

- To detect the data-dependency condition some logic should be implemented in ID stage of the pipeline.
- Structure of a logic is -

Tomasulo Algorithm :-

- This algorithm is used to detect data-dependency in a pipeline. It work on ID stage of each Instⁿ.

Eg:- $I_1: r_0 \leftarrow r_1 + r_2$

$I_2: r_3 \leftarrow r_0 * r_2$

(i) ID-Stage

(a) Instⁿ decode

S.no	Func ⁿ	Dest ⁿ	Indep- endent src1	Indep- endent src2	Depen- dent src1	Depend- ent src2	Status
I_1	ADD	r_0	r_1	r_2	—	—	0
I_2	MUL	r_3	—	r_2	r_0	—	0

(b) Access reg-file -

I_1	$r_1 \leftarrow \text{value}$ $r_2 \leftarrow \text{value}$
I_2	$r_0 \leftarrow X$ $r_2 \leftarrow \text{value}$

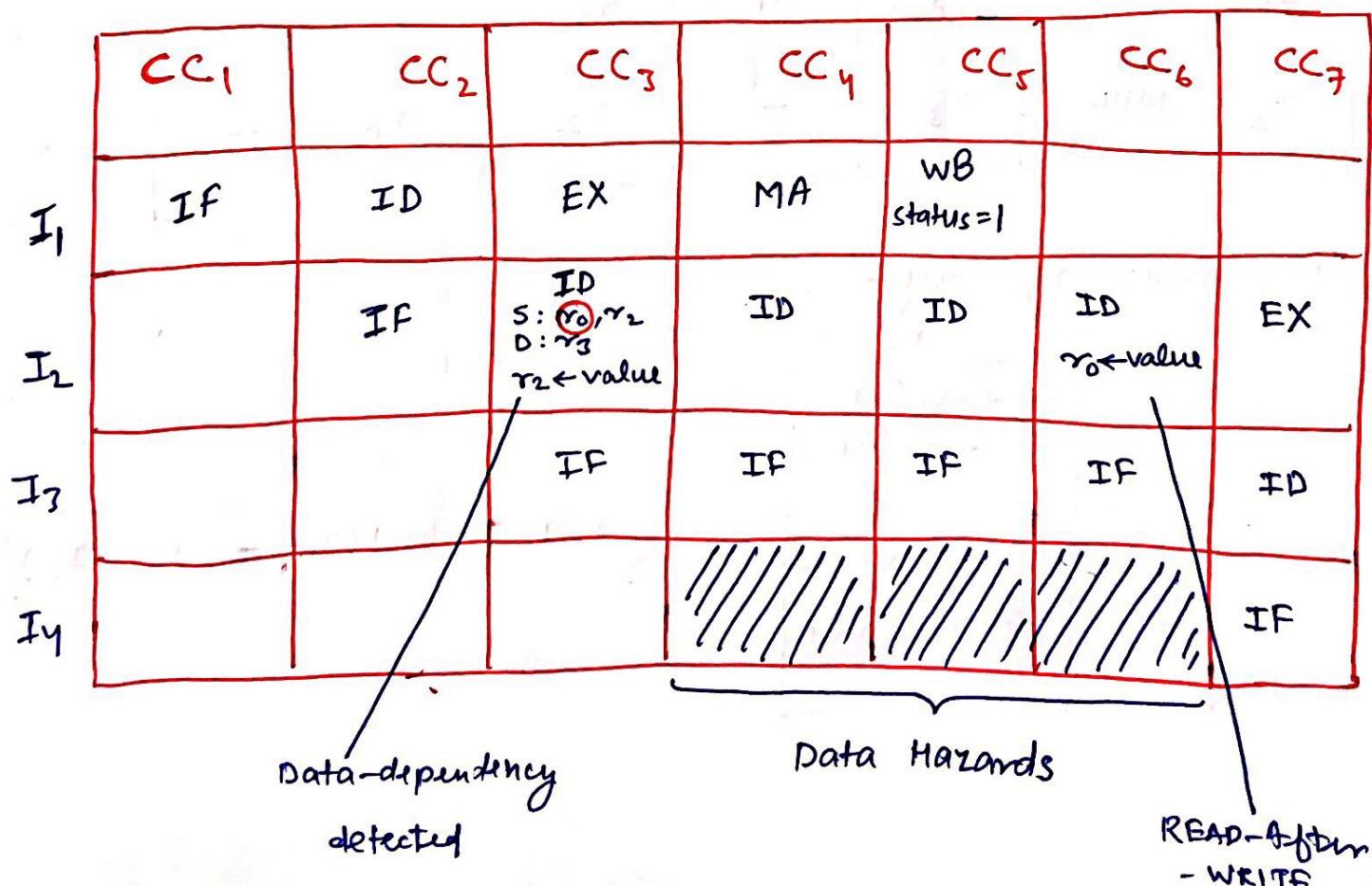
check ($\text{Dest}(I_1) = \text{Src}(I_2)$)

(ii)

Ex - stage

I ₁	EX - stage Allocated
I ₂	EX - Stage is not Allocated (I ₂ will be waiting on ID-stage until I ₁ status becomes L)

- By using this logic, we can stop the accessing of a unwanted data in the pipeline execution. Therefore, dependent Instⁿ will wait on ID- stage until the data become available in register.



- To minimize the data-hazard in pipeline, H/W technique is used i.e. operand forwarding also called as bypassing or short circuiting.
- Operand Forwarding techniques states that use the buffers b/w the stages to hold the intermediate result so that dependent instⁿ will access the new value from the buffer before updating the register file. In this process last-stage opⁿ is modified to perform the i.e. WRITE-READ opⁿ in the same cycle means updated data immediately available for READ access in the same cycle. Therefore buffer is not required at the end of last stage.
- consider the following program code, execute in the pipeline using operand forwarding technique -

Code :-

```

I1: ADD r0, r1, r2
I2: SUB r3, r0, r2
I3: MUL r4, r3, r0
I4: DIV r5, r4, r0
  
```

NOTE :- Adjacent Data-Dependency is called as TRUE - Data - Dependency

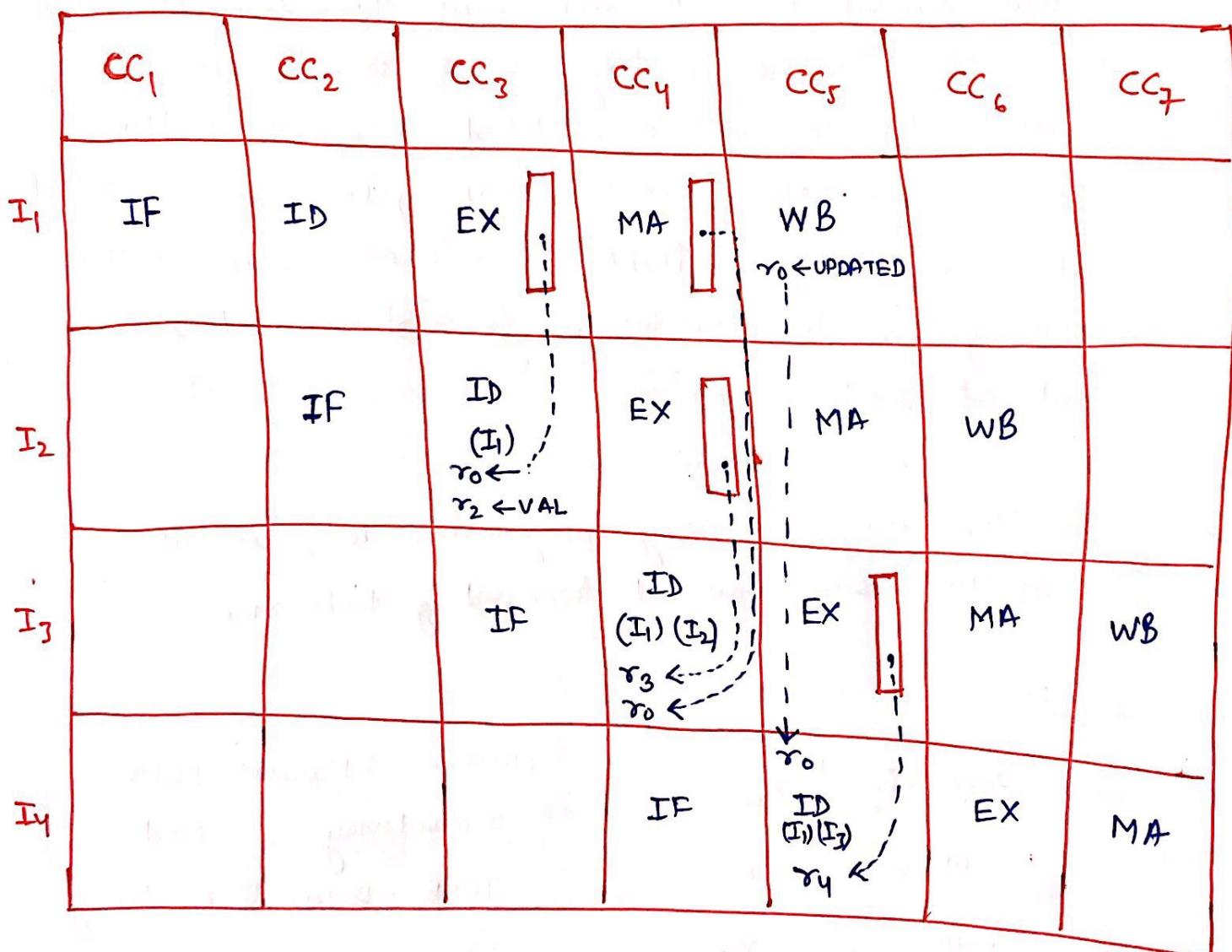
i.e.

