

## Pipeline & Vector processing multiprocessors, multi computers

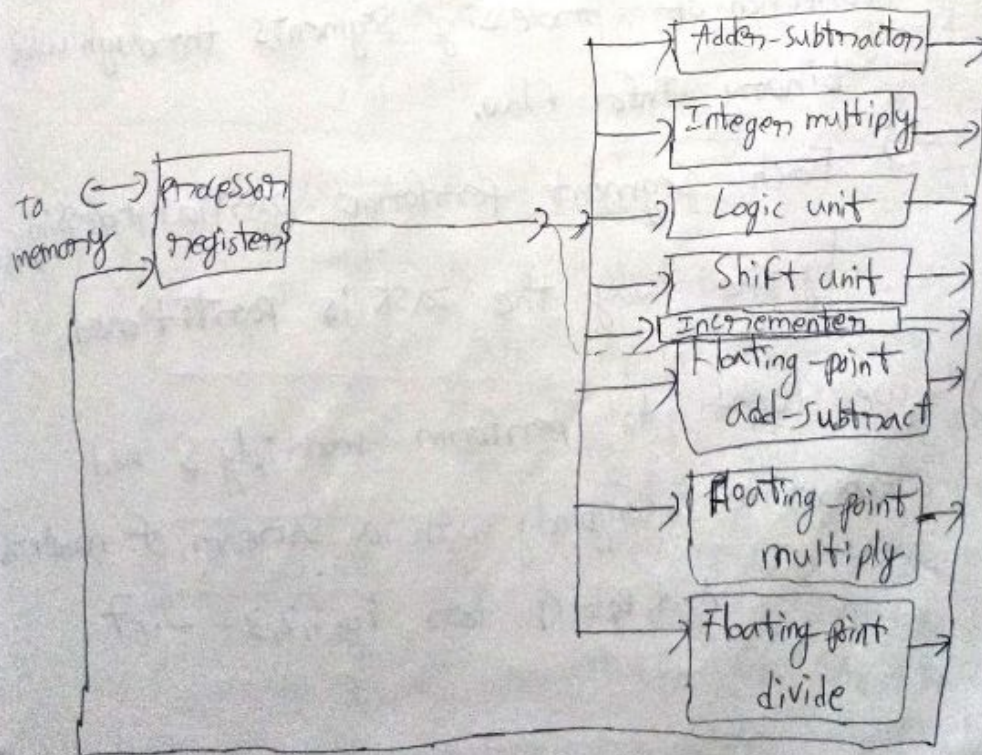
### 1) parallel processing:-

→ It is a term used to denote a large class of techniques that are used to provide simultaneous data-processing tasks for the purpose of increasing the computational speed of a computer system.

→ Instruction is being executed in the ALU, next instruction can be read from memory.

The system may have 2/more ALU's & be able to execute 2/more instructions at the same time.

→ Furthermore, <sup>the</sup> system may have 2/more processors operating concurrently.





- Single Instruction Stream & single data Stream (SISD)
- Single Instruction Stream & multiple data Stream (SIMD)
- Multiple Instruction Stream & single data Stream (MISD)
- Multiple Instruction Stream & multiple data Stream (MIMD)

⇒ pipelining:-

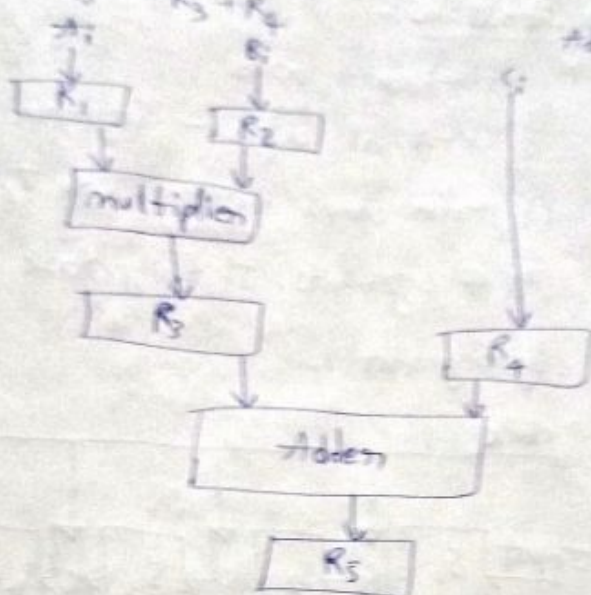
- A technique of decomposing a sequential process into sub operations, with each sub process being executed in a special dedicated segment that operates concurrently with all other segments. A pipeline can be visualized as a collection of processing segments through which binary info flows.
- Each segment performs partial processing by the way the task is partitioned.
- we want to perform multiply & add operations combinely with a stream of numbers  
 $A_i * B_i + C_i$  for  $i = 1, 2, 3, \dots, 7$



