

6/11/20

Measurement Techniques for Performance

- Latency \rightarrow (state transition \rightarrow done) time
 - Throughput \rightarrow (# operations / second)
- $\text{Throughput} \propto \frac{1}{\text{Latency}}$

Ex Car Assembly

chassis \rightarrow engine \rightarrow wheel \rightarrow producing the car

20 stages Latency \rightarrow 4 hrs

Throughput

- * 20 stages \rightarrow 4 hours \rightarrow 1 stage \rightarrow $\frac{1}{4}$ hours
- * If last stage goes 5 times per hour, so $\text{Throughput} = 5 \text{ cars/hour}$
- so $\text{Throughput} \neq \text{Latency}$.

* Comparing Performance:

"X is n times faster than Y"

$$\text{Speedup} \rightarrow n = \frac{\text{Speed}(X)}{\text{Speed}(Y)} \Rightarrow \frac{\text{latency}(Y)}{\text{latency}(X)} \text{ OR } \frac{\text{Throughput}(X)}{\text{Throughput}(Y)}$$

- Q Laptop A takes 4 hrs to compress a video.
- New laptop B takes 10 mins to compress a video. \rightarrow Execution time.
- Speedup $\rightarrow 24$

- \rightarrow Performance of a computer system / architecture is measured w.r.t benchmark programs.

— Programs & inputs which are standardized by stakeholder, academics, industry.

Benchmark Suites

- # of instructions / cycle \rightarrow Peak performance.

* Engineering workstations \rightarrow SPEC benchmarks

Not very I/O intensive \rightarrow more processor dependent

Consists of \rightarrow gcc \rightarrow software development workload compiler

BWAVES, LBM \rightarrow fluid dynamic applications

PERL \rightarrow string processing applications

CACTUS ADM \rightarrow Physics - simulators

BZIP \rightarrow Compression

XALAN BMK \rightarrow Parser for XML

	Comp X	Comp Y	Speedup
App A	9s	18s	2
App B	10s	7s	0.7
App C	5s	11s	2.2
Avg time	8s	12s	$\rightarrow 1.5$
Cube roots	7.66	11.15	$\rightarrow 1.45$

$$\sqrt[3]{9 \times 10 \times 5}$$

Geometric mean for the ratio average

Q machine A \rightarrow machine B

App 1 \rightarrow speedup of 2
App 2 \rightarrow speedup of 8
Overall avg speed $\rightarrow \sqrt{2 \times 8} \rightarrow 4$

Iron Law of Performance

CPU Time \rightarrow # of inst/program \times CPI \times clock cycle time

$$\left| \frac{\# \text{inst}}{\text{prog}} \times \frac{\text{clock cycle}}{\text{inst}} \times \frac{\text{time}}{\text{cycle}} \rightarrow \frac{\text{time}}{\text{prog}} \right|$$

Q Program executes 3 Billion (3×10^9) inst, CPI = 2, clock is 3GHz

$$\text{CPU Time} \rightarrow \frac{3 \times 10^9 \times 2}{3 \times 10^9} \rightarrow 2s$$

8/11/12 RISC Management - SD \rightarrow Cont { TODO }

9/11/12

Q SOR instructions:

CPI	
10B \rightarrow 4	
15B \rightarrow 2	
5B \rightarrow 3	
20B \rightarrow 1	

clock \rightarrow 4GHz
exec time.

$$\text{CPU Time} \rightarrow \left(\sum \left(\frac{\# \text{inst}}{\text{prog}} \times \text{CPI}_i \right) \right) \times \frac{1}{\text{clock}}$$

AMDAHL's Law:

$$\text{Speedup} = \frac{1}{(1-f) + \frac{f}{S}}$$

f \rightarrow fraction enhanced \Rightarrow part of code to be enhanced
 \rightarrow taken w.r.t original exec time

S \rightarrow speedup enhanced

Q) Program with 50B instructions, clock freq in 2GHz.

Improve Branch instr. CPI_{branch} goes from 6 to 3

Inst type	#inst of overall instructions	CPI
Int	40%	1
branch	20%	4
Load	30%	2
Store	10%	3

* cannot apply Amdahl's law w.r.t given %'s as they are in terms of # of instructions, not in terms of original exec time.

# inst	Execution
20	10
10	20
15	15
5	7.5

$$\Rightarrow \frac{20}{52.5} = f \quad \frac{1}{\frac{32.5}{52.5} + \frac{10}{52.5}} \Rightarrow \frac{52.5}{42.5} \Rightarrow \frac{105}{38} \boxed{\frac{21}{17}} \Rightarrow \underline{1.24}$$

Ex Enh1 \rightarrow Speedup of 20 on 10% of time $\Rightarrow \frac{1}{0.9 + \frac{0.1}{2}} \Rightarrow \frac{20}{18.1} \Rightarrow 1.105$
 Enh2 \rightarrow Speedup of 1.6 on 80% of time $\Rightarrow \frac{1.6}{0.32 + 0.8} \Rightarrow \frac{1.6}{1.12} \Rightarrow \underline{1.43}$

* Even if we had ∞ speedup on 10% of time then speedup $\rightarrow \underline{1.11}$

Inst type	% of time	CPI
Int	40	1
branch	20	4
Load	30	2
Store	10	3

Possible improvements

1. Branch CPI 4 \rightarrow 3
2. Clock freq \uparrow from 2 to 2.3 GHz
3. Store CPI 3 \rightarrow 2.

1. Branch 20% \rightarrow Speedup $\approx \frac{4}{3} \Rightarrow \frac{1}{0.8 + \frac{0.2 \times 3}{4}} \Rightarrow \frac{1}{0.8 + 0.15} \Rightarrow \frac{1}{0.95} \Rightarrow \underline{1.05}$

2. here speedup $\Rightarrow \frac{2.3}{2} \Rightarrow \underline{1.15}$

3. Store 10% \rightarrow Speedup $\approx \frac{3}{2} \Rightarrow \frac{1}{0.9 + \frac{0.1 \times 2}{3}} \Rightarrow \frac{1}{0.7 + 0.2} \Rightarrow \frac{3}{2.9} \Rightarrow \underline{1.035}$

Improvement of 90% of time. Speedup ≈ 2

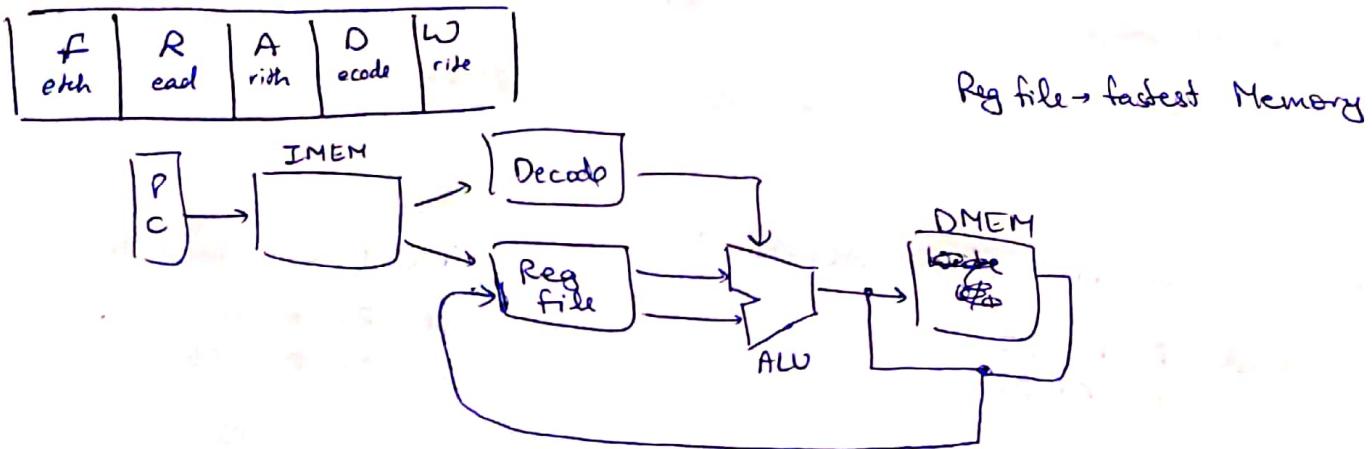
Slowing down the rest by 10x.

$$\text{Speedup} \Rightarrow \frac{1}{\frac{0.1}{0.1} + \frac{0.9}{2}} \Rightarrow \frac{1}{1+0.45} \Rightarrow \frac{1}{1.45} \Rightarrow 0.7$$

* If speedup of 90% of time is ∞ then too speedup $\rightarrow 1$

HADMA's Law → Do not mess the uncommon case too badly.