$$\begin{bmatrix} I_2 - I_1 \\ I_3 - I_2 \\ I_4 - I_3 \end{bmatrix}$$

NOTE - Non-Adjacent Data dependency analysis is data-dependency only. But it is eliminated through True-data-dependency

i.e.

$$\begin{bmatrix} I_3 - I_1 & (r_0) \\ I_4 - I_1 & (r_0) \end{bmatrix}$$



$$t_p = \max \left(\begin{array}{c} \text{stage} \\ \text{- Delay} \end{array} + \begin{array}{c} \text{Buffer} \\ \text{Delay} \end{array} \right)$$

↓ ↓

$4ns$ $1ns$

Q: Consider following Program , execute on a RISC pipeline. How many cycles are required to complete the program. Assume that all the instⁿ are spending 1 cycle on all the stages but LOAD instⁿ takes 3-cycles on 4th stage of a Pipeline.

I₁ : LOAD r₀, 3(r₁)

I₂ : ADD r₂, r₀, r₁

I₃ : LOAD r₃, 4(r₄)

I₄ : SUB r₅, r₃, r₄

	IF	ID	EX	MA	WB	
I ₁	1	1	1	3	1	LOAD
I ₂	1	1	1	1	1	ADD
I ₃	1	1	1	3	1	LOAD
I ₄	1	1	1	1	1	SUB

Now, generate reservation table -

	cc_1	cc_2	cc_3	cc_4	cc_5	cc_6	cc_7	cc_8	cc_9	cc_{10}	cc_{11}	cc_{12}	cc_{13}	cc_{14}	cc_{15}
I_1	IF	ID	EX	MA	MA	MA	MA	MA	WB						
I_2	IF	ID	ID	MA	MA	MA	MA	MA	WB						
I_3	IF	ID	IF	EX	MA	MA	MA	MA	MA	MA	MA	MA	MA	MA	WB
I_4	IF	ID	IF	ID	EX	MA	MA	MA	MA	MA	MA	MA	MA	MA	WB

* Since due LOAD opⁿ data (dependent) will be present at 'MA' stage.

Q. Consider the following program executed on 4-stages (IF, ID, EX, WB) pipeline. Assume that all instⁿ are spending 1-cycle on all the stages but 'MUL' Instⁿ takes 4-cycles on EX stage and 'DIV' Instⁿ takes 5-cycles at EX stage. Operand Forwarding mechanism is used as an optimization technique in the pipeline to handle the data dependency.

- (a) How many cycles are required to complete prog.
- (b) How many cycles are saved using optimization over non-optimization.

I₁: MUL r₀, r₁, r₂

I₂: DIV r₃, r₁, r₂

I₃: ADD r₄, r₂, r₁

I₄: SUB r₃, r₁, r₂

	IF	ID	EX	WB	
I ₁	1	1	4	3	MUL
I ₂	1	1	5	2	DIV
I ₃	1	1	1	1	ADD
I ₄	1	1	1	1	SUB

With optimization (Operand - Forwarding)

	cc_1	cc_2	cc_3	cc_4	cc_5	cc_6	cc_7	cc_8	cc_9	cc_{10}	cc_{11}	cc_{12}	cc_{13}	cc_{14}	cc_{15}
T_1	IF	ID	EX	EX	EX	EX	EX	WB							
T_2	IF	ID	ID	ID	ID	ID	EX	EX	EX	EX	EX	EX	EX	WB	
T_3	IF	IF	IF	IF	IF										
T_4															

(Total cycle = 14)

without optimization (no - operand Forwarding)

	cc_1	cc_2	cc_3	cc_4	cc_5	cc_6	cc_7	cc_8	cc_9	cc_{10}	cc_{11}	cc_{12}	cc_{13}	cc_{14}	cc_{15}	cc_{16}	cc_{17}	cc_{18}	cc_{19}
T_1	IF	ID	EX	EX	EX	EX	EX	EX	EX	EX	WB								
T_2	IF	ID	ID	ID	ID	ID	ID	ID	ID	WB									
T_3			IF	IF	IF	IF	IF	IF	IF	IF									
T_4																			

(Total cycles = 18)

$$\begin{aligned} \text{Cycle - saved} &= 18 - 14 \\ &= \underline{\underline{4 - \text{cycles}}} \end{aligned}$$

(iii) Control Dependency :-

- It occurred in the pipeline when the transfer of control & instⁿ are executed in the pipeline.

Code:-

1000 : I ₁	↓	Falling/through path
1001 : I ₂		
1002 : I ₃ (JMP 2000)	↓	
1003 : I ₄		
:	:	
2000 : BI ₁	↓	Taken path
2001 : BI ₂		
:	:	

Expected o/p sequence -

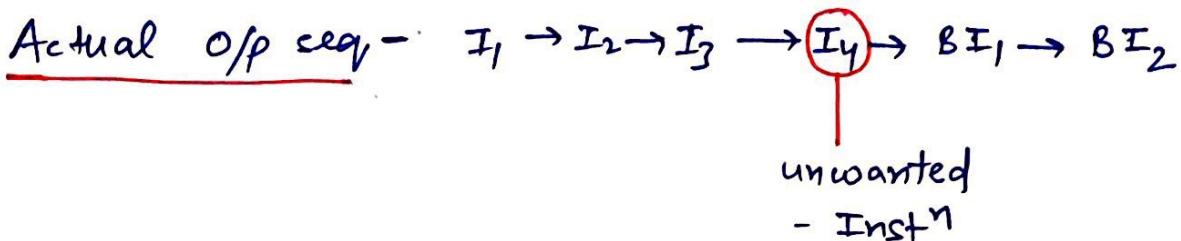
(I₁ → I₂ → I₃ → BI₁ → BI₂...)

'IF' mprog is -

T ₁ : PC → MAR
T ₂ : M[MAR] → MBR PC + I → PC
T ₃ : MBR → IR

PC:1000	cc ₁	cc ₂	cc ₃	cc ₄	cc ₅	cc ₆
I ₁	IP PC:1001	ID	EX	MA	WB	
I ₂		IF PC:1002	ID	EX	MA	WB
I ₃			IF PC:1003	ID PC:2000	EX	MA
I ₄				IF PC:1004	ID	EX
BI ₁					IF PC:2001	ID
BI ₂						IF PC:2002

IF of I₄ make
PC ← 1004 but
during ID of I₃
PC ← 2000 (JMP instⁿ)



- In the above exec. sequence, unwanted Instⁿ is executed in the pipeline. So program output is missing (unwanted behaviour). This situation is known as control-dependency.
- To handle such problems, FLUSH opⁿ or FREEZE opⁿ is used. This opⁿ states that, insert the NOP instⁿ after the JMP instⁿ in the program to suspend the unwanted instⁿ fetch.

Code:-

1000: I_1
 1001: I_2
 1002: I_3 (JMP 2000)
 1003: ~~I_4~~ (NOP)
 1004: I_5
 1005: I_6
 ; ;
 2000 : BI_1
 2001 : BI_2
 ;

	cc ₁	cc ₂	cc ₃	cc ₄	cc ₅	cc ₆
PC: 1000						
I_1	IF PC: 1001	ID	EX	MA	WB	
I_2		IF PC: 1002	ID	EX	MA	WB
I_3			IF PC: 1003	ID PC: 2000	EX	MA
NOP				IF PC: 1004	ID	EX
PC: 2000					IF PC: 2001	ID
BI_1						IF PC: 2002
BI_2						

Actual - O/P - seq : $I_1 \rightarrow I_2 \rightarrow I_3 \rightarrow \boxed{NOP} \rightarrow BI_1 \rightarrow BI_2$

unwanted Instⁿ
+ stall

- Number of stalls created in the pipeline due to a branch instn is called as branch penalty.
- It depends on the availability of target address in the pipeline i.e branch penalty is equal to & at what stage target Address available - 1.

