

- (4) Real performance - flops for specific operation
(5) Sustained performance - performance on an application

03/01/21

- * Sustained Performance < Real Peak
Continued on 09/01

High Performance Computation (HPC)

Terms

(1) flops - floating point operation per sec.

→ HPC is teraflops or more (10^{12})

→ petaflops = 10^{15} flops. (quadrillion)

(2) → Mips - Mega instructions per sec.

if there is only one instruction " " = MHz.

(3) Theoretical Peak Performance - max flops a machine can reach
= CPU speed in Gigahertz x No. of CPU cores (proxim.)

(1) Portability -

can be achieved by containerization

(storing all kind of libraries, models used

→ docker

for particular process

→ future-proofing

(2) Modular and scalable

benchmarked and universal
so that others can also use.

(3) Grid-ready → follows grid-computing

for a bigger server, efficient utilization of resources.

Applications -

(1) Climate modelling - very big real time data requiring HPC.

(2) Physics - larger calculations for FT, calculus

(3) Financial - stock prediction, currency fluctuations.

(4) Banking sector - credits

(5) Image Segmentation
(X-ray, MRI)

(5) Business sector - product prototyping

(6) Fluid dynamics

(7) Gaming

(7) Drug discovery

Challenges

- (1) cost increases
- (2) skilled professional
- (3) Integration

In HPC, latency should be low, speed should be high and it should match with the memory.

- (4) speed depends on memory. There should be no mismatch.
- (5) Power consumption.

Parallel computation

Divide the instruction in multiple units and each unit is performed by different processor.

- Motivation for parallel computation in ICs
- (1) See no. of transistors for trying the computation power
 - (2) improvement in storage technology (in disks and memory)

~~Study key notes now~~
Locality of Reference (Technique)

- (3) improvement in networking devices and system for data communication (dependent on internet)

Time to execute a given program: $T = n_i \times t_c$

$T = n_i \times (n_w) \times t_c$

$= n_i \times CPI \times t_c$

t_c = time per cycle
 n_w = total no. of instructions
 CPI = Avg. cycles per instruction
 Continue at a JAM

Motivation:

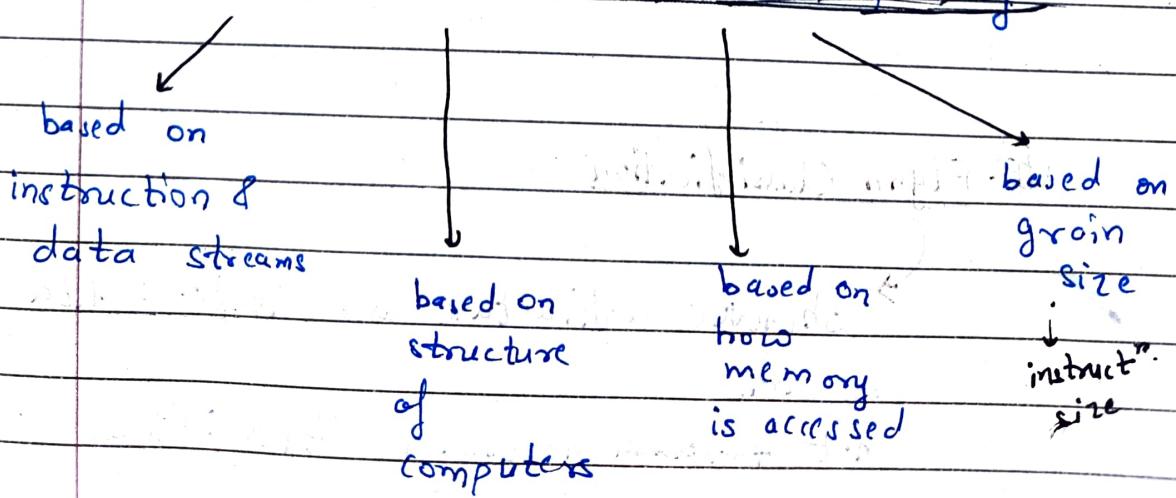
- ① Parallel computation
 - 2.) Use no. of transistors in ICs for raising the computer power
 - 3.) improvement in storage technology (in disks & memory)
- eg: L1 cache, L2 cache, main memory etc.

Locality of Reference (Technique):

- ④ Improving the networking devices & system for data communication (dependent on internet)

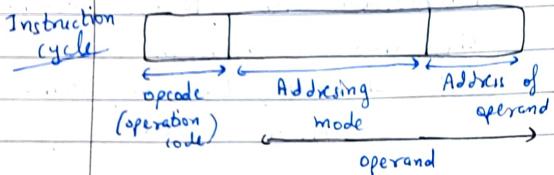
51/2u

classification based on parallel computing

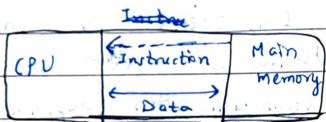


(i) Stream → Sequence or flow.
Instructions

Grain size → size of instruction of the program.



Instruction stream - Instructions come from main memory to CPU.



