



POORNIMA

COLLEGE OF ENGINEERING

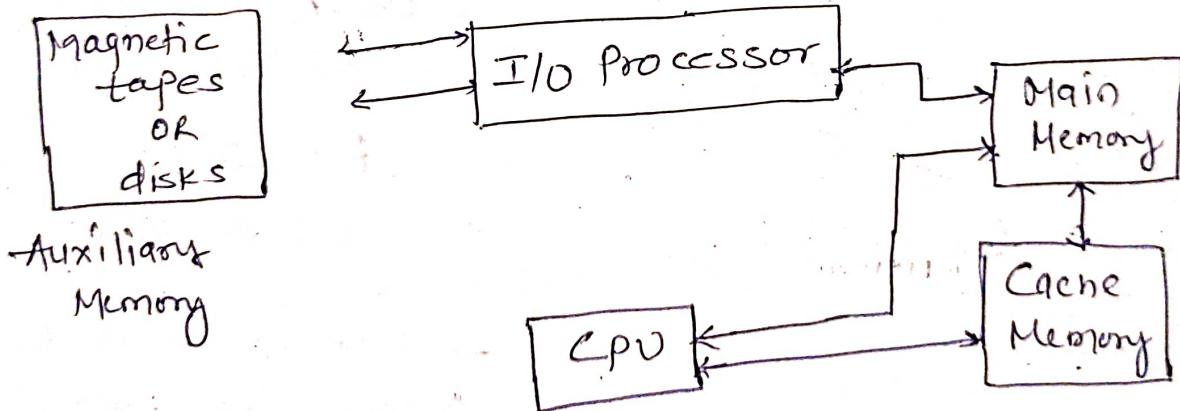
DETAILED LECTURE NOTES

UNIT-5

PAGE NO.

Memory Hierarchy

The memory unit that directly communicates with ~~memory~~ CPU is called main memory / Primary memory
Devices that provide backup storage are called secondary memory / Auxiliary memory.
- Magnetic Disk, or tapes are most common auxiliary memory



Cache memory: A special very high speed memory

Called a cache,

- It is B/W the CPU & Main memory
- This is used to compensate the speed difference B/W the main memory & CPU.

- Cache Memory hold the address of the frequently used data in Main memory
- The purpose of memory hierarchy is to obtain the highest possible average access speed while minimizing the total cost of the entire memory system.

Main memory

- It is the central storage in the computer system.
- main memory is based on Semiconductor Integrated Circuits

main memory

RAM

Random Access Memory

STATIC RAM - Made up of flop flops that store binary information

- volatile

DYNAMIC RAM - Made up of capacitors that accumulate charge to store information

- volatile

- need to refresh after particular interval

ROM

Read only memory

Types

PROM

EPROM

EEPROM

- ROM is need to store the initial program called Boot strap loader
- to start computer POST is run to check the hardware
- POST → Program on self test

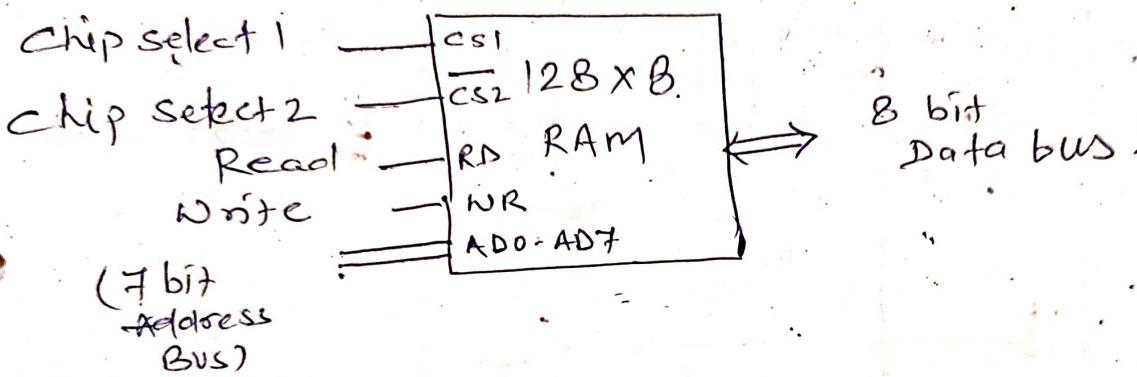


POORNIMA

COLLEGE OF ENGINEERING

DETAILED LECTURE NOTES

PAGE NO.