Classic 5 stage pipeline



13/1/20

f	D	A	M	W
ADD R2, R1, I	LW R1, w			
X	LW R1, w			
X	LW R1, w			

If command ADD R2, R1, I was moved to next step A, then it would be working on wrong data as LW R1, w haven't written value.

Assumption:

→ We can write & read in the same clock cycle



so we get stalls (X). Earlier it would take n cycles for n instructions (bubbles)

but now it takes $n+s_n$ cycles where $s_n \rightarrow \#$ of stalls.

F	D	A	M	W
Jump				
Stall	Jump			
Stall	Stall	Jump		

↳ flush the pipeline
Control dependency.

Q 5 stage processor, branch gets resolved in the 3rd stage.

* 20% of instructions are branch instructions.

* slightly more than 50% of these branches are taken

$$CPI = CPI_{\text{base}} + \underbrace{0.2 \times \frac{1}{2} \times 2}_{\substack{\text{depth} \\ \text{of pipeline}}} \Rightarrow 1.2$$

\Rightarrow # of flushed instructions

Q 25% of all instr are branch / Jump.

10-stage pipeline

Correct target for branch / jump gets computed in the 6th step

$$CPI_{\text{avg}} \rightarrow 1 + \frac{1}{4} \times 5 \Rightarrow 2.25 \quad \# \text{ all are taken}$$

RAW (True / flow dependencies)

Add R1, R2, R3

Sub R7, R1, R8

I1 MUL R1, R2, R3

I2 ADD R4, R4, R1

I3 MUL R1, RS, RL

I4 Sub R4, R4, R1

WAR

WAW (name / false dependencies)

{ Add R1, R2, R3

{ MUL R1, RS, RL

Sub R7, R1, R8

MUL R1, RS, RL

I1 → I2 RAW

I2 → I3 WAW

I1 → I4 RXW

I2 → I3 \Rightarrow WAR

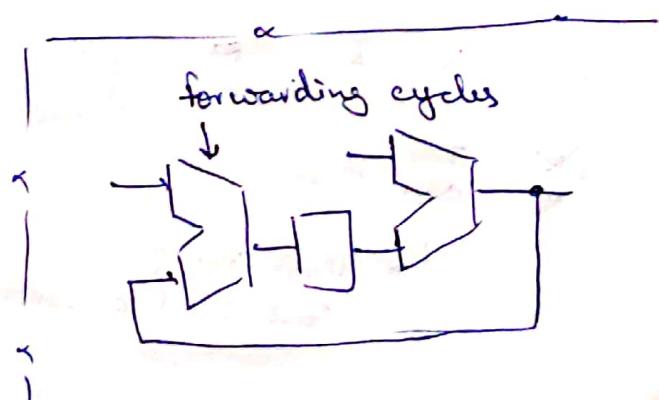
(Broken by I₃)

15/1/20

Dependencies And Hazard

Properties of the
program and oblivious
of the processor details

↳ properties of both the program
and also depends on the pipeline



I ₁	ADD	R1, R2, R3
I ₂	MUL	R7, R4, RS
I ₃	SUB	R4, R6, R7
I ₄	DIV	R10, R4, R8
I ₅	XOR	R11, R1, R7

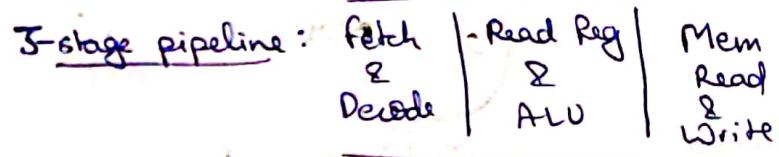
Hazards?

- true dependency (RAW)
- false dependency

- * for 5 stage pipeline → If → ID → E → MEM → WB
- * Among the two dependencies → I₁ → I₅ & I₂ → I₅ only one of them is a hazard (I₂ → I₅) as for the other one we have 3 instructions between them (This is okay due to the 5-stage pipeline)

I ₁	ADD	R1, R2, R3
I ₂	SUB	RS, R1, R4
I ₃	DIV	R6, R1, R7
I ₄	MUL	R7, R1, R8

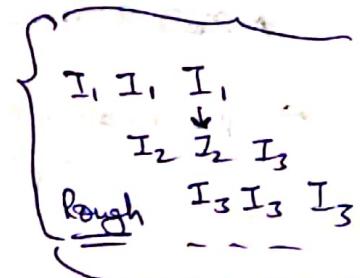
3-stage pipeline:



→ 3x dependencies : (I₁ → I₂) (I₁ → I₃) (I₁ → I₄)

<u>Dependence</u>	<u>Hazard</u>
I ₁ → I ₂	Yes (RAW) → Not a hazard if same cycle read & write is allowed.
I ₁ → I ₃	No
I ₁ → I ₄	No

if same cycle
read & write
is allowed.



Handling Hazards:

- * Control dependency → Control hazards (Branches) : flush the pipeline
- * Data dependency → stall (Introducing bubbles in the pipeline)

↓

forwarding

Example (Control dependency)

BEQ	R2, R2, label
ADD	--
SUB	--

⇒

SUB	ADD	BEQ
-----	-----	-----

in pipeline

↓ If Branch is taken

SUB	ADD	BEQ
-----	-----	-----

↓ flush the pipeline

label : { --
 --
 -- }

→ Data dependencies:

ADD R2, R2, R3

SUB R4, R2, RS

Read	ALU
ADD	
SUB	

↳ SUB instr reads old value.

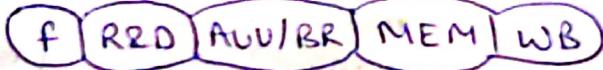
* Stall cycles can be avoided with suitable forwarding techniques
(valid for certain cases)

Example:

		flush	stall	forward
Control dependency	1) ADD R4, RS, R6 BNF R1, R4, label SUB RS, R4, R3	✓	✗	✗
	2) MUL R1, R2, R3 LW R2, 0(R1)	✗	✗	✓
	3) ADD R2, R2, R1	✗	✓	✓

* Stall & forward have no place in control dependency.

→ The organization of the pipeline follows:



Ques: Main dependencies are indicated by:

2) 8	F	R2D	ALU/BR	MEM	WB
3)			MUL ↓ LW		

(forwarding takes place in this stage)

* Forwarding is always a feedback loop in the processor data path.

