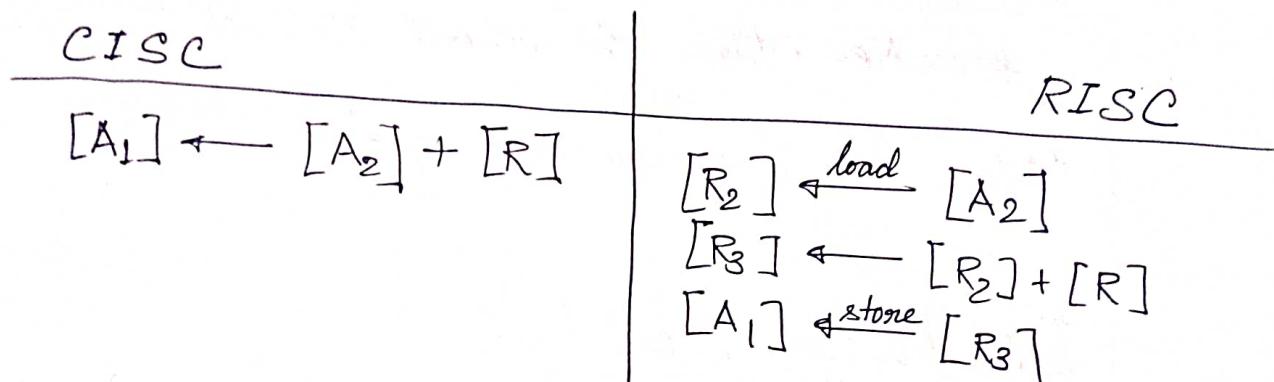


HIGH PERFORMANCE COMPUTER ARCH!

Microprocessor: A single chip containing the entire CPU.

RISC-V (32-g.p. and 32-f.p.)

- * Classic RISC design
- * Load-Store design
- * Plenty of General Purpose Registers
- * Extensible for advanced features



How to achieve high performance?

Parallelism

- * System-level (e.g. scale-up servers)
- * Individual processor level (e.g. pipelining)
- * Design level optimizations (e.g. faster clk, larger cache, e.t.c.)

Rule of Thumb: An executable spends ~~90%~~ 90% of its execution time in about 10% of its code.

"Make the common case fast"

Principle of Locality:

Spatial Locality: Instructions (or data) whose addresses in memory are close to each other are likely to be executed (or accessed) close in time.

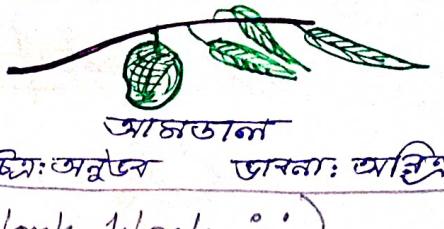
(1 out of 7 instructions are branches/function calls/jumps, e.t.c). —— 15%

Temporal Locality: Recently accessed data (or executed instruction) are likely to be accessed (or executed) again.

(about 6 instructions are in one basic block, i.e. without jumps/function-calls/branches. Processors use branch prediction for better performance)

Amdahl's Law [IBM - 1967]

- * "Law of Diminishing Return"
- * "Inherently pessimistic law" --- (Work Work :)



Definitions

- Speedup after enhancement, (S_e)
- Fraction of program which is parallelizable (f_e)

Speedup |
overall

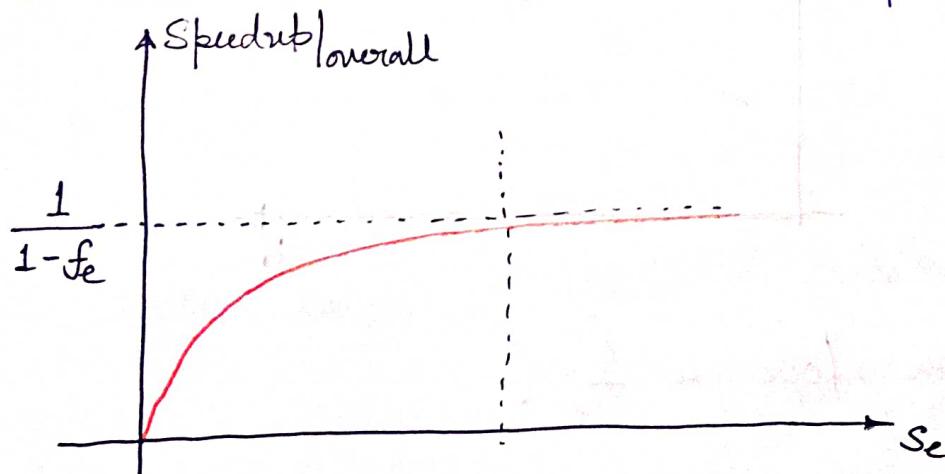
$$= \frac{\text{Execution time of entire task without enhancement}}{\text{Execution time of entire task with enhancement}}$$

$$= \frac{\text{Ex.-time}_{\text{old}}}{\text{Ex.-time}_{\text{new}}}$$

$$= \frac{\text{Ex.-time}_{\text{old}}}{\text{Ex.-time}_{\text{old}} \left[(1-f_e) + \frac{f_e}{S_e} \right]}$$

$$= \frac{1}{(1-f_e) + \frac{f_e}{S_e}}$$

$$\therefore \lim_{S_e \rightarrow \infty} \left(\text{Speedup}_{\text{overall}} \right) = \frac{1}{1-f_e}$$



Gustafson's Law [1988]

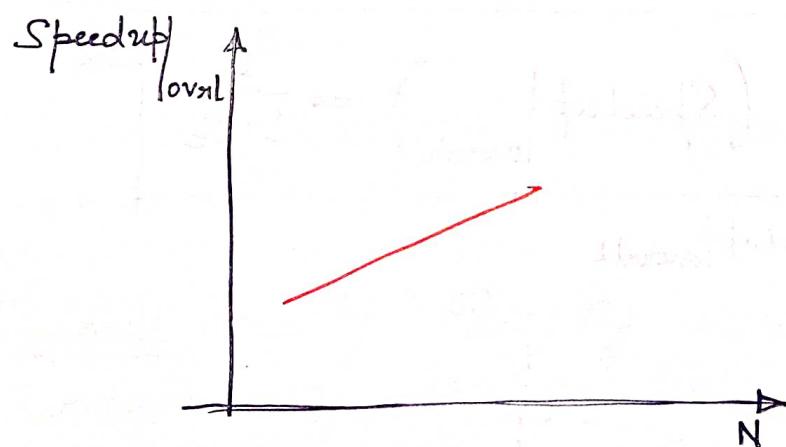
- * "Logic is the following: we scale up the problem size as well, as we are capable of scaling up the system; Problem-size \propto # processors"
- * "Extremely ~~Pessimistic~~" Optimistic."
- * "The serial part of the problem stays the same even when we scale up."

$$\frac{\text{Ex-time}_{\text{old}}}{\text{Ex-time}_{\text{new}}} \leftrightarrow \# \text{processors} = 1$$

$$\frac{\text{Ex-time}_{\text{old}}}{\text{Ex-time}_{\text{new}}} \leftrightarrow \# \text{processors} = N$$

$$\therefore \text{Ex-time}_{\text{old}} = \text{Ex-time}_{\text{new}} \left[(1-f_e) + N * f_e \right]$$

$$\therefore \boxed{\text{Speedup}_{\text{overall}} = \frac{\text{Ex-time}_{\text{old}}}{\text{Ex-time}_{\text{new}}} = (1-f_e) + N * f_e}$$



(Max slope is 1).

Comparing Results:

Example:

Sequential part, $1 - f_e = 0.2 \Rightarrow f_e = 0.8$

M1 — single core

M2 — 80 core $\Rightarrow N = 80 = S_e$

$$\therefore [\text{Amdahl}] \text{ Speedup}_{\text{AMDL}} = \frac{1}{(1-f_e) + \frac{f_e}{S_e}} = \frac{1}{0.2 + (0.8)/80} = 4.76$$

$$\therefore [\text{Gustafson}] \text{ Speedup}_{\text{GSTFSN}} = (1-f_e) + N \cdot f_e = 0.2 + 80 \cdot 0.8 = 64.2.$$

"In early 2000's, Gustafson's Law was completely trashed because of extremely optimistic results" (Moye Moye)

Processor Performance Equations:

$$0. \text{ CPU time} = (\# \text{CPU clk cycles}) * (\text{clk cycle time}) \\ = (\# \text{CPU clk cycles}) * t_e = (\# \text{CPU clk cycles}) \cdot \frac{1}{f_e}$$

$$1. T_e = N_e \cdot t_e$$

$$2. \text{ } \cancel{N_e} = T_e \cdot f_e.$$

tile world analysis 2 cases of scenarios.

Case-I: (Each instruction takes the same number of clock cycle).

A bit unrealistic case (delusion).

\therefore Clock cycles per instruction (CPI) = $\frac{1}{IPC}$

$$= \frac{\# \text{clk cycles for a program}}{\# \text{instructions in program}}$$

$$T_{\phi} = CPI \cdot n_{\phi} \cdot t_{\phi}$$

(A)

n_{ϕ} and t_{ϕ} are hard to decrease; so the natural choice is reducing CPI. A computer architect's primary goal here should be to try reducing T_{ϕ} by reducing CPI.

Case-II: (There are several classes of instructions, each with a different CPI)

n_i instructions get CPI_i

$$T_{\phi} = \left[\sum_i (n_i \cdot CPI_i) \right] \cdot t_{\phi} \quad \dots \dots \quad (B)$$

Equivalently, if ~~CPI_i~~ $f_i = \frac{n_i}{\sum n_i}$, then

$$CPI_{\text{effective}} = \sum_i f_i \cdot CPI_i$$

There are further generalizations (due to page fault/disk access)

A further generalization of (B) designed to handle penalty would be something like:

$$T_f = \left[\sum_i (n_i \cdot CPI_i) + n_{\text{penalty}} \right] \cdot t_f$$

and, $CPI_{\text{effective}} = \sum_i f_i \cdot CPI_i + \frac{n_{\text{penalty}}}{n_{\text{penalty}} + (\sum_i n_i \cdot CPI_i)}$

where $n_{\text{penalty}} = \# \text{clk cycles wasted for penalty.}$

Memory Hierarchy - I

/* slides : C */

"n-way associative" \rightarrow #entries in one set is n.

$$\# \text{sets} = \frac{\# \text{index_bits}}{\# \text{blocks}} = \frac{\# \text{blocks}}{\text{associativity}}$$

Corollary: For ~~for~~ fully associative $\rightarrow \# \text{index_bits} = 0$.

"4-way set associative is the most common"

• 4-way set associative cache has 4 sets
• each set has 4 blocks
• associativity is 4
• 4 sets \times 4 blocks = 16 blocks

Write through: Always write to main memory whenever "store" operation happens.

Write Back: Store dirty bit & write back to main memory only when needed.

Write Buffer: (in CPU) Holds data to be written back to memory; CPU continues (can be used together with Write Back).

