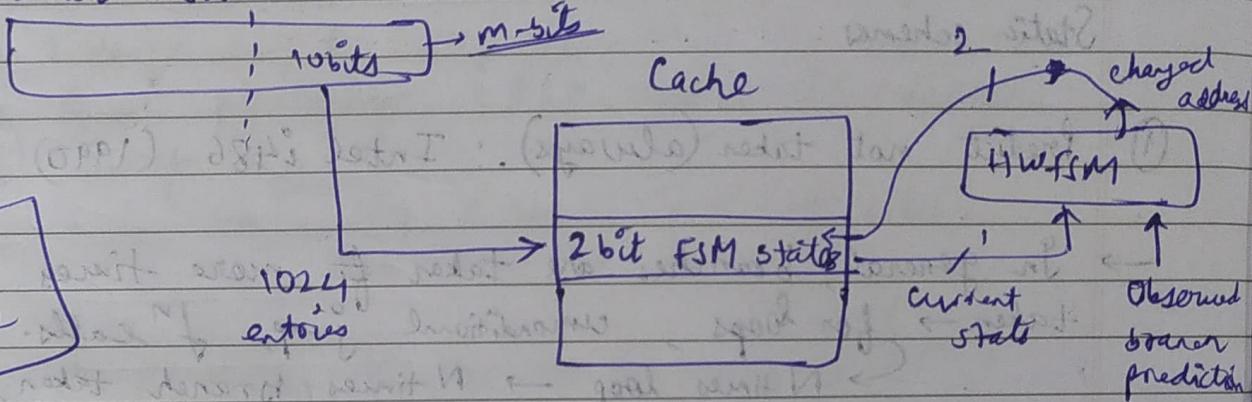


- * Dynamic Branch Predictions in follow up : Slides 4
- Observe branch outcomes, predict based on observed outcome
 - learn from mistakes.

Branch instruction address



No tag to save space

2^{M+1} bits \rightarrow size of cache

2-bit fsm state \Rightarrow predicted taken

11 \rightarrow strongly taken

10 \rightarrow weakly taken

Prediction = MSB of state

for 01 \rightarrow weakly not taken

for 00 \rightarrow strongly not taken

This cache has different names

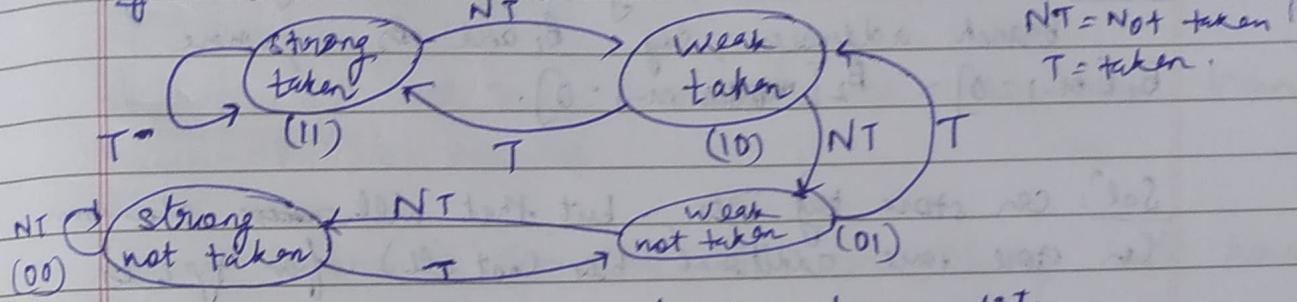
① Branch Prediction Buffer

② Pattern history table

① Local 1-level branch predictor

Branch predictor FSM.

Branch predictor FSM.
eg: 2-bit FSM has two states \leftarrow can also be 3-bit, 4-bit etc.



