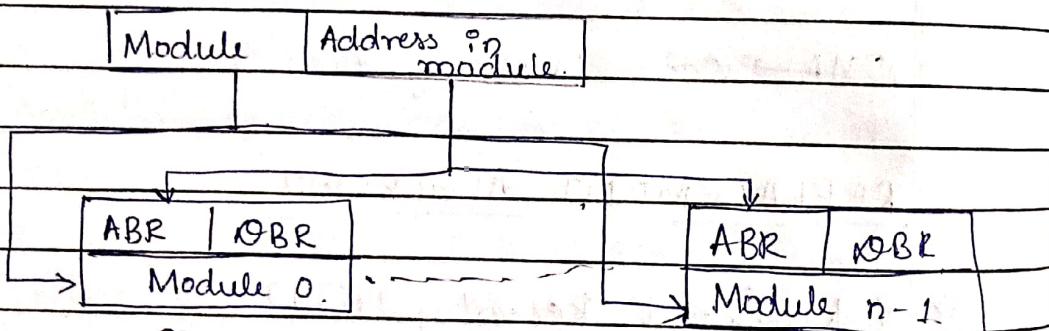
(2) oldest Block is replaced

In oldest block algorithm, the block that is there in the cache for the longest time is removed.

The simplest algorithm is to randomly choose the block to be overwritten.

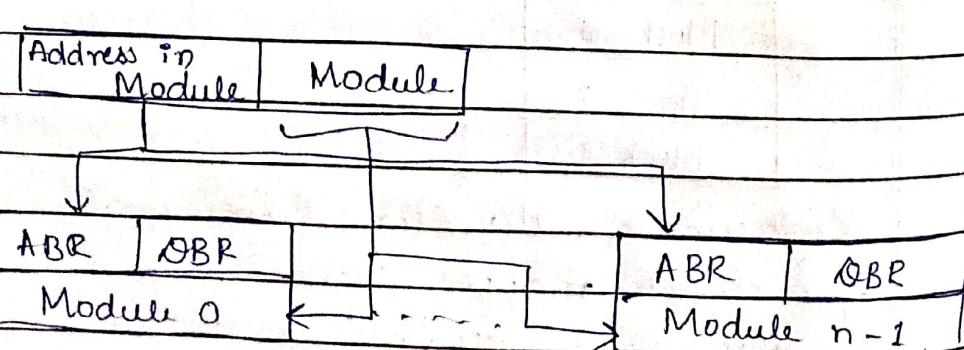
Performance Considerations:-

18 Interleaving



(a) Without interleaving.

→ Consecutive words in a Module



(b) With interleaving

Consecutive words in consecutive modules.

ABR → Address Buffer register.

DBR → Data Buffer register.

→ In the diagram main memory is structured as a physically separate module.

Example for (a).

1000i - 1st word. }
1004i - 2nd word } 1 module.

Example for (b).

1000i - 1st word 1004 : 2nd word
⇒ Module 1 Module 2

→ In memory interleaving consecutive words are located in successive modules so that several modules are kept busy at any given time as all the consecutive locations will be accessed parallelly.

Illustrate the benefit of interleaving with an example.

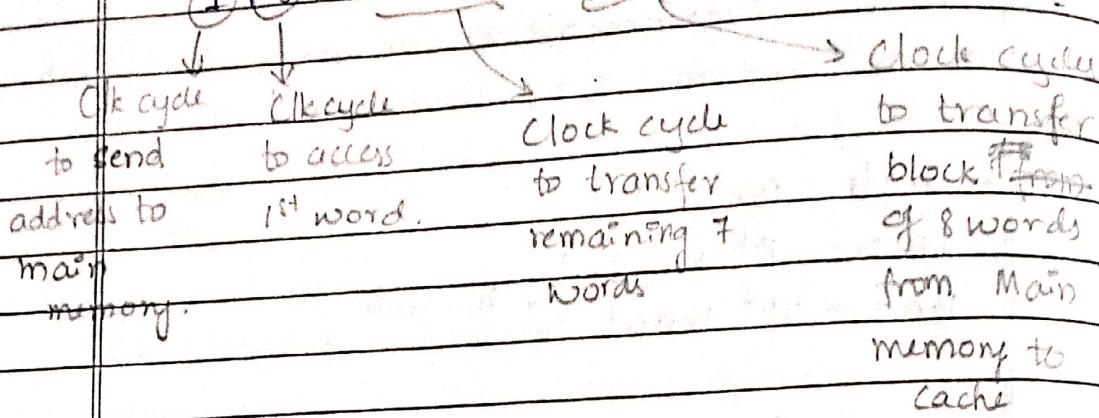
Assume cache with 8 word blocks. 1 clock cycle is to send address to main memory 1st word accessed in 8 clock cycles.