$$R_1 \leftarrow A_i, R_2 \leftarrow B_i$$

$$R_3 \leftarrow R_1 * R_2, R_4 \leftarrow C_i$$

$$R_5 \leftarrow R_3 + R_4$$



Input  $A_i, B_i$

multiply  $R_1$  & input  $C_i$

Add  $C_i$  to product

### Example of pipeline processing

Clock pulse number	Segment 1		Segment 2		Segment 3
	$R_1$	$R_2$	$R_3$	$R_4$	$R_5$
1	$A_1$	$B_1$	—	—	—
2	$A_2$	$B_2$	$A_1 * B_1$	$C_1$	—
3	$A_3$	$B_3$	$A_2 * B_2$	$C_2$	$A_1 * B_1 + C_1$
4	$A_4$	$B_4$	$A_3 * B_3$	$C_3$	$A_2 * B_2 + C_2$
5	$A_5$	$B_5$	$A_4 * B_4$	$C_4$	$A_3 * B_3 + C_3$
6	—	—	$A_5 * B_5$	$C_5$	$A_4 * B_4 + C_4$
7	—	—	$A_6 * B_6$	$C_6$	$A_5 * B_5 + C_5$
8	—	—	$A_7 * B_7$	$C_7$	$A_6 * B_6 + C_6$
9	—	—	—	—	$A_7 * B_7 + C_7$



## Contents of Registers in Pipeline Example

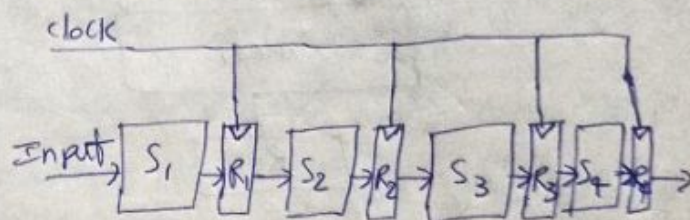
⇒ General considerations:-

Horizontal axis — clock cycles,

Vertical axis — Segment number,

Diagram shows tasks  $T_1$  to  $T_6$  executed

in 4 segments. Initially, task 1 is handled by Segment 1.



Four-segment pipeline

$$S = \frac{nt_n}{(K+n-1)t_p}$$

Space-time diagram for pipeline

		1	2	3	4	5	6	7	8	9	
											→ clock cycle
Segment	1	$T_1$	$T_2$	$T_3$	$T_4$	$T_5$	$T_6$				
	2		$T_1$	$T_2$	$T_3$	$T_4$	$T_5$	$T_6$			
	3			$T_1$	$T_2$	$T_3$	$T_4$	$T_5$	$T_6$		
	4				$T_1$	$T_2$	$T_3$	$T_4$	$T_5$	$T_6$	

As no. of tasks ↑,  $n$  becomes much larger than  $K-1$ , &  $K+n-1$  approaches the value of  $n$



line Example

under this condition, the speedup becomes

$$S = \frac{(k+1)t_n}{(k+1)t_p} = \frac{t_n}{t_p}$$

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

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

### ⇒ Instruction pipeline:-

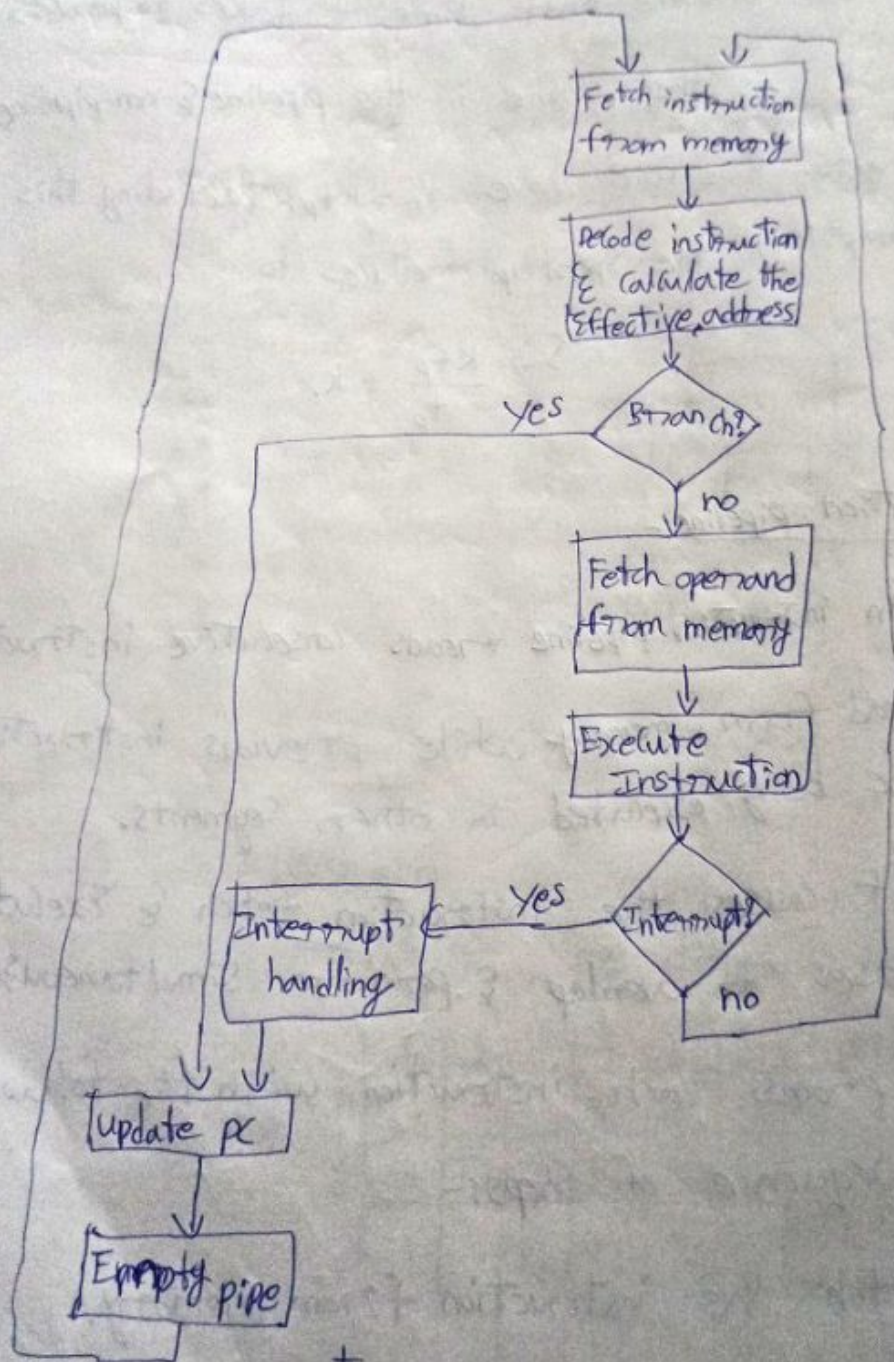
- An instruction pipeline reads consecutive instructions read from memory while previous instructions are being executed in other segments.
- It causes the instruction fetch & execute phases to overlap & perform simultaneous operations.

