

## L8-addressing modes:-

five common addressing modes

register names or memory locations for operands.

### 1) Immediate:

addi \$t0, \$t1, 5

### 2) Register:

add \$t0, \$t1, \$t2

### 3) Base address:

lw \$t1, 4(\$t2)

address in  $t_2 + 4$

is our required!

lw \$t1, (\$t2)

### 4) PC-relative:

beqz \$t0, goEnd

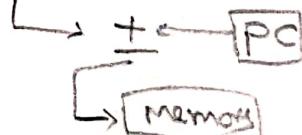
isn't appended?

assembler will calculate the rel. pos.

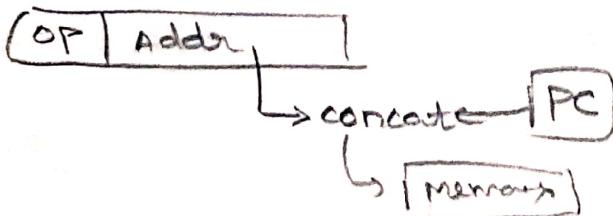
Nah!

$PC += \underline{Addr}_{18}$

[OP | rs | rt | address]



### 5) Pseudo-direct:



#### j forLoop

assembler will calculate this label's address.

J  
JAL  
JR } will not give relative to PC.  
will give only

$$PC = PC_{31:28} :: Addr_{28}$$

BEQZ

BEQ

BNE

BGTZ

} will give relative to PC

$$PC += OFF_{18}$$

## L9 - ISA and microarchitecture:

- \* Microarchitecture is implementation of ISA operations.  
Like (for "add 11111"; should we) ripple (or) carry-look ahead?
- company chips differ in their microarchitecture; but respect the ISA (x86, ARM, MIPS, RISC-V)...
- \* ISA provides:
  - Instructions
  - Opcodes
  - Addressing modes
  - Types & format
  - Registers
  - Access Control : user/OS memory
  - Address space
  - Addressability
  - Alignment
- \* ISA must satisfy requirements of software
  - assembler
  - compiler
  - OS.
- \* rest of CS305 after ISA will be caches, memory controllers, Branch predictors, prefetchers.
- \*
  - # registers → ISA
  - # cycles to register access → Microarchitecture
  - # register width → ISA
  - # instruction to access memory → ISA
  - # cycles to read memory → Architecture.

## \* Where to place our ISA specifications?

- Close to high level language? → small semantic gap,  
complex instructions (an opcode for quicksort)
- CISC → complex instruction set computer.  
like x86. (x86 is CISC with RISC mysteries)
- Closer to hardware? → large semantic gap.  
Simple instructions.
- RISC → reduced instruction set computer.  
like MIPS

## L10: world of ISAs

x86: Intel, AMD : PCs, Servers

ARM: Qualcomm, Apple, Samsung : Mobiles  
↳ Advanced RISC

RISC-V: opensource ISA. meaning, no patent filed.

## \* Subtle differences in x86 ISA:-

1) x86 arithmetic/logic instructions:

one operand should act as source and destination.

add \$t0, \$t1

to ← to + t1

2) one operand can be in memory! not allowed for MIPS

Programmer: Awesome!

Instruction: Noo! Size is too large

add \$t0, (\$t1)

3) no fixed length for instructions. (Booo!)

Now code is smaller, but decoding is tough!

<trivia in slides>

## Fallacies:

- CISC instructions provide higher performance.  
things are not that straight forward
- Assembly code provides higher performance.
- \* RISC was motivated over CISC by memory stalls.  
when fetching memory, no other arithmetic logic work can be done in complex instruction
- But RISC enables the compiler to optimize the code by fine-grained parallelism to reduce stalls.

\* ARM has some thumb instructions to reduce code density.

normal (32bit) → Thumb (16 bit)

less code written

(50%) more

35% smaller code  
50% more opcodes

more opcodes

more memory access overhead 2x-3x

more complex

instructions need more time

more overheads

more time

more power

more time

→ better compilation (compile time vs. execution time)

Code size A vs. B

Code size A vs. B

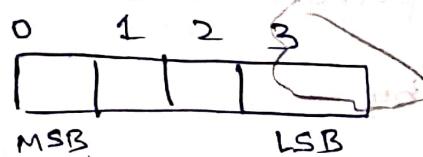
(compared to A) Code size A vs. B

## L11- Endianess and Alignment:-

Endianess : Byte ordering within a word.

1) Bigendian : word address = MSB of word

MIPS



Big end

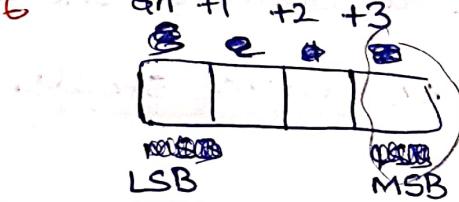
int A = 0xABCD1234

EA: [AB] [CD] [12] [34]

"The usual"

2) Little Endian: word address = LSB of word

X86



Little endian

EA: [34] [12] [CD] [AB]

Eg:

unsigned int i = 0x12345678;

char \*c = (char \*) &i;  
1 byte length  
word address.

printf ("%d", \*c);

Littleendian : 78 ✓

Bigendian : 12

→ Alignment:-

MIPS does not allow unaligned access.

lw can't use. but lh, lb okay!  
load word  
load half  
load byte

X86 allows.

Alignment for fast data fetching!  
otherwise 2-times memory requests

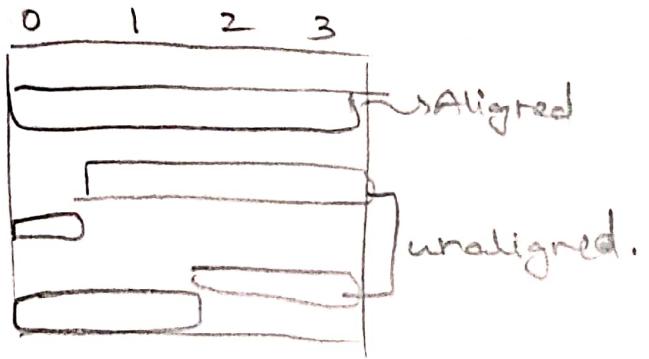
\* compiler generates aligned/unaligned access....

if "lw A" then  $A \% 4 == 0$

if "lh A" then  $A \% 2 == 0$

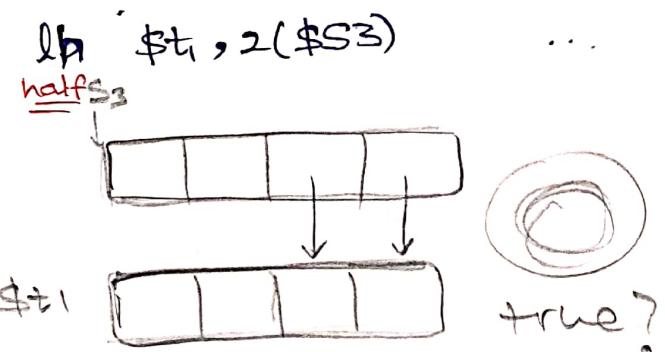
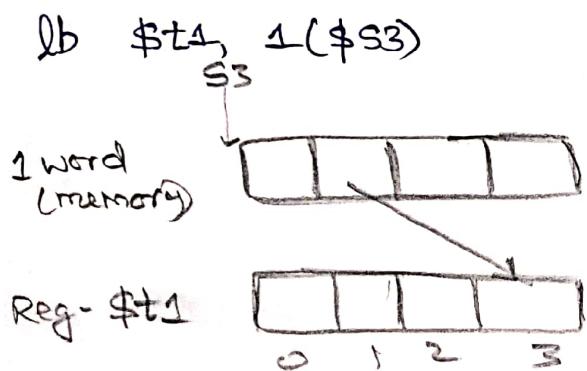
if "lb A" then  $A \% 1 == 0$  (always true)