(branch penalty \rightarrow stage in which target Address available - 1)

NOTE:- In RISC pipeline branch penalty is 1 as target address is available in 2nd (ID) stage.

NOTE:- For a hypothetical K-stage pipeline , branch penalty is $K-1$.

- if stage number of TA :- given

$$(\text{Branch - penalty} = \text{stage number} - 1)$$

- if stage name of TA :- given

$$(\text{Branch - penalty} = \text{corresponding stage number} - 1)$$

- until the Instⁿ is completed (all Instⁿ are proceed through all stages) -

$$\text{Branch penalty} = \text{Last stage in pipeline} - 1$$

NOTE:- Total no. of stalls created from the branches during the prog execⁿ is -

$$\left(\frac{\text{Branch frequency}}{\text{Branch Inst}^n / \text{prog}} * \frac{\text{Branch penalty}}{\text{stalls} / \text{branch inst}^n} \right)$$

NOTE :- To minimize the control hazard (branch penalty), h/w technique is used i.e. branch prediction buffer also called as branch target buffer or loop buffer.

- Branch - prediction buffer is a high-speed buffer present in IF(stage) used to hold the predicted target address. when target address present in first stage then there is not stall in pipeline.
- **NOTE :-** to minimize the control hazard, S/W mechanism is used i.e. delayed branch.

Delayed Branch :- It is a compiler technique, so compiler will rearrange the code if possible, otherwise substitute the NOP instⁿ after JMP instⁿ to preserve the execution path.

User-code

```

1000: I1
1001: I2
1002: I3 (JMP 2000) ] SWAP
1003: I4
|
2000: BI1
2001: BI2

```

COMPILER - M/C code

(a) Rearrangement

Addr: I₁

```

[ I3 ( JMP BI1)
  I2
  I4
  :
  BI1
  BI2
]

```

	CC ₁	CC ₂	CC ₃	CC ₄	CC ₅
PC: I ₁	IF	ID	EX	MA	WB
I ₁	PC: I ₃				
I ₃	IF	ID	EX	MA	
	PC: I ₂	PC: BI ₁			
I ₂		IF	ID	EX	
		PC: I ₄			
BI ₁			IF	ID	
BI ₂			PC: BI ₂		
				IF	
				PC: BI ₃	

O/P-sequence = I₁ - I₃ - I₂ - BI₁ - BI₂

- o/p sequence is different but program result will be same , and there is no STALL .

NOTE :- Instⁿ rearrangement swapping may be possible in unconditional branch opⁿ but no rearrangement can be taken place in case of ^{conditional} branch Instⁿ.

(b) NOP - substitution -

MIC generated -

Addr : I₁

I_2 I_3 (JMP ~~2~~ BI₁)

NOP

Iy

6

2

81

BI_2

O/p sequence -

$$I_1 - I_2 - I_3 = NOP - BI_1 - BI_2$$

$$-BI_2$$

Instruction Scheduling :-

- CPU always execute the program in a sequence from I₁ to I_n called as in-order execution.
 - In the in-order execution sequence, if any instⁿ is dependent then remaining instⁿ are also sharing the STALL cycle even they are independent.

$I_1 : ADD \quad r_0, r_1, r_2$
 $I_2 : SUB \quad r_3, r_0, r_4$
 $I_3 : MUL \quad r_4, r_5, r_6$
 $I_4 : DIV \quad r_3, r_7, r_8$

In-order execution:-

$I_1 - I_2 - I_3 - I_4$
stalls

- In the above execution sequence, I_2 dependent on I_1 so I_2 waiting until the I_1 execution is completed. Waiting creates stall. These stalls are shared by I_3 and I_4 instⁿ even they are independent.
- To handle the above problem instⁿ scheduling concept is used. This concept states that execute the independent instⁿ first called out of order (re-order) in instⁿ execution.

$(I_1 - I_3 - I_4 - I_2)$

- Re-order execution creates two more dependencies in the pipeline.

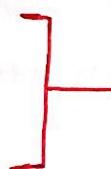
- (i) Anti-data dependency
- (ii) Output data-dependency

Instⁿ Instⁿ

READ - before - WRITE

WRITE - before - READ

WRITE - before - WRITE



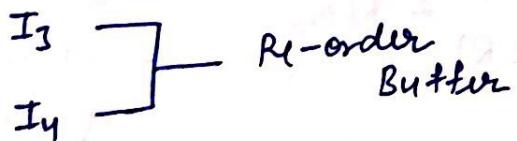
three types of hazard

- Anti data dependency will be occurred when Instⁿ J WRITE the data before Inst I reads it. eg - I₃ exec before I₂, so I₃ updates the register r₄ before I₂ reads it ∴ I₂ incorrectly access new value from reg (r₄) .
- O/p data -dependency will be occurred when the Instⁿ J tries to WRITE the data before Instⁿ I WRITE it. eg - I₄ exec before I₂, so I₄ updates the register (r₃) before I₂ WRITES it ∴ destination (r₃) updates with a old value (Data Loss)

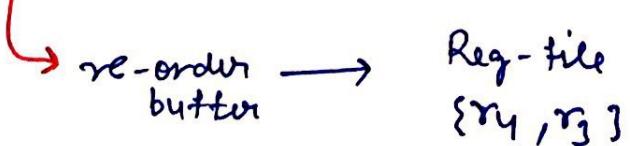
NOTE:- To minimize the anti & o/p data -dependency & stalls h/w mechanism is used i.e register renaming.

- Register renaming states use re-order buffers to store the out-of-order instⁿ o/p. Later update the reg - file with a re-order buffer contents after completion of dependent Instⁿ execution.