To process each instruction with the following sequence of steps:-

- a) Fetch the instruction from memory.
- b) Decode the instruction.
- c) Calculate the effective address.



- d) Fetch the operands (more) from memory.
- e) Execute the Instruction.
- f) Store the result in the proper place.



Four-Segment CPU pipeline



- Segments
- a) FI  $\rightarrow$  fetches an instruction.
  - b) DI  $\rightarrow$  Decodes the instruction & calculates the effective Address.
  - c) FO  $\rightarrow$  fetches the operand.
  - d) EX  $\rightarrow$  Executes the instruction.

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

we can't fetch any other instruction until branch is executed.

(True) | If 3 - Branch then FO & FI can't done at same time in steps.

Another delay may occur in pipeline, if EX Segment needs to store the result of operation, in data memory while FO segment needs to



Fetch an operand.

→ In that case, Fo must wait, until Ex finished its operation.

⇒ Vector processing:-

→ There is a class of computational problems that are beyond the capabilities of a conventional computer.

→ These problems are characterized by the fact that they require a vast no. of computations that will make a conventional computer days/even weeks to compute.

⇒ Applications:-

→ Long range weather forecasting

→ petroleum explorations

→ Seismic data analysis

→ medical diagnosis

→ Aerodynamics & Space flight simulations.

→ AI & expert systems.

→ mapping the human genome

→ Image processing

⇒ vector operation

→ many scientific  
on large  
usually  
floating

→ A vector  
array

→ I  
th



## → vector operations

→ many scientific problems require arithmetic operations on large arrays of numbers. These numbers are usually formulated as vectors & matrices of floating-point numbers.

→ A vector is an ordered set of a one-dimensional array of data items. A vector  $V$  of length  $n$  is

represented as a row vector by  $V = [V_1, V_2, V_3, \dots, V_n]$ .

→ It may be represented as a column vector if the data items are listed in a column. A conventional sequential computer is capable of processing operands one at a time. Consequently, operations on vectors must be broken down into single computations with subscripted variables.

(20)  $\begin{bmatrix} 20 \\ 20 \end{bmatrix}$   $I = 1, 100$

(20)  $\begin{bmatrix} 20 \\ 20 \end{bmatrix}$   $C(I) = B(I) + A(I)$

This is a program for adding 2 vectors  $A$  &  $B$  of length 100 to produce a vector  $C$ . This is implemented in machine language by the following sequence of operations.