Keywords:

- out of order execution
- instruction level parallelism
(defined w.r.t to ideal computer)

ILP: Instruction level parallelism

- Briefly → Tells exactly how many instructions are executed in one cycle

Cycle 1, Cycle 2	
✓	✓
✓	✓
✓	✓
✓	✓

↳ because of the dependency
these two instructions
cannot execute in the
same clock cycle

Ideal world
 $CPI = \frac{S}{\infty} \rightarrow 0$
of instructions



In real life → there are dependencies which leads to
an increase in CPI

for a classic pipeline \Rightarrow $\lim_{n \rightarrow \infty} \frac{d+n-1}{n} \rightarrow 1$

Q S-stage pipeline. Allows forwarding. If no dependencies, 10 instructions in each stage.

		Ex	WB
MUL	R2, R2, R2	2	4
ADD	R2, R2, R2	3	5
MUL	R3, R3, R3	2	4
ADD	R2, R2, R3	4	6
MUL	R4, R4, R4	2	4
ADD	R2, R2, R4	5	7

* Physical register \Rightarrow P0, P1 ... (actual location of the register value)

Q

ADD	R1, R2, R3
SUB	R4, R1, R5
XOR	R6, R7, R8
MUL	R5, R8, R9
ADD	R4, R8, R9

$$\left\{ \begin{array}{l} R_0 \rightarrow P_0 \\ R_1 \rightarrow P_1 \rightarrow P_{17} \\ R_2 \rightarrow P_2 \\ R_3 \rightarrow P_3 \\ R_4 \rightarrow P_4 \rightarrow P_{18} \rightarrow P_{21} \\ R_5 \rightarrow P_5 \rightarrow P_{20} \\ R_6 \rightarrow P_6 \rightarrow P_{19} \\ R_7 \rightarrow P_7 \\ R_8 \rightarrow P_8 \end{array} \right.$$

* Inside RAT table, we write the physical address of the corresponding architectural register.

* On write, rename the register

\Rightarrow Even after renaming, the RAW dependency remains (true dependency)
WAR & WAW dependencies are removed
(also called false dependencies or name dependencies)

16/1/20

MUL R2, R2, R2
 ADD R1, R1, R2
 MUL R2, R4, R4
 ADD R3, R3, R2
 MUL R2, R6, R6
 ADD RS, RS, R2

→ true dependencies (RAW)
 → false dependencies (WAW, WAR)

RAT	
R1	P1
R2	P2
R3	P3
R4	P4
RS	P5
R6	P6

MUL P7, P2, P2
 ADD P8, P1, P7
 MUL P9, P4, P4
 ADD P10, P3, P9
 MUL P11, P6, P6
 ADD P12, P5, P10

false dependencies are gone.

Register renaming can only remove false dependencies.
 These can be done parallelly in 2 cycles. So CPI = $\frac{2}{6} \Rightarrow \frac{1}{3}$

Instruction Level Parallelism (ILP)

- On Ideal Processor : Processor can execute all instructions in one cycle.
 Processor can execute any number of instructions in 1 cycle.
 Processor has to obey ~~true~~ true dependencies.

- * ILP holds for only ideal processor. IPC → Instructions per cycle.
 (Property of the program & not the processor)

ADD P10, P2, P3
 XOR P6, P7, P8
 MUL P5, P8, P9
 ADD P4, P8, P9
 SUB P11, P10, P5

- * Ignore false dependencies
- * 1 true dependency
- * ILP $\Rightarrow \frac{5}{2} \Rightarrow 2.5$

ADD R1, R1, R1
 ADD R2, R2, R1
 ADD R3, R2, R1
 ADD R6, R7, R8
 ADD R8, R3, R7
 ADD R1, R1, R1
 ADD R1, R7, R7

ADD P9, P1, P1
 ADD P10, P2, P9
 ADD P11, P10, P9
 ADD P12, P7, P8
 ADD P13, P11, P7
 ADD P14, P9, P9
 ADD P15, P7, P7

$$\text{ILP} \Rightarrow \frac{7}{4}$$

ILP for control & structural dependencies

- No structural dependencies (no limit to hardware)
- Assume there is an ideal branch predictor (correct prediction everytime)

Add R1, R2, R3
 Mul R1, R1, R9
 BNNE RS, R1, label
 ADD RS, R1, R2

label: MUL RS, R7, R8

- We are assuming perfect BPU - which means even before branch, i.e. executed in the fetch cycle itself we knew the outcome of the branch
- Branch itself maybe delayed (due to true dependences) but is not causing any delay in the following instructions.

ADD R1, R2, R3
 SUB R4, R1, R5
 XOR R6, R7, R8
 MUL RS, R8, R9
 ADD R4, R8, R9

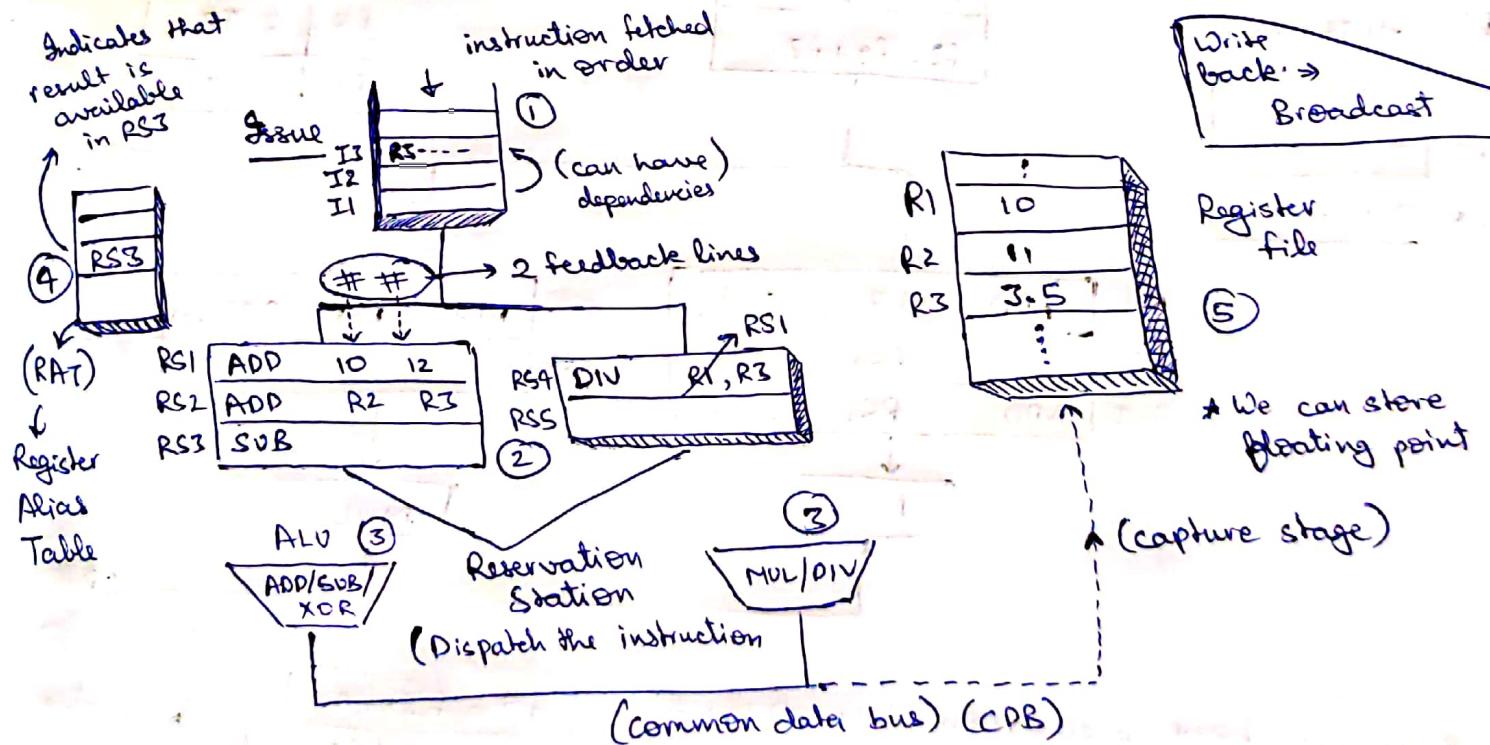
P10, P2, P3
 P4, P10, P5
 P12, P7, P8
 P13, P8, P9
 P14, P8, P9