$$I_1 \rightarrow r_0$$

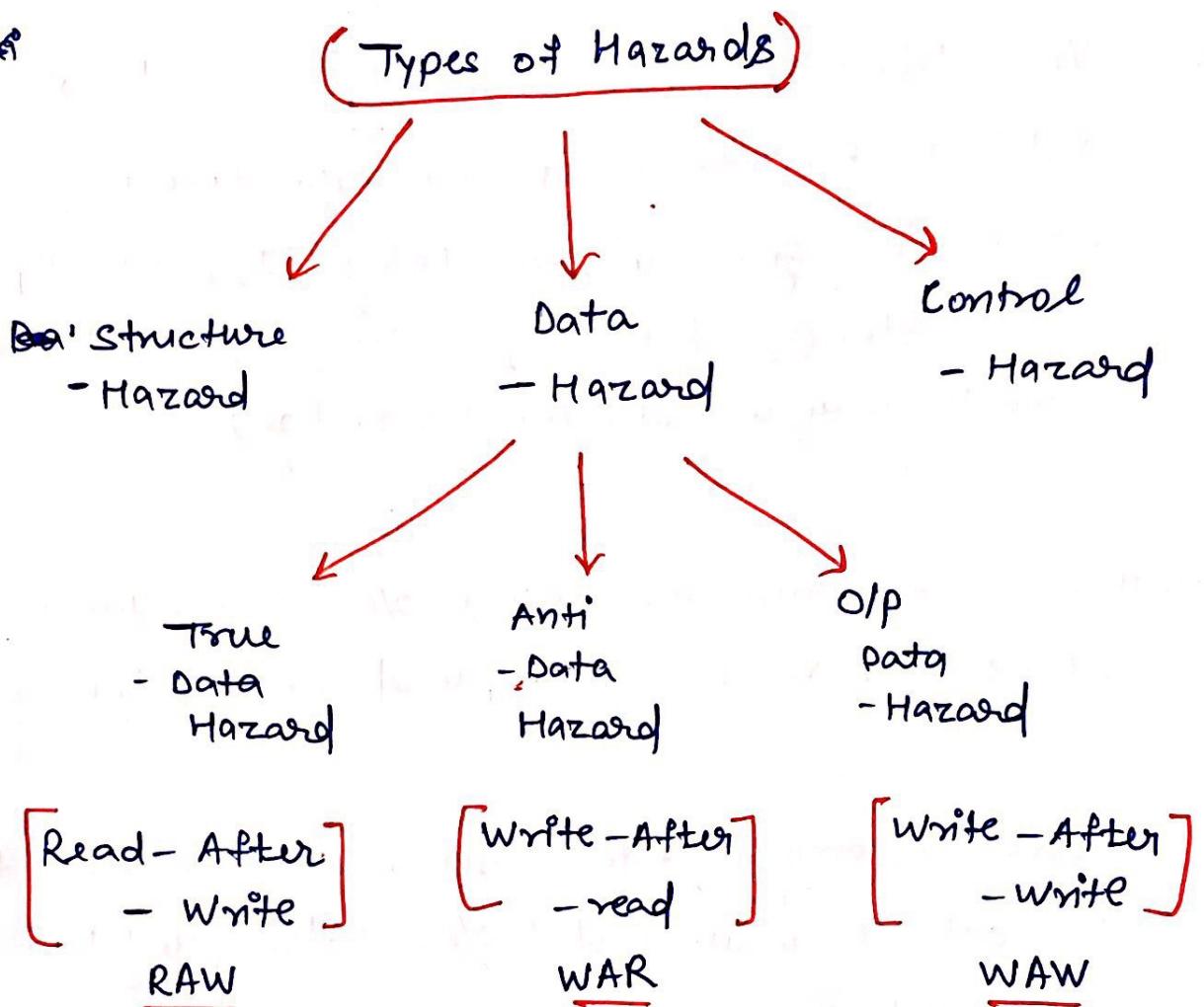


$$I_2 \rightarrow r_3 \text{ (STATUS=1)}$$



Hazard :-

- Hazard is a delay.
- Delay is present in the pipeline due to a dependency problem.
- ~~Data~~

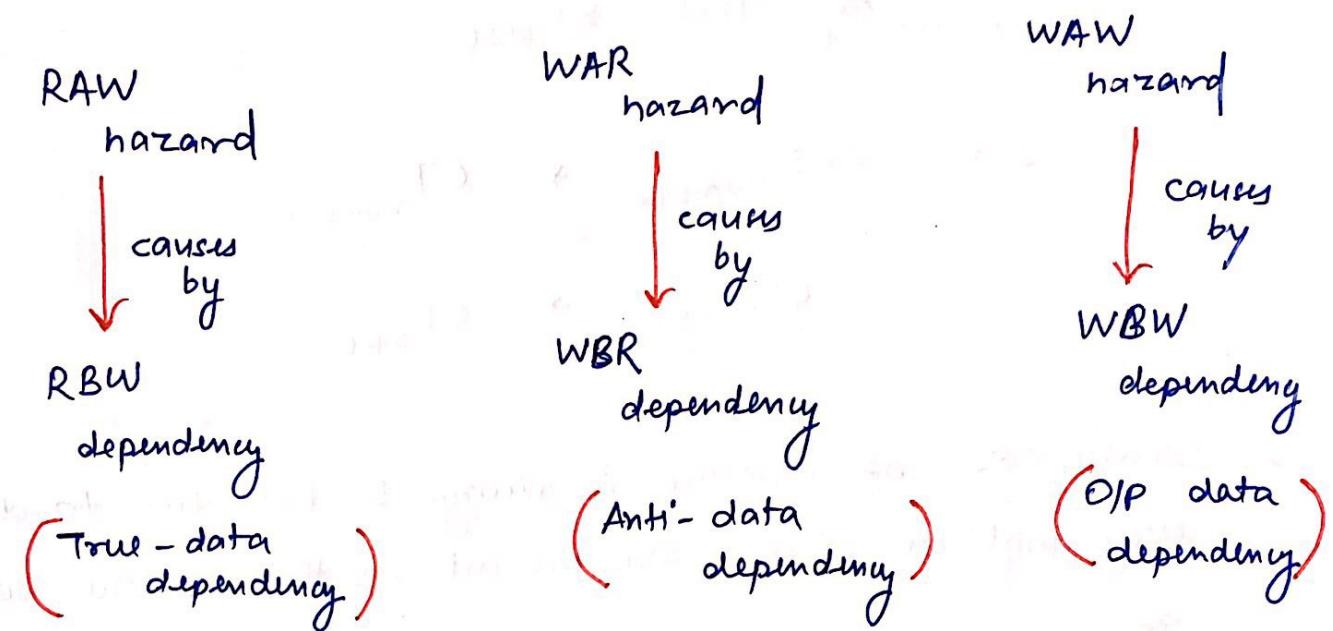


{ it is created when
the instⁿ J tries to
READ the data before
instⁿ I WRITES it.
(RBW).

{ it is created
when instⁿ J tries
to WRITES the data
before instⁿ I READS
it (WRB).

{ it is created
when instⁿ J
tries to WRITE
data before instⁿ
I WRITE it
(WBW)

NOTE:- Read- After -Read (RAR) is not a hazard as
 Read- before- Read (RBR) is not any dependency problem.



Inst ⁿ J	Inst ⁿ I	Type of Dependency	Type of Hazards
I/P reg	O/P reg	True Data	RAW (in-order)
O/P reg	I/P reg	Anti- Data	WAR (out-order)
O/P reg	O/P reg	O/P Data	WAW (out-order)
I/P reg	I/P reg	—	—
ALU	ALU	Structural	Structural Hazard
JMP		control	control Hazard

Performance Analysis with Stalls :-

$$S = \frac{\text{Avg. Inst}^n \text{ ET}_{\text{non-pipe}}}{\text{Avg. Inst}^n \text{ ET}_{\text{pipe}}}$$

$$S = \frac{\text{CPI}_{\text{non-pipe}} * \text{CT}_{\text{non-pipe}}}{\text{CPI}_{\text{pipe}} * \text{CT}_{\text{pipe}}}$$

- Ideal CPI of pipeline is always 1 but due to dependency problem stalls are present in the pipeline so

~~Stalls~~

$$\left[S = \frac{\text{CPI}_{\text{non-pipe}} * \text{CT}_{\text{non-pipe}}}{\left(1 + \frac{\text{stalls}}{\text{Inst}^n} \right) \text{CT}_{\text{pipe}}} \right]$$

- when the pipelines stages are perfectly balanced, then 1-task ET is also equal to no. of stages in the pipeline (i.e $t_n = k \cdot t_p$). under this condition

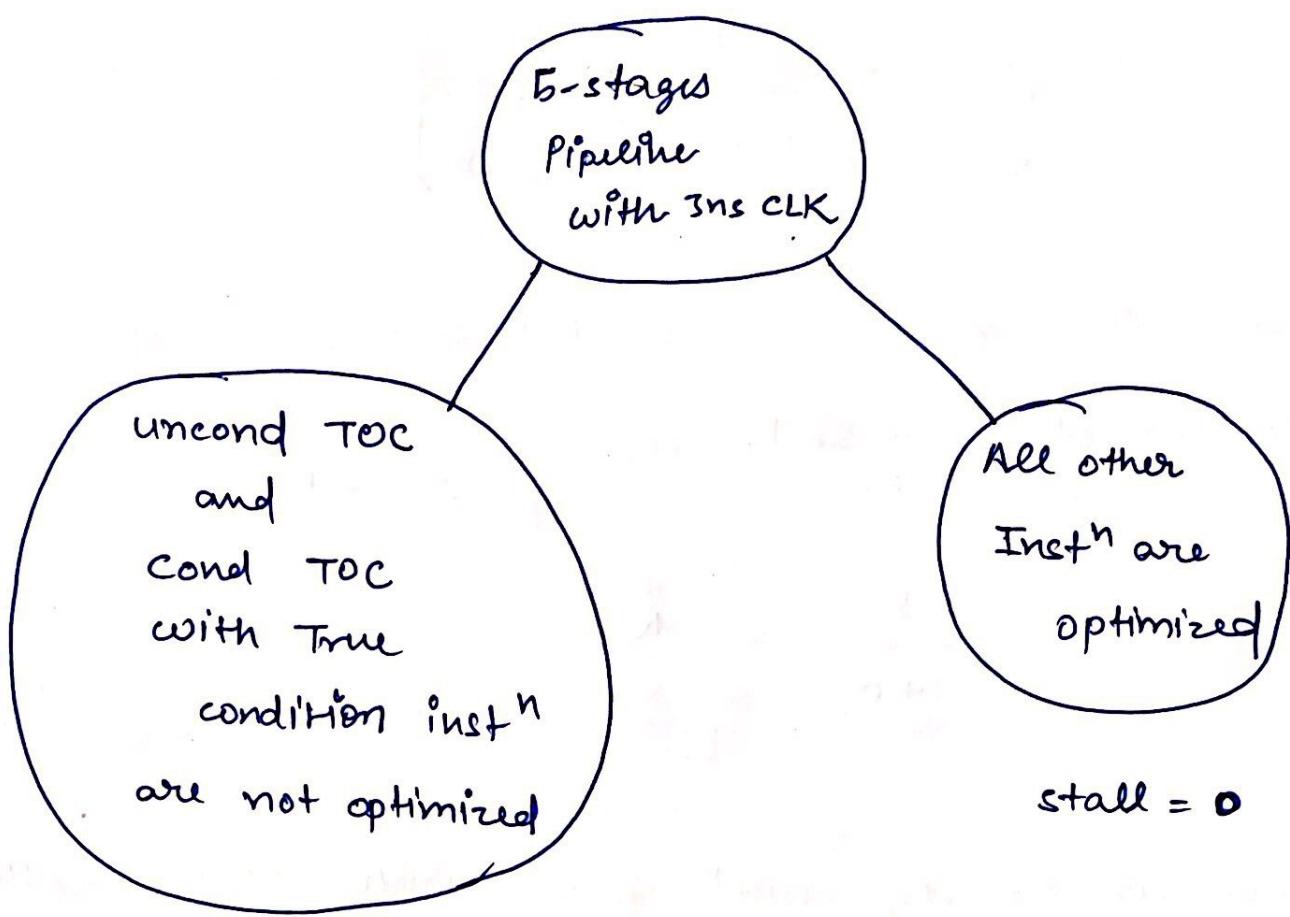
$$S = \frac{k \cancel{t_p}}{\left(1 + \frac{\text{stalls}}{\text{Inst}} \right) \cancel{t_p}}$$

