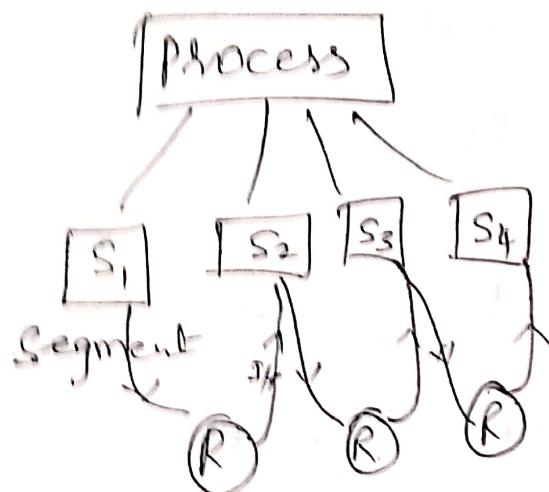


UNIT - IV

Pipelining: A process is divided into several sub-operations



We can execute all the segments concurrently

Why it is called Pipelining?

→ passing information from $S_1 \rightarrow S_2$

$$\sum A_i + B_i + C_i \quad \text{for } i=1, \dots, 7$$

$$R_1 \leftarrow A_i, R_2 \leftarrow B_i \Rightarrow \text{segment 1}$$

$$R_3 \leftarrow R_1 * R_2, R_4 \leftarrow C_i \Rightarrow \text{segment 2}$$

\Rightarrow segment 3

$$f_5 \leftarrow R_3 + R_4$$

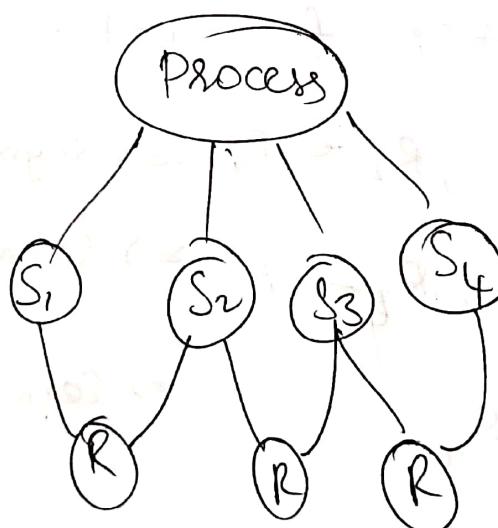
Pipelining

→ A process is divided into several sub operations.

Where each suboperation is represented with segments.

→ The output of each segment is stored in a register.

→ The output of each segment is passed as an input to next segment.

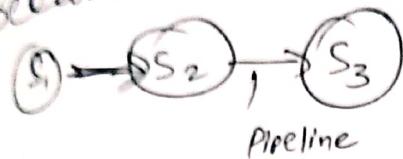


→ Here all the segments executes independently.

→ We can execute all the segments concurrently.

→ Why it is called Pipelining.

because we are transferring information



Ex

$$A_i * B_i + C_i \quad \text{for } i = 1, \dots, 7$$

① $A_i * B_i$

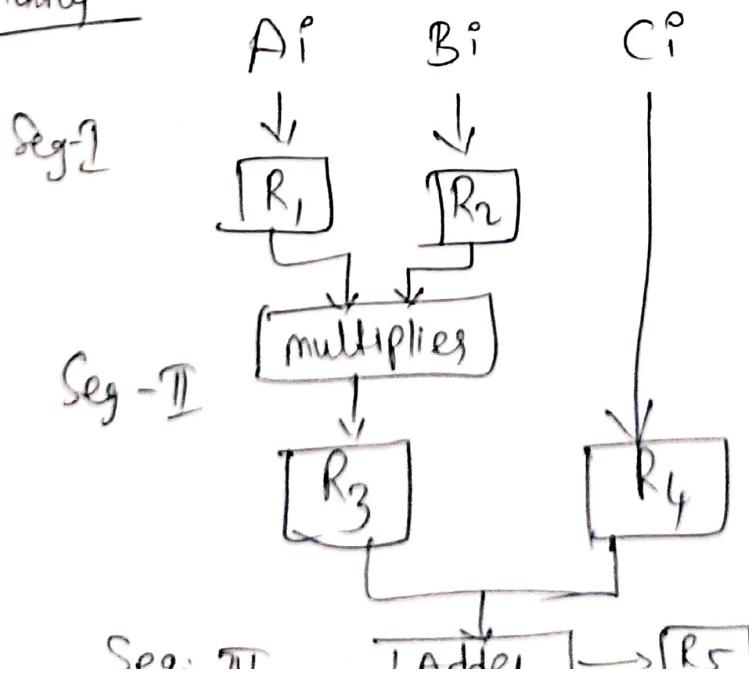
② $A_i * B_i + C_i$

Segment I $\Rightarrow R_1 \leftarrow A_i, R_2 \leftarrow B_i$

Segment II $\Rightarrow R_3 \leftarrow R_1 * R_2, R_4 \leftarrow C_i$

Segment III $\Rightarrow R_5 \leftarrow R_3 + R_4$

Pictorial



Clock Cycle	Seg I R ₁ , R ₂	Seg II R ₃ R ₄	Seg III R ₅
1	A ₁ , B ₁		
2	A ₂ B ₂	R ₁ *B ₁ CP	
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
9	A		A ₇ *B ₇ +C ₇