Data stream - bidirectional

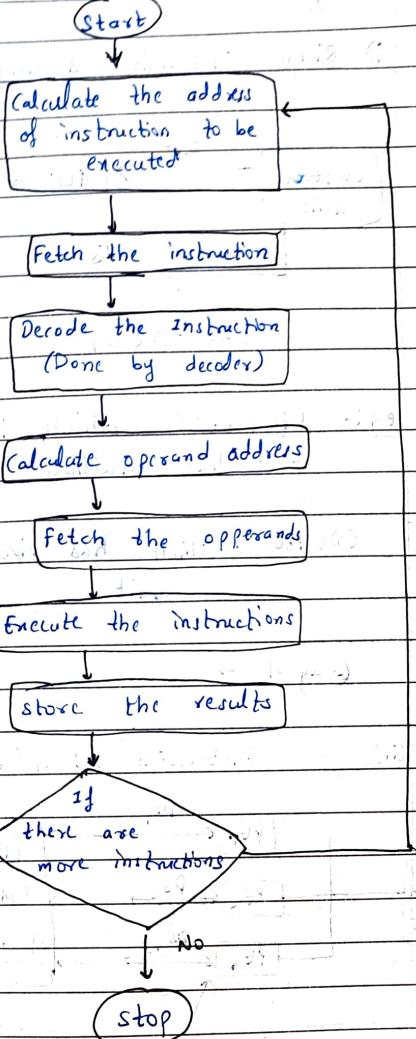
Instruction stream - Unidirectional

Flynn classification

whenever architecture uses instr. & data stream.

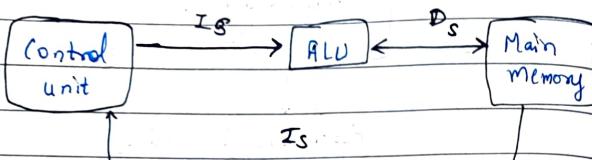
not necessarily necessary that it was parallel architecture.

Instruction
Fetch stage



Flynn's classification:

i) Single instruction & single data stream (SISD)



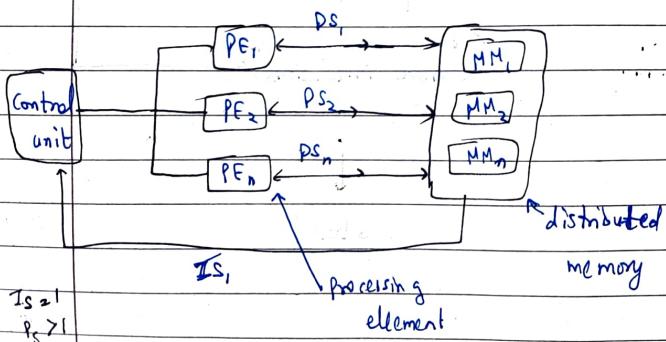
$$I_S = D_S = 1$$

e.g.: CDC 6600 which is uniprocessor but has multiple functional units

CDC 7600 → which has pipelined arithmetic unit.

Gray - 1 →

(ii) single Instructions and multiple data streams



Data streams coming from multiple memory locations

One instruction is divided & given to multiple processing elements.

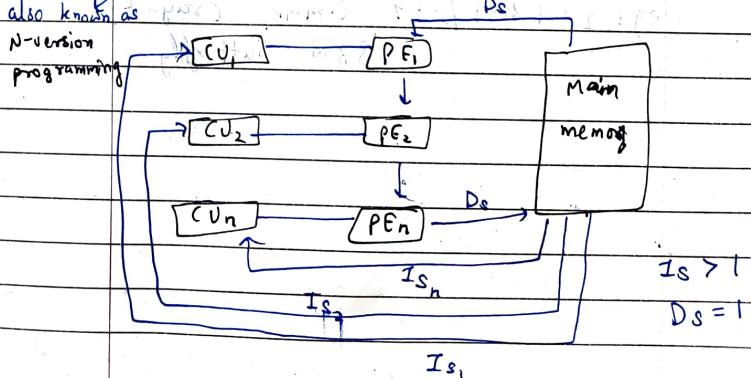
Each PE can have its own memory (e.g.: printer / fax machine can have own memory)

⇒ Global communication & another main memory communication

agar comp. ke sathe
pointer / fax machine
hui toh unki bhi
comm' overall, ho
skti hai

e.g.: BSP, MPP, DAP, CNT (connection machine).

(iii) Multiple instructions & single data stream.