Initialize  $I = 0$

DO (20) Read  $A(I)$

Read  $B(I)$

Store  $C(I) = B(I) + A(I)$

Increment  $I = I + 1$

If  $I \leq 1000$  to (20) DO

continue

$$C(1:100) = A(1:100) + B(1:100)$$

$\Rightarrow$  matrix multiplication:-

$$\begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix} \times \begin{bmatrix} b_{11} & b_{12} & b_{13} \\ b_{21} & b_{22} & b_{23} \\ b_{31} & b_{32} & b_{33} \end{bmatrix} = \begin{bmatrix} c_{11} & c_{12} & c_{13} \\ c_{21} & c_{22} & c_{23} \\ c_{31} & c_{32} & c_{33} \end{bmatrix}$$

$\Rightarrow$  mem

The product  $C$  is  $3 \times 3$ , elements  $(A \in B)$  by the inner product.

$$c_{ij} = \sum_{k=1}^3 a_{ik} * b_{kj}$$

For  $i=1$  no. in the first row & first column of matrix  $C$  is calculated by letting  $i=1$ ,  $j=1$ , to obtain

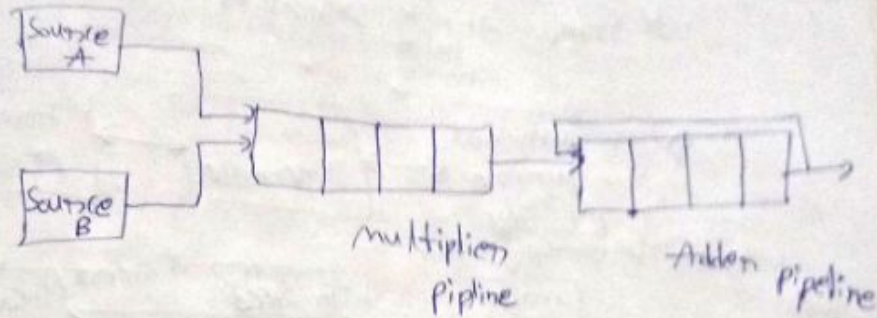
$$c_{11} = a_{11}b_{11} + a_{12}b_{21} + a_{13}b_{31}$$

Operation Code	Base address Source 1	Base address Source 2	Base address destination	Vector length
----------------	-----------------------	-----------------------	--------------------------	---------------

Instruction format for vector processor

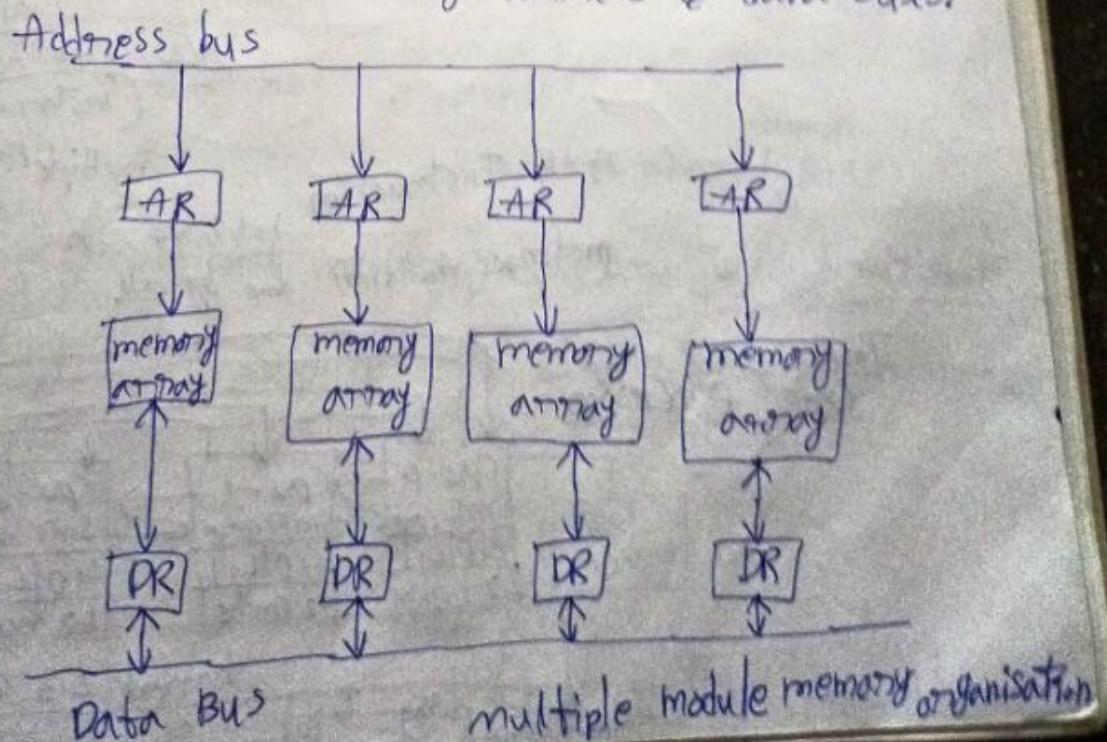


$$C = A_1 B_1 + A_5 B_5 + A_9 B_9 + A_{13} B_{13} + \dots \\ + A_2 B_2 + A_6 B_6 + A_{10} B_{10} + A_{14} B_{14} + \dots \\ + A_3 B_3 + A_7 B_7 + A_{11} B_{11} + A_{15} B_{15} + \dots \\ + A_4 B_4 + A_8 B_8 + A_{12} B_{12} + A_{16} B_{16} + \dots$$



⇒ memory Interleaving:-

→ Instead of using 2 memory buses for simultaneous access, memory can be partitioned into a no. of modules connected to a common memory address & data buses.