(9 × 3) segments

It takes 27 clock cycles to give the result

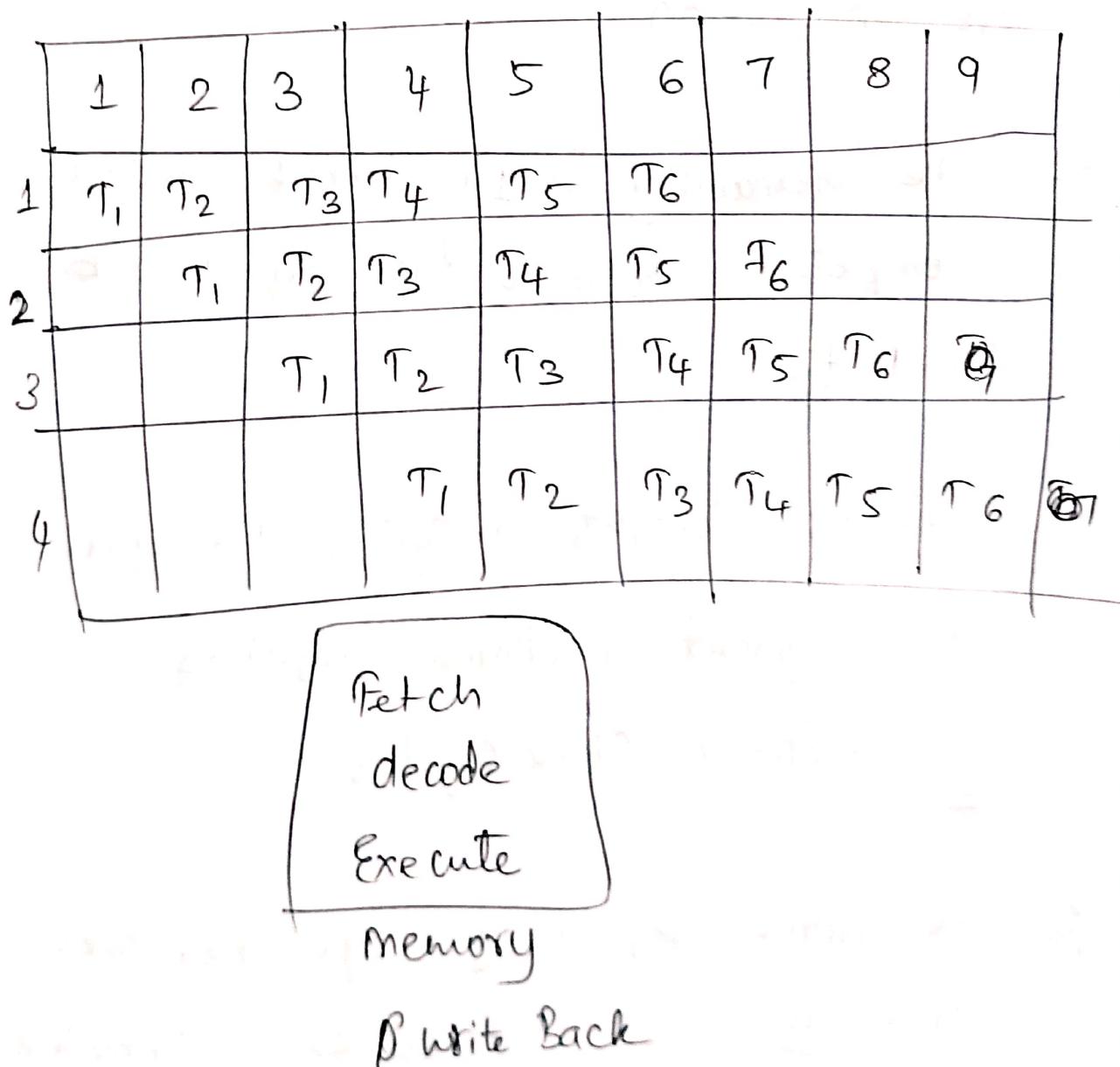
3 clock cycles for one output

$$\Rightarrow 9 \times 3 = 27 \text{ clock cycles}$$

- To Increase performance
- To decrease Time

Performance of pipelining

The Space-time diagram of a four-segment pipeline is demonstrated in following figure



Performance of Pipelining

- Inhere a k Segment pipeline with a clock cycle time t_p is used to execute n tasks
 - The first task T_1 requires a time equal to $k t_p \Rightarrow k \cdot t_p$ to complete its operation.
 - The remaining $(n-1)$ tasks will be completed after a time equal to $(n-1) t_p$
- ∴ To complete n tasks using a k -Segment Pipelining requires $\underline{k + (n-1)}$ clock cycles

for a non-Pipelining operation takes time $\underline{t_n}$ (i.e.) total task completes at $\underline{n \cdot t_n} \Rightarrow n \times t_n$