$$\left[S = \frac{K}{1 + \text{stalls/Instn}} \right]$$

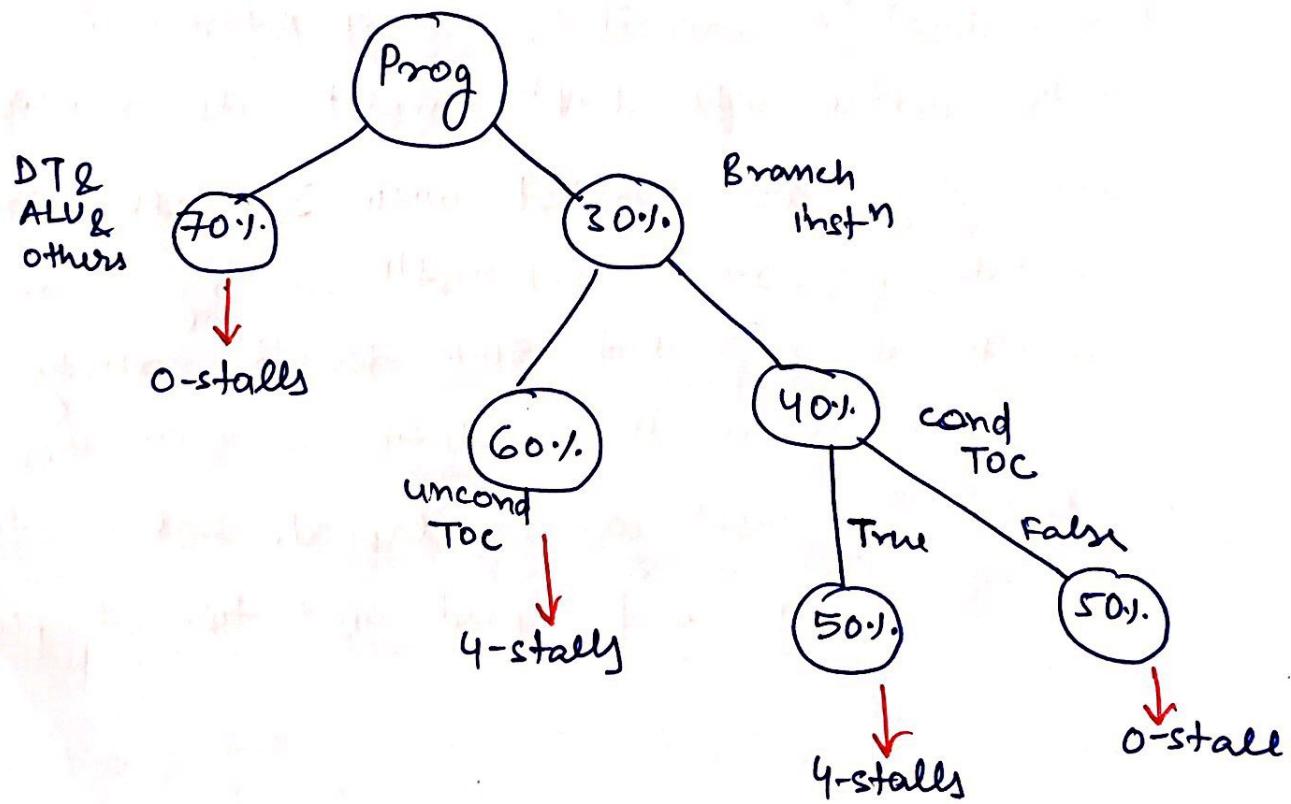
- when the system is operated with 100% efficiency then no stalls present. So,

$$S = \frac{K}{1+0} = K$$

Q. Consider 5-stage "instn" pipeline which allows all the "instn" except branch "instn". Processor stops the fetching of a following "instn" after the branch "instn" until the branch "instn" is completed. Target Address is available in the pipeline when "instn" completes its execution. Pipeline stages are operated with 3ns clk. Program contains 30% branch "instn" among them 40% are conditional in - which 50% doesn't satisfy the condition when the condition is false then the following "instn" are overlapped. Calculate the av. "instn" ET and speed-up factor of pipeline.



$$\text{stalls} = K-1 = 5-1 = 4$$



$$\text{stalls / Inst}^n = (0.7 * 0) + (0.3 \times 0.6 \times 4) + (0.3 \times 0.4 \times 0.5 \times 4) + (0.3 \times 0.4 \times 0.5 \times 0)$$

$$= \underline{\underline{0.96}}$$

(a) Avg. instⁿ ET_{pipe} = $(1 + \text{stall}/\text{Inst}^n) \cdot T_{\text{pipe}}$

$$= (1 + 0.96) (3 \text{ns})$$

$$= \underline{\underline{5.88 \text{ns}}}$$

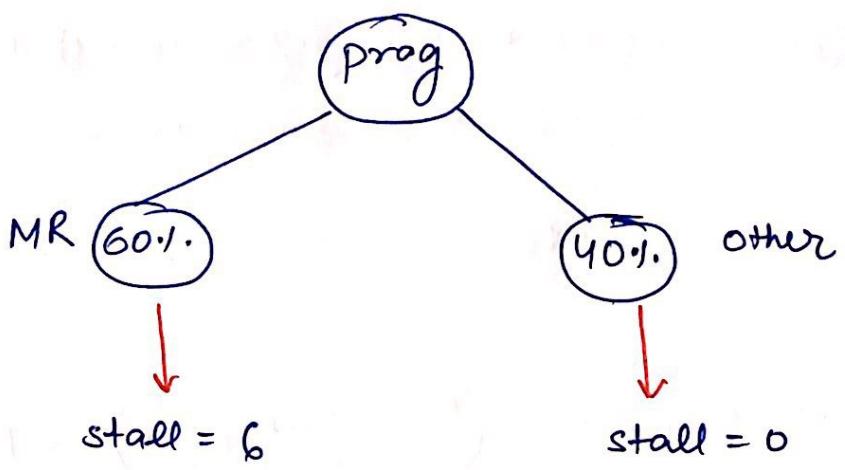
(b)

Performance gain $s = \frac{k}{1 + \text{stall}/\text{Inst}^n}$

$$= \frac{5}{1 + 0.96}$$

$$(s = 2.55)$$

- Q. Consider 8-stage pipeline operated with 4ns CLK, which allows all the instⁿ except the memory reference instⁿ. MR instⁿ penalty is 6-cycles. Program contains 60% MR instⁿ what is avg. instⁿ execution time.