MSB → decides prediction of branch predictor

MSB=1 \Leftrightarrow taken predicted

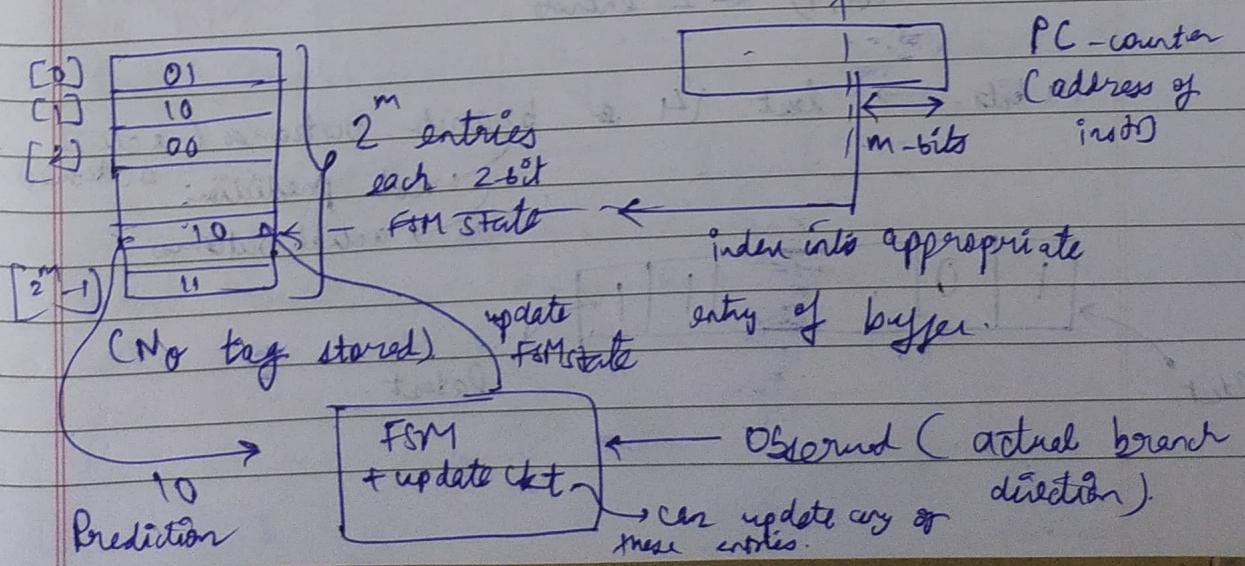
MSB = 0 \Leftrightarrow net taken predicted

Moore machine \rightarrow output is just a function of current state.

There are multiple branch instructions in a program. We maintain a buffer to store the state of FSM for each of these branch instructions.

Branch prediction buffer / Pattern history table (confusing names)
(BPB) (PHT).

The FSM can initialize to anything based on branch address.

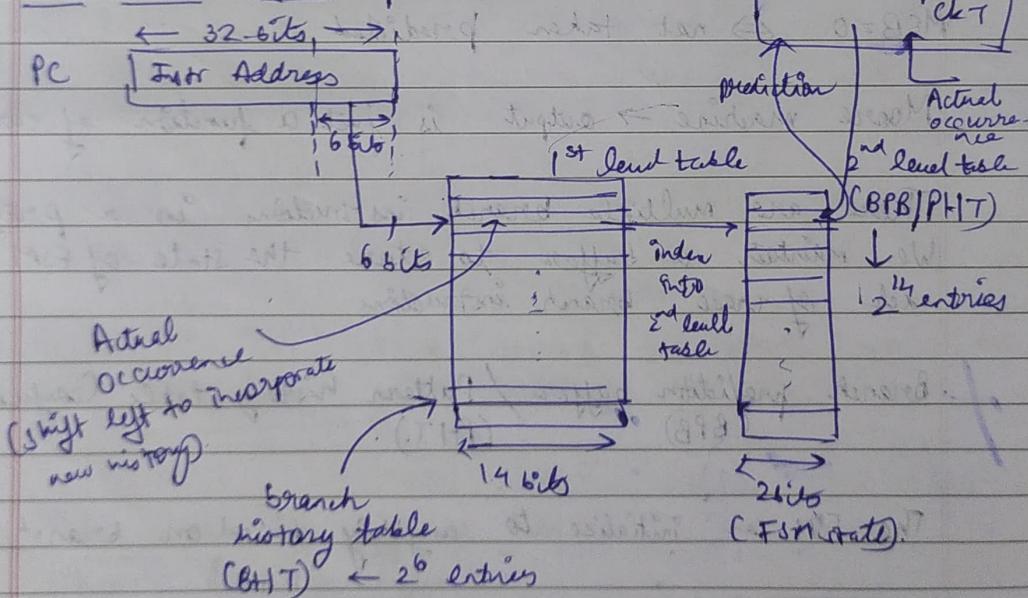


Branch prediction buffer supports only distinct 2^m branch instructions.