)  
Load byte



→ memory instructions & Alignment network :-

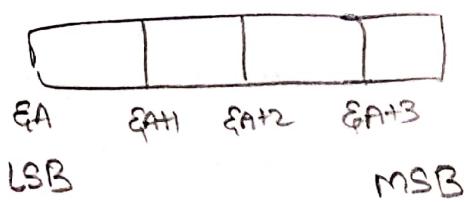
LOADS & STORES need a Alignment network that makes word/half/byte Scattered data is written aligned.



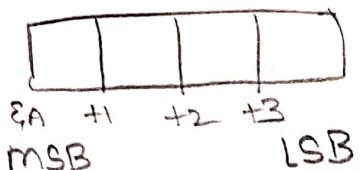
\* Quiz question based on endianess:-

cool now!

in little endian:-



Big endian:-



## L12: Single cycle CPU (processor 101)

→ so, instructions need to be ~~size~~ enough.

LIB

- \* All operations in single cycle each.

- \* Hence, clock cycle will be dependent on the longest instruction delay.  
(critical path)

- Datapath & Controlunit.

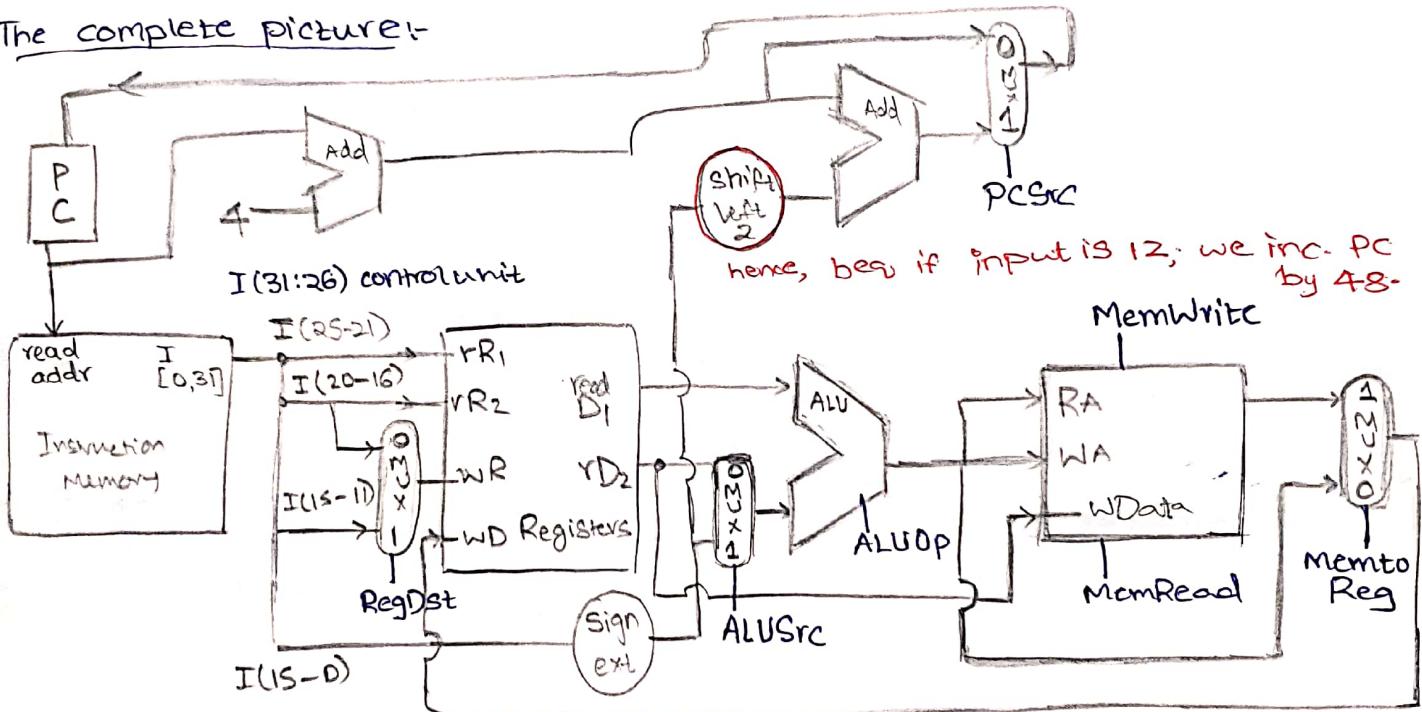
rate-determining  
path in  
a reaction

"Anything that Stores data or operates on data".

"manipulates data" ☹

"Guides" data flow.

- ## \* The complete picture!



- \* ALUsrc asserted: lower 16 bits of instructions

non-asserted! output from 2<sup>nd</sup> reg.



\* PCsrc: asserted: branch  
                  deasserted: PC+4

- \* RegDST: asserted: rd  
deasserted: rt

MIPS :- add \$dest, \$r1,\$r2.

this order different from  
MIPS.

DP	RS	Rt	Rd	Shamt	fun
----	----	----	----	-------	-----

Source → RS      temporary → Rt      destination → Rd

\* Shamt - Shift amount

in MIPS, all R instructions have opcode 0 (P&H book ref)  
but different funct.

Sll is R type!!!

(shift left  
logical)

so, ALUSrc not used.

Sll \$t0,\$t0,1 } ALUSrc is 0.

Shamt. quiz error

## L-14: Iron law:-

$$\text{Time/program} = \frac{\text{Instructions}}{\text{program}} \times \frac{\text{Cycles}}{\text{Instruction}} \times \frac{\text{Time/cycle}}{\text{microArchitecture}}$$

- Source code
- compiler
- ISA
- ISA
- microArch.
- microArch.
- Technology.
- Technology.

## L-15: 16: 17 : empirical evaluation:-

\* we see performance as,

- Latency (execution or response time)  $\rightarrow$  additive
- Bandwidth (not additive)

$$\boxed{\text{performance} = \frac{1}{\text{latency}}}$$

no. of tasks / unit time

\* latency lags bandwidth. Bandwidth hurts latency.  
(but not decreases!)

\* better latency helps bandwidth. B.W. usually hurts latency!!

## → Amdahl's law:- (common case; make it fast)

$$\text{Speedup} = \frac{\text{Overall Time Execution}}{\text{Enhanced Time Execution}} = \frac{f_{\text{enhanced}}}{(1-f_{\text{enhanced}}) + \frac{f_{\text{enhanced}}}{\text{Speedup}_{\text{enhanced}}}}$$

- Small Speedup on large fraction  $>$  Large speedup on small fraction

## → Evaluation:-

\* how to compare proc. A & B by running programs, which programs?