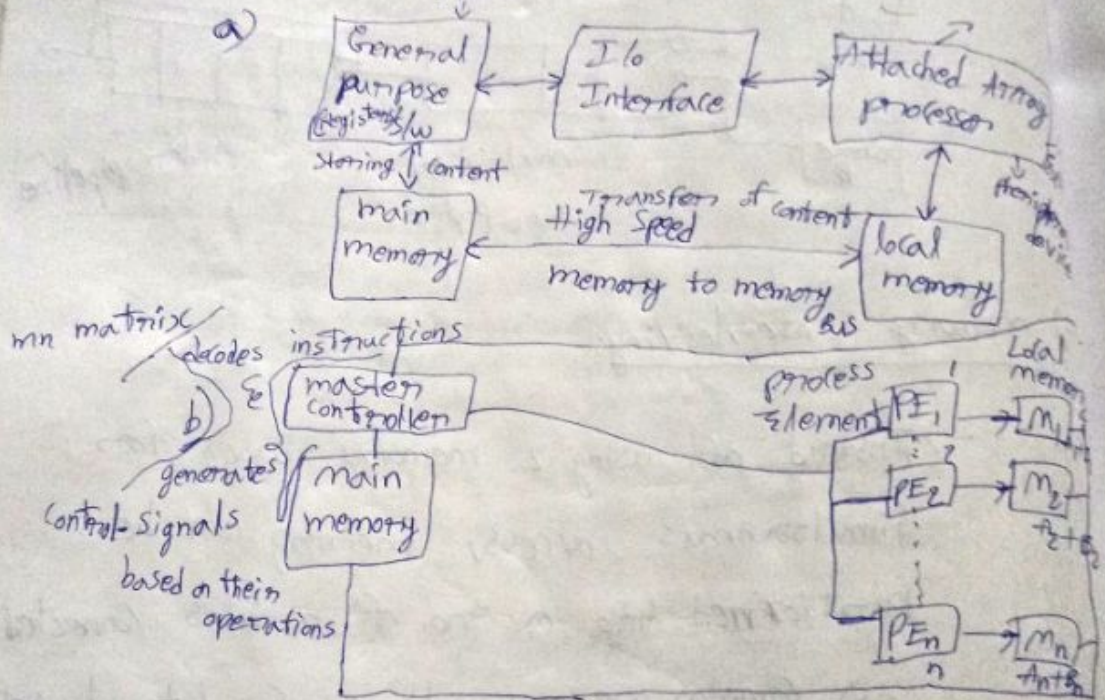
## ⇒ Array processor:-

To improve performance of the computer.

Single processor

a) Attached Array processor

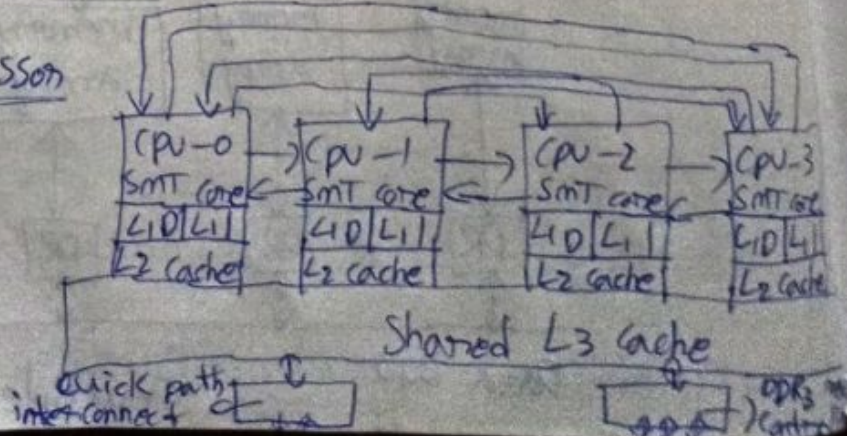
b) SIMD Array processor



Registers  
ALU, Pipeline (instructions) → multiple data

$C = A_i + B_i$  Instruction  $i = 1, 2, 3, 4, \dots, n$   
(Matrix Addition) done by all.

## Intel i7 core processor





## Characteristics of multiprocessors:-

Inter connection of 2/more (PUs with memory and input-output equipment).

Tightly coupled, Loosely coupled  
Shared memory / Distributed memory

MIMD  $\rightarrow$  multiple Instruction multiple data stream.  
computation can proceed in parallel in one of 2 ways.

$\rightarrow$  multiple independent jobs can be made to operate in parallel.

$\rightarrow$  A single job can be partitioned into multiple parallel tasks.