ILP \rightarrow 2.5
 IPC \rightarrow 2.5
 (1 MUL, 2 ADD/SUB/XOR)

IPC = 1.25

Jan 20

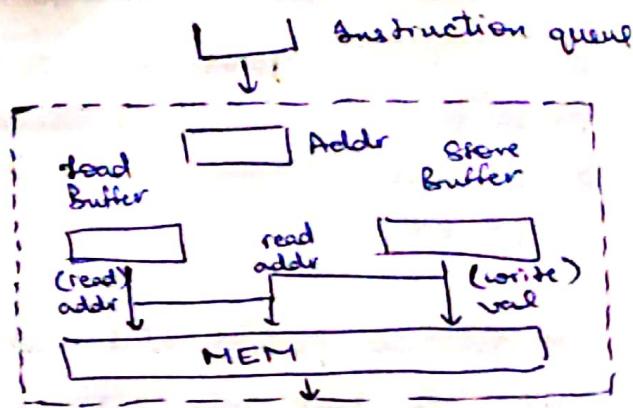
Tomasulo's Algorithm



- * for the instructions we have two operands, hence two feedback lines to update both.
- * for a value in ~~reg~~ reservation station we check if the register is ~~reg~~ under some operation (stored in RAT). If yes then take that alias otherwise take value from register file.

Load / store unit:

- * In the classical dominoes algorithm, the load/store is done in order
- * consider only one piece of memory.
- * Load/store instructions dependency is a bit tricky.

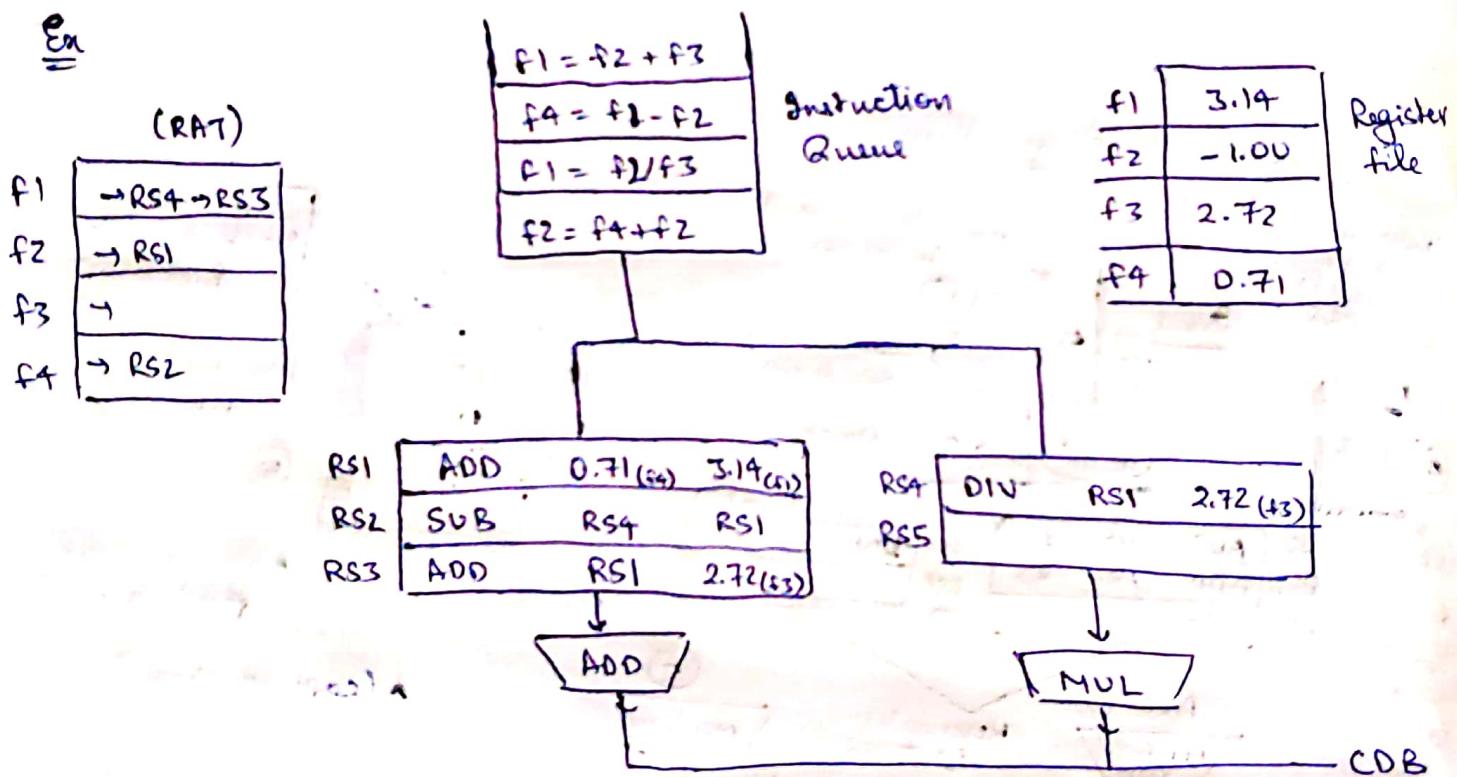


Issue:

- i) Take next instruction in program order from instruction queue.
- ii) Determine where the inputs will come from.
- iii) Get free reservation stations and of correct kind. [APP for ADD/SUB]
- iv) Put instruction in reservation station.
- v) Tag destination to register \Rightarrow RS4 | DN R3, R1, R3
 $R3 \rightarrow R1/R3$



Ex

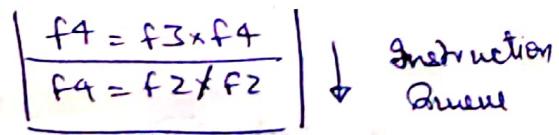


- * If we have two more instructions \Rightarrow

→ MUL/DIV Reservation station has

space for just two instructions & already had one instruction.

→ As there is no space for last MUL instruction \Rightarrow it has to wait for any instruction to complete in MUL/DIV Reservation Station.