### Benchmarks:

- SPEC CPU 2017
- Cloudsuite
- PARSEC
- MobileBench
- Geekbench
- ML perf
- Graphperf

### Pitfalls of benchmarks:-

- not representative of all features.  
if workload is I/O bound, CPU SPEC is useless.
- Too old benchmarks.

> The ones that you care & ought to occur in processor usage.

\* Okay! let's create a chip & run these benches.

Pfft! lol!  
use simulators.

- Functional Simulator:-

- Used to verify correct.
- Worthless to check performance.

- Performance Simulators:-

i) Trace-driven:

ii) Execution-driven:

\* So, pick a relevant Suite of benchmarks

1) measure clock rate of each program

2) Summarize performance using

AM      (AM)      Avg { }

AM if 100  
0.01  
is present;

→ Energy and power:-

gives unfair results

\* Power efficiency = performance / watt

Energy efficiency = performance / Joule

\* Dynamic power; static power

during  
running  
instructions even  
when  
transistor is off

## L18- Multi Cycle CPUs

\* in single cycle why every process should wait for longest proc-time!  
So load design!

### MultiCycle processor!

Now, even though longest =  $T_{proc}$ , now CPI increases.  
others  $\leq T_{proc}$ . even though cycle time decreases drastically.

\* can we have both

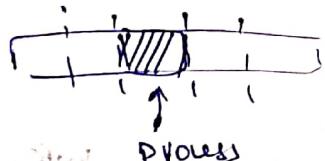
CPI = 1, C = 2ns?

Pipelining! yea.

Then, we are leaving some part of processor empty, when running on instruction.

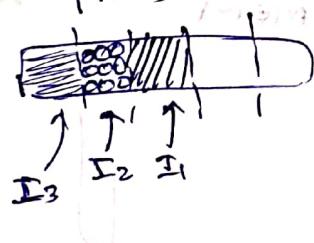
fill those too.

rm:



process

pipelined:



I<sub>3</sub> I<sub>2</sub> I<sub>1</sub>

Eq:

normal covid schedule

S<sub>1</sub>: 15 min      one after another completes  
 S<sub>2</sub>: 10 min  
 S<sub>3</sub>: 5 min  
 S<sub>4</sub>: 20 min  
 S<sub>5</sub>: 10 min

1 person per 1 hr

pipelined:

S<sub>1</sub>: (20 min) Q      waiting after S<sub>1</sub>  
 S<sub>2</sub>: (20 min) Q      waiting after S<sub>2</sub>  
 S<sub>3</sub>: (20 min) Q      waiting after S<sub>3</sub>  
 S<sub>4</sub>: (20 min) Q      waiting after S<sub>4</sub>  
 S<sub>5</sub>: (20 min) Q      waiting after S<sub>5</sub>

in a line proceed to S<sub>1</sub> S<sub>2</sub> S<sub>3</sub> S<sub>4</sub> S<sub>5</sub>

3 per hour.

even though Single person latency = 100min.

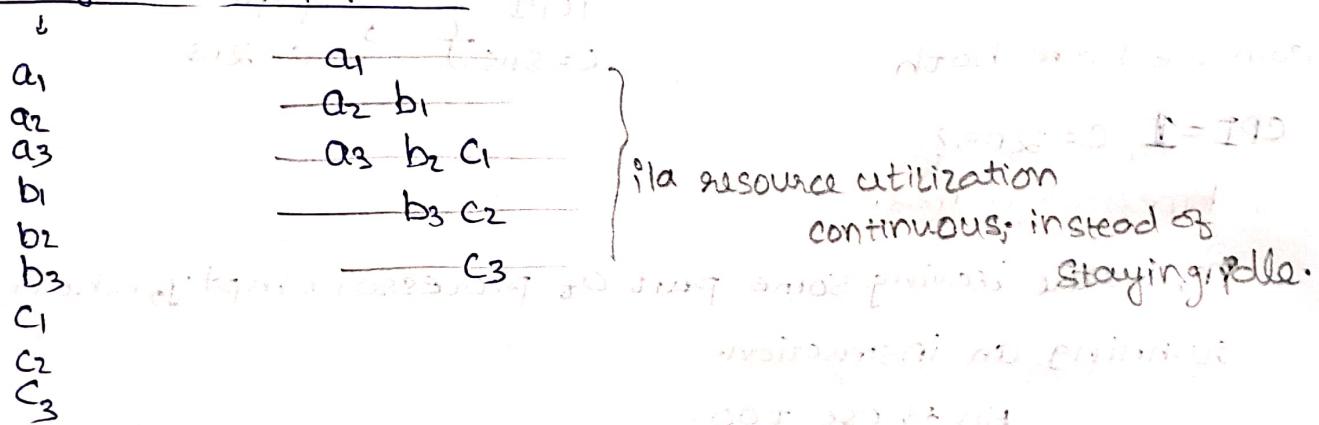
## L19, 2021 - Pipelining:-

\* single cycle CPI: 1 cycle time: long

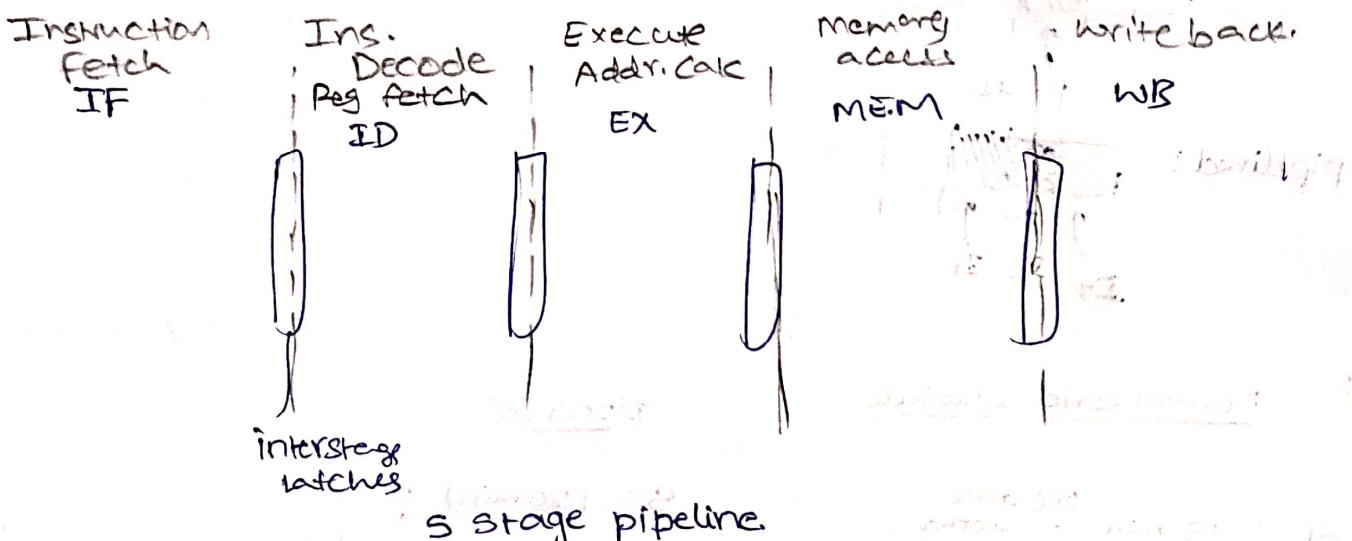
Multicycle CPI  $\geq 1$  short cycles

Pipelined CPI: 1 short cycles (but throughput increases latency worsens)