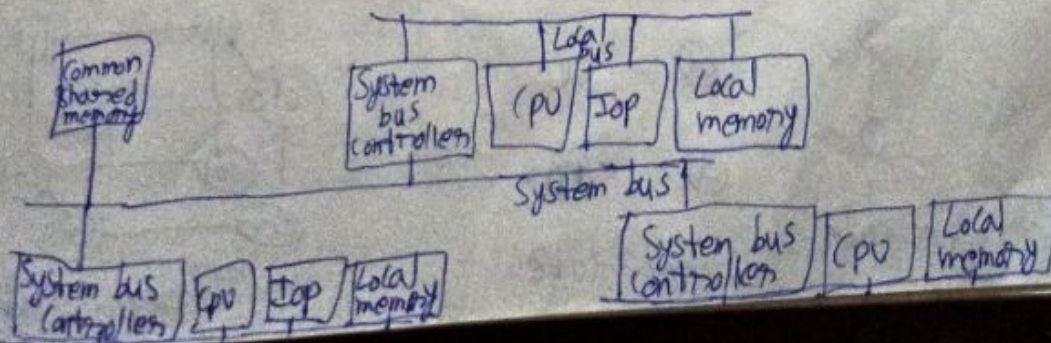
## Inter connection Structures:-

It has components:-

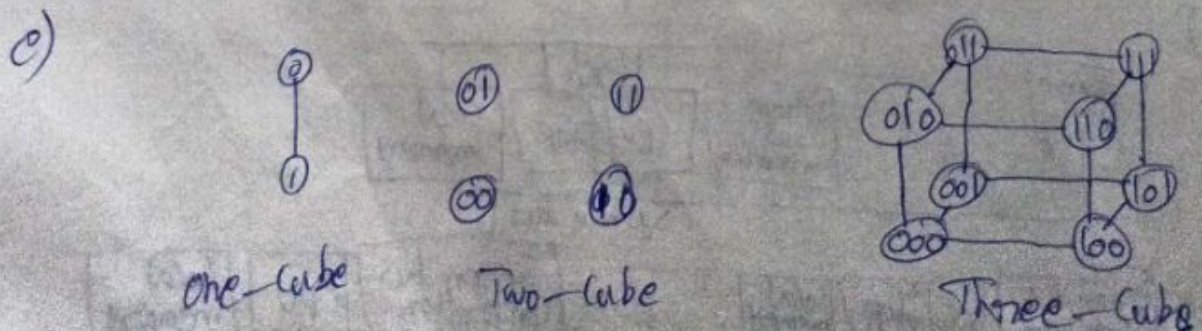
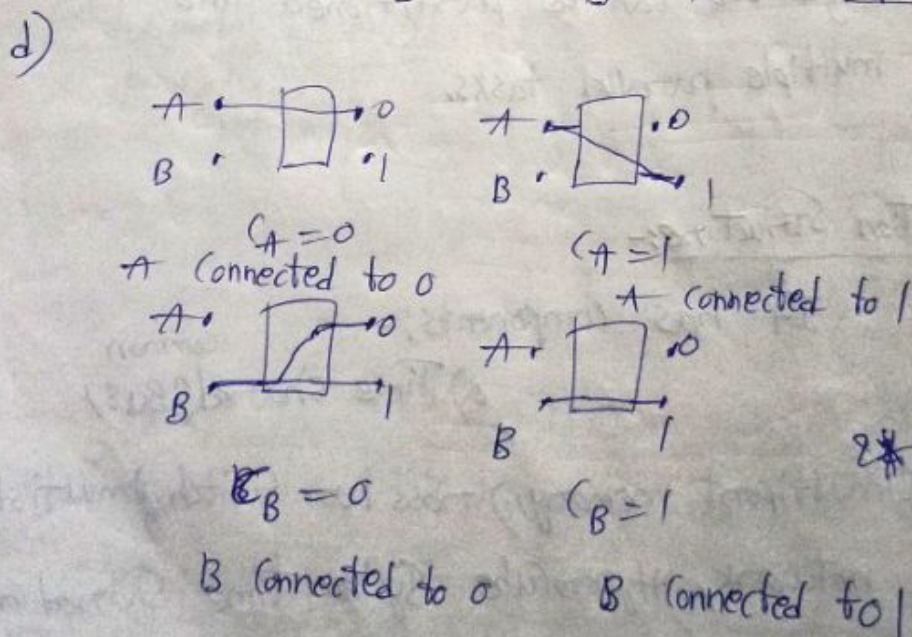
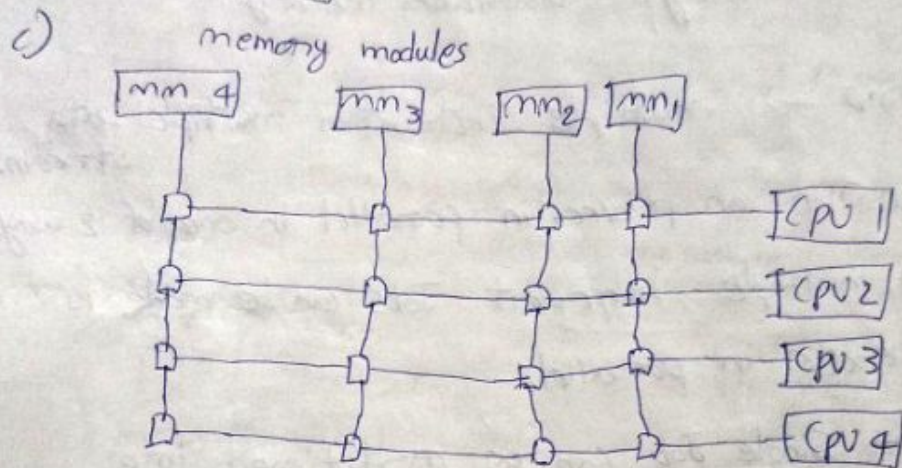
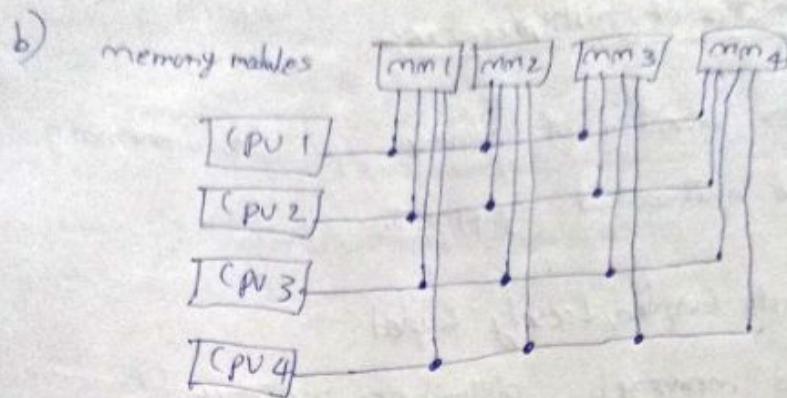
a) Time shared <sup>Common</sup> Bus /