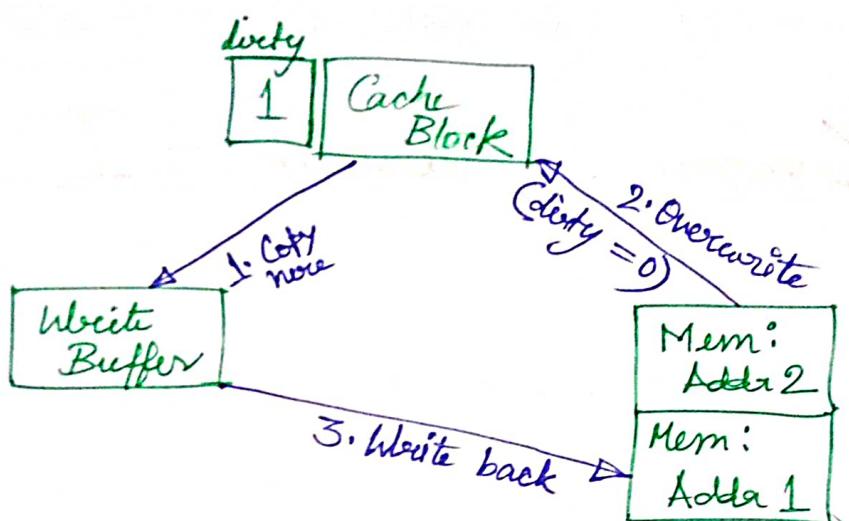


Fig: Replace data of Addr 1 in Cache (dirty) with data of Addr 2

Reads get more preference over writes

Write Allocate:

When write miss happens, swap the ~~too~~ thing on memory that is to be written (handle dirty-ness as needed)

No-write allocate:

When write miss happens, don't even update cache, write directly to memory.

"In a modern day single-core CPU, more than 100 instructions are on flight at a given moment"

Similar hierarchy is used in Hard disk-main memory boundary.

Translation Lookaside buffers are used (TLB):

- * 16-512 PTE's are brought
- * 0.5-1 cycle for hit
- * 10-100 cycle for miss
- * 0.01% - 1% miss rate
- * Fully associative.



TLB Misses:

- Case 1: In PTE but not in TLB
- Case 2: Not in main-memory.

"Tag comparison and getting out data happen in parallel. Whether the CPU will pay attention to this data will depend on the tag comparison."

Estimation of size of Multi-Level Page Table

[Worst Case analysis] Entire 48-bit VM space is utilized.

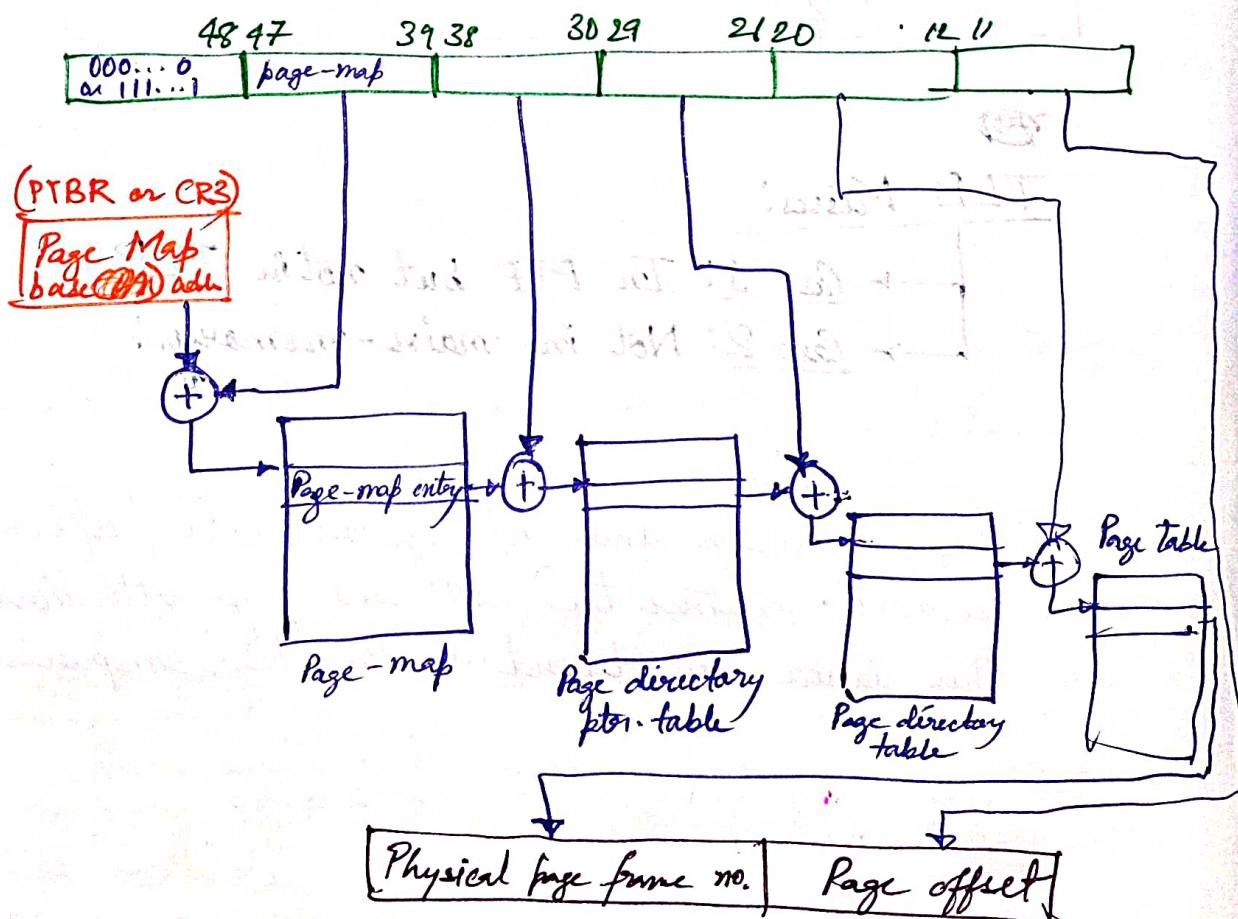
Page size - 4 kB
 $= 2^{12} \text{ B}$, Page table entry size - 8 B
 $= 2^3 \text{ bytes}$

Option 1: Single Level (Flat) P.T.

$$\# \text{ pages} = \frac{2^{48}}{2^{12}} = 2^{36}; \text{ Size of PT} = 2^{39}$$

Option 2: 4-level P.T.

Outer level P.T. = Page Map ($= 1 \text{ page}$)



$$\therefore \# \text{pages} = 1 + \underbrace{2^9}_{\substack{L_4 \\ \text{P.T.}}} + \underbrace{2^{18}}_{\substack{L_3 \\ \text{P.T.}}} + \underbrace{2^{27}}_{\substack{L_2 \\ \text{P.T.}}} + \underbrace{2^{27}}_{\substack{L_1 \\ \text{P.T.}}}$$

$\therefore \text{Memory used for P.T.} = 2^{12} (1 + 2^9 + 2^{18} + 2^{27}) \text{ bytes.}$

--- হালকা হালকা সুস্থিতা; শুক্রা ক্ষয়ের গুরুত্ব!

Look into VIPT Cache (L1 & L2) diagram!

Aliasing:

"Aliasing is basically the mapping of multiple virtual addresses being mapped into the same physical address"

— We will study aliasing in the context of Cache design.

This can happen if cache index collects some bits from V.A. which is not part of P.A. বা V.A. এর কোন অংশ

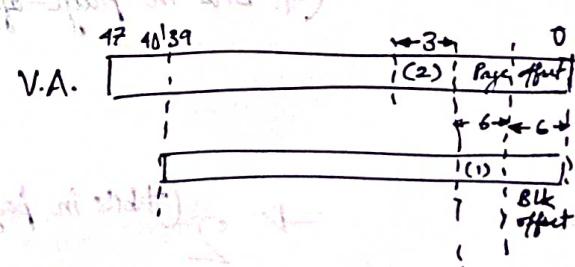
Example:

* [Opteron I Cache]:

- * 64 kB, 4 kB pages
- * 12-bit page offset
- * 2-way set associative
- * 48 bit V.A.
- * 40 bit P.A.
- * 64 B blocks

$$2^{(\# \text{index-bits})} = \frac{64 \text{ kB}}{(64 \text{ B}) \times 2_{\text{assoc.}}} = 2^9$$

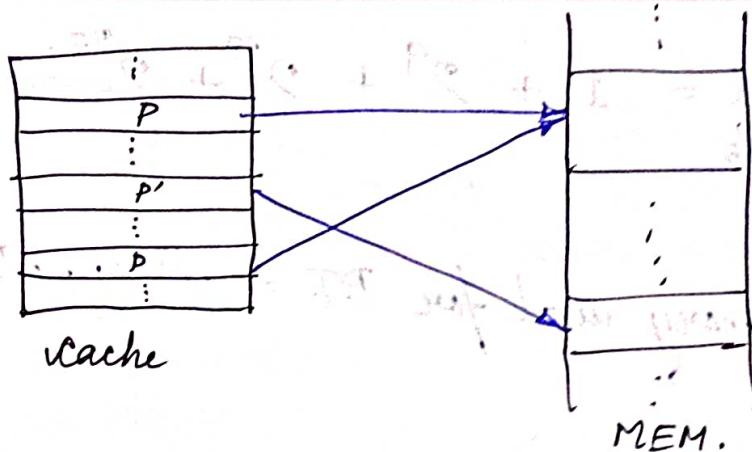
$$\Rightarrow \# \text{index bits} = 9$$



(1) 6-bit portion of cache index \Rightarrow both in P.A. & V.A.

(2) 3-bit portion of cache index \Rightarrow only from V.A.

can be used to introduce ambiguity; i.e. 2³ cache lines can point to the same physical address.



This is bad because one virtual memory might be overwritten whereas the other thinks that it is not written since the last time.

* One solution is to mark every other cache entry which can be mapped to ~~every~~ the same physical address, everytime something is written.

What we want then is the following:

$$\underbrace{(\# \text{bits in page-offset})}_{\text{virtual address}} \geq \underbrace{(\# \text{bits in block offset}) + (\# \text{bits in set index})}_{\text{physical address}}$$

$$\Rightarrow 2^{(\# \text{bits in page-offset})} \geq 2^{(\# \text{bits in blk offset})} \cdot 2^{(\# \text{bits in set index})}$$

$$\Rightarrow \text{Page-Size} \geq \text{Cache Size}.$$

Page Size < Cache Size \Rightarrow Aliasing

Average Memory Access Time (AMAT):

$$\boxed{\text{AMAT} = (\text{Hit Time}) + (\text{Miss rate}) * (\text{Miss penalty})}$$

overall goal: decrease AMAT

- * The aim is to decrease any of the three
- * Things become difficult to analyse; maybe an attempt to decrease 'Hit Time' might actually increase Miss rate.