Speed up :

The Speed up of a pipeline processing over an equivalent non-pipeline processing is defined by the Ratio

$$S = \frac{n t_n}{(k+n-1) t_p} \Rightarrow \frac{n \cdot k t_p}{(k+n-1) t_p} \Rightarrow \frac{n \cdot k}{(k+n-1)}$$

As the Number of tasks increase the speedup becomes

$$S = \frac{t_n}{t_p}$$

\therefore If we assume that the time it takes to process a task is the same in the Pipeline and non-pipeline Circuits

$$(i.e) \quad t_n = k t_p$$

The speed up reduces to

$$S = \frac{k t_p}{t_p} = k \Rightarrow S = k$$

This shows that the theoretical maximum speedup that a pipeline can provide is k , where k is the number of segments in the pipeline.

$$\text{Cycle time } (t_p) = \frac{1}{k+1} \times \text{Total time}$$

$$\text{No. of segments } k = 4$$

$$\text{Number of tasks } n = 100$$

The pipelining system will take

$$k + (n-1) t_p \Rightarrow 4 + (100-1) 20 \text{ ns}$$

$$\Rightarrow 2060 \text{ ns}$$

$$\text{Assuming } t_n = k t_p \Rightarrow 4 * 20 = 80 \text{ ns}$$

A. Non pipelining system requires $n k t_p \Rightarrow 100 * 80 * 20 \Rightarrow 8000 \text{ ns}$

Speedup Ratio is $\frac{8000}{2060} = 3.88$

Throughput :

Throughput of a processor is the rate at which operations get executed

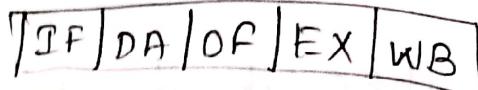
- Latency is the amount of time that a single operation takes to execute
- In non-pipelined Computer
throughput = $1 / \text{latency}$
- In Pipelined Computer
throughput $> 1 / \text{latency}$
since execution of one instruction is overlapped.

Throughput = $\frac{\text{No. of Instructions}}{\text{Total time to complete Instructions}}$

$$\therefore \text{Throughput} = \frac{n}{(k+n-1)tp} \quad \text{for Pipeline}$$
$$= \frac{1/n}{\dots} \rightarrow \text{for non-} /$$

* Pipeline Hazards

Instruction Pipeline



→ The Execution of Instructions blocked / stopped due to some problem in the Pipeline is called Pipeline Hazard.

* Situation that prevent the next instruction from being Executing during its designated Clock Cycle.

- ⇒
- ① Structural Hazard / Resource Conflict
 - ② Data Hazard / Data dependency
 - ③ Control Hazard / Branch Difficulty

① Structural Hazard

2 different segments try to use Same Resource at same time

(or)

If a Resource is needed by 2 Instructions simultaneously

② Ex

MUL \Rightarrow 2 cycles

ADD \Rightarrow 1 cycle

MUL

ADD

IF	DA	OF	EX	WB
IF	DA	OF	-	EX

Structural Hazard

Can't Execute two Instructions at a time

Solution : ① Increase Number of resources
② Wait

② Data Hazard :

Result of an instruction is used as input in next

Ex

$$P = R_1 \leftarrow R_2 + R_3$$

$$P+1 = R_5 \leftarrow R_1 * R_4$$

IF	DA	OF	EX	WB
IF	DA	OF	EX	WB

* Here in the first instruction until the value of R_1 is updated in the memory by WB (Write Back).

→ The value of R_1 in second instruction can't be fetched (i.e) (OF: Operand fetch)

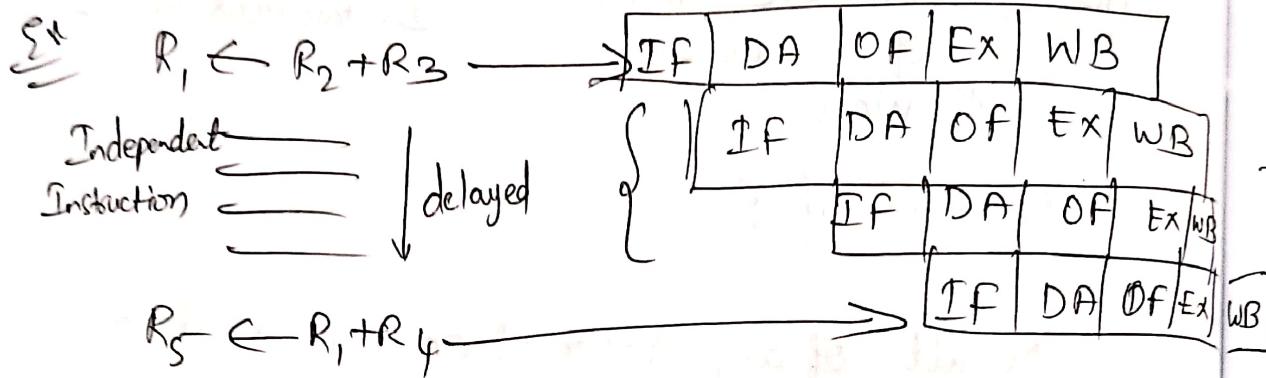
Solution : A general Instruction Pipeline cannot detect data dependency