1 clock cycle needed to send 1 word to cache
Compute time (in clock cycles) needed to transfer a block of data from main memory to Cache with / without interleaving.

without interleaving'

$$1 \text{ cache block} = 8 \text{ words}$$

$$1 + 8 + (7 \times 4) + (1 \times 8) = 45$$



with interleaving:

$$1 + 8 + \cancel{(7 \times 4)} + 1 = 11$$

⇒ Interleaving reduces the block transfer time by more than a factor of two.

2) Hit rate & Miss-penalty:

$$\text{Hit rate} = \frac{\text{Number of hits}}{\text{Attempted Memory Accesses}} \geq 0.9$$

$$\text{Miss rate} = \frac{\text{Number of misses}}{\text{Attempted Memory Accesses}}$$

⇒ Performance is affected by hit rate & miss rate

⇒ Miss-penalty indicates extra time needed to bring the desired data into the cache memory



Average access

$$\text{time experienced by processor} = hC + (1-h)M.$$

where, $h \Rightarrow$ Hit rate

$C \Rightarrow$ Time to access information in Cache

$M \Rightarrow$ Miss penalty.

- * Suppose 1 Memory Read access = 10 clock cycles.
- Time needed to load a block into cache = 17 cycles.

Assume 30% of instructions perform read/write operations

\Rightarrow 130 Memory accesses

h for instructions = 0.95

h for data = 0.9

Compute performance improvement due to cache.

$$\Rightarrow \text{Time without cache.} = \frac{130 \times 10}{\text{No. of memory access for each access}}$$

$$\text{Time with cache} = \frac{\text{No. of memory access for instruction}}{\text{No. of memory access for data}} \times \frac{(hC + (1-h)M)}{(hC + (1-h)M)}$$

$M = 17$ cycles.

$h =$

$(=)$

No. of memory access for data.
30.

$$hC + (1-h)M + hC + (1-h)M$$

$$= 0.95 \times 1 + (1 - 0.95) \times 17 + 0.9 \times 1 + (1 - 0.9) \times 17$$

$$= 0.95 + 0.05 \times 17 + 0.9 + 0.1 \times 17$$

$$= 4.4$$

$$\Rightarrow \text{Answer} = \frac{130 \times 10}{11.4} = \frac{130 \times 10}{100 \times (0.95 + 0.05 \times 17) + 30 \times (0.9 + 0.1 \times 17)} \\ = 5.04.$$

+ Calculate performance improvement of memory with & without cache given the following data:-

1 cache blocks = 8 words

Hit rate for instructions = 0.95

Hit rate for Data = 0.9

Number of processor cycles to transfer 8 words Block to / from main memory = 64

Number of processor cycles to access single word in main-memory without cache = 36

Number of instructions accessed = 100 {out of }

Number of Data accessed = 30

$\ell = 1$

BO memory access 30% are read/write transfers

130.

$$\text{Time without cache} = \cancel{100 \times 60} \times 36 \\ = \cancel{2160} \quad 4680$$

$$\begin{aligned} \text{Time with cache} &= 100 \times (0.95 + (1-0.95) 60) \\ &\quad + 30 \times (0.9 + (1-0.9) 60) \\ &= 100 \times 3.95 + 30 \times 6.9 \\ &= 602 \end{aligned}$$

$$\frac{\text{Time without cache}}{\text{Time with cache}} = \frac{4680}{602} = 7.77$$

This indicates that the memory with cache performs 7.77 times faster than memory without cache

Assume system bus clock is 4 times slower than processor clock

Hit rate improved by:

- 1) Cache larger \Rightarrow Expensive
- 2) Transfer large blocks b/w main memory & cache.