Example: Analyzing 2-level caches

$$\text{AMAT}_{\text{2-level}} = \underbrace{(\text{Hit Time})}_{\text{HT}} + \underbrace{\text{Miss rate}}_{L_1} * \underbrace{\text{Miss penalty}}_{L_1 \text{ MP}}$$

$$\Rightarrow \text{AMAT}_{\text{2-level}} = \underbrace{\text{HT}}_{L_1} + \underbrace{\text{MR}}_{L_1} \cdot \left[\underbrace{\text{HT}}_{L_2} + \underbrace{\text{MR}}_{L_2} \cdot \underbrace{\text{MP}}_{L_2} \right]$$

quite low very large

MR_{L_2} is low as:

- Larger size of L_2
- Set associativity, e.t.c.

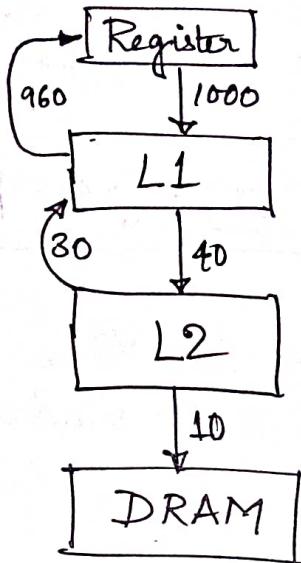
OK sorry MR_{L_2} is actually high 😞

Typically,

$$\text{HT}_{L_1} < \text{HT}_{L_2}$$

$$\text{MR}_{L_2} > \text{MR}_{L_1}$$

$$\text{MP}_{L_2} \approx 100 \text{ cycles.}$$



The local miss rate of L2 is large because most of the queries are satisfied by L1 itself.

He will call this the Local Miss Rate.
 Here $MR_L = \frac{10}{40}$

$$\therefore \text{Global Miss Rate} = MR_L * MR_D$$

"Don't disturb the memory system unnecessarily, there is a huge penalty for the same."

Example (1)

20% \rightarrow store
30% \rightarrow load

$\Rightarrow 100$ instructions $\rightarrow 150$ memory accesses.

Eg. 2

Among 1000 memory accesses

40 misses L_1

~~20~~ misses L_2

$$M.P. |_{L_1} = 10, HT |_{L_1} = 1$$

$$M.P. |_{L_2} = 200, HT |_{L_2} = 10$$

$$A \cdot MAT = ?$$

$$MR|_{L_1} = \frac{40}{1000} = 0.04$$

$$MR|_{L_2} = \frac{20}{40} = 0.5$$

$$\begin{aligned}
 \therefore AMAT &= HT|_{L_1} + MR|_{L_1} [HT|_{L_2} + MR|_{L_2} \cdot MP|_{L_2}] \\
 &= 1 + 0.04 [10 + 0.5 \times 200] \\
 &= 1 + 0.04 \times 110 \\
 &= 5.4 \text{ cycles.}
 \end{aligned}$$

Now, 1.5 memory accesses on average per instruction.

$$\Rightarrow \# \text{insts} = \frac{1000}{1.5}$$

∴ Average memory stall per instruction

=
 Memory stall can happen if L1 misses, L2 hits
 or L1 misses & L2 misses.

Advanced Cache Optimizations

1. Way Predictions:

- Enables faster S.A. cache (comparable to D.M. cache)

(Previously the last layer OR-gate creates delay)

Exact के लिए यह लाभों का; लिया चाहिए तो
नाकि आगे AMD/Intel - २३ लोकप्रिय आगे
बाकी.

- * Some kind of prediction hardware is used

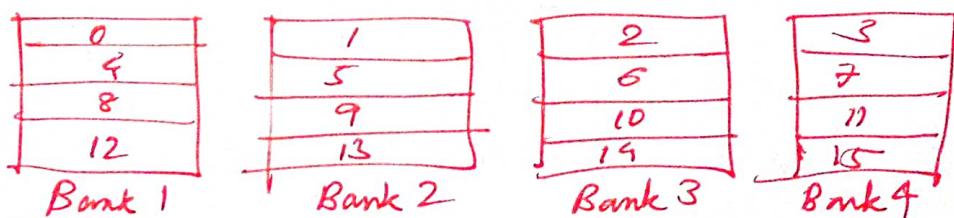
2. Pipelined Caches & Multibalanced Cache

Superscalar processors: Run multiple instructions

(will study in more detail)

- * Cache can have a different operating frequency than the microprocessor operating frequency. Therefore it is nice if the CPU can continue executing unrelated instructions rather than waiting for cache to return. [Pipelined cache]

- * Break a cache into banks



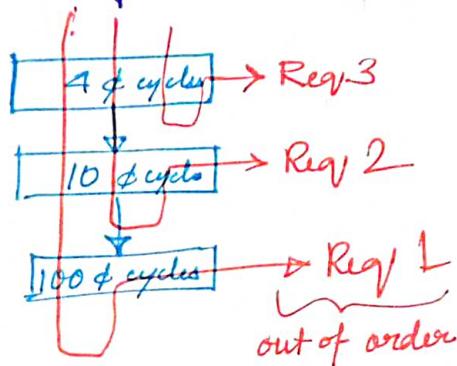
This allows "parallel" access in some sense. If two consecutive cache accesses are from different banks, the second access can be done even before first one is incomplete.

3. Non blocking Cache (Clockup-free cache)

When one memory request has missed the cache, the cache moves on to the next memory request (while remembering that one request exists)

- * Earlier caches could keep at most one request miss pending (hit-under-single-miss). Others are known as hit-under-multiple-miss.

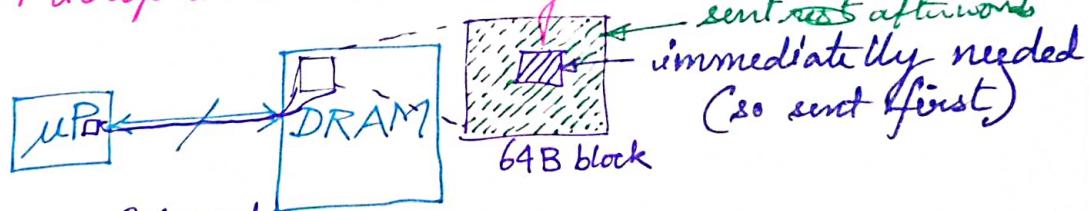
However, ordering is problematic.



Superscalar processors take care of it.

4. Early Restart / Critical Word First

Observation: Microprocessor needs only 1 word in a block.



Bus size \sim 2-4 words

Block size \sim 8/64 word

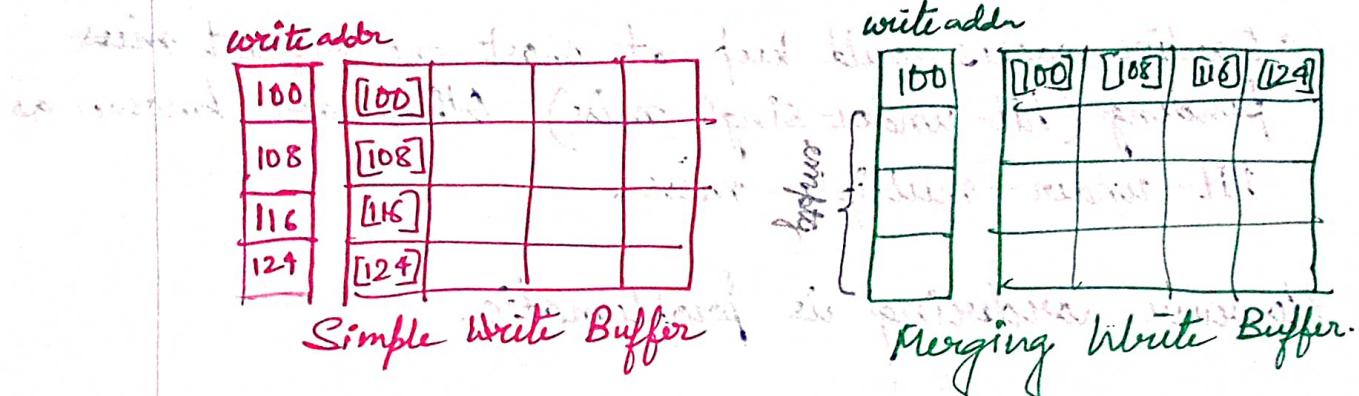
Fig. : Critical Word first

(Immediately needed word sent first, CPU works on it while other parts are sent later after that)

Early Restart: DRAM sends sequentially, when CPU finds its initially sent immediately needed block, it starts executing while the rest of the block is still being brought in.

5. Write Merging

- * If 'close-by' write operations queue up, then put them 'together' at write buffer & write them back at one write operation.



• If the writes which are close by makes write

Compiler-based scheduling \rightarrow static scheduling
 Hardware-based scheduling \rightarrow dynamic scheduling (runtime)

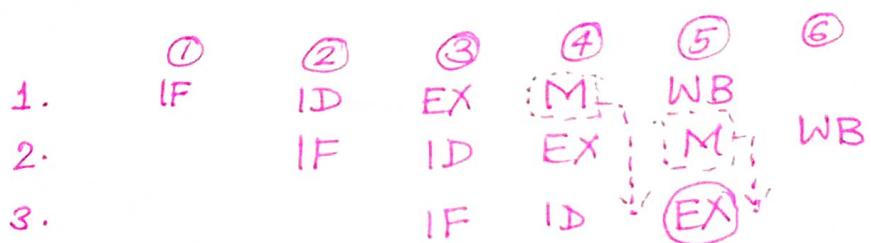
What we are going to study today is different from RISC-V.

Motivation:

for ($i = 0$; $i < N$; $++i$)
 $z[i] = x[i] + y[i];$

Each iteration gives us:

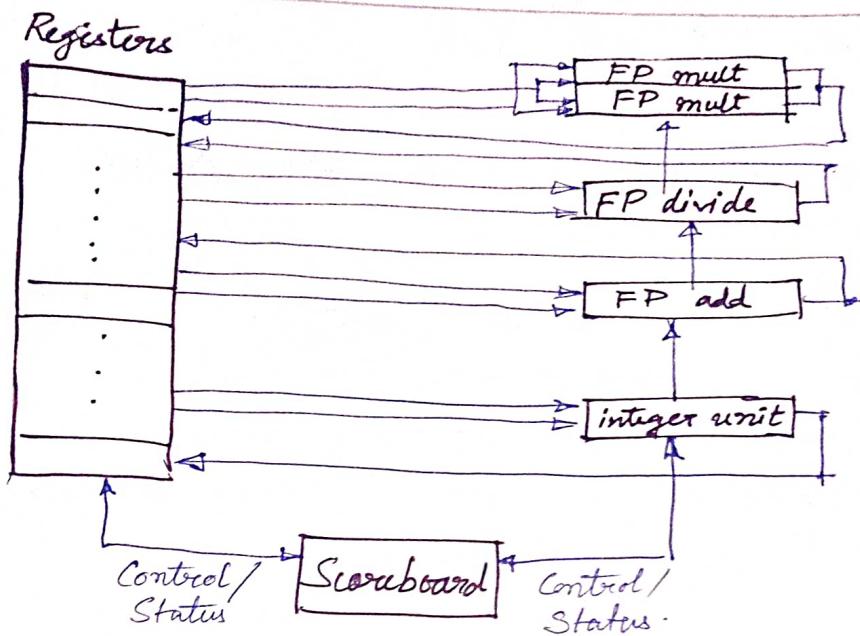
1. $R1 \xleftarrow{\text{load}} x[i]$
 2. $R2 \xleftarrow{\text{load}} y[i]$
 3. $R3 \xleftarrow{} R1 + R2$
 4. $R3 \xrightarrow{\text{store}} z[i]$
- } "Use after load"
 } \neq cycle stall!