b) multipoint memory, c) cross bar switch, d) multistage switching network, e) Hypercube System time shared common bus.

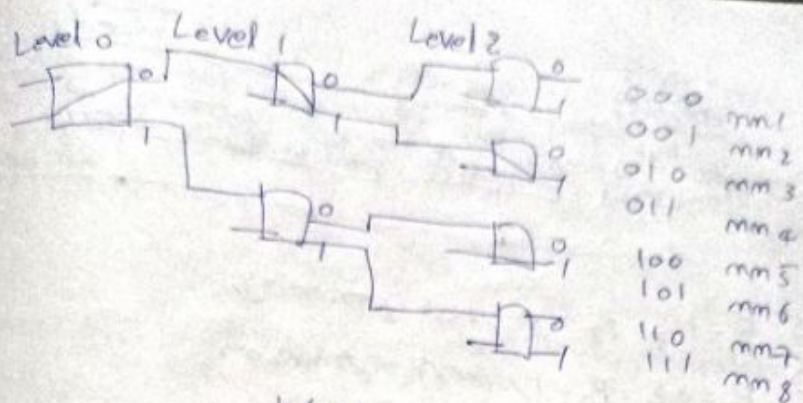
a)











1 to 8 way switch using 2x2 switch

### ⇒ Inter Connection Arbitration:-

A bus that connects major components in a multi-processor system, such as CPU, I/O & memory, is called a system bus. The processor in a shared memory (There will be) multiprocessor system request access to common (system) memory / other common resources through the system bus. If other processor is using the system bus, it must wait for some time. If many more, then arbitration logic must be performed / part of system bus controller placed b/w local bus & system bus.

### ⇒ Static Arbitration:-

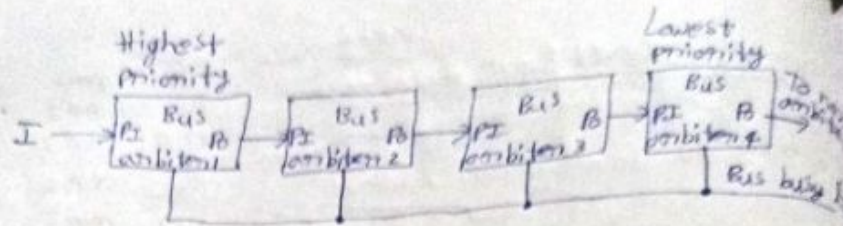
a) Daisy chain / Serial Arbitration.