→ Then solution of Data dependency is provided by compiler, (i.e) delayed Load

Ex $R_1 \leftarrow R_2 + R_3$

↓ ≡ | Independent Instructions Executed

Wait $R_5 \leftarrow R_1 + R_4$ until Instruction get Executed



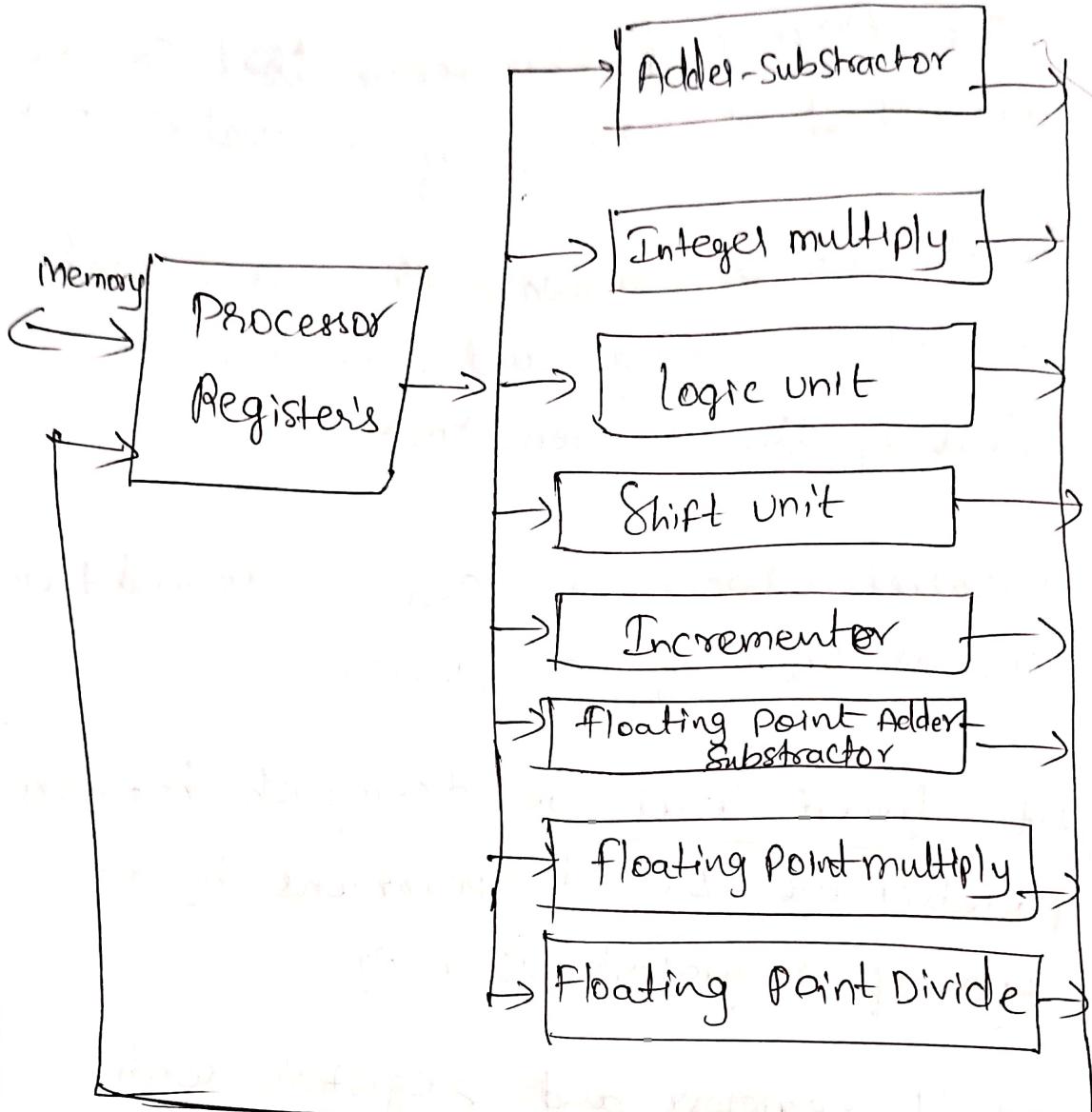
Parallel Processing

* Parallel Processing 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 computer system.