Requires operands at the beginning of the 5th \neq cycle \rightarrow stall.

Goal:

- 1: Start execution of 1 instruction (if possible) every \neq cycle
- 2: Finish execution of as many instructions as possible every \neq cycle.



example:

fmul	$f_1, \boxed{f_2}, f_3$	10¢
fdive	f_4, f_5, f_6	10¢
faddl	$\boxed{f_2}, f_7, f_8$	

Decode — (?) —

Decode — (?) —
Issue (sending an instruction to the execution unit, after fetching & decoding)
happens in-order..... apparently for each unit only

fmul f1, f2, f3 Issued to FP Mul
fadd f1, f4, f5 Not even issued to FPAdd
↓ STALL

- For WAR \rightarrow The write is issued but the permission for the ~~so~~ later instruction to write is only given after the earlier instruction reads & finishes
- For Structural & WAW \rightarrow The issue itself stalls!

গত এক ঘণ্টা বেলে, কিছু ইসলাম আলু
হাসরাও বুলেছে বেলে বিশ্বাস হচ্ছে না। কিন্তু
কি হাসরাই মুক্তি না দেবা যাবে অন্যান্য
বেলে গানে ইব্ব না।

কিছু শ্রেষ্ঠতা রাখিল, হচ্ছে, ইয়ে চলেছে।
যা শীতার মাঝে শীতার হচ্ছে, চতুর্থ বহনে
আর্কিটেকচার; শীতার হচ্ছে গত-জানুয়ার কোন
শীতার শীতার কুকুর, কুকুর, তা রেসার্চে কা
শীতার জান ও জানোকা!

Correction: It tries to (? ! ? ! ?) issue
in order.

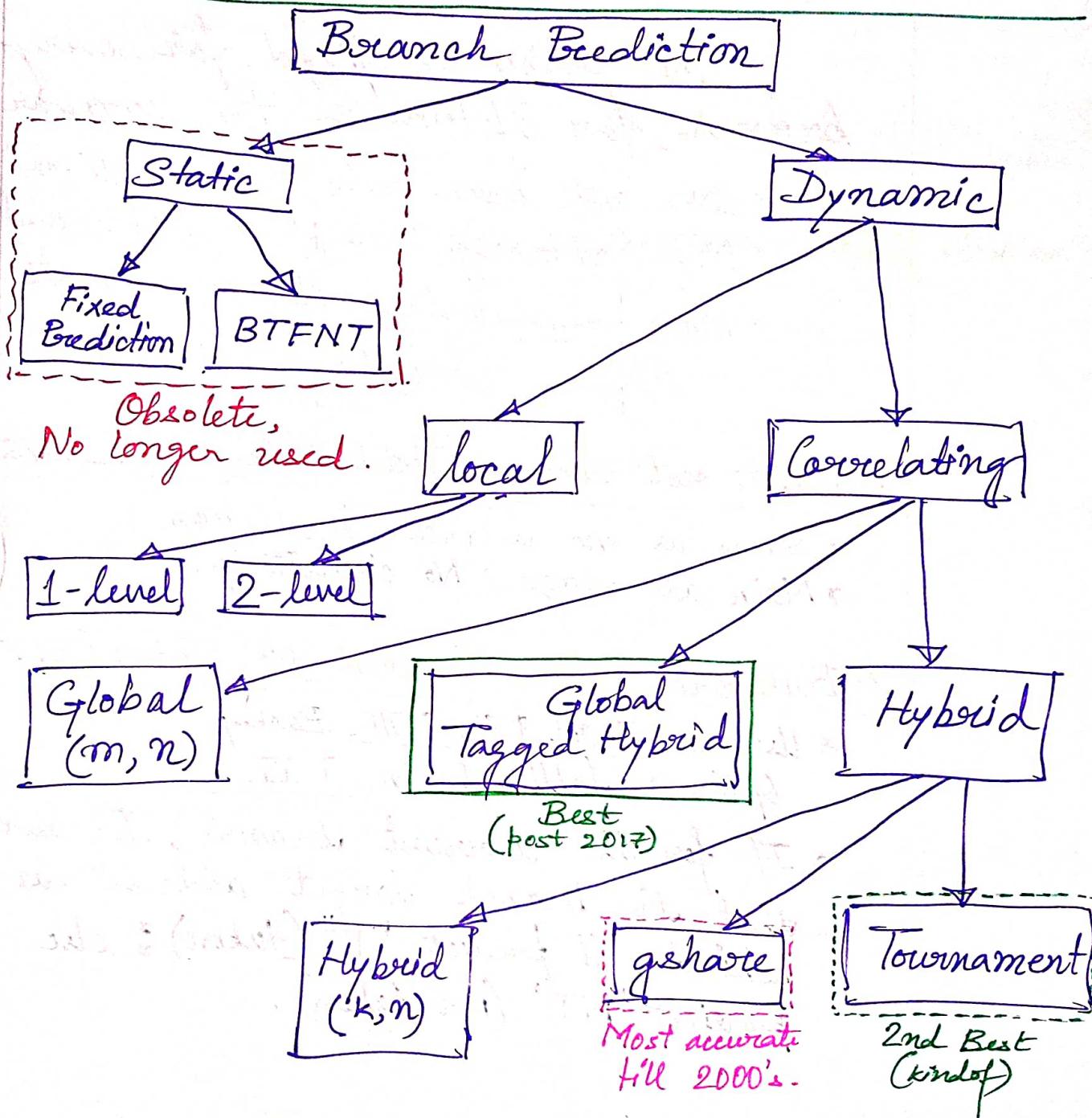
"জানো আজার হাসেনীতে রেক!"

More correction: Issue does stop. Issue
happens in order!

RAW — instruction starts, but operands
are not allowed to be read.

- * Once issued, a scoreboard can only restrict read-from/write-to Register file
- * Scoreboard is a pure controller
- * As a consequence, instructions mapped to the same execution unit are executed in order.

Branch Prediction: Hardware



Mainly two types?

- Static branch predictors (used earlier)
 - Fixed prediction
 - BTFTNT (Backward Taken Forward Not Taken)
- Dynamic branch predictors:
 - Is an FSM (state machine); changes state on misprediction.

Static Branch Prediction:

The policy is fixed for every branch, for lifetime of the processor

- Does not learn from mispredictions
- Main advantage: Simple to implement.
- Main disadvantage: Not so accurate
(often control hazard causes stalls)

1. Predict not taken (Intel i486) — ALWAYS!

- * Same as no branch prediction
- * Main advantage: No circuitry.

2. Backward Taken, Forward Not Taken (BTFTN)

- * Used in Intel PII, PIII, Backup
Optional fallback in PIV.
- * If for the current branch, I observe
that the branch target address is
behind, I predict "T" (taken); else
predict "NT" (not taken).

Dynamic Branch Prediction

Observe branch outcome, predict
based on observed outcome, learn from
mistakes.

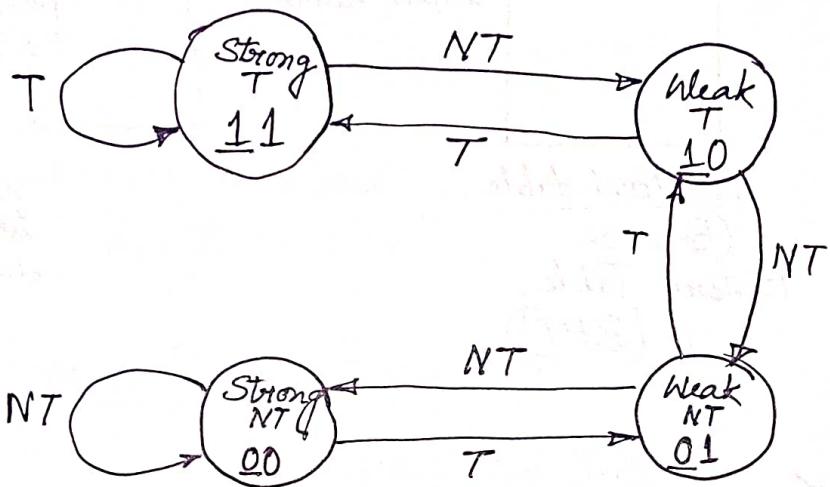
- 1 bit: predict last outcome
- 2 bit: change prediction only on 2
consecutive mispredictions.

Dynamic Branch Predictors:

1) Local-1-level:

Branch prediction FSM. e.g. 2-bit FSM/3-bit FSM or higher.

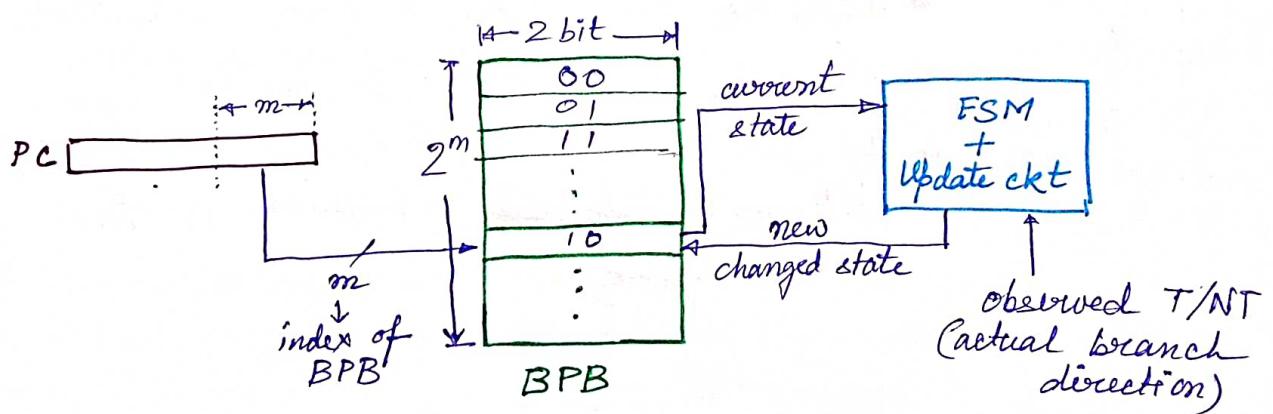
e.g. 2-bit FSM: (MSB is the prediction), Moore Machine



Branch Prediction Buffer (BPB):

(Confusingly it is also called Pattern History Table (PHT))
although this is not exactly pattern history.