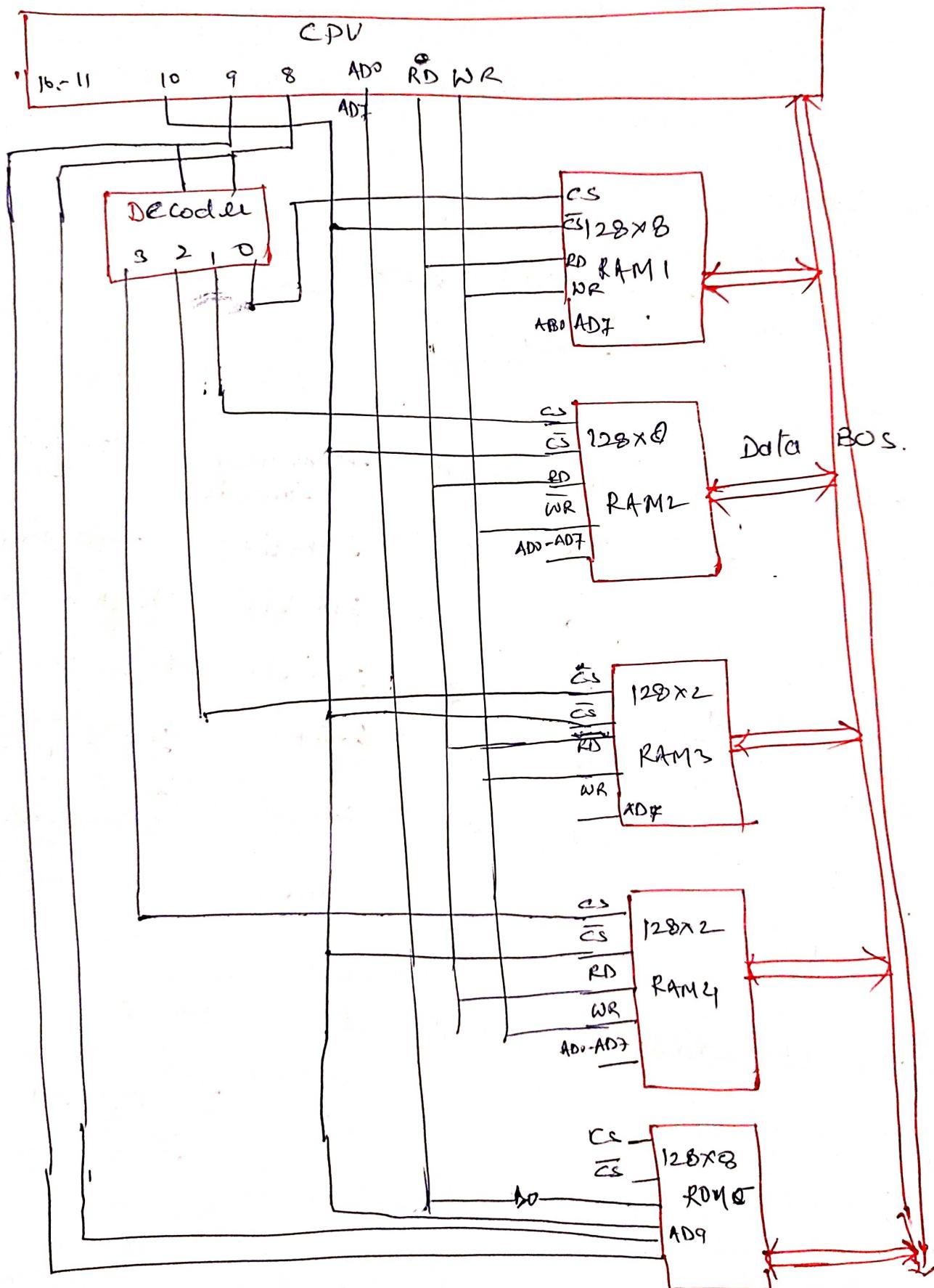
CS1	CS2	RD	WR	memory function	state of Data bus
0	0	x	x	Inhibit	high impedance
0	1	x	x	Inhibit	high impedance
1	0	0	0	Inhibit	high impedance
1	0	0	1	Write	low impedance
1	0	1	x	Read	Input Date to RAM
1	0	x	x	Inhibit	Output date from RAM
1	1	x	x		high impedance

No. of Addline is calculate $\log_2 (\text{size})$

$$\log_2 (128) = 7$$

Auxiliary Memory

Memory Connection to CPU





POORNIMA

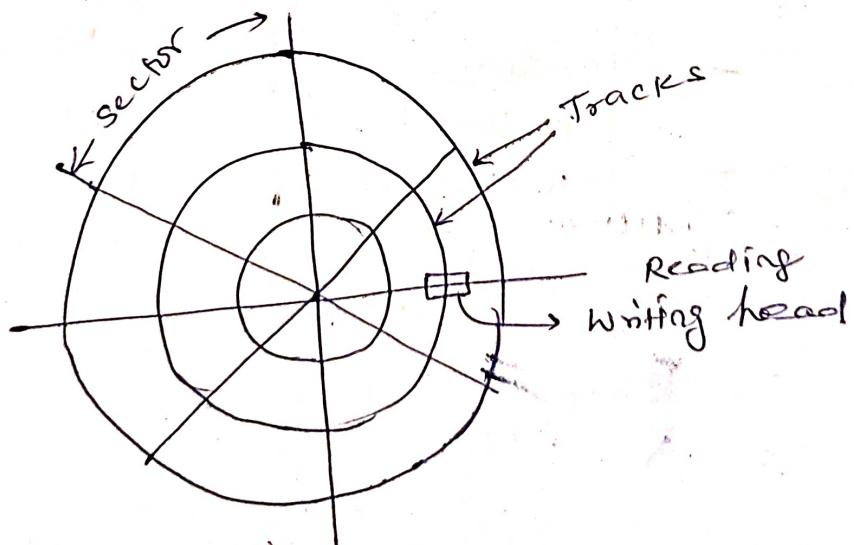
COLLEGE OF ENGINEERING

DETAILED LECTURE NOTES

PAGE NO.

Auxiliary Memory

- Most Common memory used in computer systems is magnetic disk & tapes
- The Average time required to reach a storage location in memory & obtain its contents is called the access time.
- Bits are recorded as magnetic spot on the surface as it passes a stationary mechanism called a write head.