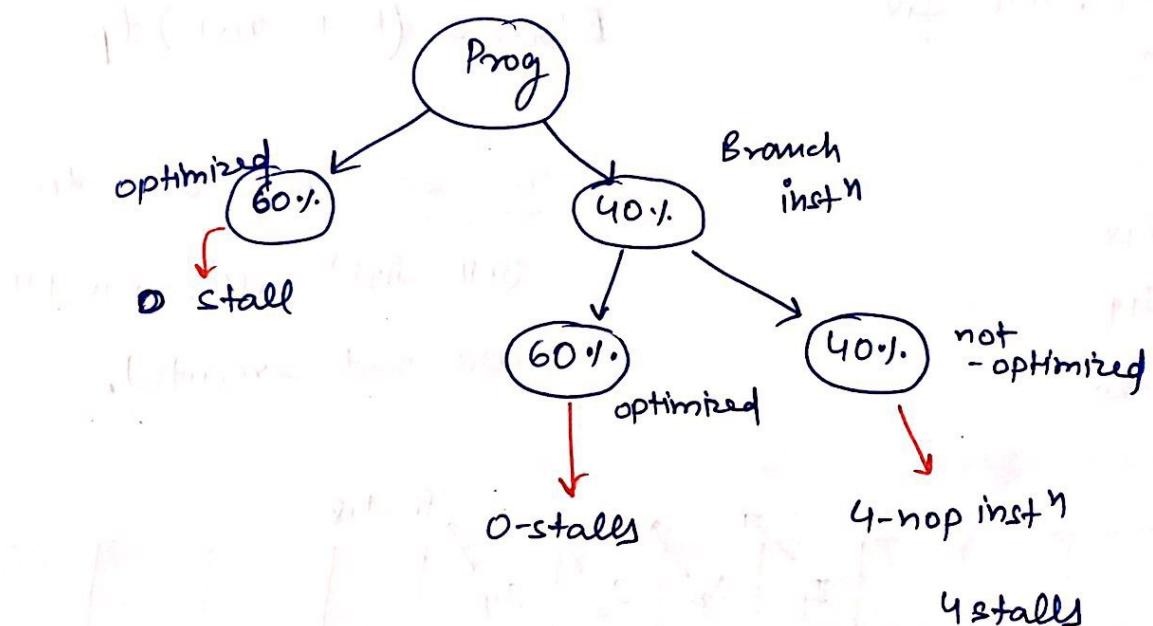
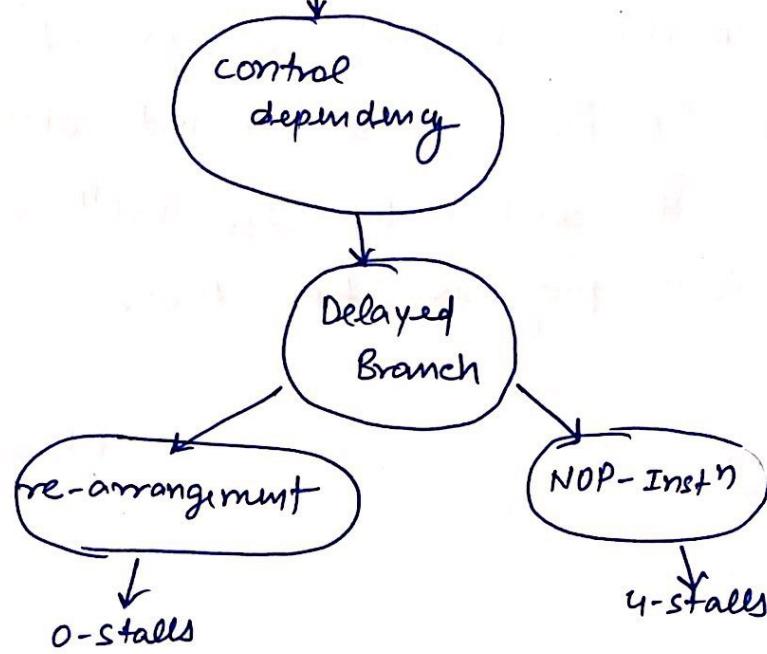
$$CT_{\text{pipe}} = 4 \text{ ns}$$

Avg. stalls / $Inst^n$ = $(0.6 \times 6) + (0.4 \times 0)$
 $= \underline{\underline{3.6}}$

$$\begin{aligned} \text{Avg. } Inst^n \text{ ET}_{\text{pipe}} &= \left(1 + \frac{\text{stalls}}{Inst^n}\right) CT_{\text{pipe}} \\ &= (1 + 3.6)(4 \text{ ns}) \\ &= \underline{\underline{18.4 \text{ ns}}} \end{aligned}$$

Q. Consider 6-stage pipeline with 2ns cblk uses, delay branch concept to optimize the control dependency stalls. In the delayed branch re-arrangement proc. is used to avoid the stalls, if not possible to re-arrangement then 4 NOP instⁿ are substituted after the branch instⁿ in the program. Program contain 40% branch instⁿ among them 60% are optimized, what is the avg. instⁿ ET.

6-stage
pipeline with 2ns
clock



$$\text{stalls/instr}^n = (0.4 \times 0.4 \times 4)$$

$$= \underline{\underline{0.64}}$$

$$\begin{aligned}
 \text{avg. instr}^n \text{ ET}_{\text{pipe}} &= (1 + \text{stalls/instr}^n) \text{ t_p} \\
 &= (1 + 0.64)(2 \text{ ns}) \\
 &= \boxed{(3.28 \text{ ns})}
 \end{aligned}$$

Q. Consider 5-stage pipeline without branch-prediction, use to execute the prog which contain 20 instⁿ ($I_1 - I_{20}$).

Pipeline stage delays are (2ns, 4ns, 8ns, 3ns, 1ns).

In the pipeline all instⁿ are proceed through all stages. If instⁿ is a uncond JMP instⁿ, it transfer the control to I_{18} instⁿ during its execution. What is the prog execution time.

Program

$$(K = 5, t_p = 8\text{ns})$$

I_1

I_2

I_3

$I_4 : \text{JMP } I_{18}$

I_5

:

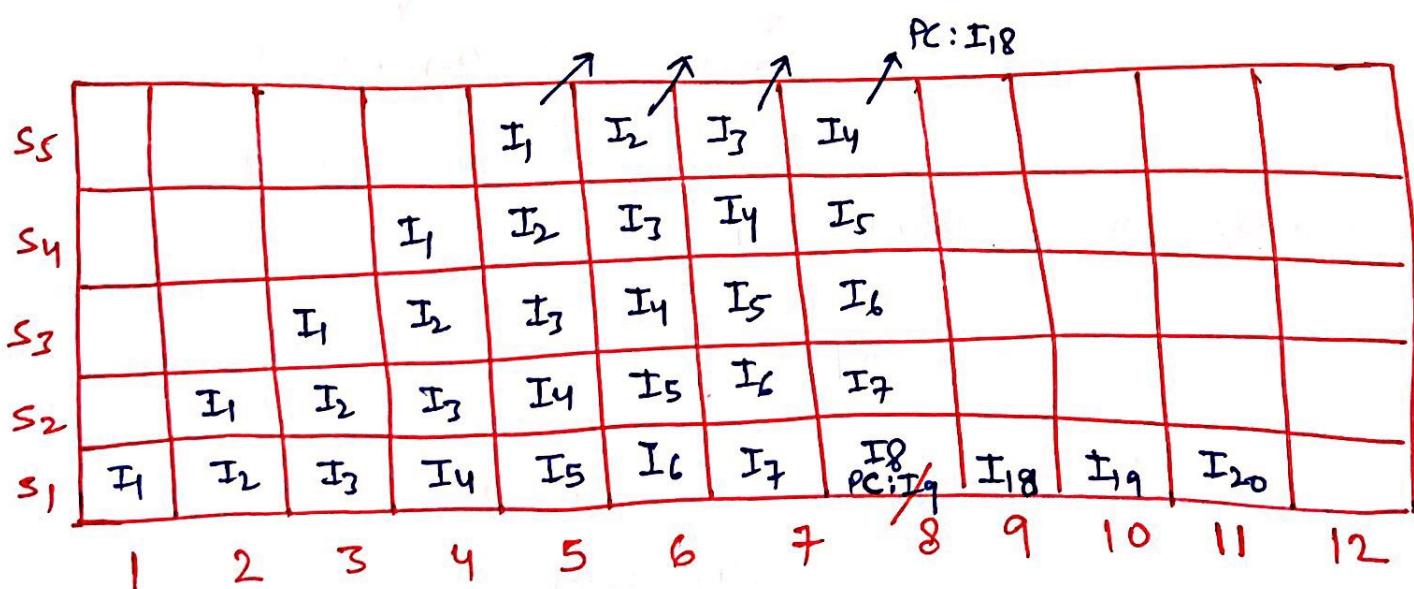
I_{18}

I_{19}

I_{20}

$$ET_{\text{pipe}} = (K + n-1) t_p$$

$n = 20$, but due to
JMP instⁿ all 20 instⁿ
are not executed.