* The purpose of parallel processing is to Speedup the computer processing capability and increase its throughput.

that the Amount of processing that can be accomplished during a given interval of time

- * The Amount of hardware increases with parallel processing and with it, the Cost of the system increases.
- * parallel processing can be viewed from various levels of complexity
 - 1) At lowest level we distinguish between parallel and serial operations by the type of register's used e.g.
 - shift registers and register's with parallel load.
 - 2) At the higher level, it can be achieved by having a multiplicity of functional units that perform identical (or) different operations simultaneously
- * A multifunctional organization is usually associated with complex control unit to coordinate all the activities among various components



In the Above figure we can see that the data stored in the processor register's is being sent to separate devices based on the operation needed on the data.

- If the data inside the processor register is requesting for an arithmetic operation then the data will be sent to the arithmetic unit and if in the same time another data is requested in the logic unit then the data will be sent

to logic unit. for logical operations.

* Now in the same time both arithmetic operations and logical operations are executing in parallel.

This is called parallel processing

* The computers are classified into ④ types based on the
① Instruction streams and
② Data Stream

These ④ classifications are called as ~~Instruction Stream~~

Flynn's classification of Computers.

① Instruction Stream: The sequence of instructions read from the memory is called as an Instruction Stream

② Data Stream: The operations performed on the data in the processor are called as Data Stream

M.J. Flynn's Classification

M.J. Flynn Consider's the organisation of a Computer System by the Number of instructions and data items that are manipulated Simultaneously.

- ① Single Instruction Stream, Single data Stream (SISD)
- ② Single Instruction Stream, Multiple data Stream (SIMD)
- ③ Multiple Instruction Stream, Single data Stream (MISD)
- ④ Multiple Instruction Stream, Multiple data Stream (MIMD)

① SISD

- Represents the organization of a single computer Containing a control unit, a processor unit, and a memory unit.
- Instructions are Executed Sequentially and the system may (or) may not have internal parallel processing Capabilities.
- Parallel processing may be Achieved by means of multiple functional units (or) by pipeline processing

SIMD :

- Represents an organization that includes many processing units under the supervision of a common control unit
- All processors receive the same instruction from the control unit but operate on different data items.
- The shared memory unit must contain multiple modules so that it can communicate with all the processors simultaneously.

MISD :

MISD Structure is only of theoretical interest since no practical system has been constructed using this organisation.

MIMD : This organization refers to a computer system capable of processing several programs at the same time.

Ex Multiprocessor and multi Computer System

One type of parallel processing that does not fit flynn's classifications is Pipelining

We consider parallel processing under the following main topics.

①

Pipeline processing:

It is an Implementation technique where arithmetic sub operations (or) the phases of a Computer instruction cycle overlap in execution

②

Vector processing:

Deals with Computations involving large Vectors and Matrices.

③

Array processing:

Perform Computations on large arrays of data



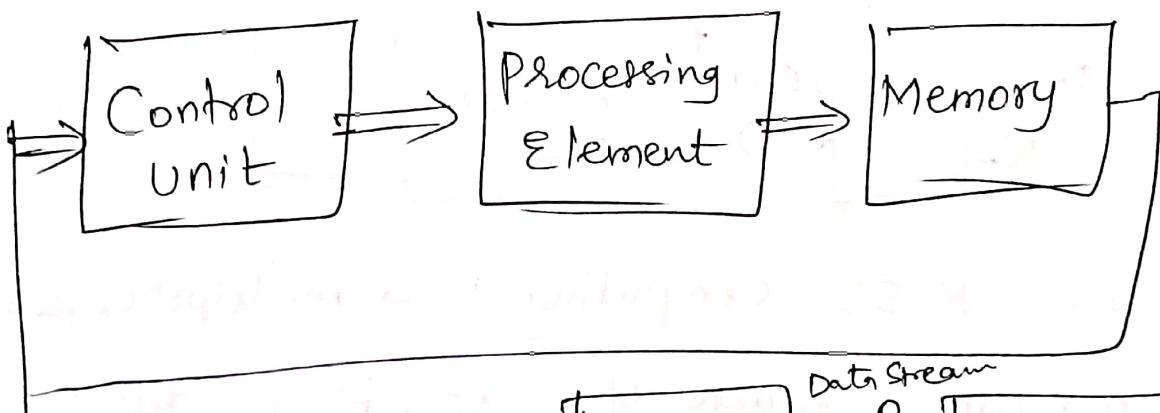
SISD

Single Instruction Stream

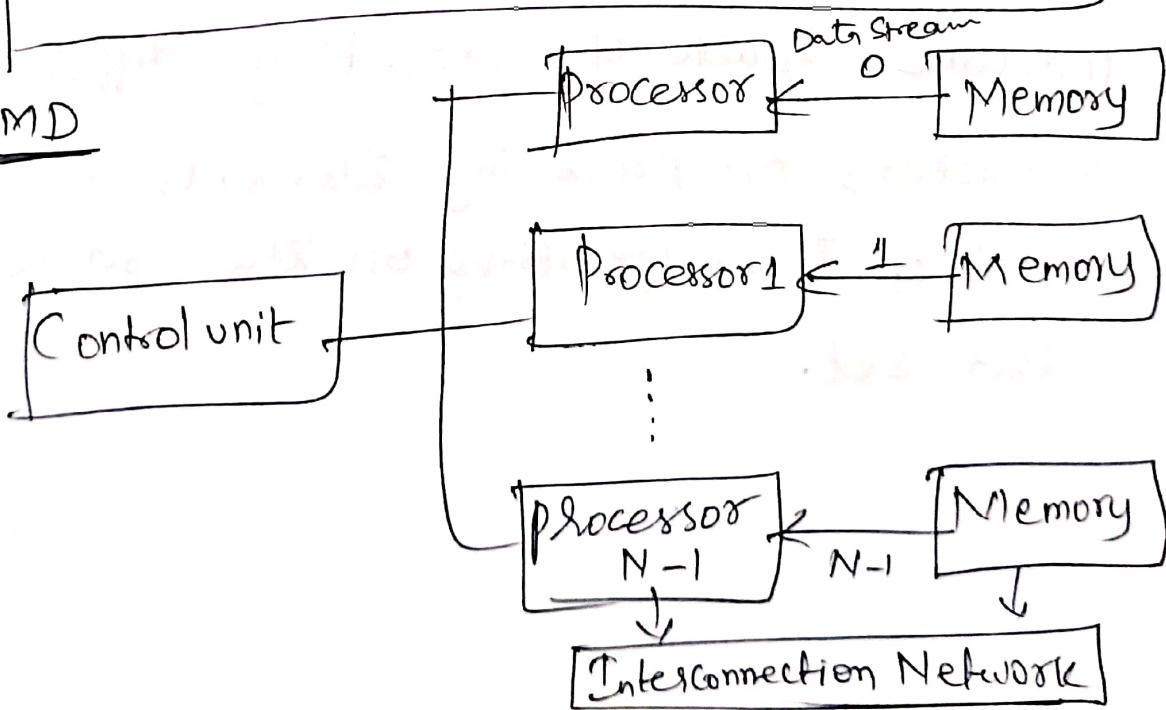
- Only one instruction stream is being acted
 (or) Executed by CPU during one clock cycle.

Single data stream: Only one data stream is used as input during one clock cycle.

Most conventional computers have SISD architecture where all the instruction and data to be processed have to be stored in primary memory.

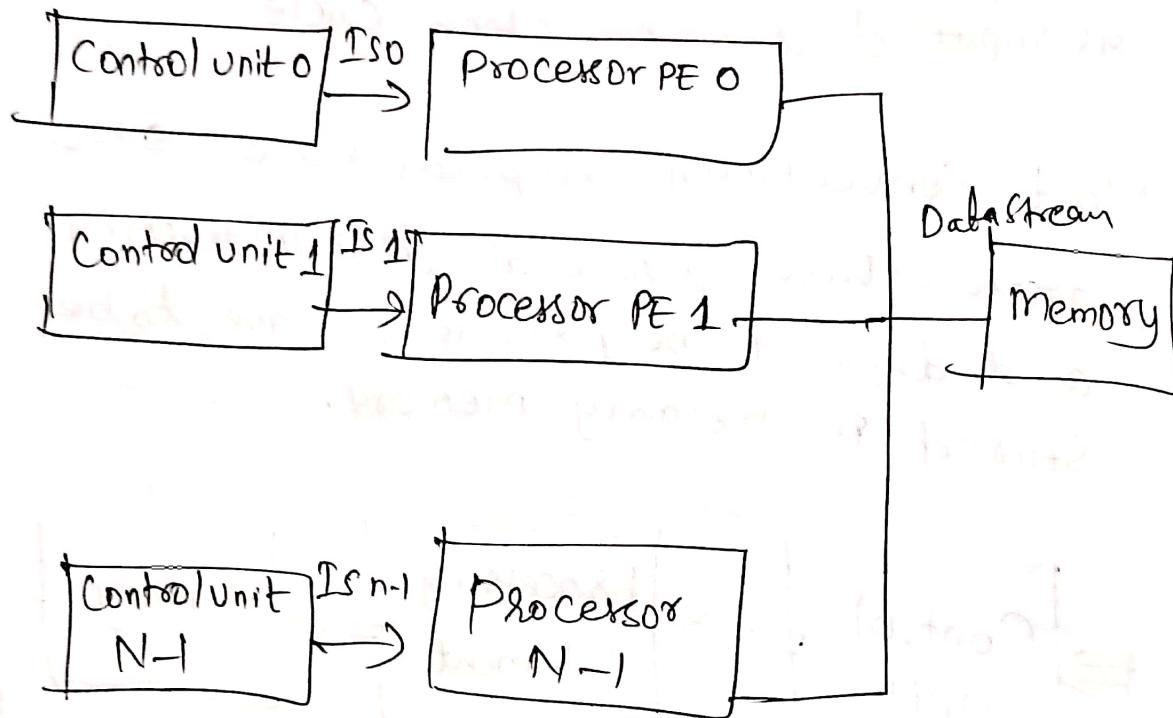


SIMD



A SIMD system is a multiprocessor machine capable of executing the same instruction on ~~the~~ all the CPU's but operating on the different data stream.