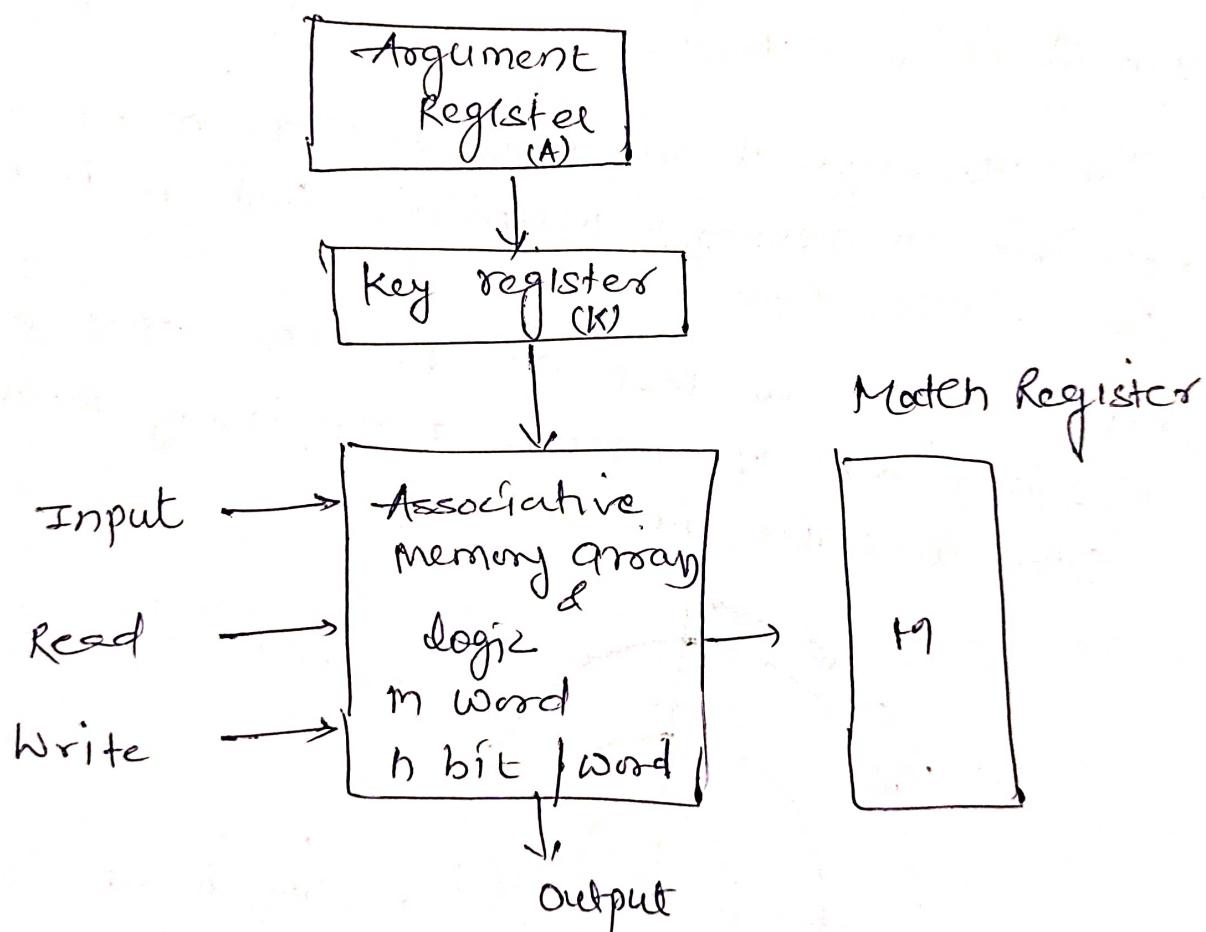


Magnetic disk

- DISKS permanently attached to the unit assembly and cannot be removed by occasional user are called Hard disk. & removable is called floppy disk

Associative Memory

- A memory unit accessed by Content is called the Associative memory or Content Addressable memory
- This memory is accessed simultaneously and in parallel on the basis of data content rather than specific address or location.



A

101 1111 00

K

111 0000 00

Word1

100 1111 00 no match

Word2

101 0000 01 Match.



Poornima

COLLEGE OF ENGINEERING

DETAILED LECTURE NOTES

PAGE NO.

Match logic:

The Match logic for each word can be derived from the Comparison algorithm for two binary numbers.

$$x_j = A_1 F_{ij} + A_2 F_{i'j}$$

Cache Memory:-

- Analysis of a large number of typical programs has shown that the references to memory at any given interval of time tend to be confined within a few localised memory areas. This is known as Locality of Reference.
- The performance of Cache memory is frequently measured in terms of a quantity called hit ratio.
- When CPU refers to memory and finds the word in Cache, it is said to be hit. If the word is not found in the Cache, it is in main memory, it is counted as miss.
- The ratio of the number of hits divided by the total CPU references to memory is the hit ratio.

The transferring of data from main memory to Cache memory is referred to as mapping process.

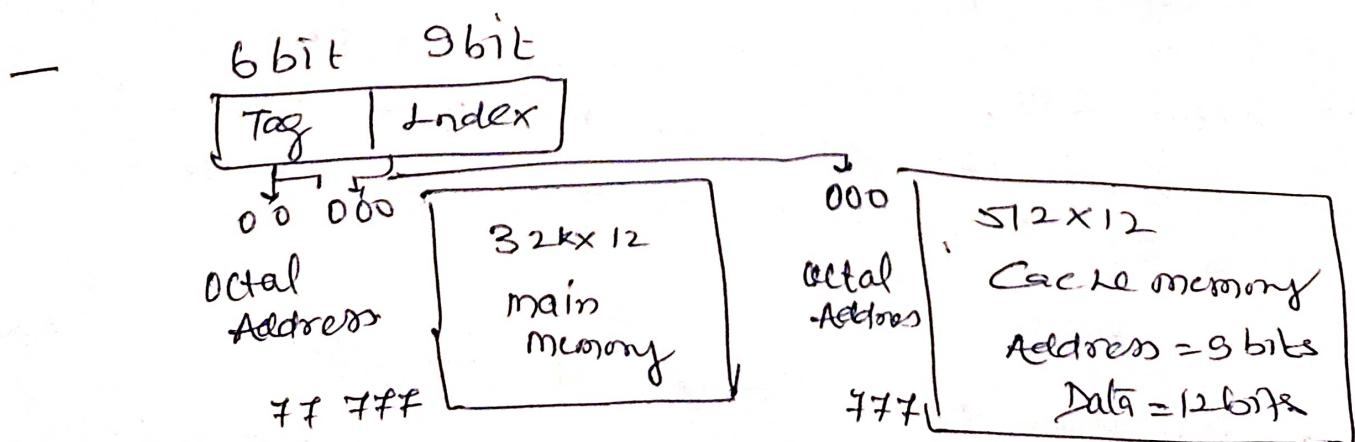
1. Associative mapping
2. Direct mapping
3. Set Associative mapping

Associative mapping

- The fastest & flexible Cache organization uses an associative memory.
- The associative memory stores the address & content of the memory word.

Direct mapping

- The direct memory mapping cache organization uses the n bit address to access the main memory and the k bit to index to access the cache.
- When CPU generates the ~~tag~~ memory request, the index field is used to for the address to access the cache.





POORNIMA

COLLEGE OF ENGINEERING

DETAILED LECTURE NOTES

PAGE NO.

Set Associative Mapping

It is an improvement over the Direct memory organization in that each word of Cache can store two or more words of memory under the same Index address.

Index address

- In general a set associative cache of size K will accommodate K words of main memory in each word of cache
- Common Replacement Algorithms are:
 - ① first in first out
 - ② least recently used
 - ③ random replacement
 - ④ most recently used
- The cache is initialised by clearing all the bits to 0. The valid bit of a particular cache word is set to be 1. the first time this word is loaded from main memory and stays unless the cache has to be initialised again.

Virtual Memory:

- Virtual memory is a concept used in some large computer systems that permits the user to construct programs through a large memory space were available, equal to the total of auxiliary memory.
- Each address that is used to give programmers the illusion that they have a very large memory at their disposal. even though the computer actually has a small memory
- Each Address that is referred by the CPU goes ~~to the~~ through an address mapping from so called Virtual Address

Address Space & Memory space

- The address used by the programmer will be called Virtual Address and the set of such addresses the address space
- An address in memory is called a location or physical address
- The set of such location is called memory space
- Address space N and the memory space by M , we have than have for this example $N=1024K$ and $M=32K$.