\hookrightarrow Miss penalty may increase.

- 3) Load through protocol. (Processor doesn't wait for entire block to be loaded)
 \Rightarrow Processor starts reading word immediately after it is brought into cache.

Cache on processor chip.

Due to space requirements cache size is limited and pentium processor uses L1 & L2 cache to increase the memory access speed

L1 cache } pentium, ARM
 L2 cache }

$$t_{ave} = h_1 C_1 + (1-h_1) h_2 C_2 + (1-h_1)(1-h_2) M$$

access where,

speed. $h_1 \Rightarrow$ Hit rate of L1

$C_1 \Rightarrow$ Time to access information in L1 cache

$h_2 \Rightarrow$ Hit rate of L2

$C_2 \Rightarrow$ Time to access information in L2 cache

$M \Rightarrow$ Miss penalty

in system with 2 cache \Rightarrow average access time by processor

Including L₂ cache reduces the average access time of memory by the processor.

Other Enhancement

- 1) Write buffering.
- 2) Pre fetching.
- 3) Lock-free cache.

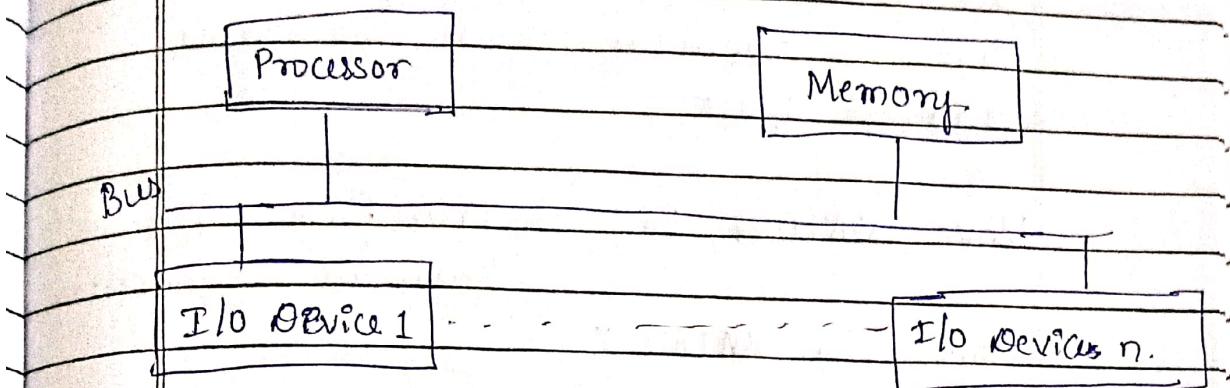
Explain the other methods for improving the performance of memory.

- Write buffering: The write buffer stores the write requests temporarily & doesn't wait for memory function complete (MFC)
↳ Read requests are high priority than write requests.
- To avoid processor waiting prefetching of data into cache is done by compiler by inserting prefetch instructions in the program.
- Lock free cache allows the processor to access the cache when a miss is being serviced.