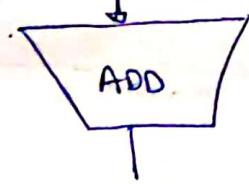


Jan 22

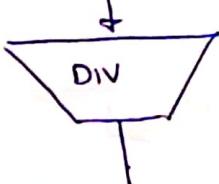
DISPATCH → Put the instruction into the execution unit from Register Station (RS)

CDB (RS1)-0.29

RS1	Add		
RS2	Sub	RS4	RS1
RS3	Div	RS1	2.72



RS4	DIV	RS1	2.72
RS5			



CDB (Common Data Bus)

(RS1)-0.29

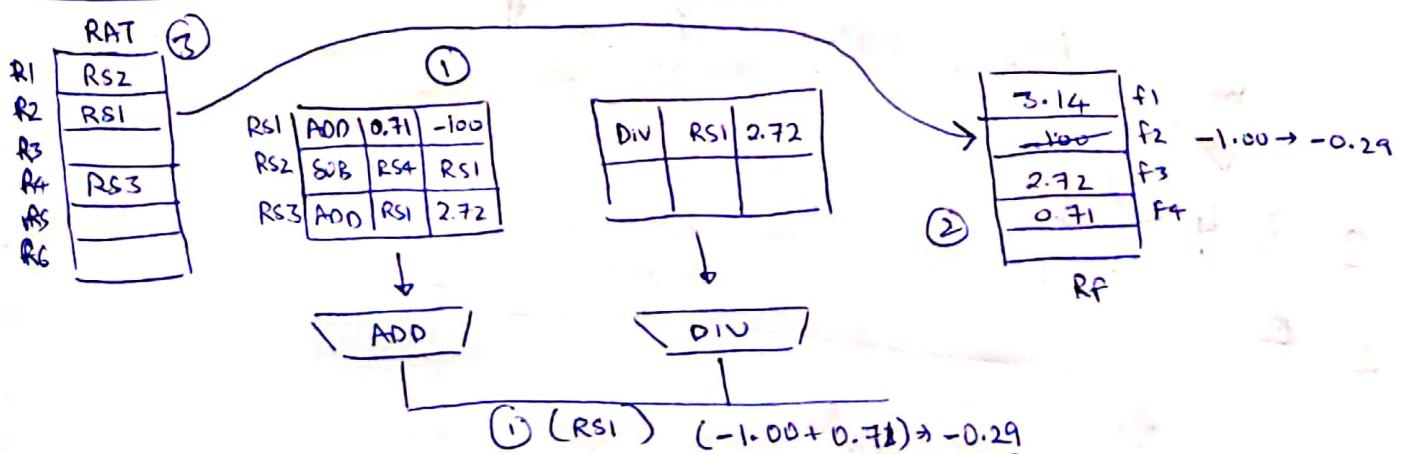
RS1			
RS2	SUB	RS1	RS1
RS3	ADD	123	2.73

DIV	RS1	2.72

Q Why has RS3 not dispatched before this cycle?

- ✓ It was issued in the previous cycle
- ✓ Another instruction was dispatched to add unit
- ✗ RS2 is older than RS3, so RS3 cannot dispatch until RS2 does.

Write Result - Broadcast onto CDB, RAT, RF



- ① Put TAG R result on CDB
- ② Write to RF
- ③ Update RAT (empty the element in RAT to point to the RF)
- ④ Free RS (invalidate a valid bit)

update the RS4 parameter in the RS. If the entries in RAT do not match the tag RS4, dont update the RAT & RF

Load → 2 cycle
 Add → 2 cycle
 Mult → 10 cycle
 Divide → 4 cycle

R2 → 100
 R3 → 200
 F4 → 25

f0	f2	f4	f6	f8	f10	RAT
ML1	✓	✗	✗	✗	✗	ML2

	Burst	Op	V _d	V _r	Q _j	Q _k	Disp	Temp tab (optional)
L01	✓ (1)	Load		134			✓ (2)	(4)
L02	✓ (2)	Load		245			✓ (4)	(6)
A01	✓ (4)	SUB	-25	7.1	✗	✗	✓ (7)	(9)
A02	✓ (6)	ADD	-9.6	2.5	✗	✗	✓ (10)	(12)
M1	✓ (3)	MUL	-2.5	2.5	✗		✓ (7)	(6)
M2	✓ (5)	DIV	-6.25	7.1	✗	✗	✓ (18)	(58)

values for 1st & 2nd operands for which the instr is waiting.

	Issue	Execute	Write
LD	f6, 34(R2)	1	4
LD	f2, 45(R3)	2	6
MUL	f0, f2, f4	3	17
SUB	f8, f2, f6	4	9
DIV	f10, f0, f6	5	58
ADD	f6, f8, f2	6	12
Capture → dispatch not in same cycle			
Issue: Dispatch not in same clock cycle.		Write back happened after execution cycle ends	Read does not dispatch in same cycle

R2 → 100 R3 → 200 F4 → 25	later f6 → 7.1 → A02 → -12.1 f2 → -2.5 f8 → -9.6 f0 → -6.25
---------------------------------	-------------------------------------------------------------------------

Rough

7



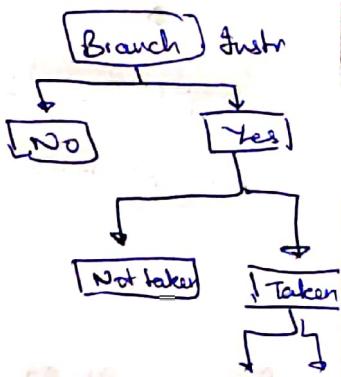
3 FebBranch Prediction Unit (BPU)

BEQ R1, R2, label (Branch if equal)

- * Refuse to predict vs predict \rightarrow choose predict
- * when we need to predict \rightarrow done at fetch stage.
- * Branch prediction works on the next value of PC

— 0110111000011 — —

- Is it a branch?
- Is it taken?
- If taken then what is target PC.

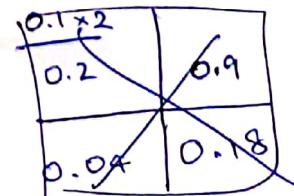
Effect of misprediction on CPI

$$CPI = CPI_{ideal} + \frac{\text{Mispredictions}}{\text{Instruction}} \times \frac{\text{Penalty}}{\text{Mispredictions}}$$

Accuracy of prediction Pipeline stage when branch is resolved
 ↓ ↓
 cycles on average per instruction that we add because of mispredictions.

	Resolve branch in 3 rd stage	Resolve branch in 10 th stage
50% accuracy for branch	$1 + 0.5 \times 0.2 \times 2$	$1 + 0.5 \times 0.2 \times 9$
100% accuracy for all other instr	<u>1.2</u>	<u>1.9</u>
90% accuracy for branch	$1 + 0.1 \times 0.2 \times 2$	$1 + 0.1 \times 0.2 \times 9$
100% accuracy for all other instr	<u>1.04</u>	<u>1.18</u>
Speedup	$\frac{1.2}{1.18} \Rightarrow 1.05$	$\frac{1.9}{1.18} \Rightarrow 1.6$

20% of all instructions are branch.



$$0.2 \times 0.1$$

$$0.2 \times 0.1 \times 9$$

Q3

5 stage pipeline

Branch is resolved at 3rd stage

fetch nothing until ~~sure~~ of what to fetch
: refuse to predict

Execute many iterations of:

LOOP	ADDI	R1, R1, -1	2 clock cycle
	ADD	R2, R2, R2	2 clock cycle
	BNEZ	R1, Loop	3 clock cycle (without prediction)

→ 7 clock cycle for one LOOP without prediction.

Simple Prediction Logic

→ not taken prediction

BNE R1, R2, label A

BNE R1, R3, label B

A:

B: using not-taken prediction

C: compute no of cycles which

are wasted on

misprediction until we

get to Y?

label A:

X:

V:

label B:

Z

Refuse to predict	not taken pred.
on branch instr penalty is 3	on branch instr penalty is 1 or 3
for other instr 2 cycle are needed	for other instr 1 cycle needed to fetch

- feb 5

Q 20% instr → branches 60% branches are taken

$CPI_{ideal} = 1$ so CPI now?