POORNIMA

COLLEGE OF ENGINEERING

DETAILED LECTURE NOTES

PAGE NO.

Virtual
Address



Virtual
Address
Register



Memory
Mapping
Table



Memory
Table
Buffer
Register

main
memory
Address
Register

Main
memory



Main memory
Buffer
Register

Address Mapping Using Pages:-

- Physical Memory is broken into group of equal size called Block.
- The term page refer to the group of address space of the same size
- The page frame is some time used to denote a block.

Associative Memory Page Table

- a system with n pages and m block would require a memory page table of n locations of which upto m blocks will be marked with block

numbers and all other will be empty.

Page Replacement

- A virtual memory system is a combination of both SW & HW.
- The memory management system software handles all the software operations for the efficient utilisation of memory space
- It must include
 - * Main memory ought to be removed to make room for a new page
 - * When new page is to be transferred from auxiliary memory to main memory
 - * Where the page is to be placed in the main memory



POORNIMA

COLLEGE OF ENGINEERING

DETAILED LECTURE NOTES

PAGE NO.

Multiprocessors

Characteristics of multiprocessors

- The SIMD systems are known as multiprocessors
 - multiprocessors can mean either a central processing unit or an input-output processor.
 - The benefit of multiprocessor organisation is an improved system performance
 - multiple independent jobs can be made to operate parallel.
 - A single job can be partitioned into multiple parallel task.
 - A multiprocessor system with common shared memory as a shared or tightly coupled multiprocessor.
 - Each system has its own private local memory. It is known as distributed or loosely coupled system.

Interconnection Structures

The components that form a multiprocessor system are CPUs, IOP's connected to input output devices and memory unit that may be partitioned into a number of separate modules.

- Several physical forms available for establishing an interconnection network.

1. Time shared Common bus.

2. Multipoint Memory

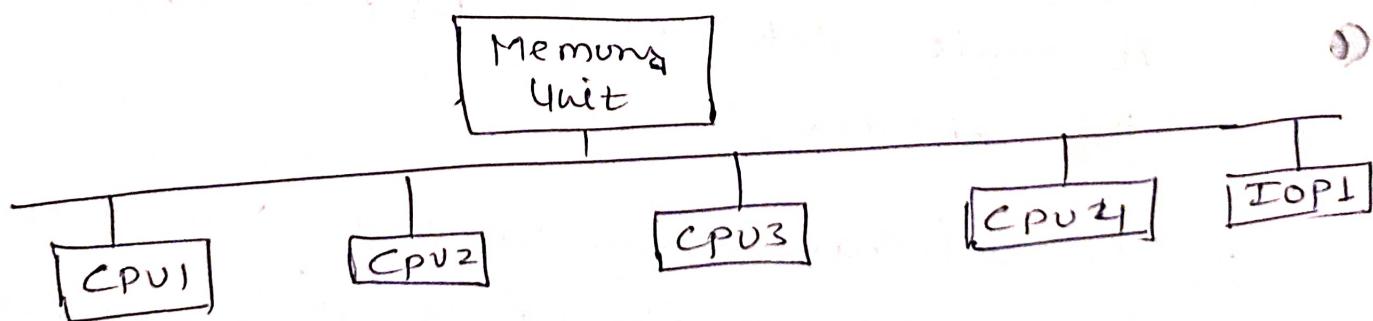
3. Crossbar Switch.

4. Multistage switching network.

5. Hyper Cube System.

① Time shared Common Bus

- Number of processors connected through a common bus.



- Only one processor can communicate with the shared memory and other common resources through the system bus at any given time.



POORNIMA

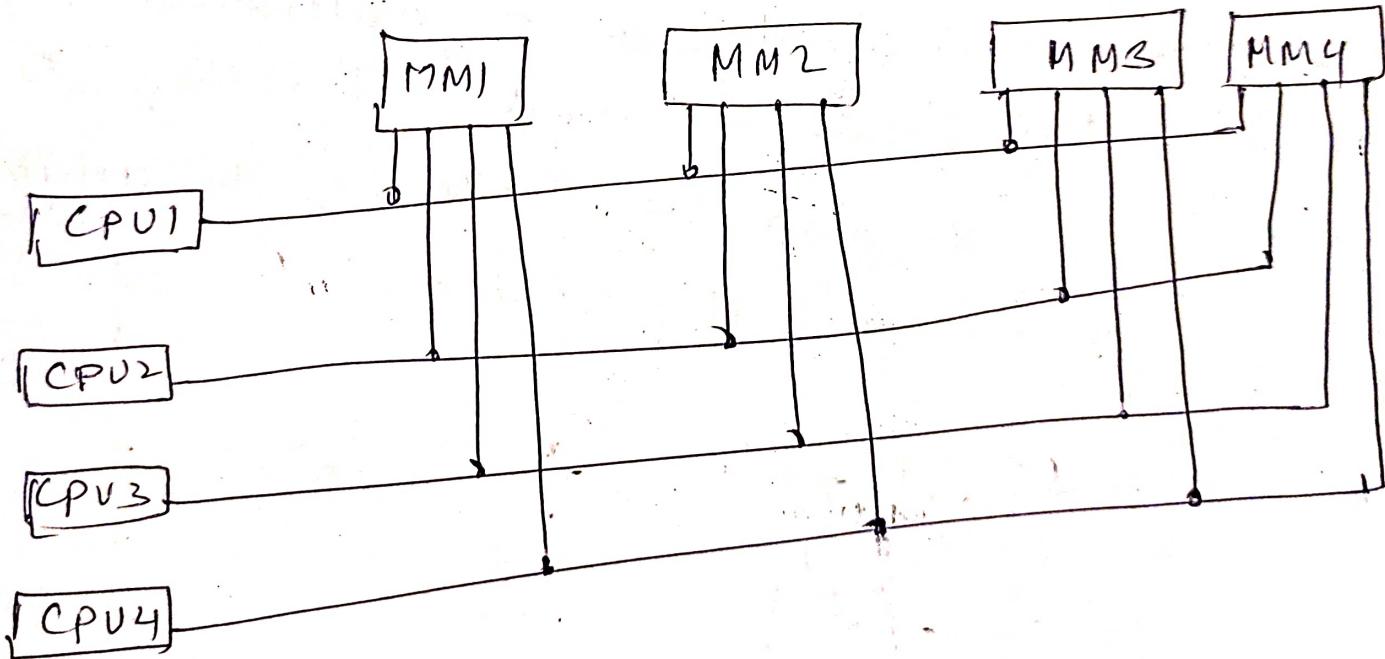
COLLEGE OF ENGINEERING

DETAILED LECTURE NOTES

PAGE NO.