\* multicycled vs. pipelined:-



→ vanilla s-staged pipeline:-



\* For a K-staged pipeline, with N instructions

1<sup>st</sup> instruction: K cycles

next instructions: 1 cycle each.

$$\text{Net} = K + N - 1$$

$$\therefore \text{avg. CPI} = \frac{K + N - 1}{N} = 1 + \frac{K-1}{N}$$

\* Single cycle design vs pipelined.

	ifetch	Decode	Execute	Memory	Writeback	Total time
Load	200	100	200	200	100	800
Store	200	100	200	200	-	700
Add	200	100	200	-	100	600
branch	200	100	200	-	-	500

in single cycle:-

$$\text{clock} = 800 \text{ ns}$$

$$1. \text{ latency for the 4 net} = 3200 \text{ ns}$$

for pipelined:-

$$\text{clock} = 200 \text{ ns}$$

$$\begin{aligned} \text{net time latency} &= 1000 + 3(200) \\ &= 1600 \text{ ns} \end{aligned}$$

$\Rightarrow 2 \times \text{speed up.}$

\* if  $10^9$  instructions:  $\frac{1000 + 10^9(200)}{800 \cdot 10^9} \approx 4 \times \text{speed up.}$

*pipeline latencies may add overhead  
200 + 10 ns.*

$\approx 3.8 \times \text{speed up}$

→ Dividing datapath into pipeline:-

$$t_{IM} = 10$$

$$t_{DM} = 10$$

$$t_{ALU} = 5$$

$$t_{RF} = 1 \quad \left. \begin{array}{l} \text{combine stages.} \\ \text{NO loss in perf.} \end{array} \right\}$$

$$t_{RW} = 1$$

unpipeline      Pipeline      Speedup

e.g.

$$10, 10, 5, 1, 1$$

4-staged pipeline

$$27$$

$$10$$

$$2.7$$

$$S, S, S, S, S$$

4-staged pipeline

$$25$$

$$5+5$$

$$2.5$$



5-Stage PP

$$25$$

$$5$$

*S// for a k-stage;*

*K is max speedup.*

## L22: Instruction Pipeline Hazards:

- Structural Hazard → Anything that prevents an instruction to move forward in a pipeline.
  - Data hazards
  - Control hazards
1. Structural: two instructions want to access the same physical structure at once.
- Eg: if we have a unified memory; rather than instruction mem & data mem.

→ what about the registers?

can we read and write in same cycle?

Yes! Since the edge-triggering property.  
whatever!

- \* Structural H兹s are highly infrequent.

## 2. Data hazards:

- \* Hazards due to data dependencies.

- \* Current instruction depends on the result(data) of previous instructions.

### a) Read after Write:-

RAW

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

True dependency

### b) Write After Read:-

WAR

add R1, R2, R3  
sub R2, R4, R6

Anti Dependency

### c) Write after write:-

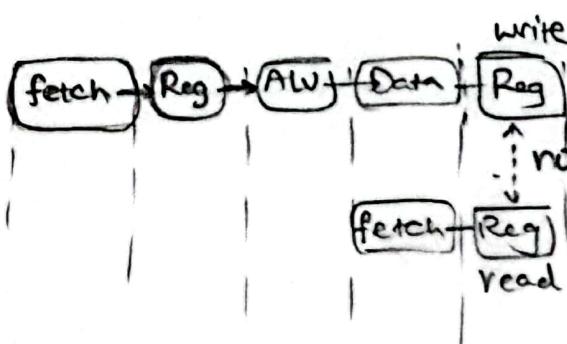
WAW

add R1, R2, R3  
or R1, R4, R5

Output dependency

not possible in  
vanilla 5-stage  
pipeline.

Slide 10:



\* If ET का point!  
Such a nuisance!

the writes happens at clock rise.  
This is reflected in read.

\* This last potential hazard can be resolved by the design of register file hardware:

We assume that write is in the first half of clock and read in the second half, so the read delivers what is written. 😊

### 3. Control hazards:-

- \* Hazards that arise from jump/branch instructions (those which change PC).

Eg:- 10: Beq r1,r3,36 goto 50

14: And r2,r3,rs X

18: Or --- X

22: And --- X

50: ✓

What to do to these 3 instructions?

- Instructions cannot move forward, & must wait for hazard to get resolved.

- Pipeline stalls.

!

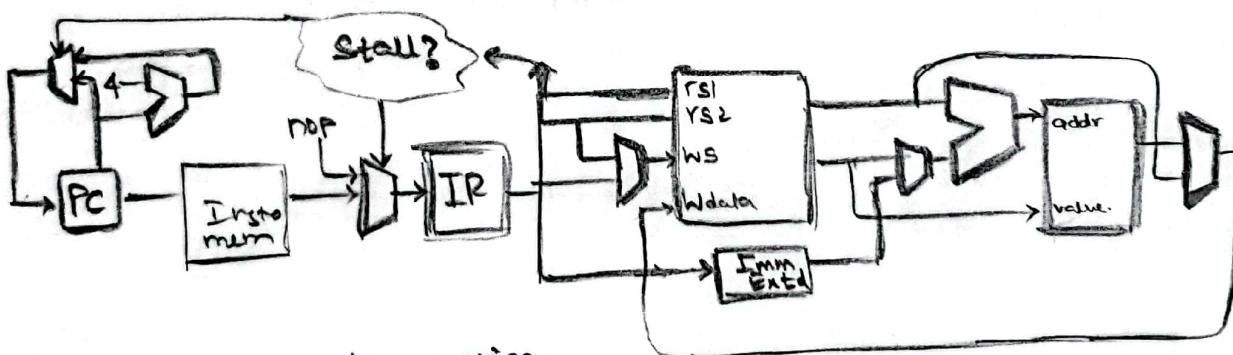
Air Bubbles in Pipeline.

- \* In a stall; PC is not updated to PC+4.

IR is not updated. (the same as prev. instruction)

deassert the control signals. ~~or~~ or ~~jmp~~ ~~for~~ ~~nop~~ instruction!

- \* Implementing a Stall:- in a stall, PC, IF/ID are frozen & ID/EX is set to all 0 controls



> don't fetch a new instruction.

> don't change PC.

> insert nop into IR. ] Compiler way  
no operations.

- \* Two ways of inserting stalls:-

1) compiler can insert nops

2) Microarchitecture has a hazard detector & gives a stall signal.