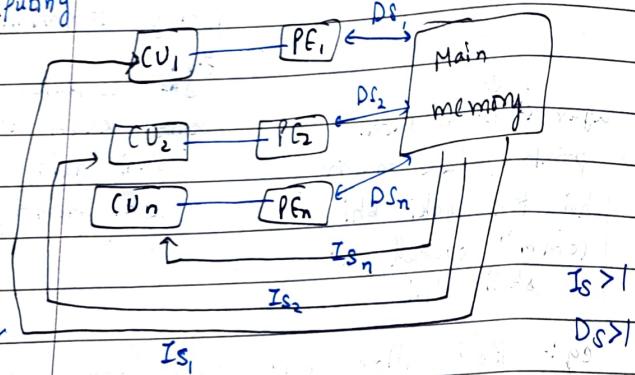


eg: Carnegie Mellon University
derived \hookrightarrow C.MMP

not used frequently } Real time computers must be fault tolerant
But here if fault in PC \rightarrow This same goes to PE \rightarrow So not fault tolerant

i) multiple instruction & multiple data streams

also known as parallel computing



This is asynchronous

eg: IBM 3081/3084, C.MMP, Cray-2, BBN - Butterfly

GJAN

To increase computer speed / speed of program

- ① reduce cycle time (t_c) by increasing clock freq.
- ② reduce number of instructions
- ③ reduce CPI \rightarrow using parallelism

- ① Instruction level Pipelining (overlapping of instructions)
- ② Internal Parallelism
- ③ External Parallelism

Temporal Parallelism :-

Time-based

\Rightarrow This method breaks up job into a set of tasks to be executed overlapped in time

\rightarrow This is also called assembly pipeline / vector processing

Conditions: (in which temporal parallelism can be applied):-

- The jobs to be carried out are identical.
- A job can be divided into many independent tasks.

(5) Moore's law

No. of transistors on computer chips doubles approximately every 2 years.

This implies CPU speed doubles every 18 months.

Time to execute a given program, $T = n_c \times t_c$

n_c = total no. of CPU cycles or

t_c = cycle time.

$$T = n_i \times \left(\frac{n_c}{n_i} \right) t_c$$

n_i = total no. of instructions

↓
avg cycles
per instruction

$$= n_i \times CPI \times t_c$$

To ↑ computer speed / speed of program

- (1) reduce cycle time (t_c) by ↑ing clock freq.
- (2) reduce no. of instructions (n_i) by using better compilers and efficient algs
- (3) reduce CPI by using parallelism.

Types of \rightarrow (1) instruction level parallelism

by pipelining - overlapping of instructions.

(2) internal parallelism

(3) external parallelism

Temporal parallelism \Rightarrow This method breaks up time-based a job into a set of tasks to be executed overlapped in time.

aka assembly-line processing, pipeline processing & vector processing.

→ Conditions in which temporal parallelism can be applied

(1) The jobs to be carried out are identical.

(2) A job can be divided into many independent tasks.

(3) Time taken for each task should be approximately same.

(4) Time taken to send a job from one processor to the next is negligible compared to the time needed to do a task.

(5) the no. of tasks are much smaller as compared to the total no. of jobs to be done.

Let $n = \text{no. of jobs}$

$p = \text{time to do a job}$

~~but~~ $k = \text{no. of tasks}$ (let each job be divided into k tasks)

$t = \text{time to do each task} = \frac{p}{k}$

* Time to complete n jobs = $n p$

with no pipeline processing
(no parallelism)

* Time to complete n jobs with pipelining = $\frac{p + (n-1)p}{k}$

~~repeat~~

→ This can be sped up by pipelining

in which the time taken will be speeded up by $\rightarrow \frac{np}{p + \frac{(n-1)p}{K}}$

= Ratio of w/o pipelining

$$= \frac{np}{\left(\frac{n+k-1}{k}\right)p}$$

~~w/ pipel~~

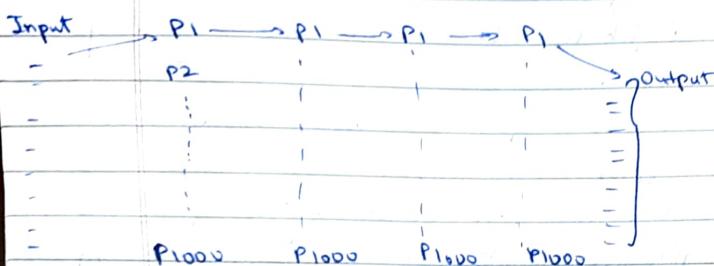
Teacher
TA1 Answer

Page No.

Date: / /

Page No.

Date: / /



Disadvantages

- ① Synchronization - for this identical time should be taken for each task in the pipeline so that a job can flow smoothly in the pipeline without holdups (ek job koi bhot hota nahi na baitha h).
- ② Bubbles in the pipeline - If some tasks are absent in a job, bubbles are formed in the pipeline. Other processor will sit idle if 03, 04 are not attempted.
- ③ Fault tolerance - If a processor crashes in plus a cycle, synchronization breaks and thus the system does not tolerate fault.
- ④ Inter-task communication - If too much time is consumed to pass the task b/w the processor in

the pipeline. This time should be smaller as compared to the time taken to execute the job.

⑤ Scalability - The no. of processors working in the pipeline cannot be increased after a certain limit. There comes a limit after which no. of processors can't be increased.