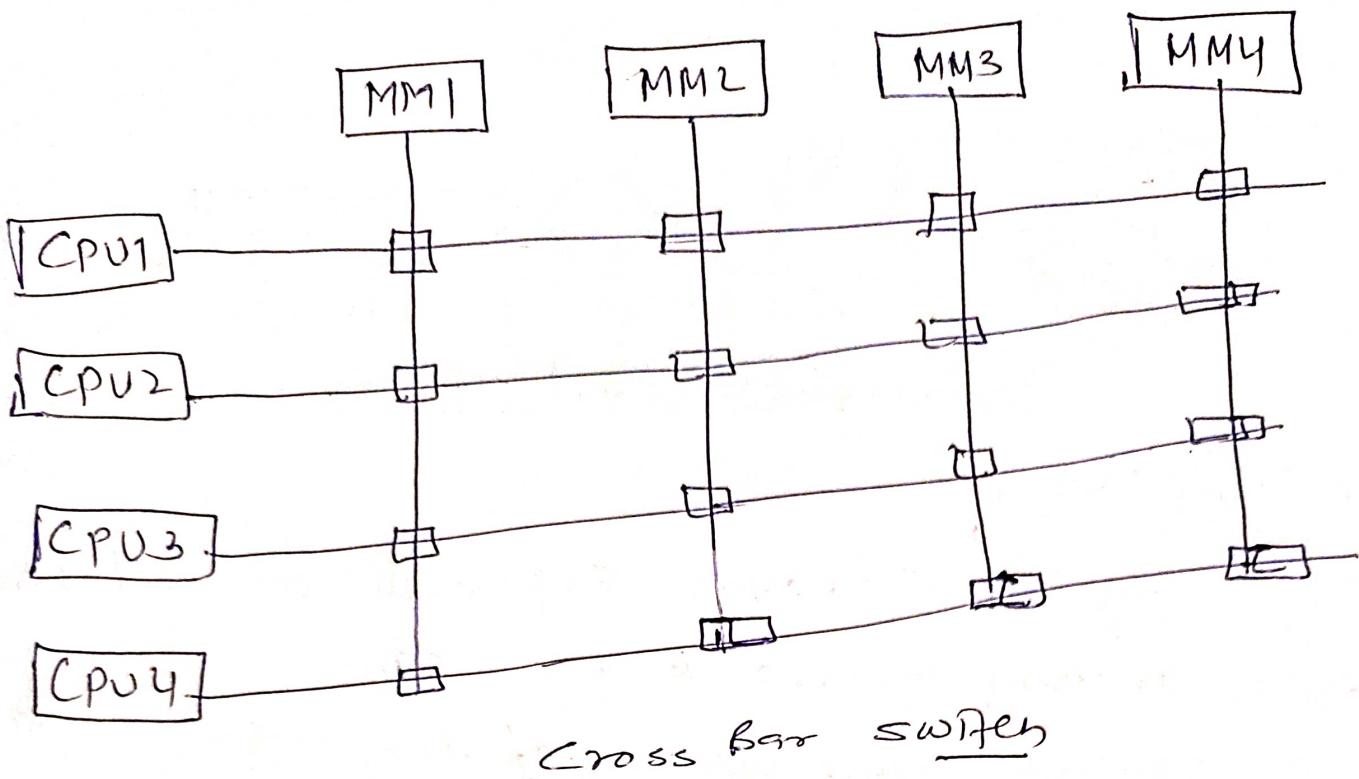
Multipoint Memory :-

- This system employs separate buses between each memory module & each CPU.
- A processor Bus consist of the address, data & control bus to communicate with the memory.



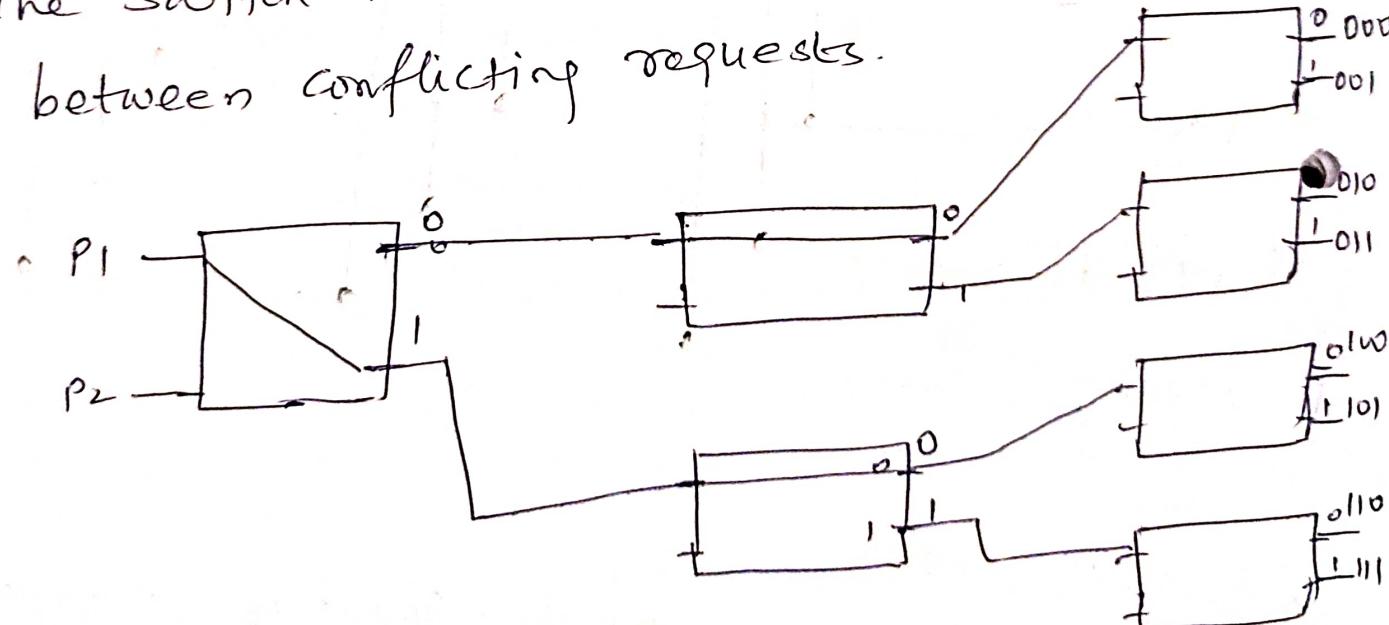
Crossbar switch

The crossbar switch organization consists of a number of cross points that are placed at intersections b/w the processor buses & memory module paths.