Hazard detector : ID stage  
forwarding unit : EX stage

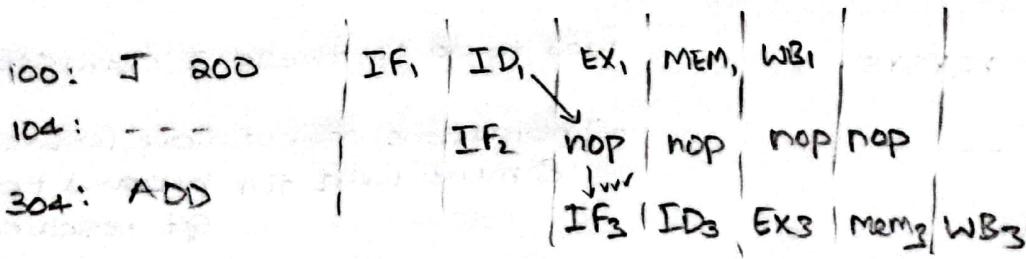
- \* Eg for NOP:

sll \$0 \$0 (in MIPS)

NOPs for control hazard.

How for data hazard?

\* Jump instruction → PC value fetched in ID stage.  
only one NOP instruction.



\* for branch instruction, 2 NOPs needed....

in EX stage we get  
next PC.

\* Hence

$$\text{Speedup} = \frac{\text{CPI unpipelined}}{\text{CPI pipelined}} = \frac{\text{CPI unpipelined}}{\text{ideal CPI} + \text{stalls/Instruction.}}$$

stalls  
affect  
speedup

④ for balanced pipeline.

(meaning; all stages  
are equal time).

## L23: Hazards detection & mitigation:-

→ Data hazard detectors :-  
conditions.

- Execute to decode:

$$\text{EX/MEM} \cdot \text{Reg Rd} = \text{ID/Ex} \cdot \text{Reg Rs}$$

or

$$\text{ID/Ex} \cdot \text{Reg Rt}$$

All

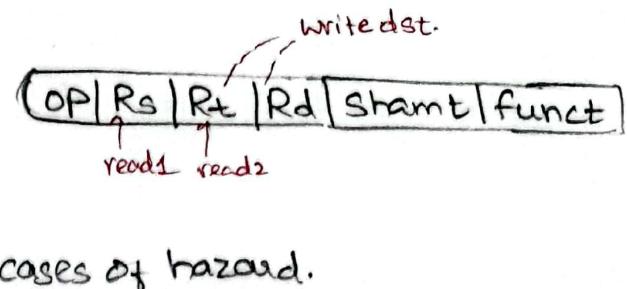
- memory to decode:

$$\text{MEM/WB} \cdot \text{Register Rd} = \text{ID/Ex} \cdot \text{Reg Rs}$$

or

$$\text{ID/Ex} \cdot \text{Reg Rt}$$

Load

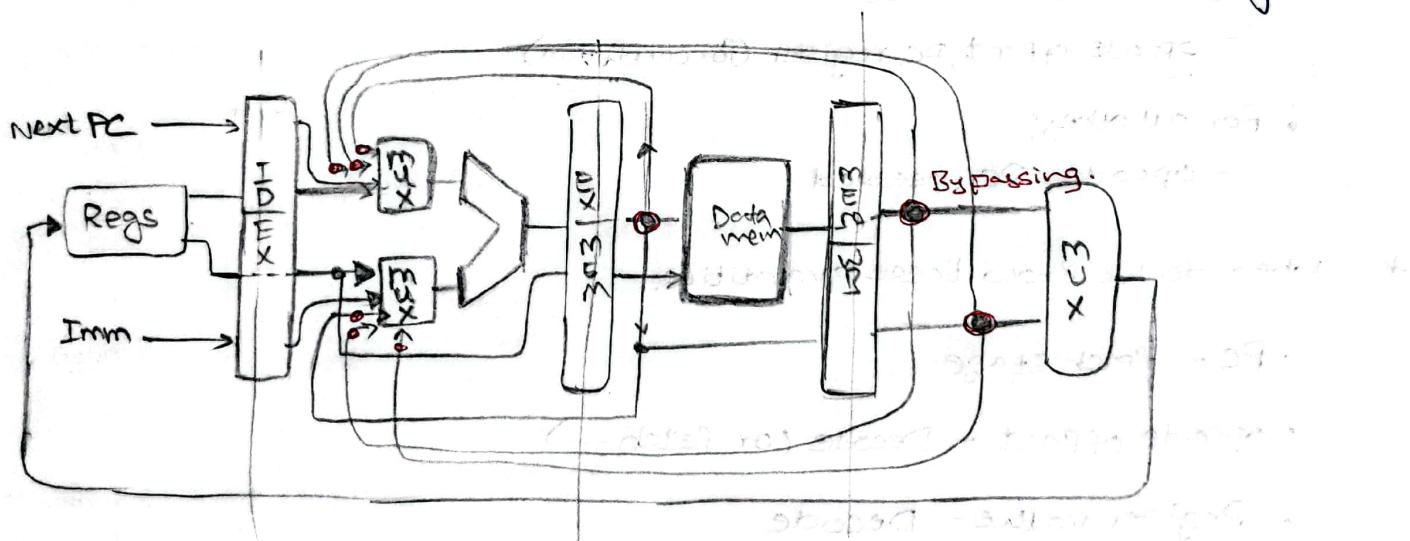


what about instructions that don't write into registers?

- Bypassing:

PPA-2: logic of hazard detector  
in 5-stage pipeline.

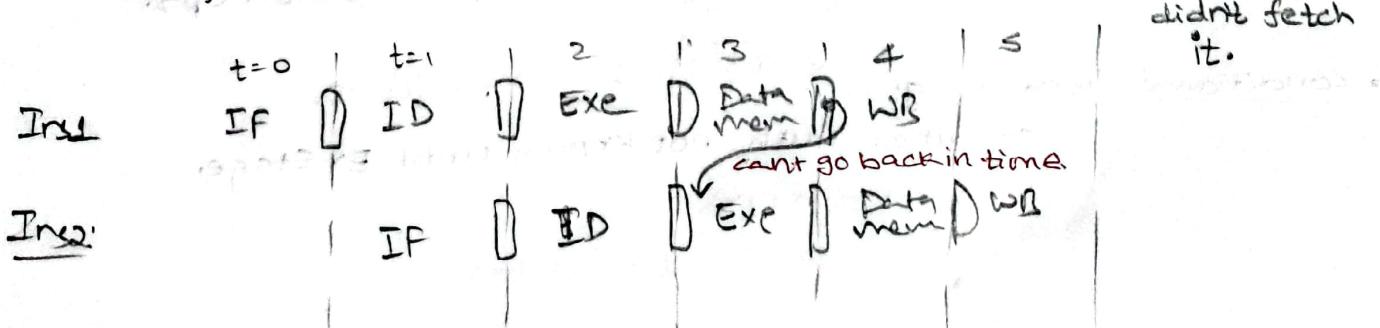
- \* Route data as soon as possible after its calculated, to an earlier pipeline stage.



Imp note!

- \* When we load something & immediately try to read; data hazard! even with bypass!

Since; when Ins2 is in Exec, Ins1 is still in DATA MEM stage.



- \* Bypass doesn't mitigate control hazards at all! only data hazards.

→ Control hazards: occurs during branch or jump.

\* When do we calculate next PC: what we need

- Jumps
  - opcode, offset, PC

- Jump Register
  - opcode, register value

- Conditional Branches
  - opcode, offset, PC, register (for condition)

- For all others
  - opcode, PC needed

\* When do we know these quantities:

- PC - fetch stage

- opcode, offset - Decode (or fetch...)