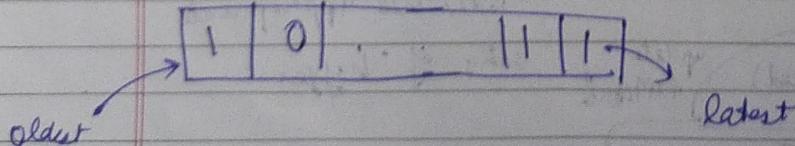
Issue: the index bits might cause aliasing among branch addresses, say B_1 and B_2 , where $B_1[m-1:0] = B_2[m-1:0]$.

Sol: can store tag bits, but that will ~~cause~~ slow down
can store some additional bits (not all) for tag bits.

② Load 2-level branch predictor



14 bits \rightarrow last 14 of branch outcomes (actual)
not prediction outcomes
actual outcomes



Initial branch

Possible that B_1, B_2 (with

14 bits change \rightarrow instruction is

Also 2 branch
same predict

B₁: if C

B₂: if C

B₃: if

Observation: if
 \Rightarrow Behaviour of

* Correlating bne

① global (m,n):

global for
all branches
on

Initial branch history table \rightarrow maybe all 0s or all 1s

Possible that 14-bit history may be of two branch instructions B_1, B_2 (with 6-bits same).

14-bit change ~~at~~ every time $\&$ the corresponding branch instruction is executed.

Also 2 branch instructions with some history will have same prediction.

Q2: $B_1: \text{if } (aa == 2)$

$aa = 0;$

$B_2: \text{if } (bb == 2);$

$bb = 0;$

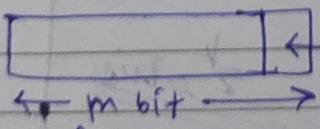
$B_3: \text{if } (aa != bb) \{ \quad \}$

Observation: if B_1 is taken and B_2 is taken $\Rightarrow B_3$ is not taken.
 \Rightarrow Behaviour of B_3 is correlated to behaviour of B_1 and B_2 .

* Correlating branch predictors:

① Global (m, n) :

Global for
all branches



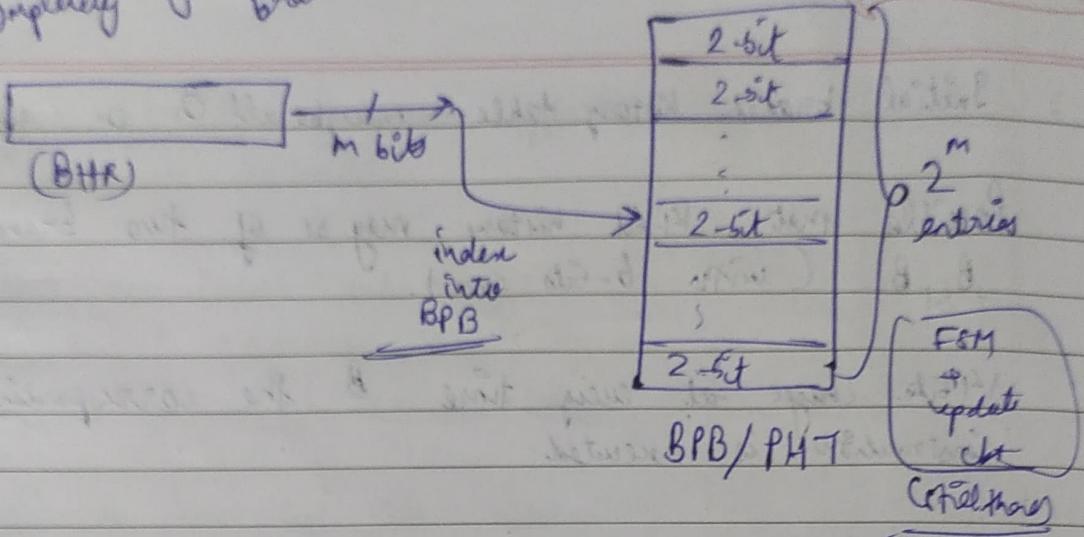
Branch history register (BHR).