* contains FSM states:



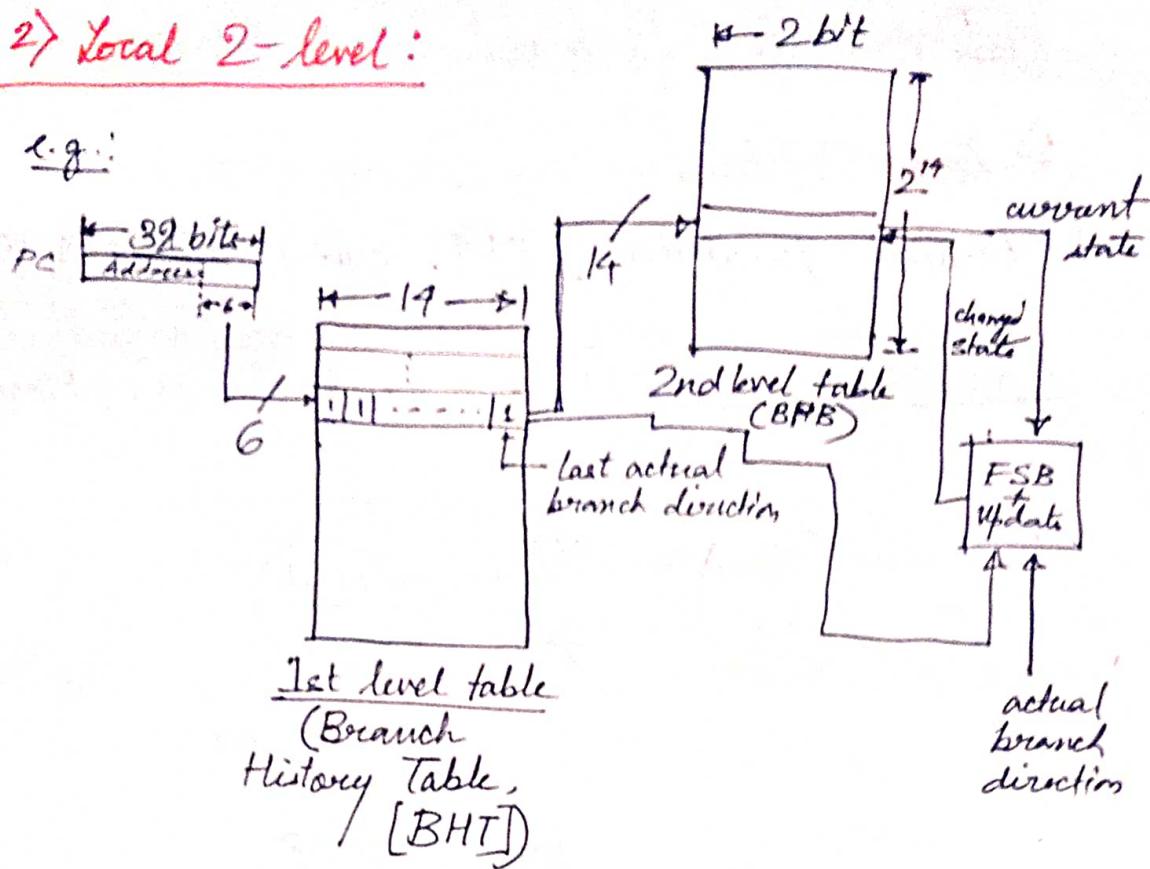
* DM Cache

* There is no tag — allows aliasing.

(The index bit might cause aliasing among branch addresses B_1, B_2 where $B_1[m-1:0] = B_2[m-2:0]$)

2) Local 2-level:

e.g.:



- * Aliasing is still not gone
- * (History can have a mixture if two branches B1 and B2 have the same last m bits).
- * Update on BPT

- depends on actual branch direction
- depends on history as well.
- If an entry of BPT is read it will get updated (possibly to the same value).

Example:

B1: if ($aa == 2$) $aa = 0$;

B2: if ($bb == 2$) $bb = 0$;

B3: if ($aa \neq bb$) {---}

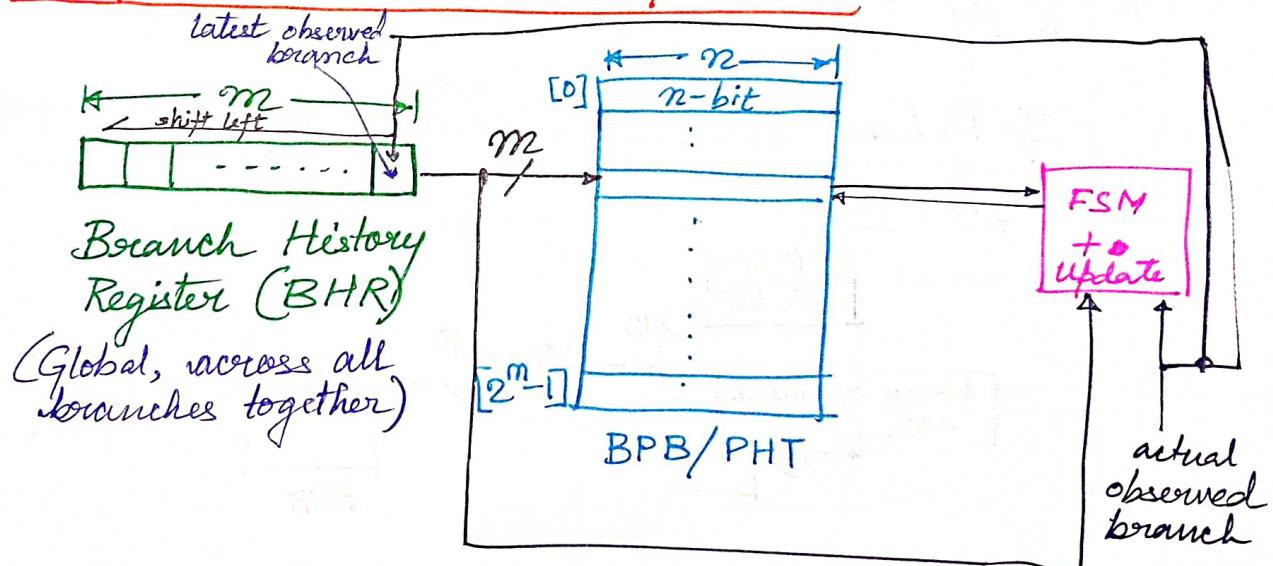
Observation: If B1 & B2 are taken \Rightarrow B3 is not taken

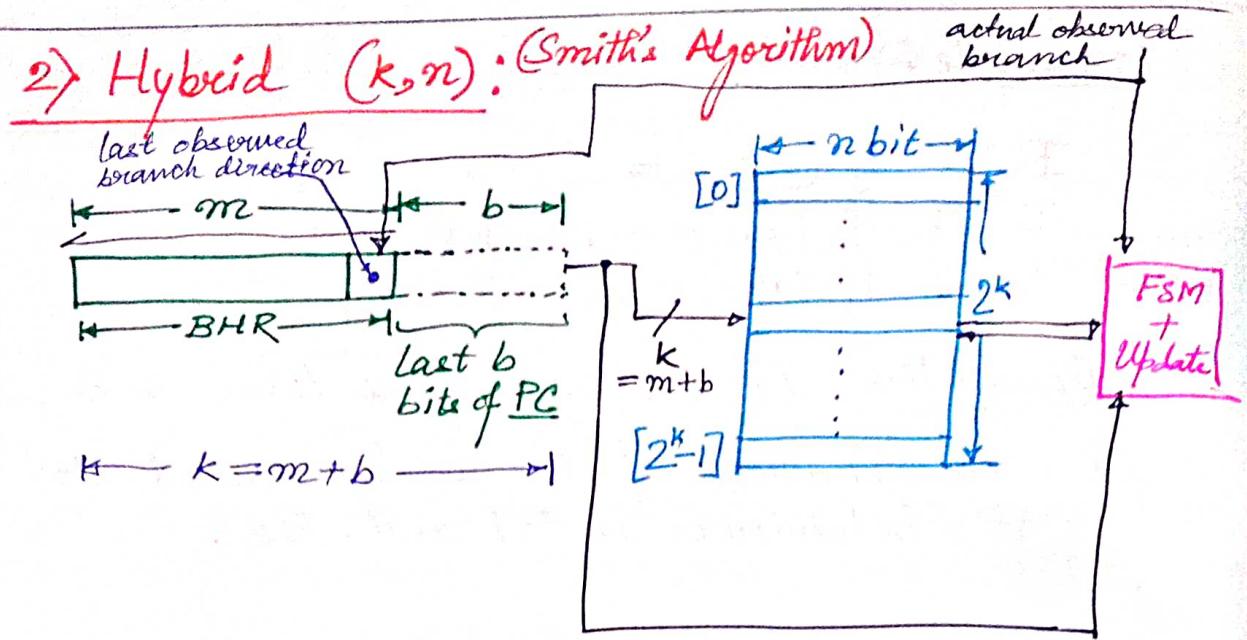
\Rightarrow Behaviour of B3 is "correlated" to the behaviours of B1 and B2.

These correlations give us the idea of correlating branch predictors.

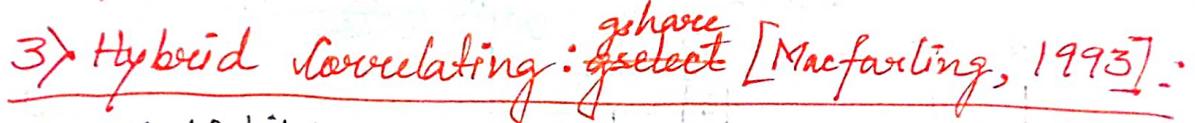
Correlating Branch Predictors:

1) Global (m, n) branch predictor:

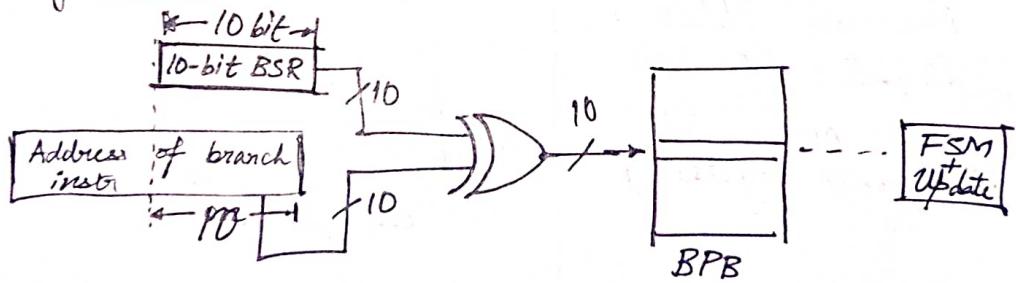




- * Unconditional jumps (jor) are always taken, then any prediction of NT would be wrong. This can be avoided by separately noting down jor instructions.
 - * If k is fixed then $m \uparrow \Rightarrow b \uparrow$.



e.g.: 10 bit:



- * Insight: LSB's of branch instruction address & BHR is uncorrelated.
 - * Best for ~20 years despite its adhoc nature