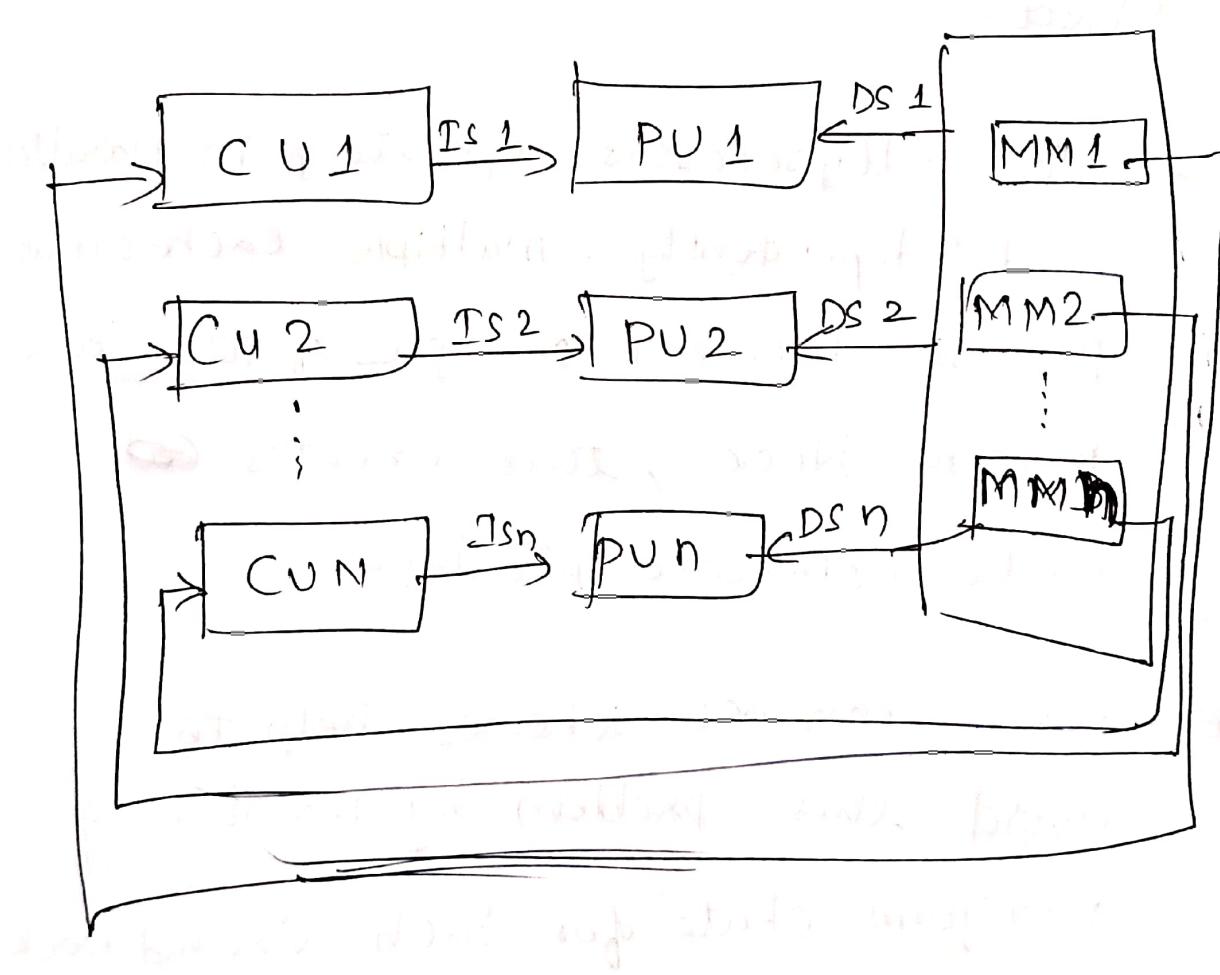
MISD



An MISD Computing is a multiprocessor machine capable of executing different instructions on processing elements but all of them operating on the same data set.

MIMD

A MIMD system is a multiprocessor machine that is capable of executing multiple instructions over multiple data streams. Each processing element has a separate instruction stream and data stream.