Multi stage switching Network -

- A basic component of a multi stage network is a two input, 2 output interchange switch.
- The switch has the capability to arbitrate between conflicting requests.



Binary Tree with 2x2 Switches



Poornima

COLLEGE OF ENGINEERING

DETAILED LECTURE NOTES

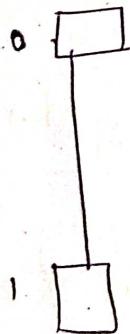
PAGE NO.

Hypercube Interconnection

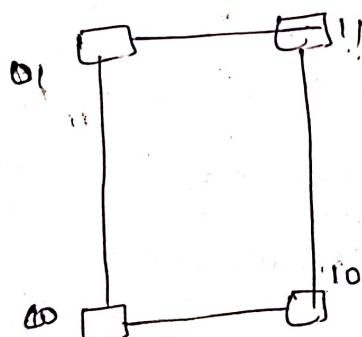
The hypercube or binary n cube multiprocessor

structure is a loosely coupled system composed of $N = 2^n$ processors interconnected in a n -dimensional binary cube.

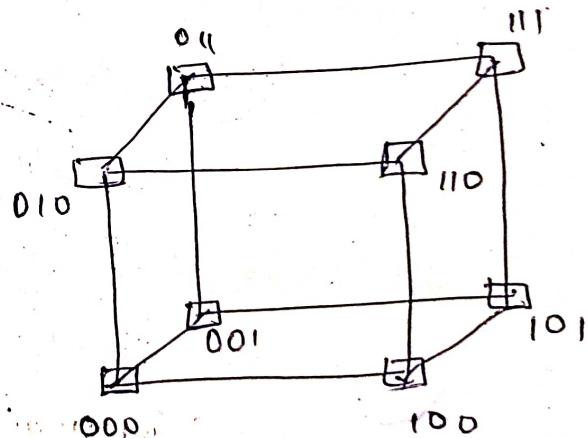
- Each processor forms a node of a cube



One Cube



Two Cube



Three Cube

Inter processor ~~Arbitration~~ Arbitration:-

- Computer system contains a number of buses at various levels to facilitate the transfer of information between components.
- A bus that connects major components in the multiprocessor system such as CPU, I/O's and memory is called a system bus.
- A system bus typically consists of 100 signal lines. These are divided into three functional groups.
 - Address Bus
 - Data Bus
 - Control Signal.
- Data transfer over the system bus may be synchronous & asynchronous.
- In synchronous bus, each data item is transferred during the time slice known in advance to both the source & destination unit. It is achieved by driving both units from common clock source.
- In asynchronous bus each data item being transferred is accompanied by handshaking control signals to indicate when the data are transferred from the source and received by destination.



POORNIMA

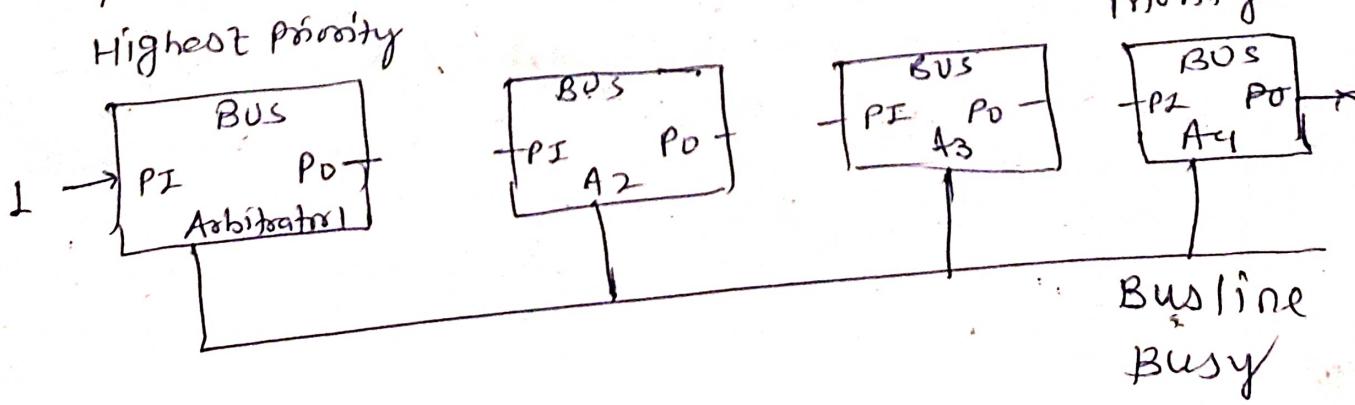
COLLEGE OF ENGINEERING

DETAILED LECTURE NOTES

PAGE NO.