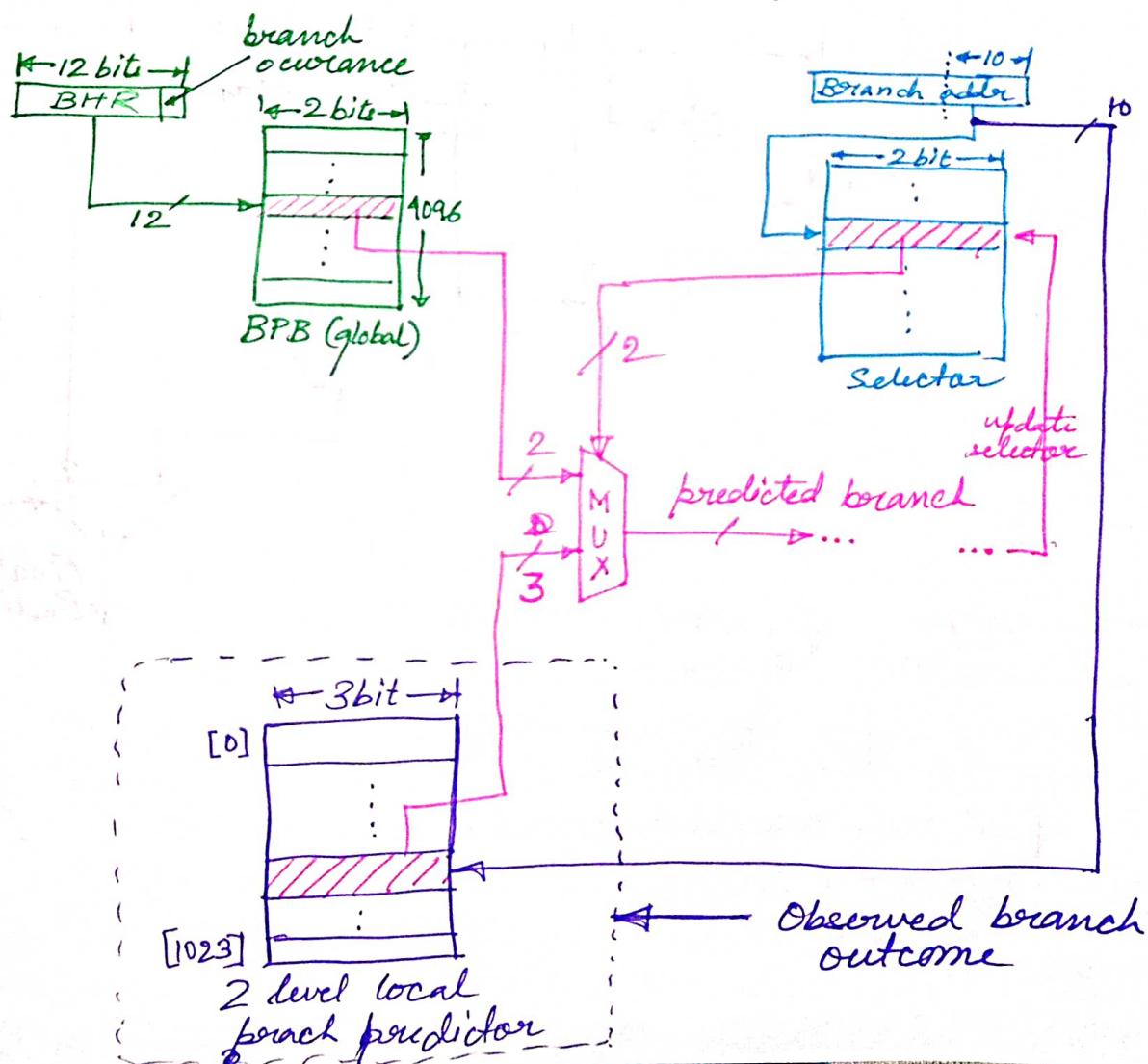
Tournament Predictor

Idea:

- (i) Have both local & global predictor
- (ii) Local predictor: 2-level
- (iii) Dynamically (during runtime), by observing the recent performance (accuracy of predicting) select either global/local predictor.

[Change predictor type after 2 mispredictions]

Example: BHR is 12-bits,
Global predictor BPB is 4096 entries,
each two bit.

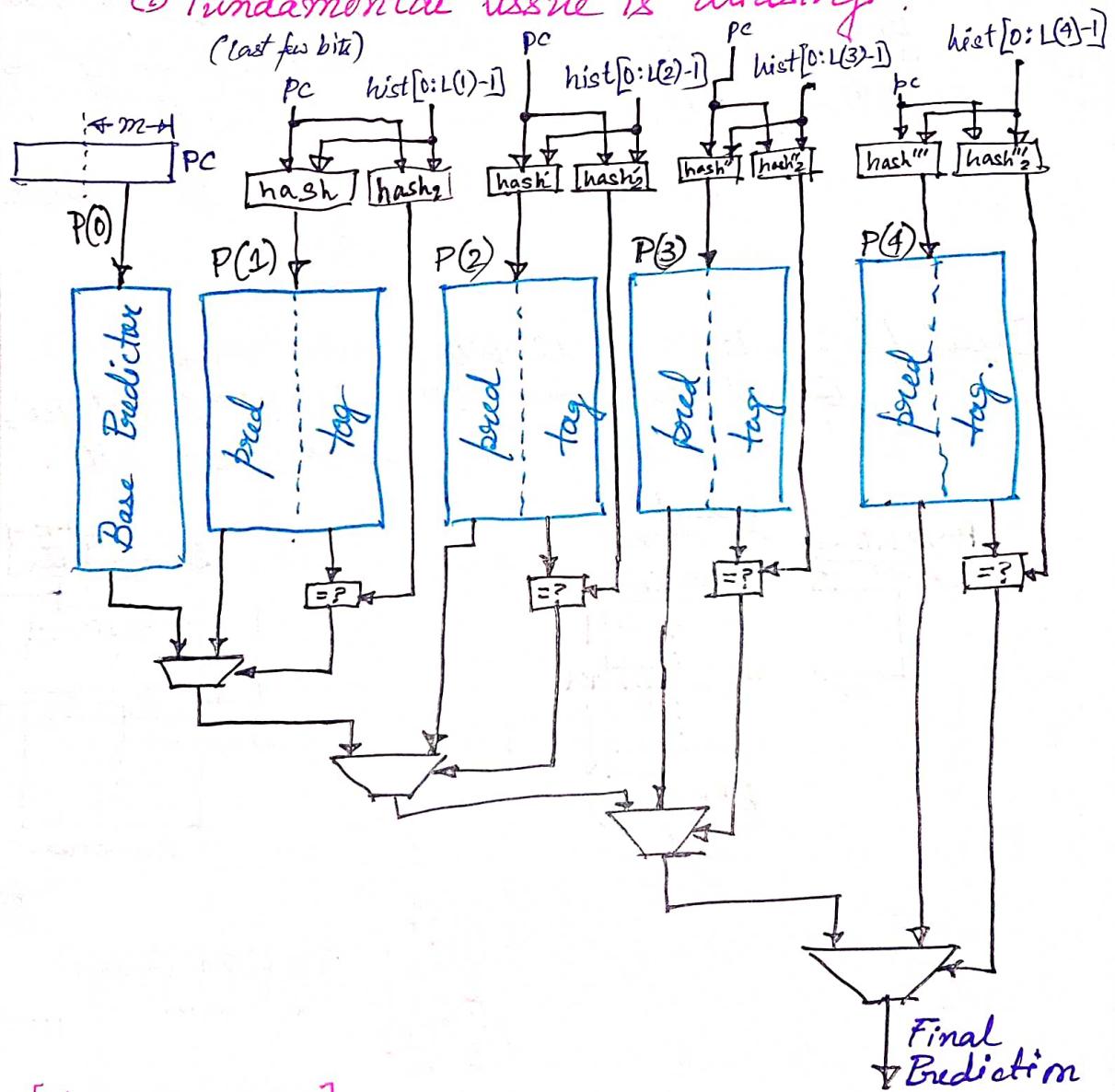


Tag Hybrid Predictor:

* Widely believed that modern processors do this.

* Motivation:

(i) Fundamental issue is aliasing.



[tag \leq 8 bits]

As compared to Tournament, with equivalent hardware overhead, achieves better prediction.
 $\approx 100\%$.

Instruction Level Parallelism: Hardware Techniques

As of now, issue is done once per clk cycle
 $\Rightarrow CPI \geq 1, IPC \leq 1$.

Tomasulo, IBM, 1967: Tomasulo's Algo.

* Differences from Scoreboard: (Better Phrased in slides)

- Communication between units are direct (can also be through Reg. File)
- No central controller (execution units can have some amount of hazard detection)
- Every unit has local buffers (once read, don't need to read again).

* Similarities

- Instruction issue in order.
- Issue stalls for structural hazard.
- Read Operand & Write to Registerfile can be out of order.

Example: [WAR]

add R1, R2, R3
sub R2, R4, R5

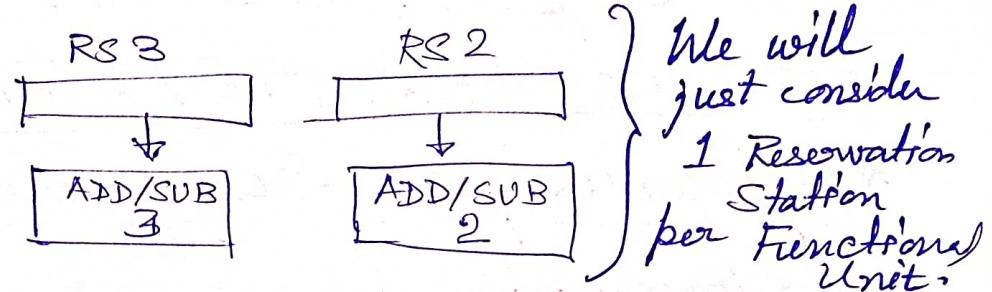
R2

add can read R2 into local buffer (at the very beginning)
& sub can write to R2 ~~but~~ anyway, anytime! (yes!).



Diagram & beyond.
→ Slides

Reservation Stations cannot communicate with each other except CDB (central data bus).



* Steps of Tomasulo! (check slides)

1. Issue (or Dispatch):

- Perform renaming to avoid WAW/WAR
- When things go to reservation stations, register names have gone.

2. Execute:

- When operand becomes available (check by monitoring CDB), store these in reservation station
- Execute when all operands are available. (This delay in execution until operands are available, saves from RAB avoids RAW hazards).

3. Write Result

:

In Tomasulo, issue stops only on structural hazard (not even when WAW hazard is possible).

More slides.

"ବେଳିଯେ ଗେହେ ହୁବି ଆଜାର କାହାର ପାଇଁ
ଆଜାର କୁଳାତ୍ମି ପାଇଁ ପତନ
କରାନ୍ତେ ପାରିବି ଆଜାରୀ!"

Branch Target Buffer (BTB)
on slides.

<u>Branch addre.</u>	<u>Branch target addr.</u>	<u>Instre @ target</u>
A1	A2	add \$1, \$2, \$3

আস্বিকে মুক্তিবাহু চেষ্টা করিব। (জনজাত্যু গু
রুমা কেন কে আনে!) শিক্ষাজ্ঞলোর উপ
ন্মুক্তিয়া এবং অনুচ্ছে আরু আলোচন
অভিযান জামিয়ে উঠিবে।

- In classical 5 stage pipeline, write back is in order.
- Commit of Tomasulo with speculation is also in order.
- lw/sw latencies are in MEM stage.

stall: issue stops (or in classic 5 stage pipeline, IF stops).