Latent branch direction
observed.

shifts left

(completely ignore current branch instruction address)

classmate
Date _____
Page _____

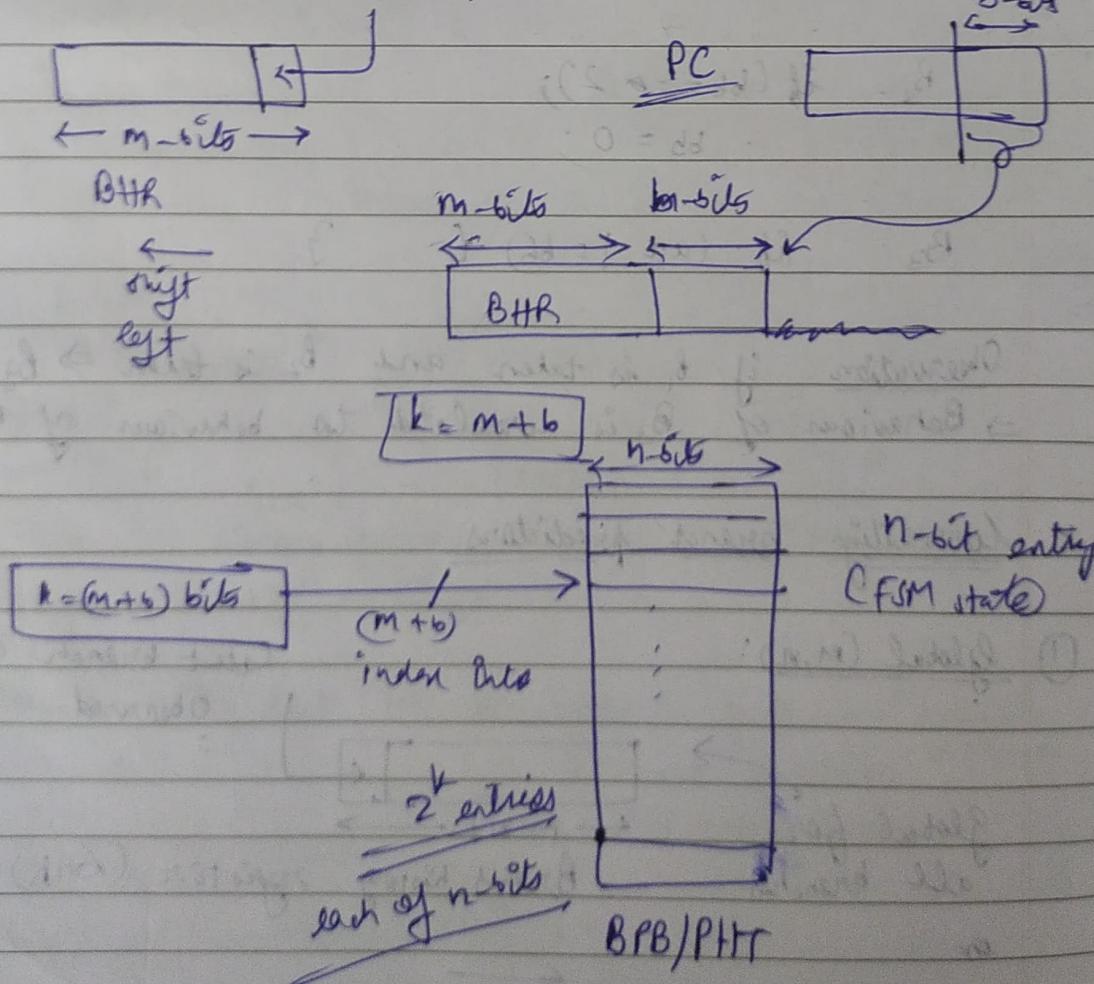


m = size of BHR in bits

n = no. of bits in FSM state (here, $n=2$)

② Hybrid (k, n) :