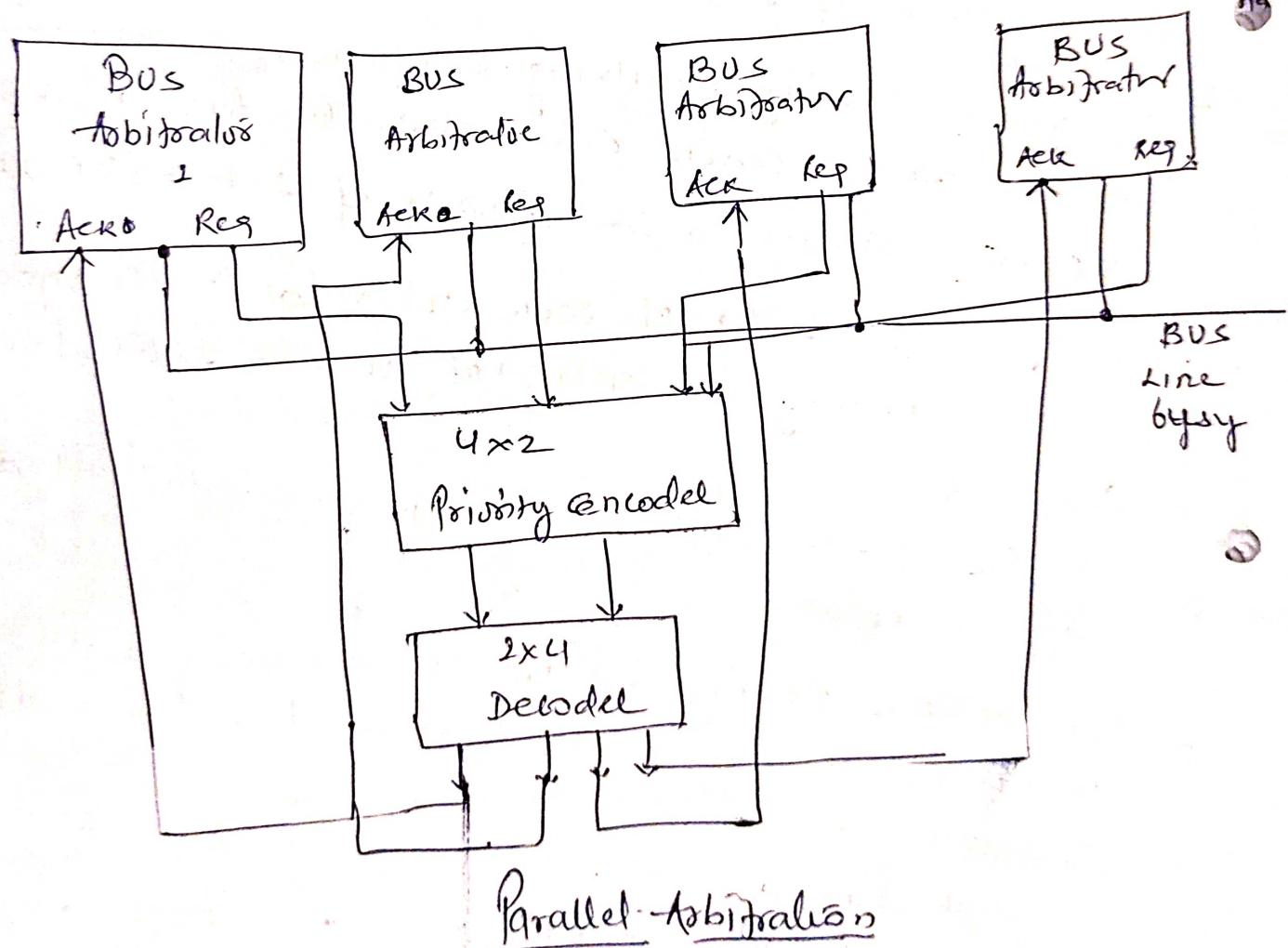
Serial Arbitration Procedure

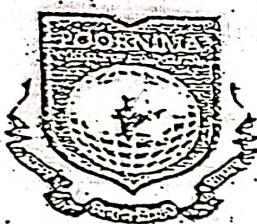
- Arbitration procedures services all processors requests on the basis of established priorities.
- The serial priority resolving technique is obtained from a daisy chain connection of bus arbitration circuits similar to the priority of interrupt logic.
- The priority out (PO) of each arbitrator is connected to the priority in (PI) of the next lower priority arbitrator.
- The PI of the highest priority is maintained at a logic 1 value.
- The highest priority unit in the system will always receives access to the system but when it requests it.



Parallel Arbitration logic:

- The Parallel Bus Arbitration technique uses an external priority encoder and a decoder.
- Each bus Arbitrator in the Parallel scheme has a bus request output line and a bus acknowledge input line.
- Each arbitrator enables the request line when its processor is requesting access to the system bus.





POORNIMA

COLLEGE OF ENGINEERING

DETAILED LECTURE NOTES

PAGE NO.

Dynamic Arbitration Algorithms

The Two bus Arbitration procedures just described use a static priority algorithm since the priority of each device is fixed by the way it is connected to the bus.

- In the dynamic priority algorithm gives the system the capability for changing the priority of the devices while the system is in operation.
- The time Slice Algorithm allocates a fixed length time slice of bus time that is offered sequentially to each processor in Round Robin fashion.
- The Least Recently Used (LRU) algorithm gives the highest priority to the requesting device that has not used the bus for the longest interval.
- The rotating daisy chain procedure is a dynamic extension of the daisy chain algorithm. In this

There is no central bus controller and the priority line is connected from the priority out of the last device back to the priority in the first device in a closed loop.

Interprocessor Communication & synchronization

- Multiprocessor system must be provided the facility for communicating with each other
- A communication path can be established through common Input Output channels.
- In a shared memory multiprocessor system the most common procedure to set aside a portion of memory to that is accessible to all processors.
- In addition to shared memory, a multiprocessor system may have other shared resources.
- for example, a magnetic disk storage unit connected to an I/O may be available to all CPU's.
- Interprocessor synchronization

The instruction set of a multiprocessor contains basic instructions that are used to implement communication and synchronization between co-operating processes.

Cache Coherence

In a shared memory system, all the processors share a common memory. In addition, each processor may have a local memory, part or all of which may be a cache. The compelling reason for having separate caches for each processor is to reduce the average access time in each processor.

The same information may reside in a number of copies in same caches and main memory. To ensure the ability of the system to execute memory operation correctly the multiple copies must be kept identical. This requirement imposes a Cache Coherence problem. A memory scheme is coherent if the value returned on a load instruction is always the value given by the last store instruction with same address.

Condition for Incoherence

Cache Coherence problem exist in multi processor with private caches because of the need of share writeable data. Read only data can safely be replicated without cache coherence enforcement mechanism.



POORNIMA COLLEGE OF ENGINEERING DETAILED LECTURE NOTES

PAGE NO.

- Communication refers to the exchange of data between different processes.
- Synchronization refers to the special case where the data used to communicate b/w processes is control information.
Synchronization is needed to enforce the correct sequence of processes and to ensure mutually exclusive access to shared writable data.

Mutual exclusion with a semaphore

- A properly functioning multiprocessor system must provide a mechanism that will guarantee orderly access to shared memory and other shared resources.
- It is necessary to protect data from being changed simultaneously by two or more processors. This mechanism has been termed as mutual exclusion.