→ Not-taken Prediction → penalty → 2

$$CPI \rightarrow 1 + \underline{0.2 \times 0.6 \times 2} \Rightarrow \underline{\underline{1.24}}$$

Q

	Not taken predictor	Better prediction 99%	Speedup
5 stage (3rd stage)	$1 + 0.12 \times 2$ ⇒ 1.24	$1 + 0.01 \times 2$ ⇒ 1.02	1.22
14 stage (11th stage)	$1 + 0.12 \times 10$ ⇒ 2.2	$1 + 0.01 \times 10$ ⇒ 1.1	2
14 stage (11th stage) (4 instr/cycle)	$0.25 + 0.12 \times 10$ ⇒ 1.4	$0.25 + 0.01 \times 10$ ⇒ 0.35	4.14

~~% branch mispred~~
~~20%~~

Not taken prediction

88%
(accuracy)

{ 20% instr * 60% }
taken

+ 12% mispredict

+ 88% accuracy

Q Program has 20% of br inst

10% of branches are mispredicted

$$CPI_I = 0.5 \quad \text{Penalty} \rightarrow 30 \text{ cycles}$$

If 2% of branches are mispredicted

$$\Rightarrow CPI_I = CPI_{I_1} + 0.2 \times \cancel{0.07} \times 30 \Rightarrow \cancel{0.07} \quad 0.5 = CPI_I + 0.06$$

$$CPI_{I_1} = CPI_I + 0.06$$

$$= 0.5 - 0.06 + \cancel{0.12}$$

$$\Rightarrow \boxed{0.44} \quad \boxed{0.56}$$

Reasons for Prediction accuracy:

① Improve CPI

② Reduce wasted instr/cycle

$$PC_{next} = f(PC_{now})$$

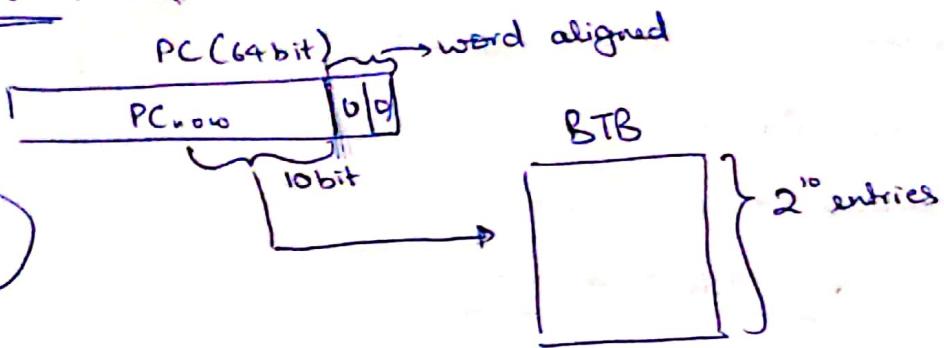
→ Assumption: Branches tend to behave the same over & over

→ History [PC_{now}] can help



→ we can make BTB very big but that won't help as then the access time will be large. Also we want no aliasing.

Solution?



BTB has 2¹⁰ entries & size 4 byte instructions (word aligned)

- instructions would begin at addresses divisible by 4
- 32 bit addresses
- which BTB entry is used for PC = 0x0000AB0C

$$\Rightarrow \text{PC} = 0x0000AB0C$$

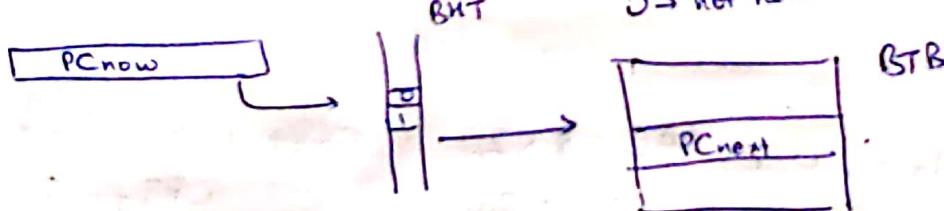
$\Rightarrow 0x_\underline{0000}\ \underline{0000}\ \underline{0000}\ \underline{0000}\ 1010\ |\underline{1011}\ \underline{0000}\ \underline{1100}$

$\Rightarrow 1011000011 \Rightarrow \underline{707}$

Direction Predictors

new table → BHT → Branch history table

0 → not taken



Q Q

0xC000	MOV	R2, 100
0xC004	MOV	R1, #
loop: 0xC008	BEQ	R1, R2, Done
0xC00C	ADD	R4, R3, R1
0xC010	LW	R4, 0(R4)
0xC014	ADD	RS, RS, R4
0xC018	ADD	R1, R1, 1
0xC01C	B	loop

Assume perfect predictor

BHT → 16 entries

BTB → 4 entries

	# BTB access	BTB row
0xC008	1	2
0xC01C	100	3
All rest	0	—

# BHT access	BHT row
1	1
1	2
101	3
100	4
100	5
100	6
100	7
100	8

* If we had non-perfect predictor then we would have BHT entries for the instructions after the loop too.

Q Let's say we have BHT with 1-bit prediction.
i.e. predict with the outcome in the previous access



Misprediction

- | (Initially D. Misprediction at end)
| ↘ (Initially mispredict)

→ This might have just 2 misprediction in total for a loop
→ But there can be some cases where we have short loops like NNNNTTNNN --- . There we will have many mispredictions.

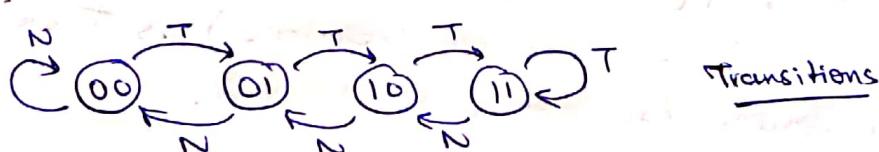
Solution?💡 use 2-bit predictor (Prediction with history)

00 strong not taken

01 weak not taken

10 weak taken

11 strong taken

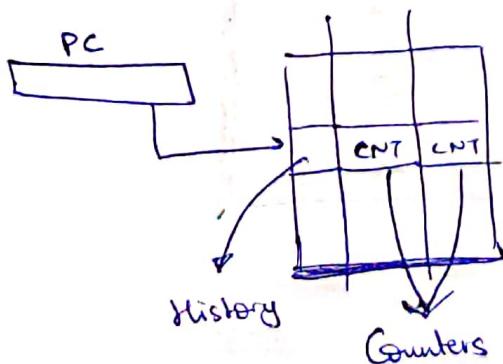


* Saturation Prediction ($2BC \rightarrow 2\text{bit counter}$) 2-bit predictors.

-feb6

1-bit predictors → 2-bit predictors (SPEC'89 benchmark → 88% to 92% accuracy.)