b) Parallel Arbitration

### a) Daisy chain Arbitration:-

Based on serial mechanism, it

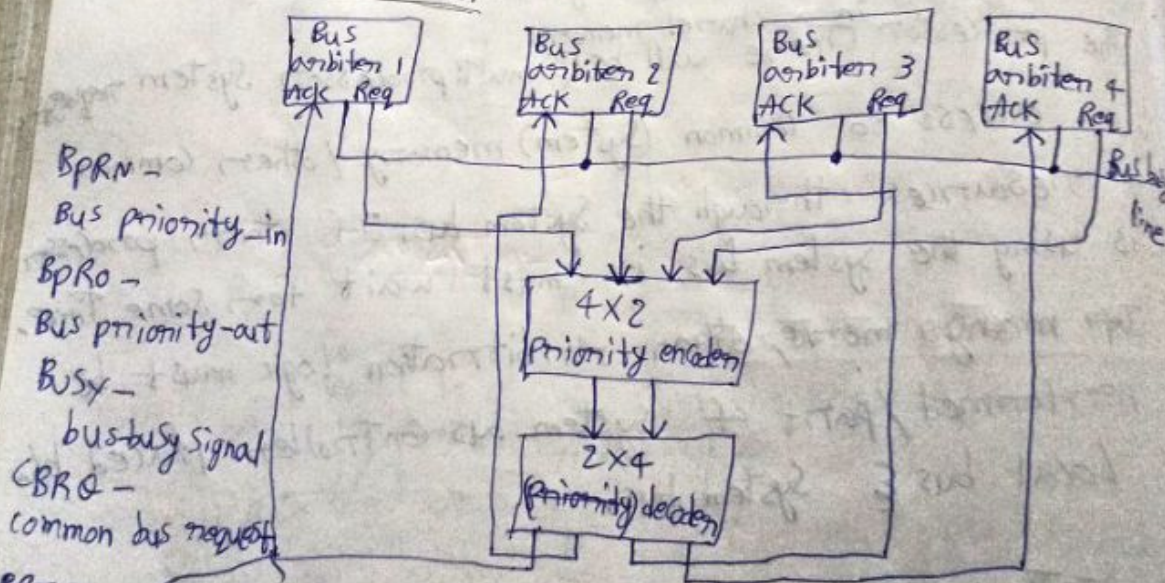




- $P_1 = 1 \ \& \ P_2 = 0$ , has priority.
- $P_1 = 0 \ \& \ P_2 = 1$ , Next processor
- $P_1 = 0 \ \& \ P_2 = 0$ , Lowest priority

when arbiter receives control of bus, it examines the busy line. If it is inactive, no other processor is using the bus. It (arbiter) activates the busy line & its processor takes control of the bus.

#### b) Parallel arbitration:-



BPRN - Bus priority-in  
 BPRO - Bus priority-out  
 BUSY - bus-busy signal  
 CBRO - common bus request  
 BREO - Bus Request  
 BCLK - Bus clock

Requests from 4 arbiters into encoder & 2 bit code will generated then decoder which enables proper acknowledge line to give to highest priority unit.



most priority  
Bus  
To next  
ambiten  
Bus busy line  
a) Dynamic Arbitration Algorithm:-

→ Time slice,

a) Time slice:-

No priority, in Round Robin fashion, Given time limit, processors will communicate with bus.

b) Polling:-

In this, bus grant signal is replaced by poll lines, which are connected to all units. Define address to each device, connected to the bus. If processor requires access, it recognizes its address & in a prescribed manner, it activates the bus busy line. Polling done by choosing a different processor, under program control.

c) LRU:-

(LRU) Least recently used processor that uses bus will have the highest priority.

d) FIFO:-

rotating daisy chain

The first processor will use the bus as first in first out principle.

Round bus connections as daisy chain & bus releasing ambiten will have the least priority.

e) Hardware performance issues:-

→ microprocessors have seen an exponential increase in performance

- Improved organization

- Increased clock frequency



- Increase in parallelism
  - pipelining, - superscalar (multi-issue)
  - simultaneous multi-threading (SMT)

→ Diminishing returns

- more complexity requires more logic
- Increasing chip area for coordination & signal transfer logic
- Harder to design, make & debug

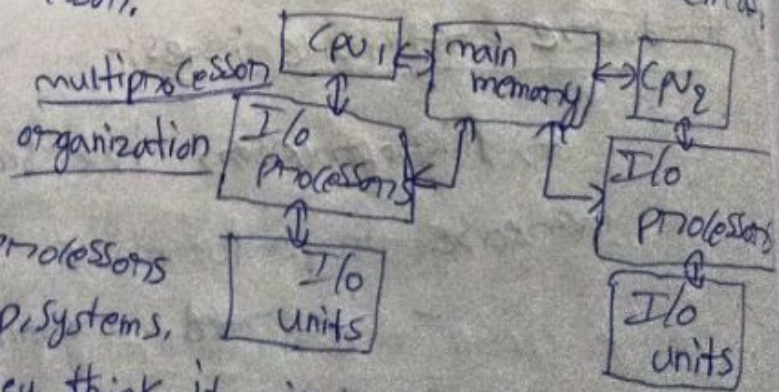
⇒ Software performance issues:-

- performance benefits dependent on effective exploitation of parallel resources (obviously).
- Even small partition impacts performance.
- many overheads of mpsoc; - (context switch, communication of work, cache coherence)
- Some apps effectively exploit multicore processor.

⇒ multicore organization:-

An arch. in which a single physical processor incorporates the core logic of more than one processor. A single <sup>integrated</sup> circuit is used to package/hold these processors. Multicore arch. places multiprocessor cores & bundles them as a single physical processor.

⇒ Intel core i7 q900:-



Used for computer processors it makes high end desktop systems, fastest processors, they think it will be used to build most powerful consumer devices.