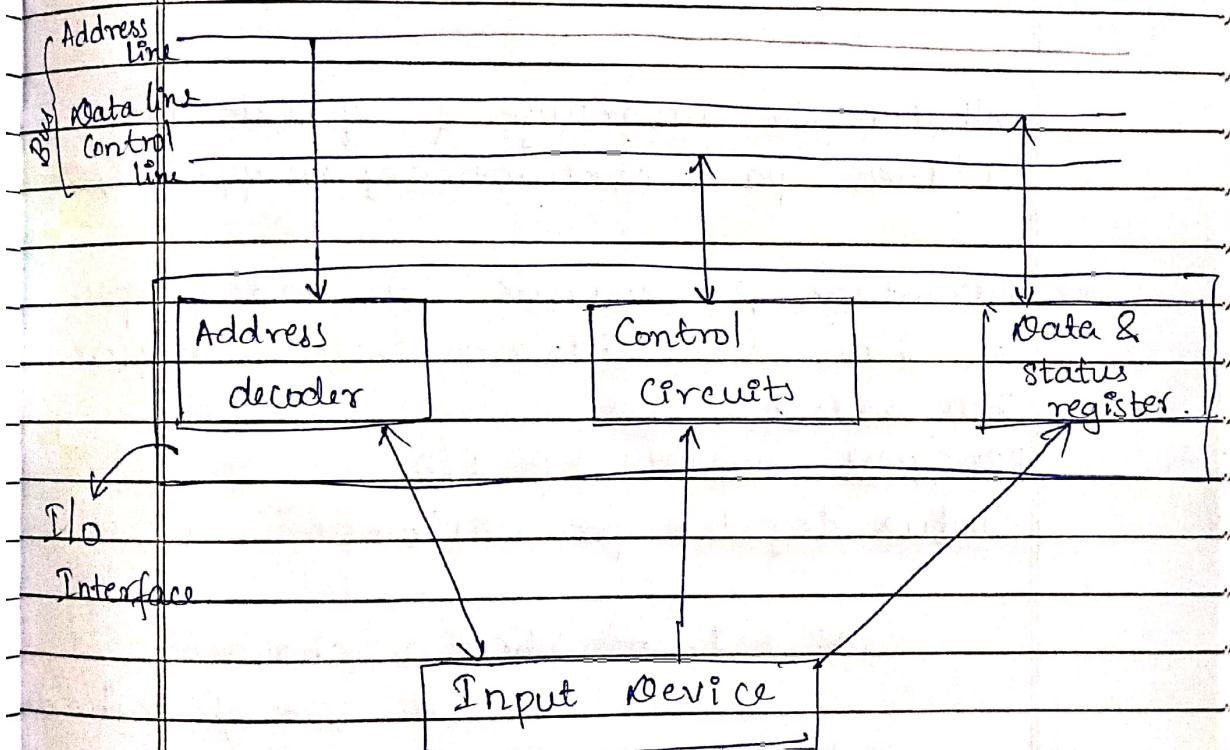
UNIT - 4

I/O Organization & Arithmetic

Accessing I/O Devices



Single-Bus structure



I/O Interface for an I/O Device

- The I/O device examines the lower ordered bits of address line to determine whether they should respond.

→ Memory mapped I/O :- Memory & I/O device share same address space.
It is an arrangement where the I/O devices & memory share the same address space.

Move $\text{RDATAIN}, R_0 \Rightarrow$ Transforms from keyboard to processor

Move $R_0, \text{RDATOUT}$

$\text{RDATAIN} \Rightarrow$ I/p buffer

What is the advantage of separate I/O address space over memory mapped I/O

→ Advantage of separate I/O address space is that I/O addresser deals with pure I/O address lines.

PROGRAM - CONTROLLED I/O.

Status Register for keyboard

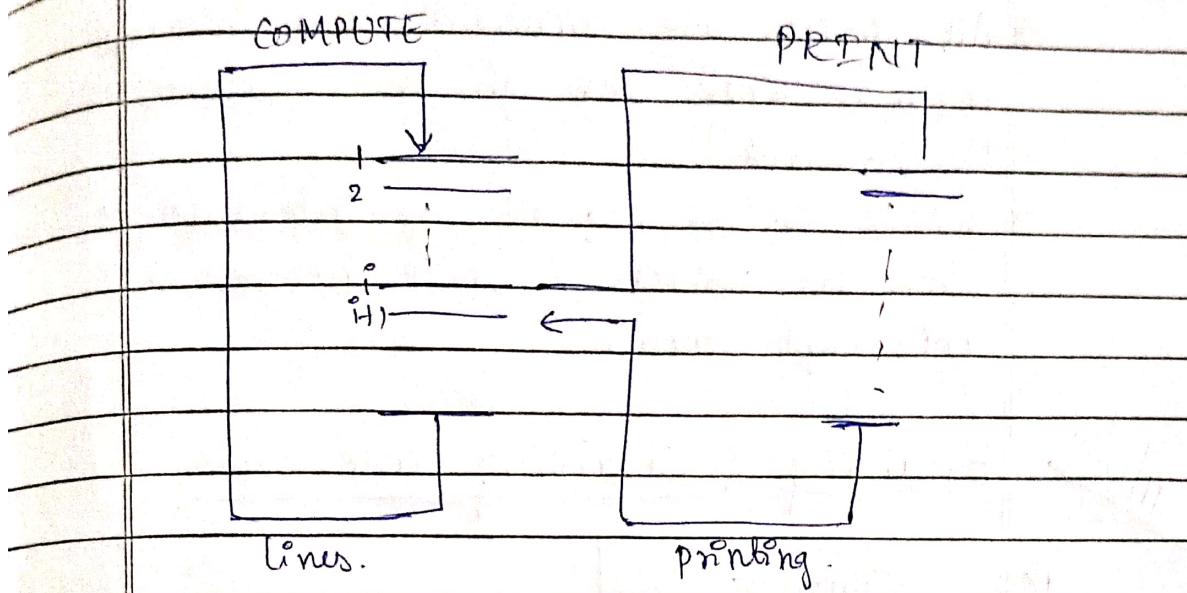
$SIN = 1 \Rightarrow$ When a character entered through keyboard