USE AFTER LOAD

We consider the following instruction.

		1	2	3	4	5	6	7	8	9
lw	x_1	$0(x_2)$	IF	ID	EX	MEM	WB	-		
add	x_4	x_3	x_1	IF	ID	stall	EX	WB		
sub	x_7	x_5	x_6	IF	stall	EX	ID	EX		
mul	x_{10}	x_{11}	x_{12}				IF	ID		

BRANCH HAZARD (RISC-V), (Branch ^{cond.} eval: EX)

Instn #	^{clock}									
	1	2	3	4	5	6	7	8	9	10
i beg x1 x2 for	IF	ID	EX	MEM	WB					
pred: i+1 add . . .	IF	ID		nop	nop	nop				
NT { i+2 sub . . .	IF		nop	nop	nop	nop				
for: i+300 add . . .		IF	ID	EX	MEM	WB				

{ only prediction

finishing in cycle 8 instead of 6

জ্ঞান আগবং, কেবল অটো কুক,
শীর রেখ করে দেখ-কুক !!

(Side again)

Multiprocessor Computers and Thread Level Parallelism

Instruction level parallelism

[Adv] throughput increase, transparent to programmer.

Disadv

* (in slides)

We must look at other sources of parallelism.

Parallelism exploited by Software

Data-level parallelism: Operating on several data items at the same time, in the same task.

Task-level parallelism: create & operate multiple tasks (independently or cooperatively)

Parallelism exploited by Hardware

Instruction level parallelism (ILP): DLP in modest level

Vector processors/ GPU's/multimedia instruction set: DLP

Thread level parallelism: DLP/TLP

Request level parallelism: DLP/TLP

Flynn's Taxonomy: (Slides)

SISD/SIMD/MISD/MIMD
single instr. stream
single data stream

(slides)

(More slides T-T)

(more slides)



If a local processor wants to update contents of a local memory, it sends a signal to others and the others invalidate in their own cache. The ~~original~~ initial local processor becomes the owner of the block.

If a processor finds a block to be invalid, it sends a signal, asking for that block. If someone (the owner) sends the block — great! Otherwise it takes from LLC (last level cache)

(MS I in slides)

"Syllabus is the entire syllabus, stress will be given on stuff taught after midsem"

Hardware Support for Synchronous Primitives

"Suppose Processor P1 does this in an infinite loop & P2 does that in an infinite loop"

P1: `while (1) { count++; }`

P2: `while (1) { count--; }`

variable count is shared between P1 and P2. Both can R/W. So basically, the following happens ($\&A$ be address of count)

P1

<code>lw</code>	X_1 , A
<code>addi</code>	X_1 , 1
<code>sw</code>	X_1 , A

P2

<code>lw</code>	X_2 , A
<code>addi</code>	X_2 , -1
<code>sw</code>	X_2 , A

Ideal case:

At time $t=t_0$, $count == 5$

At time $t_1 > t_0$, P1 executes, $count \leftarrow 6$

At time $t_2 > t_1$, P2 executes, $count \leftarrow 5$

this is what is expected; but here we may the (incorrect) assumption, that P1 executes entirely, then P2 executes.

Let us see what can happen in reality.

Possible Actual Execution Sequence:

P1: lw X1, A // $X1 = 5$
P1: addi X1, 1 // $X1 = 6$
P2: lw X2, A // $X2 = 5$
P2: addi X2, -1 // $X2 = 4$
P1: sw X1, A // count = 6
P2: sw X2, A // count = 5 !!! aaaaaa!!!

This is called a data race, where data depends on the actual sequence these set of instructions are executed.

"To avoid data races, you must incorporate synchronization in your processes."

The framework of a process:

Process P:

```
do {  
    ...  
    [entry section]  
    critical section; // Process potentially performs  
    // Read/Write operations on  
    // one or more shared  
    // variables in memory  
    [exit section]  
    ... (remainder section)  
} while (1);
```

"ও মন, কখন কুরু কখন রে শোব, কে আনে
এ রাজিকবের ঘোলা রে মন, রাজিকবের ঘোলা রে মন
যাব ঘোলা হয় কে আনে!"

Simple Hardware Based Solution

This requires atomic operation (basically assume that the instruction set implements the same thing without interruption).

"This assumption is really aggressive & difficult to achieve in modern processors."

~~do~~

The Process framework becomes:

```
do {  
    :  
    [ acquire lock ] //atomic, 'lock' is shared  
    critical section  
    [ release lock ] //also atomic  
    remainder section  
} while(1);
```

There are two methods to achieve this

1) TestAndSet Atomic Operation:

```
bool TestAndSet (bool *target) {  
    bool orig = *target;  
    *target = true; //causing it to spin  
    return orig;  
}
```

"orig probably means 'original value' and hence not $\frac{2}{3}$ rd of RVS's initials."

We can use it as follows:

```
bool lock = false; //global variable shared  
//across process
```

The 'lock' shared variable is initially false. For each process, we have the following fragment of code:

```
do {  
    while (TestAndSet(&lock)) {  
        ; //do nothing & spin  
    }  
    //execute critical section  
    lock = false;  
    //execute remainder section  
} while (1);
```

"अन्य सिर्फ अपनी, नहीं अपने फ्रेंड्स के लिए रुकावा, लेकिन अपने..."

TestAndSet returns false the first time, & hence the first process gets to execute critical section (say, P1). Others who want to enter in their own critical sections, wait by spinning until P1 releases the lock by lock = false instruction).

This TestAndSet was available in many old ISA's as a single instruction. It is no longer popular due to ~~the~~ the fact that two memory accesses are required.

2) Atomic Swap/Atomic Exchange : (EXCH)

```
void swap (bool *a, bool *b){  
    bool temp = *a; //load, a in mem  
    *a = *b; //store, b is in reg.  
    *b = temp;  
}
```

We use this in the following way:

```
bool lock = false;
```

Each process executes the following code segment:

```
do {  
    bool key = true; //local variable  
    //mapped to register  
    while (key) {  
        swap (&lock, &key);  
    }  
    //execute critical section  
    lock = false; //unlock  
    //remainder section  
} while (1);
```

Here basically, the key (of at most one function) captures the lock; others wait, till the key releases the lock.

Again, this is too complicated a design to implement as a single instruction for modern day processor.

We will look into an alternative which kind of takes care of the atomicity assumption. We will look into a special load and a special store. The design is such that if these two instructions are not executed in a ~~non-interrupting~~ the non-interruptive way, there will be a mark. We can understand that something went wrong and we will consider this to be 'lock not acquired'.

"This is a defensive measure & there is nothing better that we can do."

We will have a look how such a pair of instructions can simulate EXCH.

Special Set of Instructions

1) Load Linked (LL) / Load Reserved (LR) :

lrd or. lri.w.

$lrd\{d/w\} \langle Reg \rangle, \langle$ Address of memory contained in a register \rangle

example: lrd.w R1, \$R2
or, lrd.w X2, (X1)

There is a special (non-architectural) register "Link Register" / "Reserved Register" (R₁)

∴ for lrd.w X2, (X1), the following are done

$X2 \leftarrow \text{MEM}[X1]$ and $R_1 \leftarrow X1$.

In other words, "loc has created a reservation on $\text{MEM}[X1]$ ". This is because of the following reason.

If an interrupt (including context switch) happens after loc, for the cache block to which $\text{MEM}[X1]$ is invalidated, R_{oi} is set to 0 forcibly. Because $R_{oi} \leftarrow 0$, happens by force!

2) Store Conditional (SC)

$\text{sc} \cdot \{d/w\} \langle \text{Reg} \rangle, \langle \text{Address in Reg} \rangle$

Example: $\text{sc} \cdot w X3, (X1)$

This is very complex ... T-T

[Step 1] SC checks value of R_{oi} .

~~if ($R_{oi} \neq \langle \text{Address-argument} \rangle$) { // arg 2 of sc
X3 \leftarrow Non-zero value; // sc failed
else { MEM[X1] \leftarrow X3; //
X3 \leftarrow 0; } // sc passed //
} for sc-w X3, (X1)~~

~~if ($R_{oi} \neq \langle \text{Address-argument} \rangle$) { // 2nd arg of sc
X3 \leftarrow Non-zero value; // sc failed
} else { // sc passed
MEM[X1] \leftarrow X3;
X3 \leftarrow 0;
}~~

[Example] Atomic EXCH using LL/SC

EXCH: $\text{MEM}[X_4] \leftrightarrow X_4$. atomically

LL/SC:

try:	mv	X_3, X_4	// $X_3 \leftarrow X_4$
	lr	$X_2, (X_1)$	// $X_2 \leftarrow \text{MEM}[X_1]$, $R \leftarrow X_1$
	sc	$X_3, (X_1)$	// If sc succeeds, // $\text{MEM}[X_1] \leftarrow X_3$ // $X_3 \leftarrow 0$
	bnez	X_3	try again if sc failed
	mv	X_4, X_2	// $X_4 \leftarrow X_2 (\text{:=MEM}[X_1])$

If there was an interrupt in between lr & sc, I would try again.

So in five instructions, we can 'effectively' implement an atomic exchange

"*तुम अपनी जगह बदल दूँ तब तुम रखना कहु तोह!*"

We will now look into the atomic fetch & increment operation (done using LL/SC).