10/01/24 In QPS

$$\text{speed up by } \frac{n}{(k+n-1)} = \frac{nk}{(k+n-1)}$$

If $n > k$

$$= \frac{k}{1 + \frac{k-1}{n}}$$

If $n > k$ then $\frac{k-1}{n} \approx 0$ and speed up is $\approx k$ \hookrightarrow no. of tasks in divided and each task is given to different processor

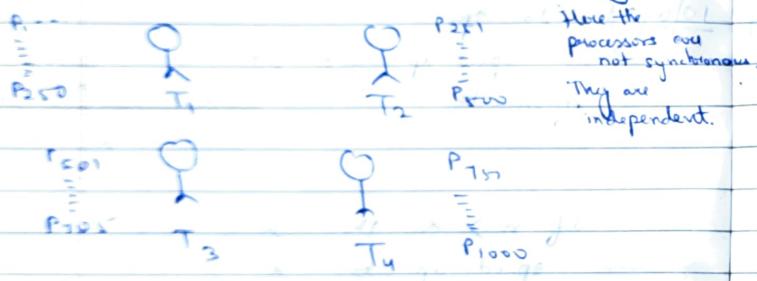
Influence
⇒

speed up is proportional to no. of processors for condition - no. of job > no. of task

Adv → ① each processor can be fine-tuned to do one kind of a job.
→ One processor can become expertise for one task.

- * If a speed up of a processor of a method is directly prop to the no. of processors then the method is said to be scale well.
→ temporal claim - scales well.

Data claim with static-assignment
The ifp data is divided into individual sets and processed simultaneously.



Let

$$n = \text{no. of jobs}$$

$$p = \text{time to do a job}$$

$$k = \text{no. of processors to do the job}$$

$kpg = kp = \text{time to distribute the job to } k \text{ processors (teachers)}$

Observe that this time is prop. to no. of processors.

∴ time to complete n jobs by single processor $\approx np$

$$\therefore \text{by } k \text{ processors} = \frac{np + kp}{k}$$

$$\begin{aligned} \text{speed up} &= \frac{np}{\frac{np + kp}{k}} = \frac{knp}{kp + np} \\ &= \frac{k}{1 + \frac{kp}{np}} \end{aligned}$$

If $k^2p < np$ then speed up = k.

efficiency = $\frac{\text{speed-up}}{k}$

If the no. of processors \uparrow ers, the time to distribute jobs to processors also \uparrow ers.

Hence, speed up is not directly prop to the no. of processors.

Advantages

- (1) No synchronisation required. - all teachers can check copies without dependency on other
- (2) Dependency due to bubbles are not present. All processor can perform complete job. If bubble arises then processor can tackle independently.
- (3) More fault tolerance. Agar los break will not have effect in large
- (4) No inter-task communication delay
↳ does not exist here

~~difficult~~ ~~basic~~ → Temporal - has processor do particular task in
kisi ekta ha pipeline.

→ Data - The jobs are divided & given to different
processor which can do complete job
independently

Disadvantages

- (1) Static assignment - The assignment of the jobs to each processor is pre-decided. Thus, if a processor is slow then the completion time of the total job will be slowed down. If there are more bubbles, that single processor will take more time.

- (2) We must be able to divide the set of jobs into subsets of mutually independent jobs, each subset should take the same time to get completed
⇒ evenly distributed

- (3) Each processor must be capable of executing all the jobs.

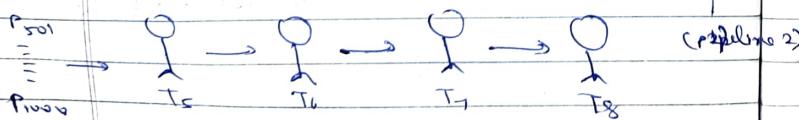
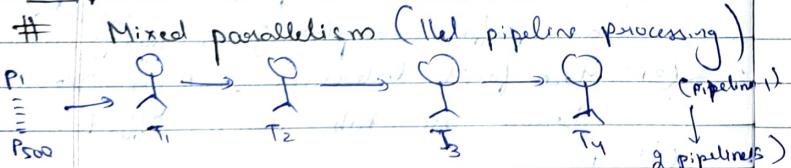
- (4) The time taken to divide the set of jobs into equal subsets of jobs should be small.

- This method is effective only if the no. of jobs given to each pipeline is much larger than the stages/segment of the pipeline (T_1, T_2, T_3, T_4)
- Both pipelines process simultaneously
- This method almost halves the time taken by a single pipeline.

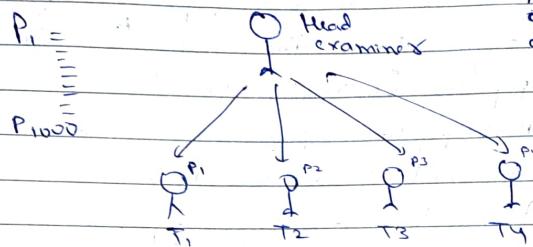
Disad → ~~both~~ disadvantages of temporal & data

Examples of computers performing this \rightarrow
Gray, NEC-SX

This method is very efficient for numerical computing in which a no. of long vectors and large matrices are used as data and could be processed simultaneously.



Data I/Osim with dynamic assignment



Here a head or processor distributes 1 job to each processor & when it finishes the job it waits in the queue & to get a new job assigned then completes it

This process continues until all the jobs are completed.