Actual/observed branch direction



Eg:
 $k=6$
 $m=4$
 $b=2$

$0 \ 1 \ 1 \ 0 \ 1 \ 0$

F_{26}



[F&M state]

[26]

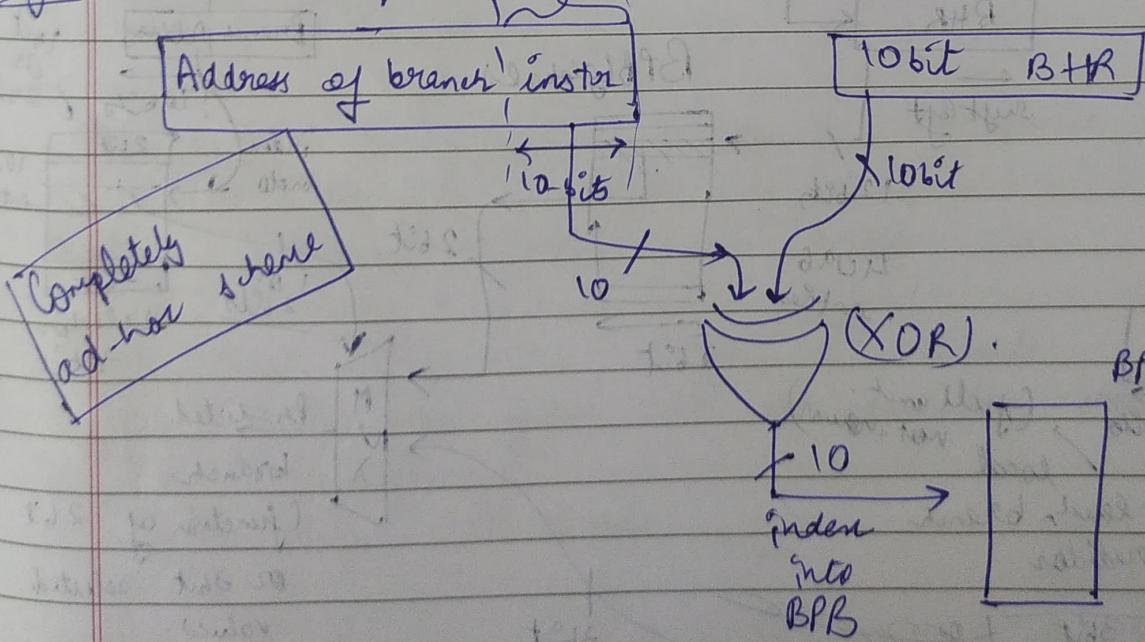
Say k is fixed.

$b \uparrow (\Rightarrow m \downarrow)$ \rightarrow aliasing probability decreases.
 \rightarrow less global history can be preserved.
 (BHR is smaller)

Hybrid (k_n) \equiv Hybrid (k, m, n) \equiv Smith's algorithm.
 (share)

③ Hybrid correlating: select [Macfarling, 1993] :-
 → Insight: LSB of branch address & BHR bit strings are
 not correlated.