$SIN = 0 \Rightarrow$ When processor reads character.

Interrupts :

→ A program controlled I/O repeatedly reads the status register which leads to

slower execution. To overcome this ~~problems~~ interrupts are used in which the I/O device sends a special signal when it is ready for a data transfer.

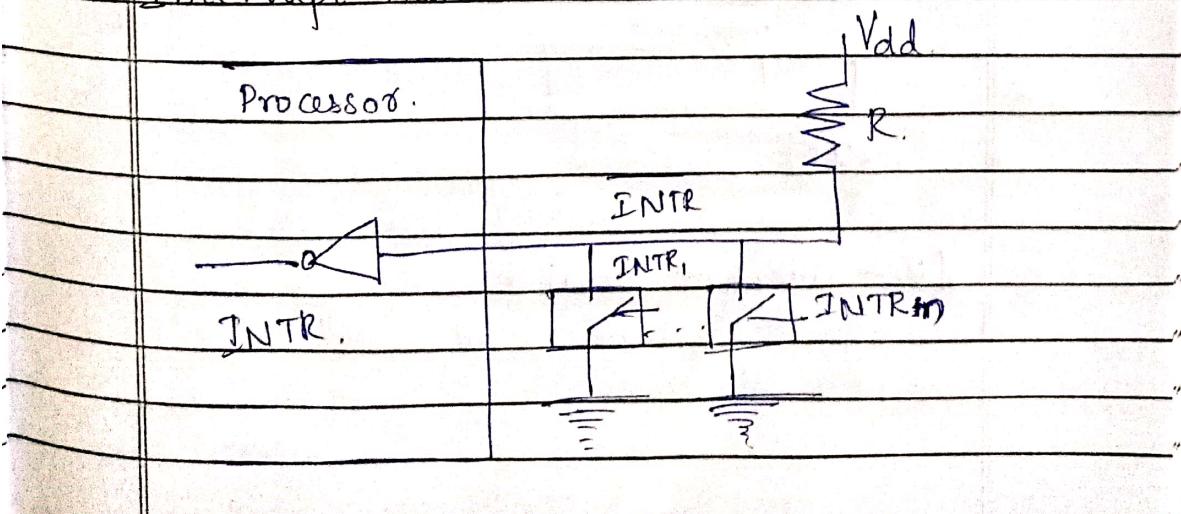


\Rightarrow ISR \rightarrow Interrupt Service Routine.

\Rightarrow ISR is a program that is executed when interrupt occurs.

Interrupt latency is the delay b/w time of interrupt request is received and the start of execution of ISR

Interrupt hardware.