Advantage

- (1) Balancing the work assigned to each processor directly dynamically as the processor progresses.
- (2) A processor that finishes the job quickly gets another job immediately and does not need to sit idle.
- (3) Time to execute a job may widely vary without creating a bottleneck.
- (4) This method is not affected by bubbles.

Disadvantages

- (1) If many processors complete their job simultaneously, they need to wait in the queue as only one will be attended at a time.
- (2) The main processor can become a bottleneck.
- (3) The main processor itself is idle b/w handing out jobs.

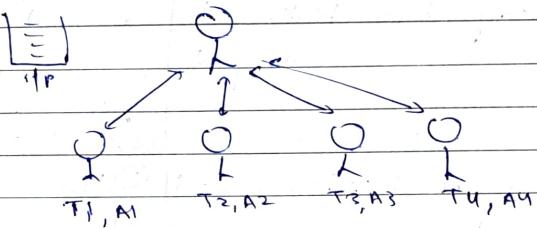
(4) It is difficult to increase the no. of processors as it will increase the prob of many processors completing their jobs simultaneously thereby leading to long queue of processors.

Data I/Osim Quasi Dynamic Scheduling

Giving each teacher / unequal sets of worksheets to correct / jobs

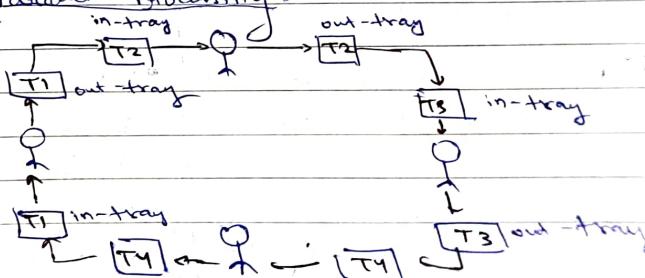
Other I/Osim

(1) Specialist I/Osim



one processor do one task of one job.

(2) Coarse Grained specialist Temporal Parallel Processing



- Instead of spending time to wait for other processors to complete their tasks, the processor takes the task from their in-tray, executes it and then places it in their out-tray.
- Initially, if there is no task in its in-tray, it will wait for it to be filled.

Temporal Vs Data Migr.

Temporal

① Job is divided into set of independent tasks and these tasks are assigned for processing.

② Tasks should contain equal time. Pipeline stages should thus be synchronized.

③ Bubbles in jobs lead to idling of processors.

④ Processors are specialized to do specific tasks efficiently.

Data

① Complete job is assigned to the processor.

② Jobs may take different times. No need to synchronize the jobs.

③ Bubbles don't cause idling of processors.

④ Processors should be general purpose and may not do all tasks efficiently.

- | | |
|--|--|
| ⑤ Not tolerant to processor faults | ⑤ Tolerates processor faults |
| ⑥ Task assignment is static (pre-decided). | ⑥ Task assignment can be static, dynamic & quasi-dynamic. |
| ⑦ efficient with fine-grained tasks. | ⑦ efficient with coarse-grained tasks and more specialized quasi-dynamic scheduling. |

16/01/24 Pipelining

Pipeline is useful when same processing is applied over multiple units.

→ It is a technique to decompose a sequential process into sub-operations

Task - one operation performed at all segments/stages.

example - $A_i \times B_i + C_i$, $i = 1, 2, \dots, 5$

$$R_1 \leftarrow A_1 \quad R_2 \leftarrow B_1$$

$$R_3 \leftarrow R_1 \times R_2$$

$$R_4 \leftarrow C_1$$

$$R_5 \leftarrow R_3 + R_4$$

Clock Pulse	Segment 1	Segment 2	Segment 3
1	R1 R2	R3 R4	R5
2	A1 B1		
3	A2 B2	A1+B1 C1	A1+B1+C1 (Conv inst executed (in 3 clock cycles) A2*B2+C2)
4	A3 B3 simultaneously fetch as well to create bcz it is pipelining	A2*B2 C2	A3*B3+C3
5	A4 B4	A3*B3 C3	A3*B3+C3
6	- -	A4*B4 C4	A4*B4+C4
7	- -	A5*B5 C5	A5*B5+C5
Total time with pipelining = 7 → without pipelining = 15 ←			
3 for 1st 1 for rest all			

Time required to perform 1st task =

$$\frac{k \times t_p}{\text{no. of segments}} + \text{time for one segment} \times (\text{no. of remaining tasks})$$

$$(n-1) \times t_p$$

$$\text{Total time} = (k+n-1) \times t_p$$

(after sequential)

- # Consider a non-pipelining system, that takes t_n time to perform a task
 \Rightarrow for n tasks = $n \times t_n$