→ 1 bit history with 2BC per history

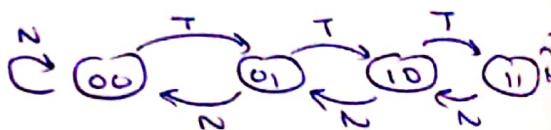


STATE	PREDICT	OUTCOME	CORRECT?
(0, SN, SN)	N	T	X
(1, WN, SN)	N	N	✓
(0, WN, SN)	N	T	X
(1, WT, SN)	N	N	✓
(0, WT, SN)	T	T	✓
(1, ST, SN)	N	N	✓
(0, ST, SN)	T	T	✓

use steps as:
 → update
 → shift
 → predict
 (USP)

SN → Strong N WN → Weak N
 ST → Strong T WT → Weak T

- STATE → (Index, Predict₀, Predict₁)
- UPDATE → on current index update based on state transition:



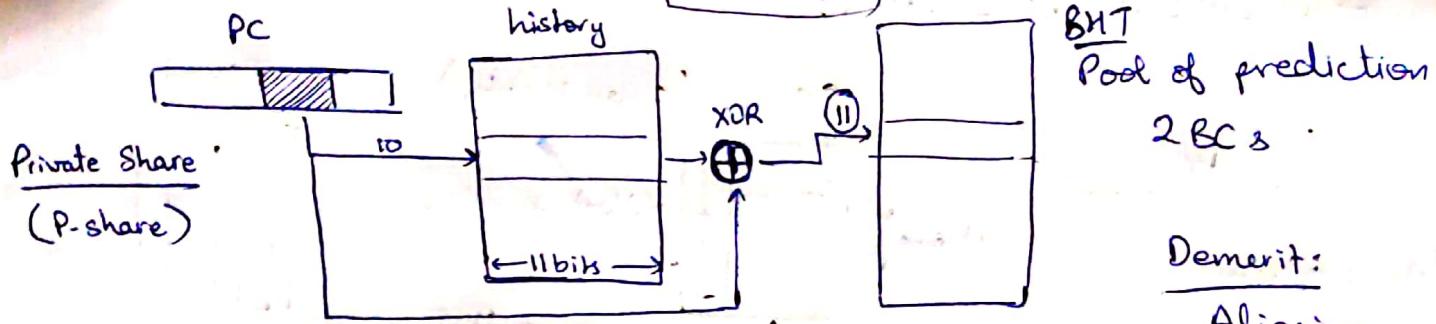
- Q for sequence $\rightarrow (NNT)^*$. We operate it for N repetitions of sequence what are the no. of mispredictions.
- Total no. of instructions $\rightarrow \underline{N \times 3}$
 Each cycle faces 1 misprediction on $\rightarrow T$ so:
 Total mispredictions $\rightarrow N$
- * In general, N-bit history predictors can learn all patterns of length $\leq N+1$. BUT we might not use some counters ever as for $(NNT)^*$ we never use (counter₃) as TT state never occurs. So the resources are not well utilized.

N-bit history | 2BC per history

Need 1024 entries for each branch to have an entry

N =	cost (bits)	(NNNT) [*] learned?	# of 2BCs used for (NNT) [*]
1	5×1024	X	Both (no wastage) 2
4	$(4 + 2^8 \times 2) \times 2^{10}$	✓	2048 2
8	$(8 + 2^8 \times 2) \times 2^{10}$	✓	2
16	$(16 + 2^{16} \times 2) \times 2^{10}$	✓	2^{17}

10 feb



$$\text{Size here} \Rightarrow (2^1 \times 11 + 2^{11} \times 2) \Rightarrow 31 \text{ Kbits}$$

→ for same size of history in previous method, it would take some more memory.

* update the PHT (Pattern history table) in a non-speculative ~~non~~ manner.

• $A[i] \in [0, 99]$

```
{
    if ( $A[i] < S_0$ )
        do task1()
    else if ( $A[i] > S_0$ )
        do task2()
}
```

(correlated branches)

second else if condition is benefitted by knowing the history of first if. So here global share helps (G-share)

→ In g-share we cannot use non-speculative way as it will introduce a lot of delay in which the commands might be executed.

→ So in g-share we update in speculative manner.

• for ($i = 1000; i != 0; --i$)
 if ($i \% 2$) $n = n + i$;

→ for good accuracy on all branches history should be:

→ PSHARE → 1 bit

→ GSHARE → 3 bits → order is (0 1 1 0 0 1 ---)

↳ instances to predict.

Loop

BEQ R1, Zero, Exit

AND R2, R1, 1

BEQ R2, Zero, Even

ADD R3, R3, R1

Even

ADD R1, R1, -1

Loop.

→ (B1)

1000
(NTNT---) ↗ T
without no history too it remains same.

→ (B2)

→ (NT)*

can be predicted with just 1 bit history.

→ (B3)

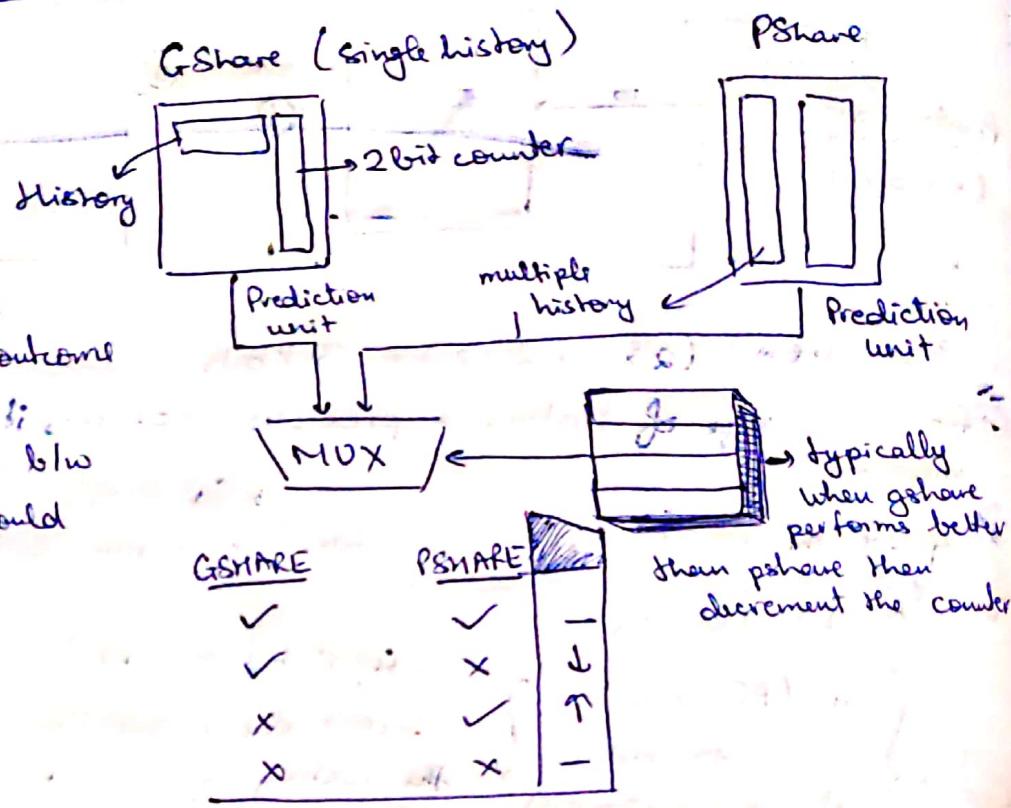
→ normal loop
no history needed.

Tournament Predictors:

Meta Predictor:

- array of 2 bit counters
- does not predict the outcome of the branch
- tries to tell which one b/w GSshare & PSshare should be selected.