- Register value - Decode

- Branch condition ( $(rs) == 0$ ) - Execute (or decode?)

\* Speculate  $PC = PC + 4$  might work ~50%.

\* If we miss-speculate, then Kill, or insert nops.

- conditional branching:  $\rightarrow$  **2 nops** 1 nop for jump.

condition truth not known until Ex Stage.

$\langle L23, S11 \rangle$

Pipeline visualization.

- Instructions b/w branch, target are called wrong-path.

I<sub>1</sub>: 96: ADD

I<sub>2</sub>: 100: BEQZ 200

I<sub>3</sub>: 104: Add

I<sub>4</sub>: 108: SL

I<sub>5</sub>: 304: Add

	t=0	1	2	3	4	5	6	7	8	9	10	11
IF	I <sub>1</sub>	I <sub>2</sub>	I <sub>3</sub>	I <sub>4</sub>	I <sub>5</sub>							
ID		I <sub>1</sub>	I <sub>2</sub>	I <sub>3</sub>	nop	I <sub>5</sub>						
EX			I <sub>1</sub>	I <sub>2</sub>	nop	nop	I <sub>5</sub>					
MEM				I <sub>1</sub>	I <sub>2</sub>	nop	nop	I <sub>5</sub>				
WB					I <sub>1</sub>	I <sub>2</sub>	nop	nop	I <sub>5</sub>			

conditional

Branch.

EX stage  
completed!

if we jump; give 2 nops.

\* Nops do no change in datamem

SLL \$0,\$0. or registers.

\* Branches instruction; when will we know if taken or not. When will we know target address.

Inst.	Taken Known?	Target Known?
J	After ID	After ID
BEQZ	After EX	After EX

→ if we add an ALU in decode stage,

BEQZ	After ID	After EX
------	----------	----------

really?

hmm...

vanilla 5-stage pipeline do in the.

EX - we calculate.

MEM stage; we write back....

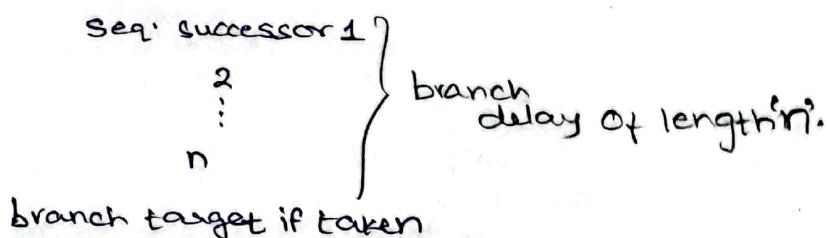
## Q24: Delayed branch and branch delay slot

old idea.

modern processors have branch predictors.

- \* Define a branch, to take place after a following instruction.

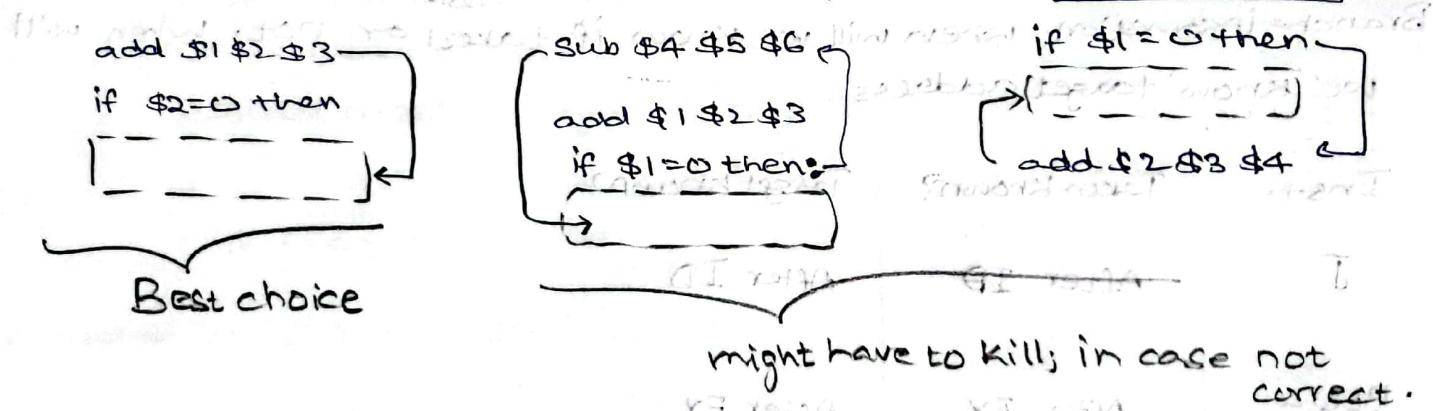
branch instruction



- \* Don't put a branch in a branch delay slot.

→ Scheduling branch delay slots:-

- a. from before branch:
- b. from branch target:-
- c. from fastthrough:-



→ Stalls and performance:-

- \* N instructions & S stalls;

$$CPI = \frac{N+S}{N}$$

- \* Pipeline Speedup <sub>new</sub> =  $\frac{\text{Pipeline depth}}{1 + \text{Stalls due to branches}}$

where Stalls (branches) = Branch frequency  $\times$  penalty.

### Data hazards

- Bypassing
- stalls (NOPs) if bypassing won't work

### Control hazards:-

- Speculate  $P_C = P_C + 4$ , Kill if wrong
- Delayed branch & branch slots.
- Branch predictors

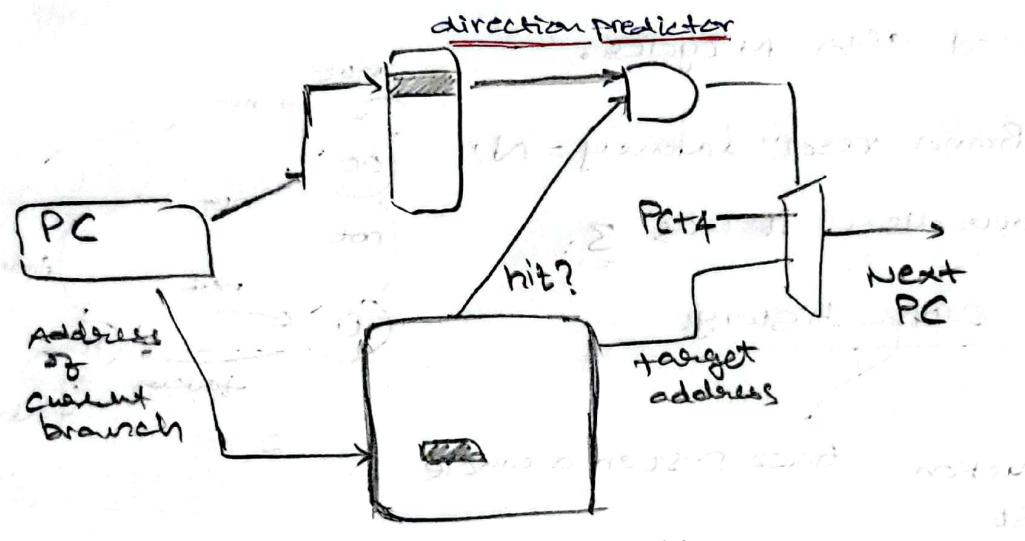
## L25: Branch Predictions:-

L26

- \* One way to mitigate control stalls. If we mispredict, kill all!
- \* Done in the IF stage itself.