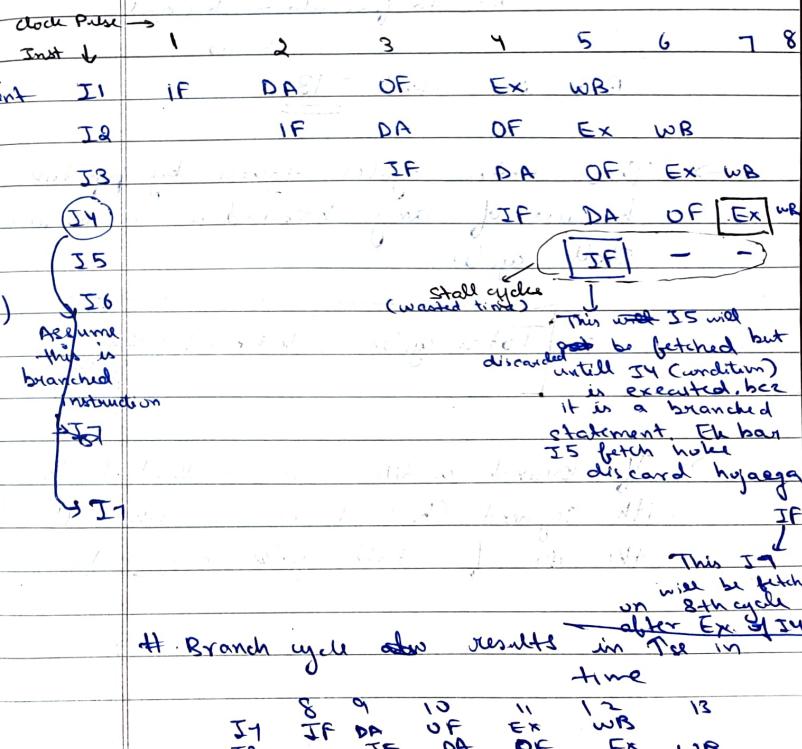
$$\star \text{ speed up} = \frac{n \times t_n}{(k+n-1) t_p} = \left(\begin{array}{l} \text{(non-pipelining)} \\ \text{(pipelining)} \end{array} \right)$$

Case 1 - If $n \gg k$, speedup = $\frac{t_n}{t_p}$ (max speedup) = ideal

Case 2 - When time req to perform one task is same in pipeline and without pipeline :-
 \Rightarrow speed up = $\frac{n}{k+n-1} = \frac{n}{k+n-1}$
 (can't be greater than k)

5-segment instruction Pipeline

- ① Instruction fetch (IF)
- ② Decode & Address calculation (DA)
- ③ Operand Fetch (OF)
- ④ Execution (EX)
- ⑤ WB (Write-Back Result)



Total instructions executed = 6 ($I_1 - I_4$, I_7, I_8)

Condition is checked in the execution stage of the branch instruction and the instructions which are fetched after the branched inst will get discarded. (these are the stall cycles which are discarded and wasted time)

$$\text{Here, } k = 5$$

$$n = 6$$

$$\text{Acc to formula} = k+n-1$$

$$= 10$$

$$\text{But in actual} = \underline{\underline{13}}$$

These extra cycles taken by pipeline because of any hazard or reason are stall cycles.

$$\boxed{\text{Total cycles} = k+n-1 + \text{stall cycles}}$$

If branch condition is evaluated in i th segment (here it was in 4th) then the no. of stall cycles are $i-1$.

delay th. inst
due to
any

Pipeline Hazards

Pipeline hazards are often seen as mis-steps in the optimization of processes within the pipelining concept.

Hazard can occur during ^{stage} ① fetch ② decode or ③ execution

17/01/24

① Hazard ~~can be~~ usually related to issues with the memory system such as latency. (conflict b/w d instruction)

Decoding ② An instruction might need to wait for completion of another instruction which currently occupies the decoder.

Execution ③ If there is a conflict in accessing a functional unit (like register, ALU)

Types of pipeline hazards

① Structural Hazard / Resource conflict

These occur when the same hardware resource is desired by multiple instructions simultaneously.

MUL takes more time for exec than ADD

MUL IF DA OF EX EX WB

ADD IF DA OF [EX] EX WB

becomes stall cycle b/c ALU is busy with MUL and is executed later

Solutions to this -

- ① incur stall cycles. (increase) in case one is dependent on other inst.
 ② ↑ no. of resources. (not much feasible)

(2) Data Hazard / Data dependency

When the execution of one inst. is dependent on completion of the other.

These are categorized into 3 types

a) RAW (Read After Write)

- aka True dependency
- occurs when an inst. depends on result of prev. inst.

$$\text{eg: } i: R_1 \leftarrow R_2 + R_3$$

$$j: R_5 \leftarrow R_1 + R_4$$

here you also have to check whether time bad dependency aa yahi hai due to hazards occur.
eg - i & j

I5 or \rightarrow j reads a source before it is dependent.
So I5 is written by i.
I5 bhi complete ho jago, no stalls would occur

(b) WAR (Write after Read)

- aka anti-dependency & false dependency
- when an inst. depends on the reading of a value before that value is overwritten by some other instruction

(possible much value than R1, me but is updated later. But you read prev val not updated one)

$$\text{eg: } i: R_1 \leftarrow R_2 + R_3$$

$$j: R_2 \leftarrow R_4 + R_5$$

\rightarrow If j first writes in R2 and i reads later.

c) WAW (Write after Write)

- aka write-dependency & false dependency
- when a value is written by the instruction before the prev. inst. writes that value.

$$\text{eg: } i: R_1 \leftarrow R_2 + R_3$$

$$j: R_1 \leftarrow R_4 + R_5$$

\rightarrow If j writes in source before i writes
 (duru satte aapne time pe fetch hui
 lekin, j get exec before i, and then
 i get exec, so this is hazard
 bcoz current value should be j but
 here it is i).