Cache Coherence

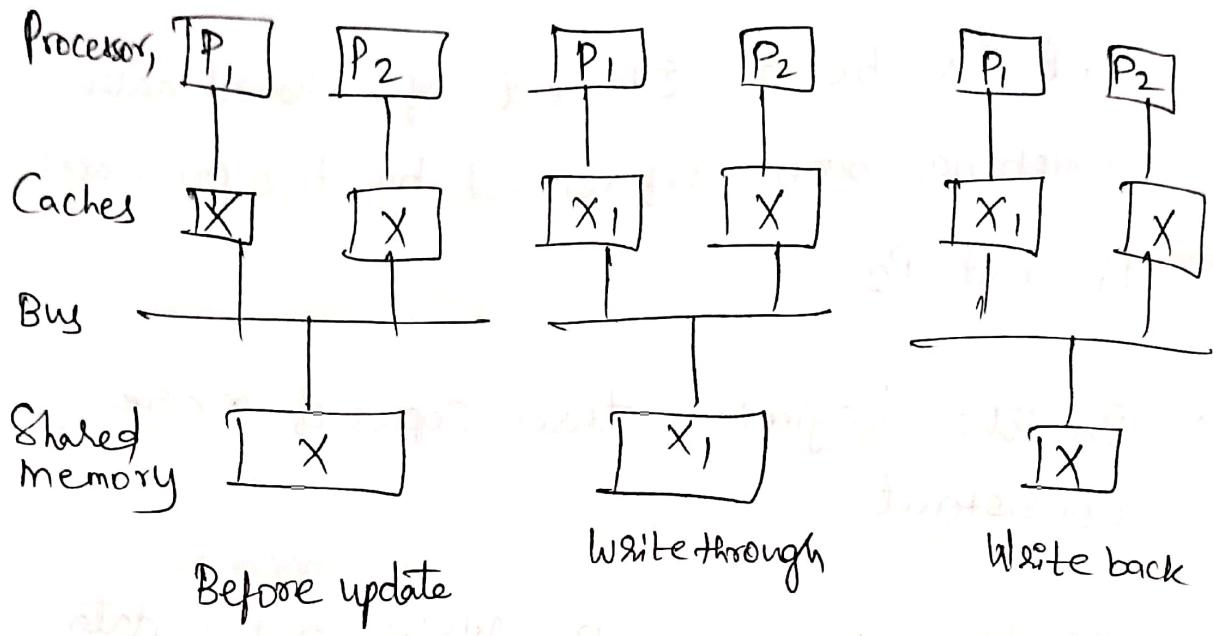
In a multiprocessor system, data inconsistency may occur among adjacent levels (or) within the same level of memory hierarchy.

For Example: the Cache and the main memory may have inconsistent copies of the same object.

- As multiprocessors operate in parallel and independently, multiple caches may possess different copies of the same memory block, this creates ~~an~~ Cache Coherence problem.
- * Cache Coherence Schemes help to avoid this problem by maintaining a uniform state for each cached block of data.

Cache coherence

- Let X be an element of Shared data which has been referenced by two processors P_1 and P_2
- In the beginning three Copies of X are Consistent
- If the processor P_1 Writes a new data x_1 into the Cache,
By Using Write-through policy , the same Copy will be written immediately into shared memory.
- In this Case, inconsistency occurs between Cache memory and the main memory
When a write-back policy is used
 - the main memory will be updated when the modified data in the Cache is replaced (or) invalidated.
 - If Cache is full then the Value is Replaced (or) modified.

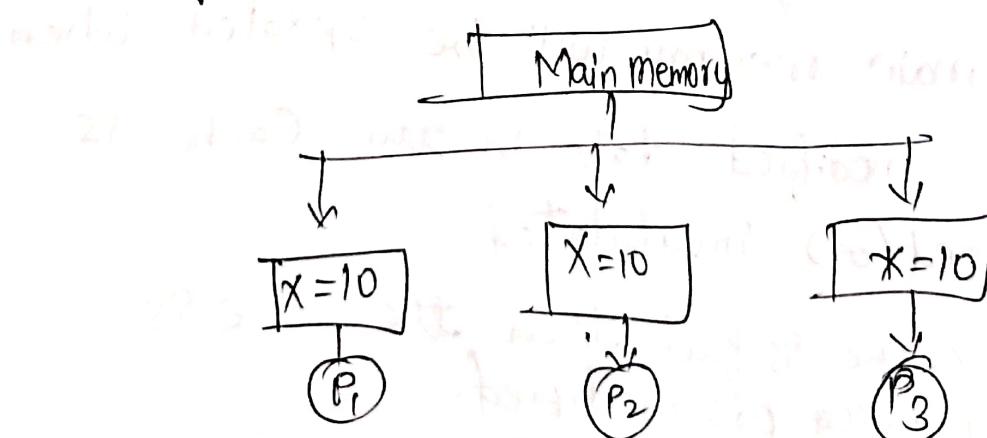


i) Write update:

Whenever the value of \textcircled{X} is changed in Cache of $\textcircled{P_1}$ then it is updated as same in $\textcircled{P_2}$ and $\textcircled{P_3}$.

ii) Write through:

When \textcircled{X} is modified in Cache and shared memory then it is updated in main memory also.



* Write Invalidate :

Whenever the value of (X) is modified
(or) changed in Cache then other
Cache of (P_2) and (P_3) Values are
invalidated.

② protocols

① Write update \rightarrow update in memory
at a time of modification

② Write Invalidate

↳ If any updating of
data done then other
~~Cache~~ Cache values will
be invalidated.