Eg: 10-bit example



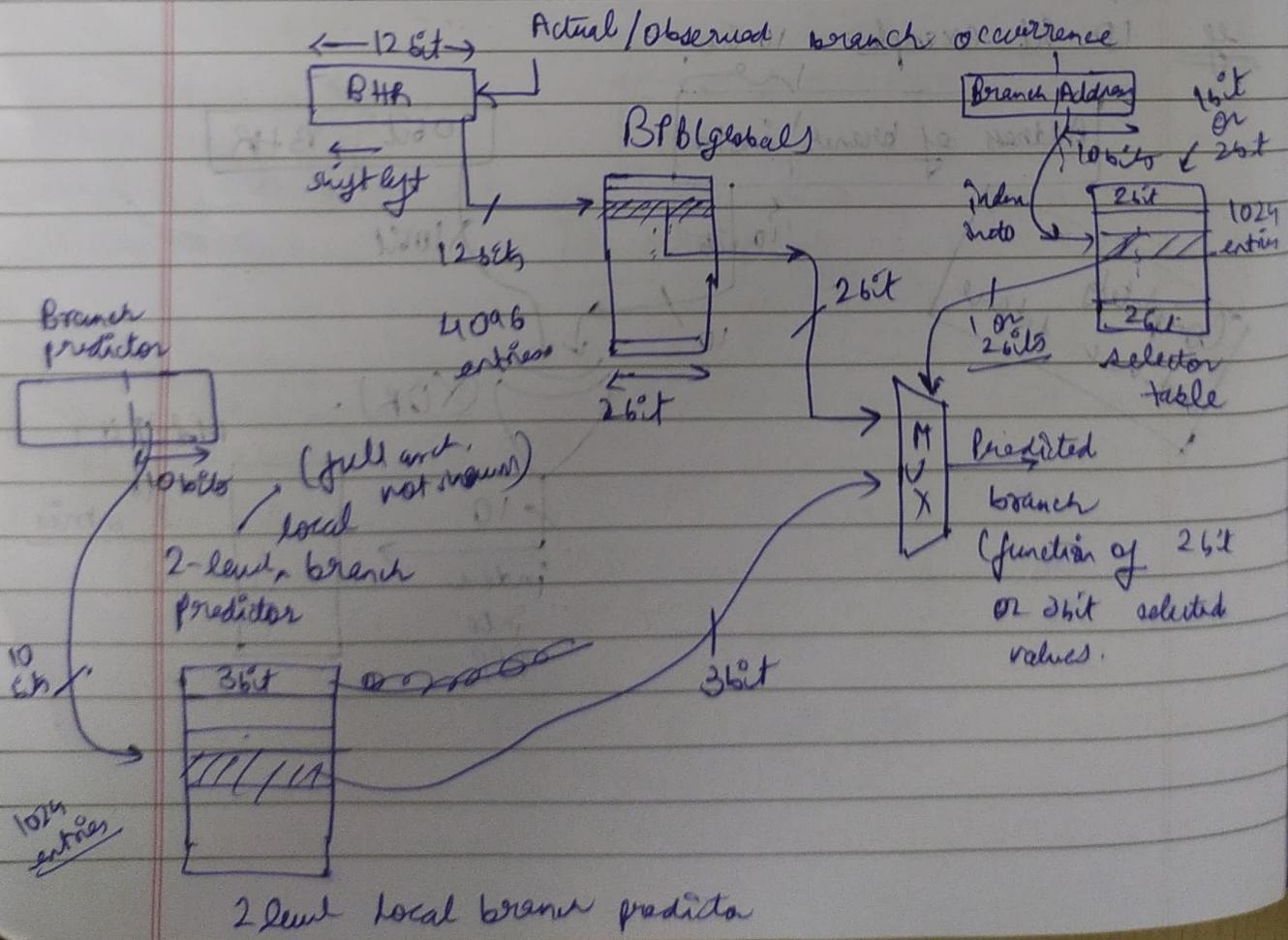
④ Tournament Predictor

Idea:

- ① Have both local and global prediction
- ② Local prediction: 2-level local predictor
- ③ Dynamically (during runtime), by observing the recent performance (accuracy of predictions), select either global/ local predictor. (Change predictor based on ~~performer~~ performance).

→ Change predictor type after 2 consecutive mispredictions.

Example: BHR \rightarrow 12-bit, global predictor, BPB: 4096 entries, each 2-bit



$$F \supset X \supset W$$

$$F \supset D$$

We can allow both local and global predictors to update ~~to~~ always. But we will use the selector table to determine and choose prediction from appropriate predictor.

$$\rightarrow L(1) < L(2) < L(3) < L(4)$$

$$P(1) \rightarrow h[0 : L(1)-1]$$

$$P(2) \rightarrow h[0 : L(2)-1]$$

$$P(3) \rightarrow h[0 : L(3)-1]$$

$$P(4) \rightarrow h[0 : L(4)-1]$$

\rightarrow hash \rightarrow only some bits of pc and RAW are used

* Tomasulo's algorithm:

1 reservation station for each add/multiplier etc.

① feed b_1, b_2, f_3 RAW hazard

② feed b_3, b_4, f_1, f_2 RAW hazard

① say EX finishes in 5th cycle. Value pushed to CDB at end of 5th cycle.

② WB stage completes in 6th cycle.