so, effective no. of input (n) = 11

$$\begin{aligned} \text{so, } ET_{\text{pipe}} &= (K + n - 1) t_p \\ &= (5 + 11 - 1)(8 \text{ ns}) \\ &= 120 \text{ ns} \end{aligned}$$

Q. Consider the following code, executed on the pipeline CPU.

$$I_1: r_0 \leftarrow r_1 + r_2$$

$$I_2: r_1 \leftarrow r_0 - r_2$$

$$I_3: r_0 \leftarrow r_1 * r_2$$

$$I_4: r_2 \leftarrow r_0 * r_1$$

$$I_5: r_1 \leftarrow r_2 - r_0$$

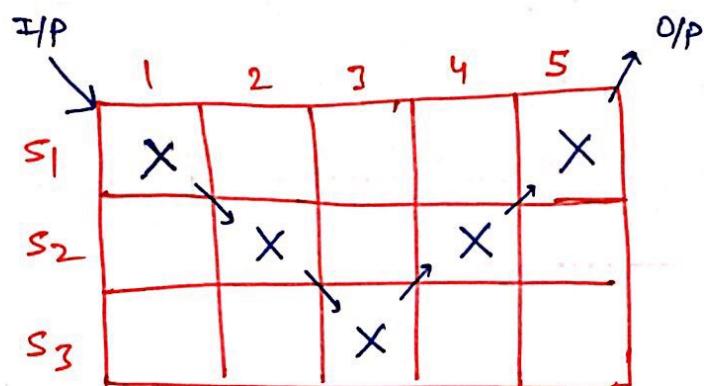
$$I_6: M[X] \leftarrow r_1$$

How many RAW, WAR & WAW Hazards possible?

RAW [In-order] (IN-OUT)	WAR [out-order] (OUT-IN)	WAW [out-order] (OUT-OUT)
$I_2 - I_1 (r_0)$	$I_2 - I_1 (r_1)$	$I_3 - I_1 (r_0)$
$I_3 - I_2 (r_1)$	$I_3 - I_2 (r_0)$	$I_5 - I_2 (r_1)$
$I_4 - I_3 (r_0)$	$I_4 - I_3 (r_2)$	
$I_5 - I_4 (r_2)$	$I_4 - I_2 (r_2)$	
$I_6 - I_5 (r_1)$	$I_4 - I_1 (r_2)$	
	$I_5 - I_4 (r_1)$	
	$I_5 - I_3 (r_1)$	
	$I_5 - I_1 (r_1)$	

Non-linear Pipeline :-

- This pipeline contains forward and backward connection, so reservation table is used to process the inputs.
- Let us consider the following reservation table to process the inputs -



path : S₁ — S₂ — S₃ — S₂ — S₁

- Multiple checkmarks in the row, shows repeated use of same stage in different cycles.
- contiguous checkmarks in the row indicates extended use of same stage.
- Multiple checkmarks in the column indicates parallel use of ^{diff} stages in the same cycle.

Latency analysis :-

- Latency means clk cycle difference b/w to successive initiation in the pipeline.
- It depends on the reservation table.
- In the non-linear pipeline, some of the latencies will causes collisions called as forbidden latency (non-permissible latency). And some of the latency doesn't causes collision called as non-forbidden latency (permissible latency)
- To detect the forbidden latency we need to take the clock cycle difference b/w any two checkmarks in the row. i.e

$$\text{row 1 (S}_1\text{)} : 5 - 1 = 4 \quad (C_4 = 1)$$

$$\text{row 2 (S}_2\text{)} : 4 - 2 = 2 \quad (C_2 = 1)$$

$$\text{row 3 (S}_3\text{)} : \underline{\quad}$$

- Collision vector is derived based on the forbidden latencies.

i.e $[c_n \dots c_2 c_1]$

c $\begin{cases} 0 & ; \text{permissible} \\ 1 & ; \text{non-permissible} \end{cases}$

length of collision vector = no. of cycles needed to execute one instn.

eg:- $c_5 \ c_4 \ c_3 \ c_2 \ c_1$
 $[0 \ 1 \ 0 \ 1 \ 0]$

- Latency sequence is derived based on collision vector w.r.t permissible latencies. i.e

(i) Latency seq w.r.t 'c₁' :-

	1	2	3	4	5	6	7	8	9	10	11	12	13
s ₁	1	2			1	2	3	4			3	4	5
s ₂		1	2	1	2			3	4	3	4		
s ₃			1	2					3	4			

Latency cycle

(1,5), (1,5), (1,5) ...

$$(2-1) = 1$$

$$5(7-2) = 5$$

$$(\text{Avg - latency}) = \frac{1+5}{2} = 3$$

$$(8-7) = 1$$

$$(13-8) = 5$$

b.

(ii) Latency seq w.r.t 'c₂' :-

	1	2	3	4	5	6	7	8	9	10	11	12	13
s ₁	1			2	1		3	2		4	3		
s ₂		1		1	2		2	3		3	4		4
s ₃			1			2			3			4	

Latency cycle

3, 3, 3 ...

$$(4-1) = 3$$

$$(7-4) = 3$$

$$(10-7) = 3$$

(Avg-latency = 3)

(iii)

Latency seq wrt 'Cs' :-

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
s_1	1				1	2	3			2	3	4	5			4
s_2		1	1			2	3	2	3			4	5	4	5	
s_3			1					2	3				4	5		

Latency -cycle

$$(6-1) = 5$$

(5,1), (5,1), (5,1), ...

$$(7-6) = 1$$

$$(12-7) = 5$$

$$(13-12) = 1$$

(Avg-latency = $\frac{5+1}{2} = 3$)

Memory Organization :-

$$\text{Performance} \propto \frac{1}{\text{Access time}}$$

$$\text{Access-time} \propto \text{Execution time}$$

$$\text{Total Time} = \frac{\text{Hit/miss}}{\text{Latency}} + \frac{\text{Access time}}{\text{Depends on Bandwidth}} + \frac{\text{Transfer time}}{\text{Depends on Bandwidth}}$$

↓
(Negligible)

Types of memory organization :-

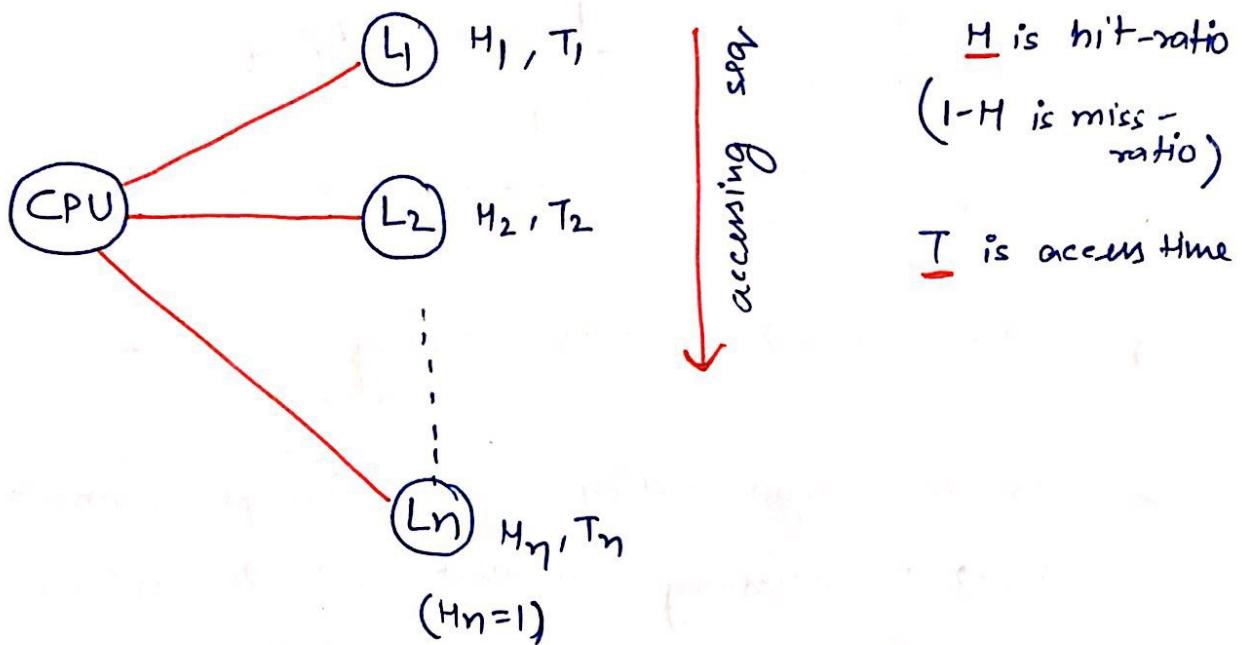
(i) Simultaneous Access mem orgⁿ.