### 3 tasks:-

1. Is the instruction a branch instruction? In ID we'll know
2. If yes, can be predict taken or not? → A hit in BTB means that Instruc<sup>n</sup> is a branch instruct!
3. If taken, Predict the target address.



BTB: branch target buffer

### → static direction prediction techniques:-

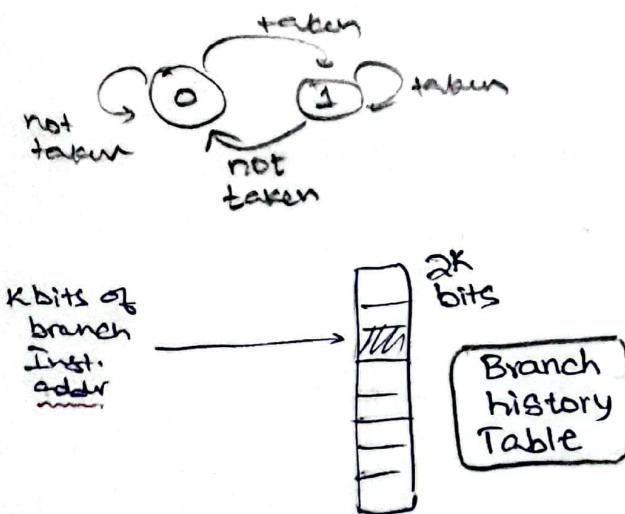
- 1) Always not - taken:
  - no need BTB
  - 30~40% accurate
- 2) Always taken:
  - need BTB
  - 60~70% accuracy

Since loops are taken most of the time.

### → Dynamic:-

micro-architectural way! not software....

#### i) Last time predictor



TTTNNNNN - 1 out of 4 times wrong....

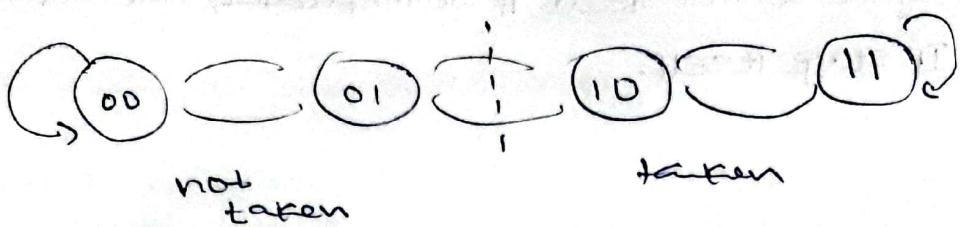
\* for loop with N-iterations

$$\text{Accuracy} = \left( \frac{N-2}{N} \right)$$

Mind you! Instruction  
each branch address has its  
own predictor.

Hence a BHT.

## 2) 2-bit bimodal predictor :-



## \* Branch resolution latency :-

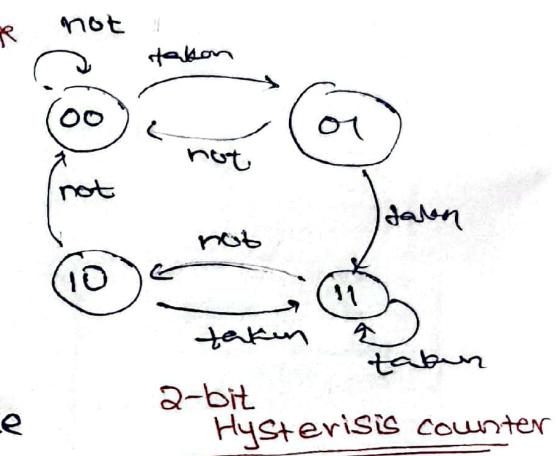
next fetch address after a control flow instruction is not determined after N cycles.

$$\text{Branch resov latency} = N.$$

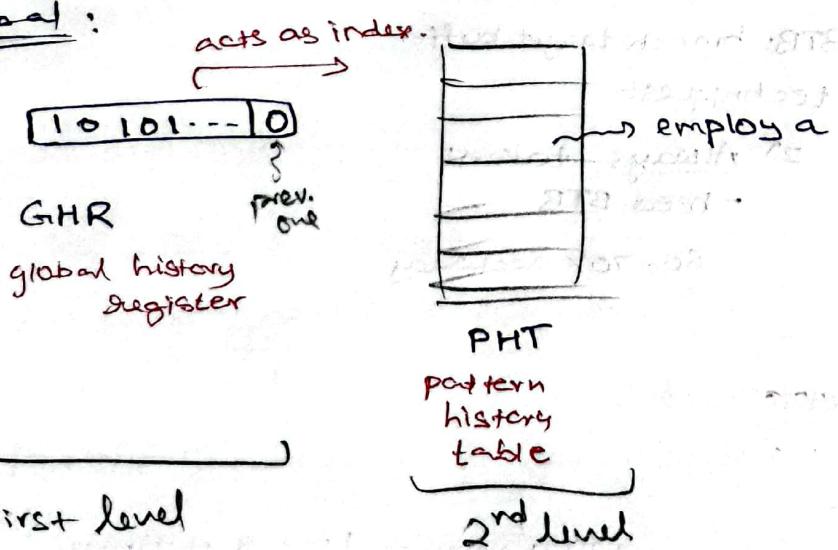
for our design it was 3.