② Cannot start EX before 7th cycle.

→ first read from CDB into reservation station in 6th cycle.

→ EX in 7th cycle.

① fadd b_1, f_2, b_3

RAW hazard

② fadd f_1, f_4, f_2

① not allowed to execute

② must write f_1 after ①.

or not allowed to update reg file.

Registers
hardware used then internally.
→ buffers
(internal, non-architectural)
Broadcast values on CDB.
Every unit is having CDB.
architectural registers
(can be used by programmers).

classmate

Date _____
Page _____

(W)
(W)
R
W

Reservation
in clock
It is a pre
cycle

Network on chip → use LAN type protocols on-chips to communicate b/w units.

→ Read TB before speculation

1 Reservation station is present for 1 adder/multiplier unit.

WAW and WAR hazards can be removed by register re-naming. but RAW hazards cannot be removed by this technique.

Register re-naming can be done by the compiler as well. (software techniques).

Branch predictive speculative execution is used in deep front pipelines and EX stage takes multiple cycles and even ID stage takes multiple cycles.

→ branch predictors are extremely accurate. Hence this allows 10-20 instruction to complete execution before knowing correct branch decision.

Address calculation is always in-order. (Simplify hazards through memory. The addresses are pushed to load buffer and store buffer only after they are resolved.)

Reservation stations are intelligent to detect RAW hazards.
No direct inter-connection b/w units.
CDB holds value does not hold values it transmits values.

1 reservation
for 1 e
Say va
cycle.

Say exe
8th &
9th &
update
adder/m
execution

→ Reservation
update
It stays
because
maintain
order

Say we
Load/stor
in-order

DMA co
allow

CDB w
units can

Reservat

Reservation station reads value into CDB in clock cycle i and updates itself. It is supposed to adder/multiplier in clock cycle $(i+1)$.

classmate

Date _____

Page _____

1 reservation units holds information for 1 instruction for 1 execution unit.

→ Not the case for modern processors like Intel i7.

Say values are available on CDB in 8^{th} clock cycle. The reservation unit

Say execution unit writes to CDB at end of 8^{th} f cycle. The value is still available during 9^{th} f cycle. The reservation station reads and updates value during 9^{th} f cycle. The corresponding adder/multiplier for this reservation station starts execution in 10^{th} f cycle.

→ Reservation station of a unit is not allowed to update till execution of this unit completes. It stays in busy state till then. This is needed because the internal data structures need to be maintained (for other units) till execution completes.

Outer execution

Say we have 2 reservation stations per adder.

Load/store buffer has only 1 memory unit. Execute in-order.

DMA controller can take hold of memory bus and not allow CPU to access memory.

CDB has to be wide enough because multiple units can write to CDB within same clock cycles.

Answers

We can reorder loads. But we cannot reorder "stores and stores" or "stores and loads".

When issuing a load to memory, we first check if there is a store before this load with same effective address as that of load. If there is such instruction, we first allow store to complete.

Once functional unit completes, it updates its reservation station and CDB in ~~simultaneously~~.

Once functional unit completes, it updates all other reservation stations, not its own reservation station. However, the status of execution is stored in parent reservation stations. This status is used by other reservation stations.

Trap: software interrupt

→ say a process calls a system call. This will raise software interrupt / trap.

Interrupt: May or may not be error.
→ Timer interrupt

Exceptions: Generated internally by hardware
→ Div by 0
→ Invalid op-code.

If exception is caused by an instruction that must not have been executed (due to speculation), we ignore this exception.