(ii) Hierarchical Access mem orgⁿ.

(i) Simultaneous Access memory orgⁿ:

- In this orgⁿ, CPU directly connected to all the levels of a memory, so that CPU will be accessing the data from any level of the memory in a sequence i.e
 - when there is miss in level 1 mem, then CPU directly access the data from level 2 mem without copying it into level 1 mem

- when there is a miss in level 2 ^{mem} then CPU access data directly from level 3 mem without copying it into level 2 and level 1 memory and so on.



- Hit-Ratio (H) = $\frac{\text{Hits}}{\text{Total Access}}$

Hit-Ratio (H) + miss-ratio = 100%.

- $T_1 < T_2 < T_3 < \dots < T_n$

Time required to access 1 word Data from the mem is called as avg. memory access Time

$$T_{\text{avg}} = H_1 T_1 + (1-H_1) H_2 T_2 + (1-H_1)(1-H_2) H_3 T_3 + \dots - (1-H_1)(1-H_2) \dots (1-H_{n-1}) H_n T_n$$

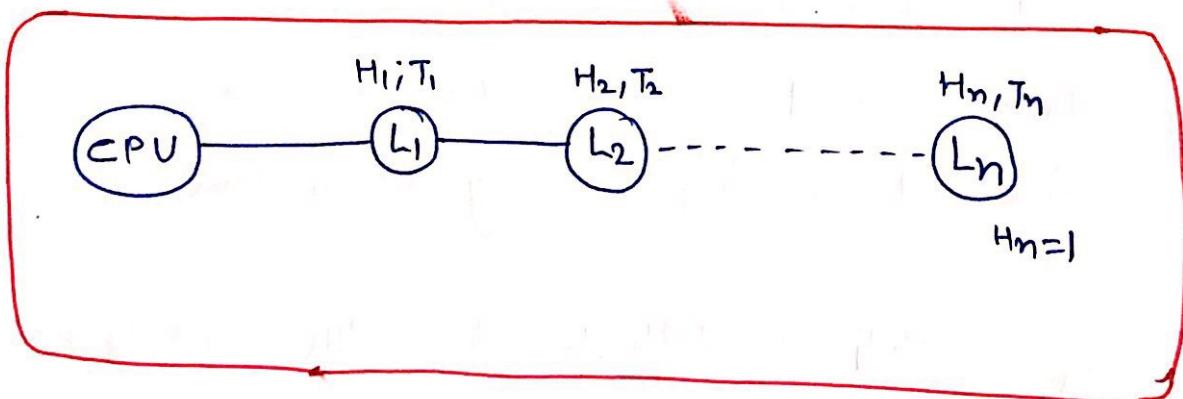
1 word $\rightarrow T_{avg}$

so in 1 sec = $\frac{1}{T_{avg}}$ words/sec

$$\eta_{mem} = \frac{1}{T_{avg}} \text{ words/sec}$$

(ii) Hierarchical Access memory orgn:-

- In this organization, CPU always connected to level 1 memory so that CPU always accessing the data only from the level 1 memory.
- When there is a 'miss' in level 1 memory then the respective data will be copied from higher levels to level 1 memory. Later CPU will access the data only from the level 1 memory.



$$T_{avg} = \underbrace{H_1 T_1}_{+} + \underbrace{(1-H_1) H_2 (T_2 + T_1)}_{+ \dots +} + \underbrace{(1-H_1)(1-H_2) H_3 (T_3 + T_2 + T_1)}_{+ \dots +} \\ + \dots + \underbrace{(1-H_1)(1-H_2) \dots (1-H_{n-1}) H_n (T_n + T_{n-1} + \dots + T_1)}$$

Analysis:-

I₁: Data (Assume that data present in L2 Mem)

I₂:

I₃: Data

I₄:

I₅: Data

I₆:

I₇: Data

Re-use the I₁ Data

Assume L₁ → T₁ (2ns)

L₂ → T₂ (10ns)

Simultaneous

I₁: M_{L₁}, T₂

miss in
L₁

I₃: M_{L₁}, T₂

I₅: M_{L₁}, T₂

I₇: M_{L₁}, T₂

Hierarchical

I₁: M_{L₁}, T₂ + T₁

I₃: T₁

I₅: T₁

I₇: T₁

Temporal
locality

access same
data word

$$\text{Time} = 4T_2$$

$$= (40\text{ns})$$

$$\text{Time} = T_2 + 4T_1$$

$$= (18\text{ns})$$

I_1 : Data (w_1)

I_2 : " (w_2)

I_3 : " (w_3)

I_4 : " (w_4)

Assume all the words
are present in L_2 , Mem

Simultaneous

I_1 : M_{L_1}, T_2

I_2 : M_{L_1}, T_2

I_3 : M_{L_1}, T_2

I_4 : M_{L_1}, T_2

Hierarchical

I_1 : $M_{L_1}, T_2 + T_1$

I_2 : $M_{L_1}, T_2 + T_1$

I_3 : $M_{L_1}, T_2 + T_1$

I_4 : $M_{L_1}, T_2 + T_1$

$$\text{Time} = 4T_2$$

$$= (40 \text{ ns})$$

$$\text{Time} = 4(T_2 + T_1)$$

$$= (48 \text{ ns})$$

So, ~~adjust~~ adjustment is done, data is transfer
in the form of blocks.

Assume block size = $16W$ so,

in one data - transfer 16 words will
be move to L_1

$$I_1 : M_{L_1}, T_2 + T_1$$

$$I_2 : T_1$$

$$I_3 : T_1$$

$$I_4 : T_1$$



Spatial
Locality



adjacent
data words

Accessing

$$\text{Time} = T_2 + 4T_1$$

$$= (18 \text{ ns})$$

NOTE :- In the hierarchical access memory org", access time is low because of the principal 'locality of reference'.

- "Locality of reference" means accessing the higher levels of a memory data from the lower-level memory (L_1). It is of two types:
 - (i) Temporal locality
 - (ii) Spatial locality
- To satisfy the above locality of reference, we need to organise the data in the memory in form of blocks. Block contain multiple words and its size is decided by designer.

NOTE:- when the question contain hierarchy word or cache memory then use hierarchical orgⁿ otherwise use simultaneous orgⁿ.

- Q.** In a two-level memory orgⁿ, Level 1 memory is 8 times faster than level 2 memory. And its access time is 30ns less than avg memory access time. Let L₁ memory access time is 40ns. What is the Hit-ratio?

$$(H_2 = 1)$$

$$T_{avg} = H_1 T_1 + (1-H_1) H_2 T_2$$

$$S = \frac{T_2}{T_1} = 8$$

$$T_{avg} = H_1 T_1 + 8(1-H_1) T_2$$

$$(T_2 = 8T_1)$$

$$(T_1 = 40\text{ns})$$

$$70\text{ns} = H_1(40) + 8(1-H_1)(40)$$

$$\text{so, } T_2 = 320\text{ns}$$

and

$$70 = 40H_1 + 320 - 8 \cdot 320H_1$$

$$T_{avg} = 40 + 320H_1$$

$$70 = -280H_1 + 320$$

$$280H_1 = 250$$

$$\left(H_1 = \frac{25}{28} \right)$$

Q. Three-level memory org' has the following specification :

Level	Access time per word	Block-size in word	Hit - ratio
1	40 ns	—	0.6
2	100 ns	2	0.9
3	200 ns	4	1

If the referred data is not in L1 memory then transfer it from L2 to L1 and if not in L2 mem then transfer it from $L3 \rightarrow L2 \rightarrow L1$. What is T_{avg} ?

$$T_1 = 40 \times 1 = 40 \text{ ns}$$

$$T_2 = 100 \times 2 = 200 \text{ ns}$$

$$T_3 = 200 \times 4 = 800 \text{ ns}$$

so,

$$\begin{aligned}
 T_{avg} &= H_1 T_1 + (1-H_1) H_2 (T_2 + T_1) + (1-H_1)(1-H_2)(T_2 + T_1 \\
 &\quad + T_3) \\
 &= 0.6(40) + (0.4)(0.9)(240) + (0.4)(0.1)(1040) \\
 &= 24 + 86.4 + 41.6 \\
 &= \underline{\underline{152 \text{ ns}}}
 \end{aligned}$$