- * Both predictors are updated.



as we go on decrementing it gets to 0.
Therefore automatically GSshare gets selected

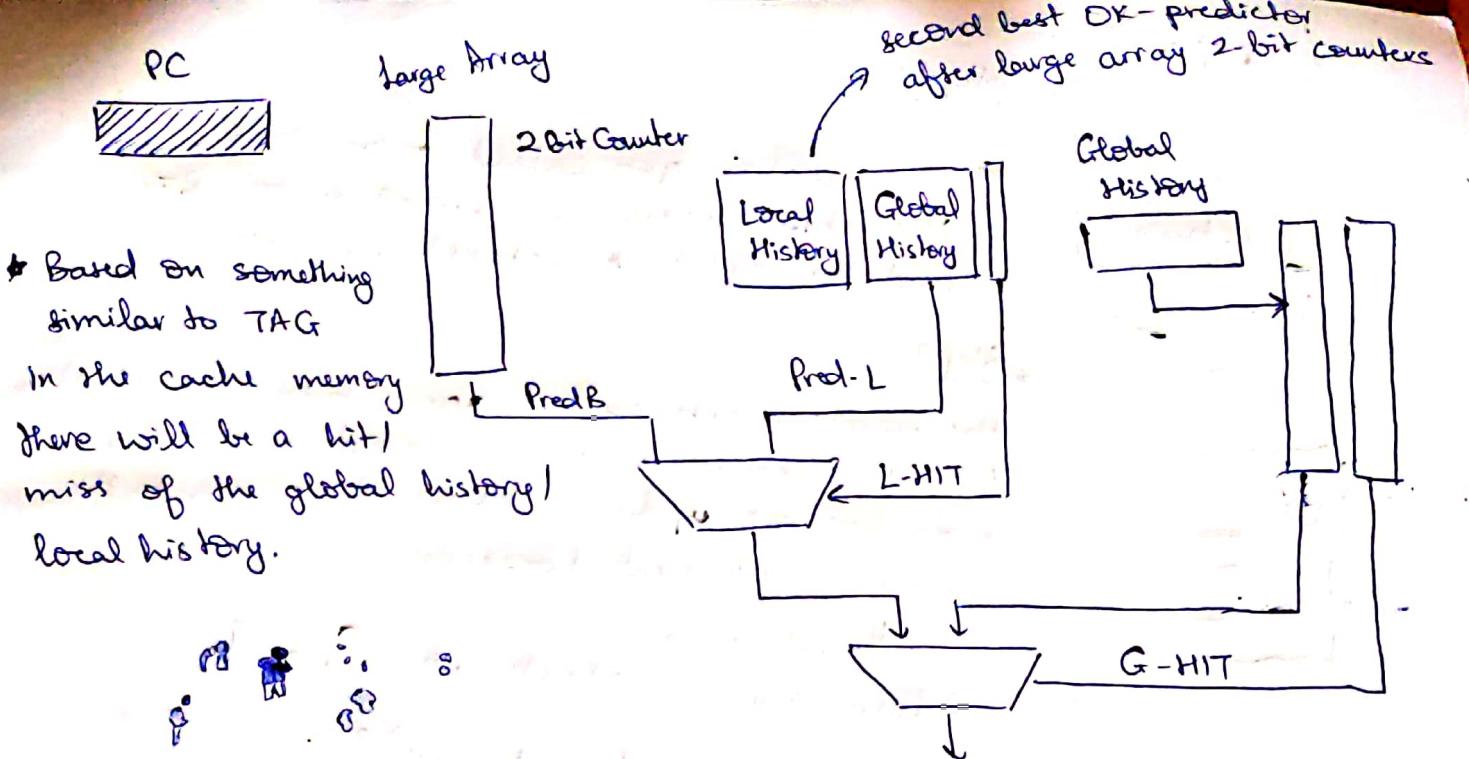
Hierarchical Predictors

- * Present in most modern day processors
- one good predictor (Good Predictor)
- one not so good/simple predictor (2Bit Counter) (OK predictor)
- * The update gets reflected on the GSshare or PSshare predictor depending upon the one which is being used.
- // Update OK predictor on each decision.
Update Good predictor only if the OK predictor fails.

Ex Consider an OK predictor which is a 2-bit counter

- GSshare (Global history)
- PSshare (Private history)

→ consider GSshare the best in this example. Hence it is used at the end if both the remaining predictors fail.

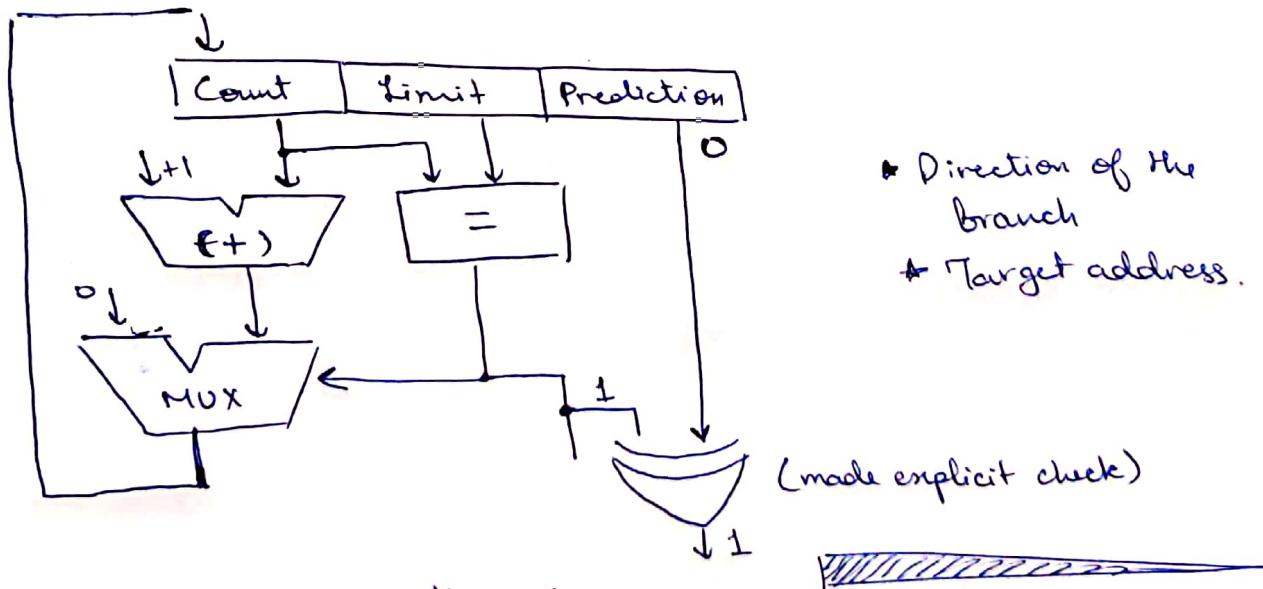


- Based on something similar to TAG
- In the cache memory there will be a hit/miss of the global history/local history.

Loop Detection

since there is no restriction on the value of N (no. of iterations) N in a loop can be very high \gg size of PHT.

- Although we have the predictors available, it is better to have a dedicated block to handle such loop structures.
- Say $(NNN \dots NT)^*$ is a long sequence.



- Consider in a program there is a function f1, which gets called from two locations L1 & L2 in the program

Hardware Dedicated unit called:

- RAS (return address stack (LIFO))
- Predressing.



Step 1:

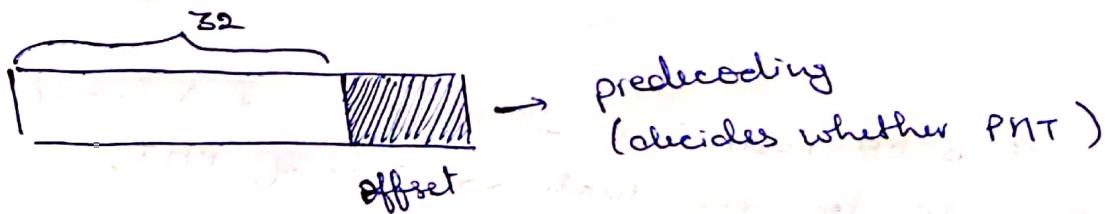
→ the function gets called from here, the return address get stored in the BTB and the function returns correctly.

Step 2:

the function gets called again from the location L2. The return address stored in the BTB does not get updated, hence it returns to the previous L1 next location.

Q. Solved by
Return Address Stack (RAS) [LIFO]

Q How do you know if it is RET (Return)?



- Instruction comes from memory to L2-cache to I-cache, maintain an extra bit to precede whether it is RET.