[Example] Atomic Fetch & Increment using LL/SC

We want to achieve the following atomically

$$\text{MEM}[X_1] \leftarrow \text{MEM}[X_1] + 1.$$

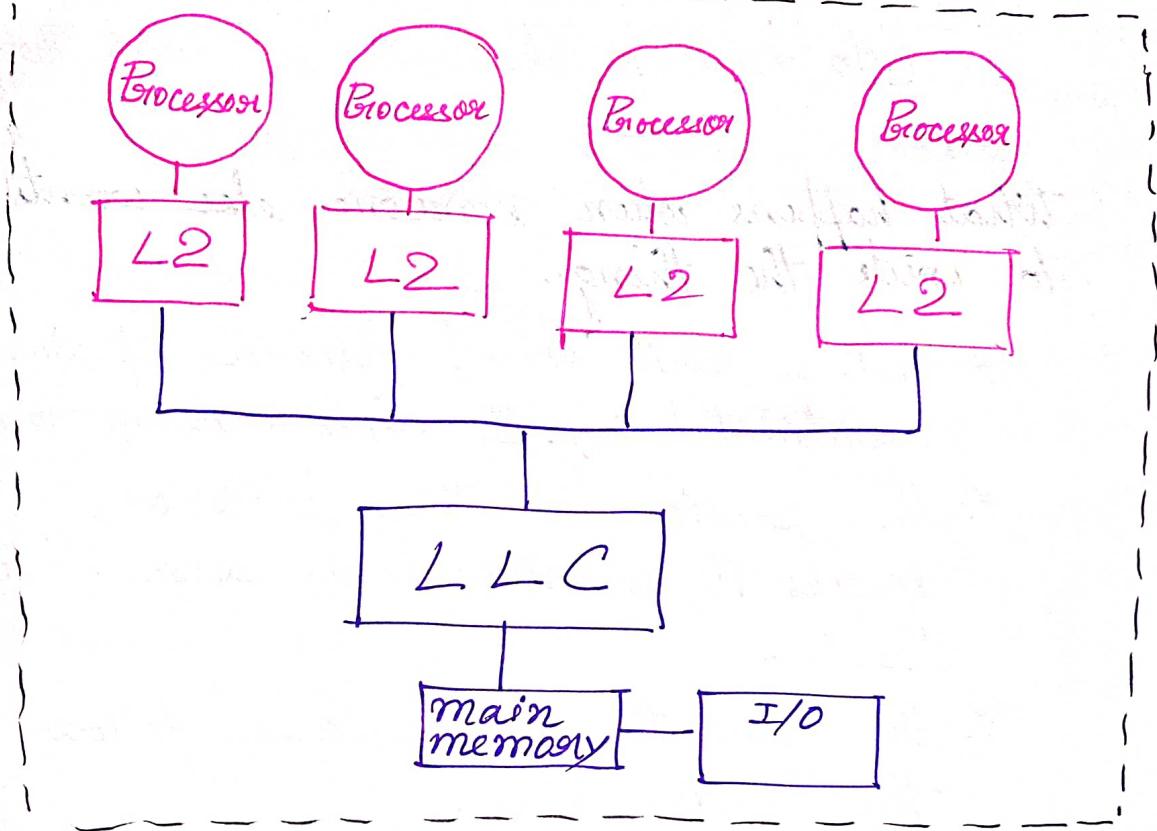
try:	lcr	$X_2, (X_1)$
	addi	$X_3, X_2, 1$
	sc	$X_3, (X_1)$
	bnez	X_3, try

Basic Rule of Thumb:

- * Completely avoid memory access between lcr & sc.
- * Try to add only simplest register operations between lcr & sc. Stuff like exceptions (divide by 0) is problematic.
- * Best case is to put nothing in between lcr and sc.

IN THE EXAM: You may be asked to write other atomic operations with these lcr/sc command given their intended behaviour

(Back to MSI slides [T-T])



When a processor wants to write something:

- * send invalidate-message to everyone (broadcast)
- * other processors invalidate the block immediately.
- * Hardware based mechanisms are there to arbitrate and break ties
- * After ~~is waiting~~, this, the writing processor marks this cache block as "modified" and claims ownership.

- * Writing processor writes to their L2 (without updating LLC or other ~~L2~~ L2 caches).

(if it immediately modifies to LLC, it is a write-through scheme).

(No modern ~~is~~ caches are write-through)

What happens when someone else wants to write the thing.

- * Gets a write miss! (because it was invalidated it). It sends a message to everyone
- * The previous writing processor, marks it as shared as sends it to LLC.
- * The present writing cache follows the previous step.

Optimization (MESI):

E: (Exclusive): only the owner has this block and is not dirty.

[LLC maintains directory ~~who~~ of which all L2 caches have the block]

"SHARED" state is Read Only.

"In the MSI state diagram (for one block), the 'exclusive' is actually 'modified'."

এখনে এত কিছি আছি, এন এন আরও কোথাই আছি

[Sidle with 'blue' & 'orange' table]

"LLC writing takes a lot of time."

"On a 'modified' block request, the owner writes to LLC, but also directly provides it to the seeker."

আবি, এত বলি আবি
এই প্রক্রিয়াটি কীভাবে হয়ে যাবে এবং
কে কে করে কীভাবে হবে

"Caches are well behaved"

"All actions on cache blocks are done by local cores in response to some messages!"

"Processors keep on listening!"

আবি আকারে শার্টিয়া কোন
ওবেছি ওবেছি আবারেই কোন,

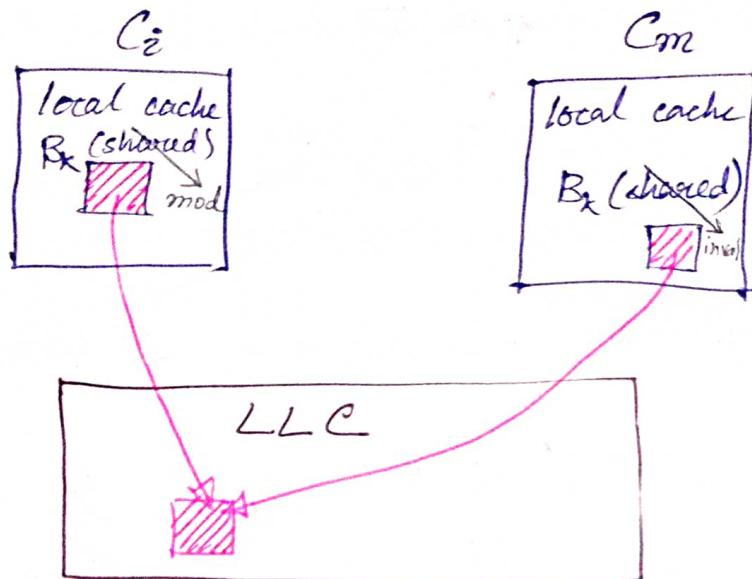
"This micro architecture is confusing!"

কোথাও আবার হারিয়ে যাওয়ার লেই আবা

বনে বনে

Coherence Miss:

সমিরণ কৰিবলাব পৰা, মুলি দিয়ে আৰু হৰি



- * C_i wants to write $W_t \in B_k$
- * C_i invalidates (requests to invalidate) B_k in C_m
- * If C_m wants $W_t \in B_k$
 - ↳ Read miss [True Sharing Miss]
- * If C_m wants $W_t \in B_k$, $W_t \neq W_t$
 - ↳ Read miss still happens.
[False Read Miss]

"Now that you've grown up
you understand there is no Santa!"

Directory Protocol

Based \rightarrow Coherent

[slides]

For each directory LLC block, store in directory:

- A bit vector with (#cores) bits.
- Denoting if the core has a private copy in the L2 cache.

• Coding logic

Go through details of this

selected mechanism

• 12 cores

• 16 LLC blocks

• 16 bits

• 16 bits for each block

• LLC block address

• 16 bits

• Cache miss after cache hit will be detected by comparing

• LLC block address with LLC block address

• Cache miss after cache hit will be detected by comparing

• LLC block address with LLC block address

Back to critical section

Solution to Critical Section Problem :

Most common solution:

do {

acquire lock // atomically (try to)
// acquire lock, if no
// cannot acquire lock
// loop ("spin").

{critical section} // only one process
// can execute its
// critical section at
// a time

release lock }

{remainder section}

} while (1)

We will assume "effectively atomic"
instruction for EXCH

We will also ~~base~~ base our discussion
considering the shared multi-core
"write-invalidate" protocol!

Solution 1. (Classic Solution)

Assume that the lock variable is located at $O(X_1)$, in memory. Assume $X_0 = \$0$, contains the value 0 always.

// entry section (can enter if lock == 0)

```
addi    X2, X0, 1
lockit: EXCH  X2, O(X1) // try (pseudo) atomic EXCH
bnez   X2, lockit // if lock not acquired, spin
```

Now, say, P1 successfully performs EXCH.

* P1 "holds" lock in its local cache, ownership.

* P2 "wants" lock.

- It will exchange X_2 to $O(X_1)$.

- lock will not be acquired, (as both X_2 and $O(X_1)$ hold value of 1)

- However, this writes 1 to $O(X_1)$ anyway

- Other caches (of P3, P4, etc. and P1) are invalidated

This causes unnecessary crowding in the common bus.

Solution: Change over - Only perform write when lock value is 1

So the issue was: Every process/thread, which performs a spin to acquire the lock, must send an invalidate message on the bus following its EXCH. However such invalidates are useless to the originating process/thread, if the lock value in the memory is already 1. This happens ~~due~~ on every spin.

Solution 2.

Do not send invalidate messages unnecessarily. First check if lock variable is zero. Only if its is zero, perform EXCH

```
lockit : ld. X2, 0(X1) //lock at 0(X1)
          bnez X2, lockit //lock not available
          addi. X2, X0, 1 //X2+1
          bnez X2, lockit
          EXCH X2, 0(X1) //same as before
          bnez X2, lockit
```

"Invalidate messages are not eliminated, the volume just decreases by a lot."

Notice that EXCH performs load (of the lock variable) anyway. So why not use that, instead of keeping another Id?

Solution-3. (Optimized Solution)

We would use the ~~lcr~~ in the ~~EXCH~~ pseudoatomic operation to read the value. Another way of looking into it is that, we slightly modify EXCH.

```
lockit : brc    X2, 0(X1)
          bnez  X2, lockit
          addi   X2, X0, 1
        bnez  X2, 0(X1)
        sc    X2, 0(X1)
          bnez  X2, lockit
```

could had been dangerous if this gets "memory access violation". (Won't happen here, cause lockit is right here, i.e. trustable)

We now hand-simulate this on 3 processor (example in text-book).

Hand-Simulation:

~~Step #~~ ~~P0~~ ~~P1~~ ~~P2~~

তার যা হিল শুরু শীর্ষ
প্রেসেসর তার নাও

শান রেখে, শান রেখে

তার হই প্রেস শীর্ষ

শান শুরু আবিষ্য দে

অবেহিল শুরু শীর্ষ

[P.T.O]

Step

P0

has lock = 1

in L2 of P0

(in critical section)

P1

spins to

check lock == 0,
get read miss

P2

spins to

check lock == 0,
get read miss

2.

lock \leftarrow 0

(Don't care)

invalidation

same as P1

but: gets the
bus, read miss
gets hold of
lock from P0
wait for bus

invalidation

invalid in P2

shared in

P0, P1

P0 gets hold of
bus, sends
read miss. P0
gives lock value (=0)

P0 waits
back to LLC

P0 sends

invalidate

P2 gets hold of
bus, sends
read miss. P0
gives lock value (=0)

P2, P0 waits
back to LLC

invalidation

shared, invalid

lock ...

lock ...

lock ...

lock == 0, but succeeds.

wait for bus

test succeeds

wait for bus

test succeeds

execute read

receive read

miss from P1

send to P1

execute read

receive read

miss from P1

send to P1

execute read

receive read

miss from P1

send to P1

execute read

receive read

miss from P1

send to P1

5.

Windows 8.1

into an iron

4. <don't care>

test succeeds
(lock == 0)

shared,
invalid, shared

meh...

5. <don't care>

writ for bus
tried to execute
swap, gets cache
miss (as shared
is read-only)

P1 got the
lock value.

6. <don't care>

Execute swap,
but got ~~lock~~
cache miss as
P2 has invalidated.

P1 invalidates
from P2.

P2: Modified
P0, P1: Invalid

Performs swap
after changing
local cache to
MODIFIED.
set lock $\leftarrow 1$

P1 sends
invalidation
over bus.

P1: modified
P0, P2: invalid

P1: modified
P0, P2: invalid
continues with acc
swaps on local copy
depends on local copy

meh...

7. <don't care>

8.

— END OF COURSE —

বিজ্ঞান বিদ্যালয়