POORNIMA

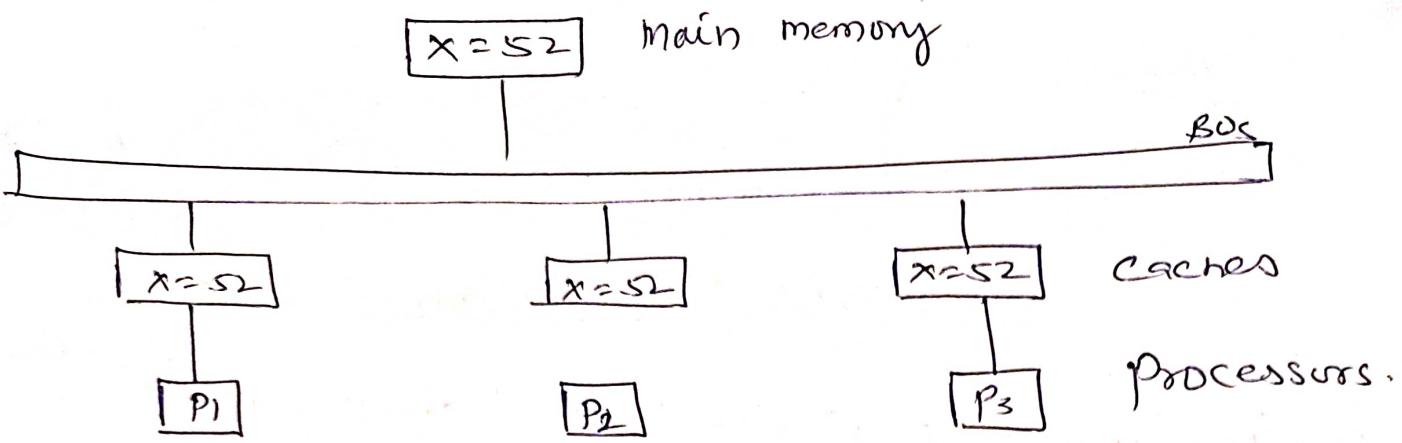
COLLEGE OF ENGINEERING

DETAILED LECTURE NOTES

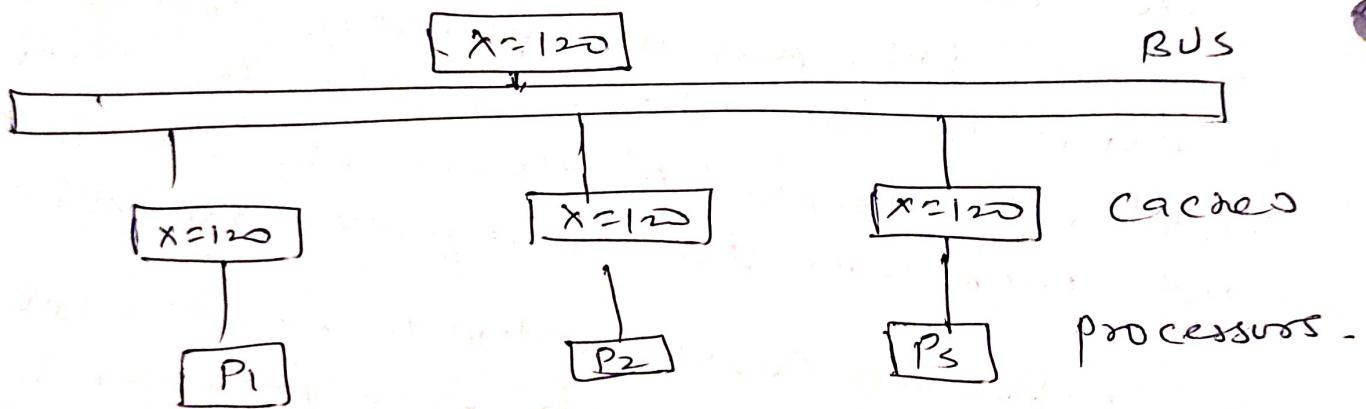
PAGE NO.

Solution to Cache Coherence Problem:

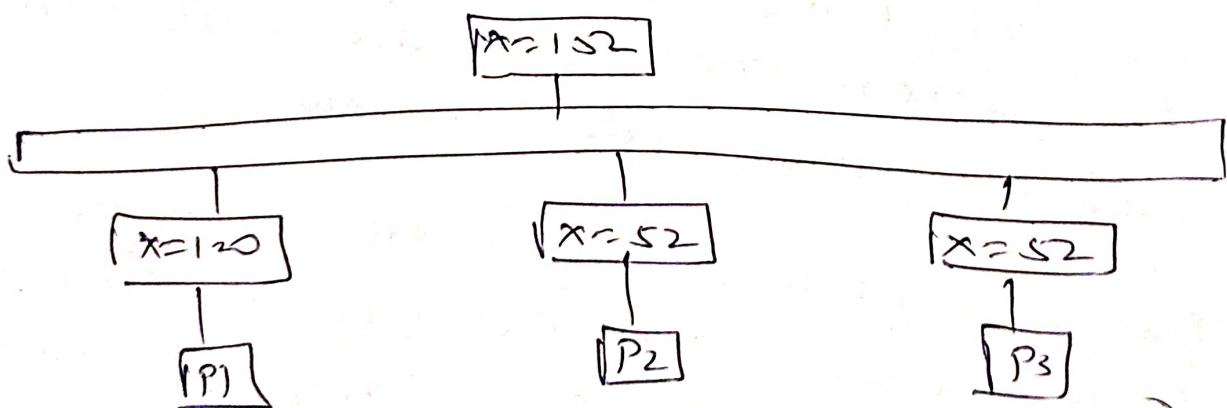
- A simple scheme is to disallow private caches for each processor and have a shared cache memory associated with main memory. Every data access is made to the shared cache. This method violates the principle of closeness of CPU to cache and increases the average memory access time.
- For performance consideration it is desirable to attach a private cache to each processor. One scheme that has been used to allow only non shared and read only data to be stored in caches. Such items are called cacheable.
- Shared writable data are non cacheable.
- The compiler must tag a data as either cacheable or non cacheable and the system hardware makes sure that only cacheable data are stored in caches. The non cacheable data remain in main memory.



Cache configuration after a load on x .



With write through Cache policy



With write back Cache policy

End