→ This needs +

There is a
This buffer
buffer

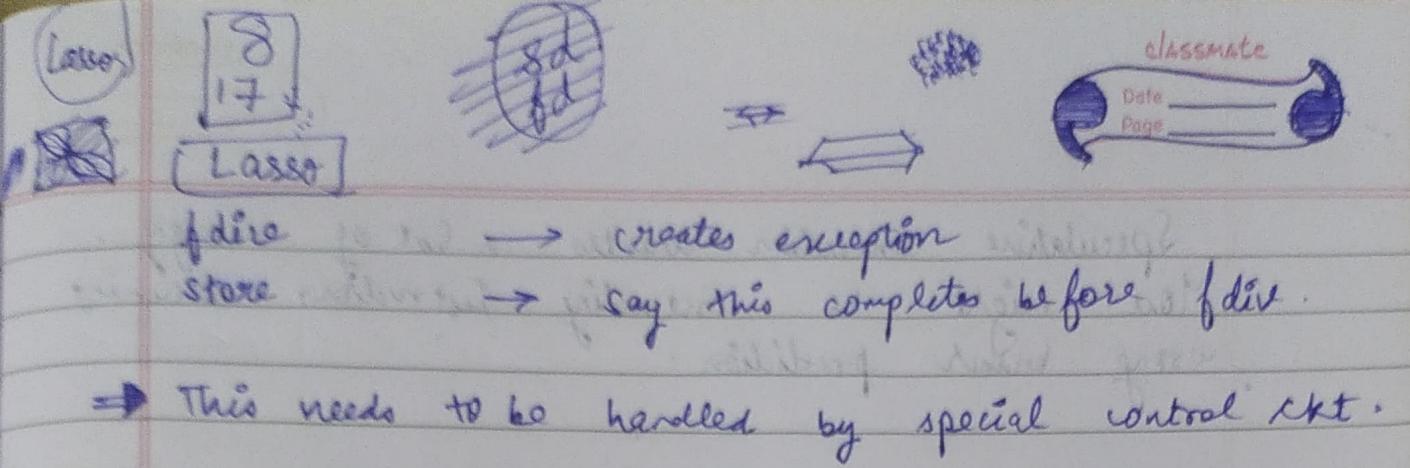
B7B
Branch address
(of instr)

It is helpful
not in ins

8 Reorder buffer →
this.

Full speculate

Precise except
the process
did speculate
an exception
an exception



ldi

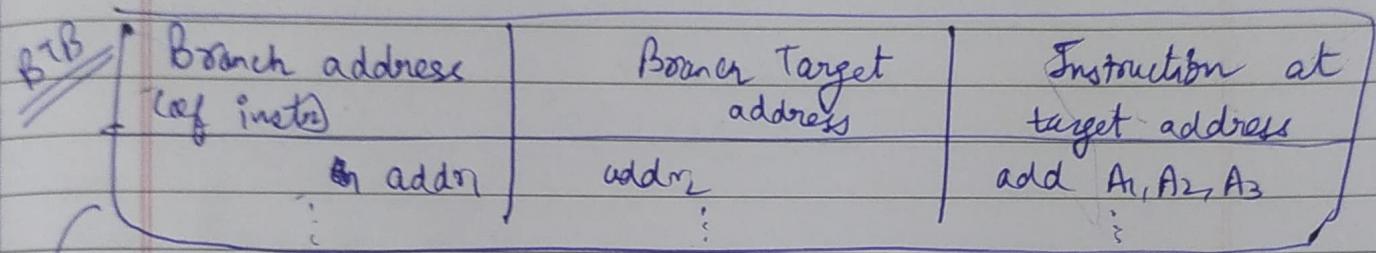
→ creates exception

store

→ say this completes before ldi

⇒ This needs to be handled by special control rkt.

There is a buffer b/w L3 cache and memory.
 This buffer brings in next cache block ~~100100~~ as ~~100100~~ was not needed.



It is helpful if instruction at address $addr_2$ is not in instruction cache.

8 Reorder buffer → stores values in program order (logic ensures this). Commit happens in program order.

Full speculation → not implemented ~~as of~~ till now.

Precise exceptions → Exceptions are not raised if the processor executes an instruction that it did speculatively. ~~or never~~ (Do not raise an exception if there would not have been an exception during perfect in-order execution)

Speculative execution wastes a lot of energy.
You might be executing instructions assuming wrong branch prediction.

The extra things get released out of which heat is

generated due to which heat is emitted in a small
area which is called noise. If noise is emitted into
the system, then it becomes a problem.

To understand

what is target

and what is noise

target noise

noise

noise

variables should be

constant

noise

extra numbers to calculate the frequency of the
noise which is in the system.

extra signal reduces memory access because when
extra memory is accessed then it is called noise.

extra noise is not helpful for the system.

extra noise is not helpful for the system.