Register Renaming
 * We can write the result in another register in case of false dependency.

Solutions to this -

- ① Compiler or language translator will generate a program where data dependency is not present (aka Delayed load) — Software soln.

This delay is produced by writing other independent or no ops instructions b/w i and j.

$$\text{eg: } i: R_1 \leftarrow R_2 + R_3$$

$$j: R_4 \leftarrow R_1 + R_5$$

~~IF DA OF~~



Page No.

Date: / /

All segments have different separate H/w

- ② H/w solⁿ - after detecting data dependency.

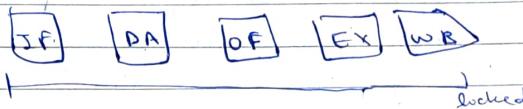
When a H/w can detect and solve the data dependency.

(i) Hardware Interlock

when an inst is executed, all segments have separate hardware. When values is stored in one hardware that particular H/w is locked until the whole inst is executed.

when it completes, then it allows other inst. to update further values.

- All prev H/w get locked until WB is done.

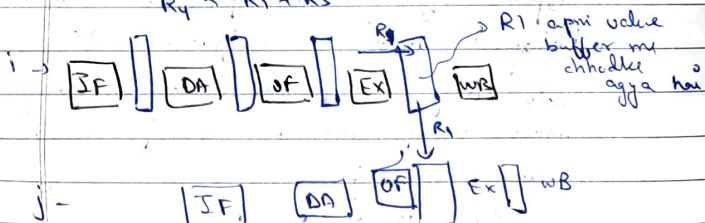


(ii) Operand Forwarding / Rypassing

All segments have own buffer to store prev. value results in it.

$$cg \quad R_1 \leftarrow R_2 + R_3$$

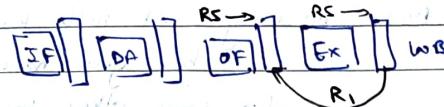
$$R_4 \leftarrow R_1 + R_5$$



Page No.

Date: / /

- R₁ is executed and stored in EX buffer and then becomes stored in OF buffer (dependent - R₁) and R₂ is not dependent.



This does not make us wait till WB and removes time delay and stall cycles.

- No stalls for ALU to ALU dependency. But for ALU to Mem or Memory to ALU stalls are required.

Condition Done inst ALU, me exec hani chahiye

③ Control Hazard (Branch Difficulty)

When the CPU decodes an inst as branch inst then CPU has to wait till branch outcome is available to fetch new instruction.

Solutions

i) Delayed Branch

Job take branch inst exec, add nope and independent inst b/w

Q) HW soln - Branch Prediction.

It predicts whether the current branch inst will jump or not.

If Pred = correct - no prob

Pred = wrong - it rolls back ~~as~~

bcz you have gone
to some other inst and
then comes back.

— X —

14/01/24

Bernstein's Conditions

Bernstein gave a set of conditions on the basis of which we check if 2 processes can execute safely or not

Here, $P_i \rightarrow$ software entity

Input $\rightarrow I_i$ (set of all if/p variables
for a process)
aka ReadSet / Domain of P_i

Output $\rightarrow O_i$ (set of all of p variables)
aka Write Set / Range

→ Consider 2 processes which are interdependent



$w(P_1) O_1 : O_2 w(P_2)$
For 2 processes
to execute safely,
they should be
independent also.

Conditions (given by Bernstein's)

(1) $I_1 \cap O_2 = \emptyset$ (anti-independent)
if p of 1st process and o/p of 2nd process
should have nothing common

(2) $I_2 \cap O_1 = \emptyset$ (flow-independent)
Read(P_2) \rightarrow should be independent of $w(P_1)$

(3) $O_1 \cap O_2 = \emptyset$ (output independent)

In general, P_1, P_2, \dots, P_k can execute in
parallel if Bernstein's conditions are satisfied
on a pair wise basis.

$\Rightarrow P_1 \parallel P_2 \parallel P_3 \dots \parallel P_k$ [iff $P_i \parallel P_j$ for
all i, j]

Note → If $I_1 \cap I_2 \neq \emptyset \rightarrow$ it does not prevent parallelism
or \Rightarrow (I_1, I_2 to be fault navi paths)

Properties of Parallelism relation

(1) Commutative ($P_i \parallel P_j \Leftrightarrow P_j \parallel P_i$)

(2) Non Transitive (if $P_i \parallel P_j$ & $P_i \parallel P_k$,
does not imply $P_j \parallel P_k$)

(3) Associative [$(P_i \parallel P_j) \parallel P_k = P_i \parallel (P_j \parallel P_k)$]