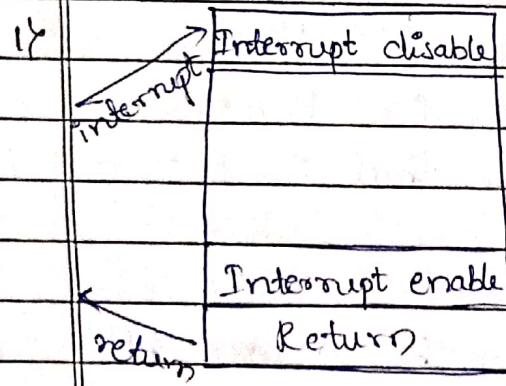
Interrupt request line implementation using open drain Bus

$$INTR = INTR_1 + INTR_2 + \dots + INTR_n$$

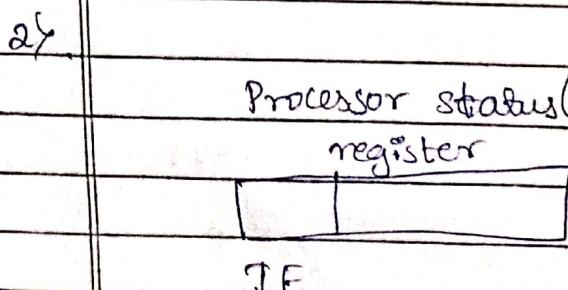
- All devices are connected to a interrupt request line through the switches to the ground.
- When device requests an interrupt it closes its switch & INTR becomes 0 & interrupt occurs.

~~Tholia~~

Enabling & Disabling interrupt



JSR



IE - Interrupt enable

- In option 1 processor uses interrupt disable instruction has 1st instruction in the TSR.
- In Option 2 the IE bit is set to 1 in the PS register to enable further interrupts.

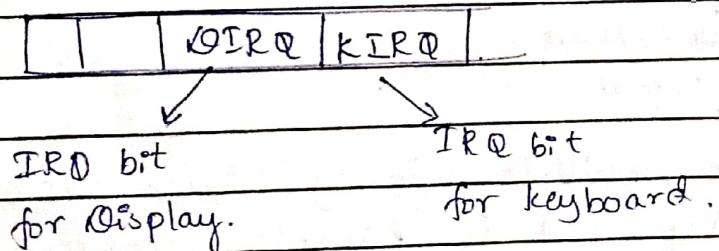
$IE = 1 \Rightarrow$ Processor enables interrupt

$IE = 0 \Rightarrow$ Processor disables interrupt

Handling Multiple devices.

- Polling -
- Vectored interrupt
- Priority (Interrupt Nesting)
- Daisy chain (Simultaneous requests).

Polling

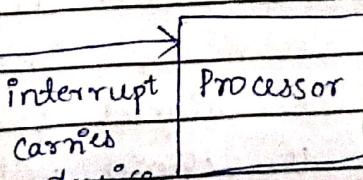


(IRQ → Interrupt Request)

* In polling Scheme.

the first device encountered with IRD bit set is the device that should be serviced

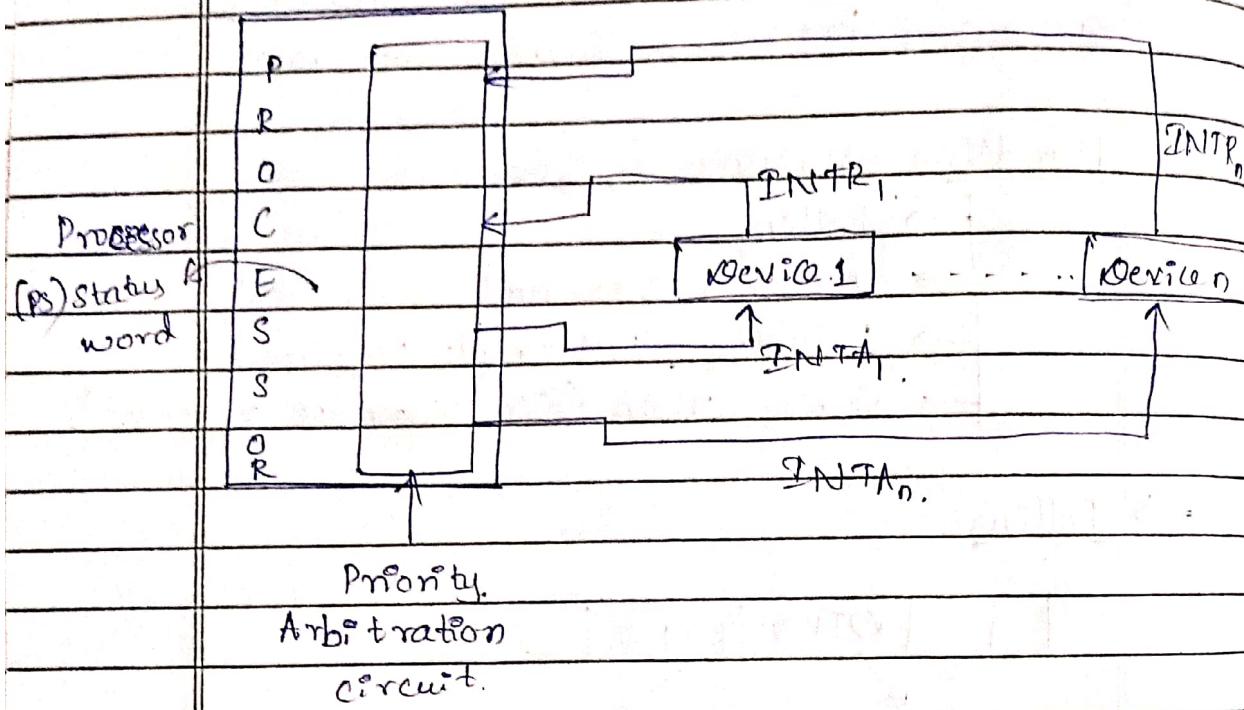
Vectored interrupt:



address & IRR address.

- * In vectored interrupt scheme device requesting interrupt identifies itself to the processor

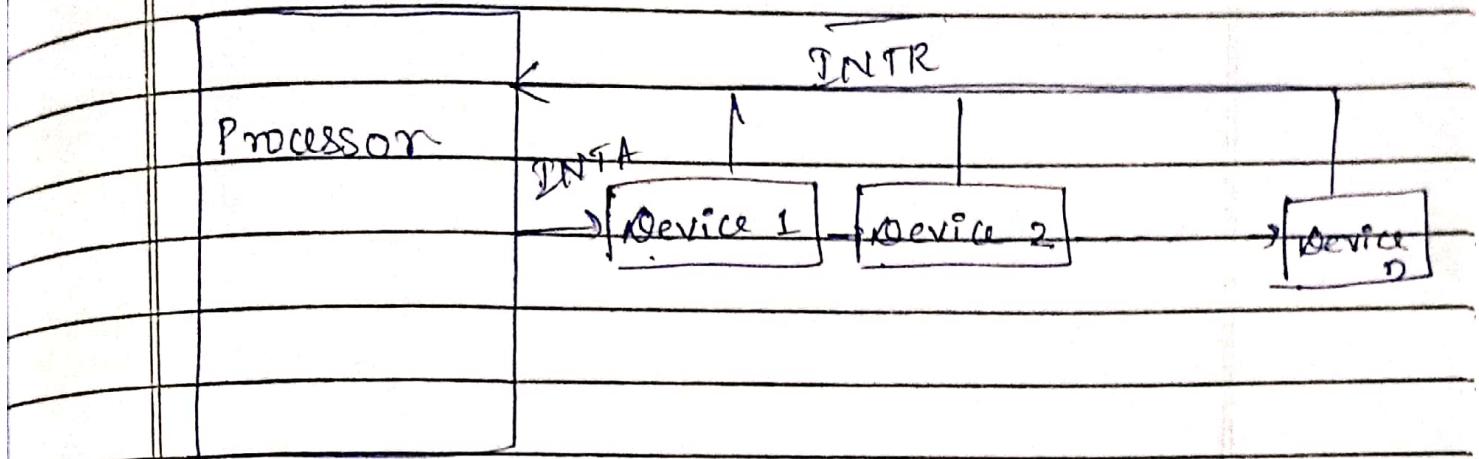
⇒ Priority (Interrupt Nesting)



- * In a multiple priority scheme processors priority level is changed according to the device priority that it is servicing.
- * In this scheme request is accepted only if it has a higher priority level than that is currently assigned to the processor.

⇒ Daisy chain

- * In daisy chain arrangement the device that electrically closest to the processor has the highest priority. The second device along the chain will have second highest priority



+ Devices can be organized into groups and each group is connected at a different priority level.

Within a group devices are connected in daisy chain.