Eg- $P_1 \Rightarrow C = D+E$

$P_2 \Rightarrow M = G+E$

$P_3 \Rightarrow A = B+C$

$P_4 \Rightarrow C = L+M$

$P_5 \Rightarrow F = G+E$, Check Meim

Starts	I_1	O_1	I_2
P_1	{D, E}	{C}	
P_2	{G, C}	{M}	
P_3	{B, C}	{A}	
P_4	{L, M}	{C}	
P_5	{G, E}	{F}	

You have to check all pairs $\Rightarrow 5C_2 = 10$ pairs

a) For $P_1 \& P_2$: ① $I_1 \cap O_2 = \emptyset$

$\{D, E\} \cap \{M\} = \emptyset$

② $I_2 \cap O_1 = \emptyset$

$\{G, C\} \cap \{C\} \neq \emptyset$

$P_1 \nparallel P_2$

b) For $P_1 \& P_3$: ②: $P_1 \nparallel P_3$

c) $P_1 \nparallel P_4$

d) $P_1 \nparallel P_5$

e) $P_2 \nparallel P_3$

f) $P_2 \nparallel P_4$

g) $P_2 \nparallel P_5$

h) $P_3 \nparallel P_4$

i) $P_3 \nparallel P_5$

As studied earlier, there can be delays in pipeline

Types of Pipeline Delay

① Uniform Delay Pipeline

All stages / segments will complete their operations by taking the same time,

$\text{cycle Time } (T_p) = \text{Stage Delay}$

If there are buffers b/w the stages, then

$T_p = \text{Stage Delay} + \text{Buffer Delay}$

② Non-uniform Delay Pipeline

All stages will complete their operations by taking different times,

$T_p = \max(\text{Stage Delay})$

eg - In 4-segment pipeline, delays are

$1, 2, 3, 4 \Rightarrow T_p = \max(1, 2, 3, 4)$

$T_p = 4 \text{ ns}$

If buffers are there, $\Rightarrow \text{Max}(\text{Stage Delay} + \text{Buffer Delay})$

eg - Assume 4-segment pipeline having stage delay as 2ns, 8ns, 3ns, 10ns. Here, we have to determine the time taken to execute 100 tasks in the above pipeline

$\text{Here, } T_p = 10 \text{ ns.}$

Page No. _____
Date: / /

Page No. _____
 no. of cycles without stall for stall
 when branch inst add stall
 # done before

$$\begin{aligned}
 \text{Execution time} &= (k+n-1) T_p = \text{no. of cycles} \times T_p \\
 &= (4+100-1) \times 10 \\
 &= 1030 \text{ ns}
 \end{aligned}$$

Formula for Performance of pipeline with stalls

(1)

$$\begin{aligned}
 \text{speedup} &= S = \frac{\text{Avg exec time non-pipeline}}{\text{Avg exec time pipeline}} \\
 &= \frac{(CPI \times \text{Cycle Time}) \text{ non-pipeline}}{(CPI \times \text{Cycle Time}) \text{ pipeline}}
 \end{aligned}$$

Note \Rightarrow In the pipelined processor, ideal CPI = 1 but because of the stalls CPI becomes greater than 1.

Then speed up becomes -

$$S = \frac{(CPI \times \text{cycle time}) \text{ non-pipeline}}{(1 + \text{no. of stalls per inst}^*) \times (\text{cycle time}) \text{ pipeline}}$$

(2)

$$\text{If } (\text{cycle time}) \text{ pipeline} = (\text{cycle time}) \text{ non-pipeline}$$

then

$$\text{speedup} = \frac{(CPI) \text{ non-pipeline}}{1 + \text{no. of stalls per instruction}}$$

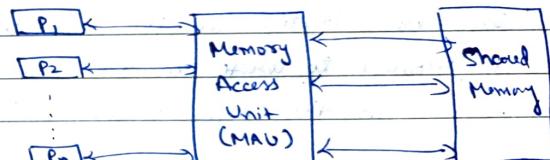
PRAM Model

(Parallel Random Access Machine)

PRAM is an abstract machine for designing an algorithm applicable to parallel computers.

A PRAM Model contains -

- ① a set of similar type of processes.
- ② all processors share a common memory unit.
processes can communicate among themselves through shared mem only
- ③ A Memory Access Unit (~~MAU~~) (MAU) connects the processors with the single shared memory



PRAM Architecture

'n' no. of processors can independent operations on 'n' no. of data in a particular unit of time.

Problem \Rightarrow This may result in simultaneous access of same memory location by different processors.

To solve this, some constraints are there:

① Exclusive Read Exclusive Write (EREW)

No 2 processors are allowed to read from or write to the same memory location simultaneously.

② Exclusive Read Concurrent Write (ERCW)

No 2 processors are allowed to read from the same memory location at the same time but are allowed to write at the same memory location simultaneously.

③ Concurrent Read Exclusive Write (CREW)

All processors are allowed to read from the same mem location at the same time but not allowed to write to the same mem location at the same time.

④ Concurrent Read Concurrent Write (CRCW)

All processors are allowed to read from and write to the same mem location at the same time.

Methods to implement PRAM Model

- ① Shared Memory Model
- ② Message Passing Model
- ③ Data Parallel Model