## → Local history & Global history :-

for each branching instruction track past on a whole  
track its past

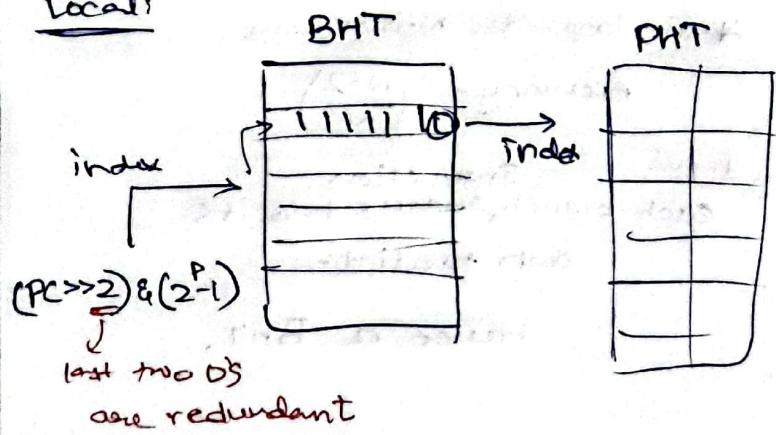


### Global :

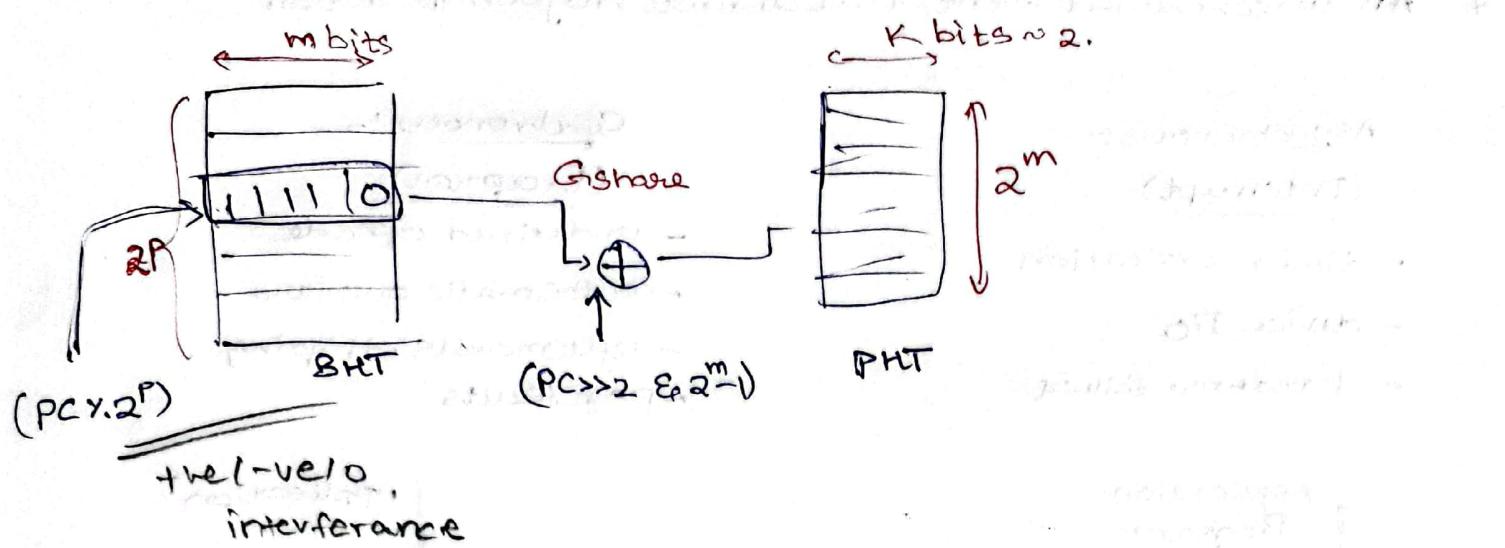


2-bit predictor.  
we mostly use a 2-bit predictor rather than 3/4/...

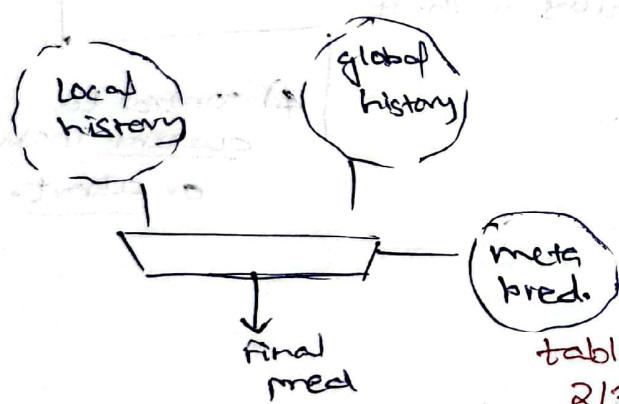
### Local :



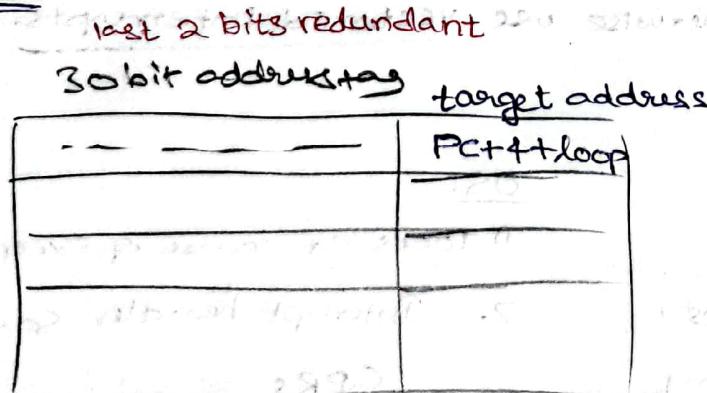
\* we also use set of branches



→ Tournament predictor:



→ BTB:



\* A hit in BTB indicates PC is branch instruction

399 words

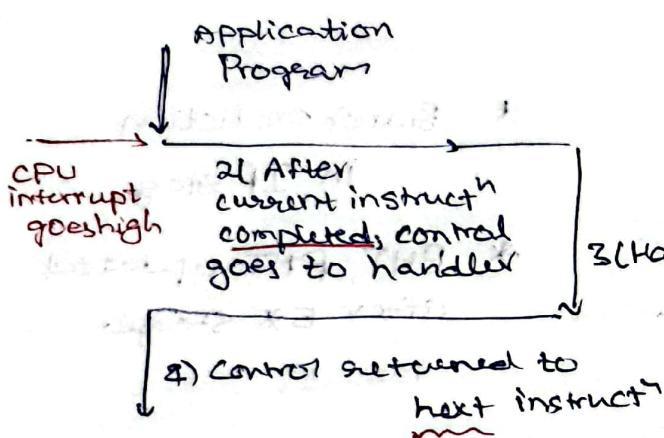
## L27:- Exceptions in Pipeline:-

- \* An unscheduled event, that disrupts program in action

### Asynchronous:-

#### (Interrupt)

- timer expiration
- device I/O
- hardware failure

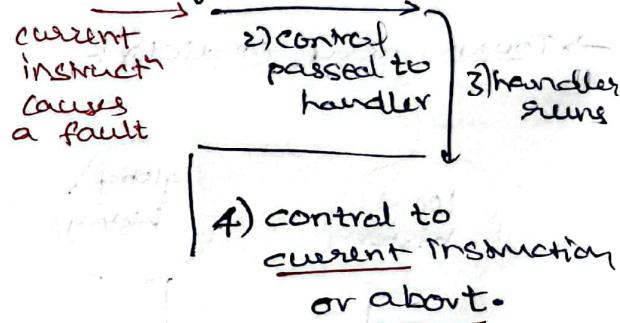


### Synchronous :-

#### (Exception)

- undefined opcode
- arithmetic overflow
- systemcalls - trapping
- page faults

### Application



### Interrupt Handlers:-

- \* 1) EPC: exception program counter.
- 2) Save EPC onto stack; before allowing nested interrupts.
- 3) Also save the Cause register.  
otherwise use vectorized interrupts.

### Processor - OS handshake:-

#### Processor:-

1. Stops the offending ins.
2. Completes the prior instructions
3. Flushes the future instruction
4. Sets cause register
5. Saves EPC
6. Disables further interrupts.
7. Jumps to predecided address.

#### OS:

1. Looks at cause of exception
2. Interrupt handler saves the GPRs general purpose reg.
3. Handles the exceptions
4. Calls RFE  
 }  
 indirect jump instruction:  
 return from except<sup>h</sup>  
 - enables interrupts  
 - goes from OS mode to user mode.

Except handling pipeline.

watch lec again.

the two basic types of communication

information exchange between machines

ATM switching fixed connection

switching connectionless

variable length packet

connectionless delivery mechanism

fixed length packets  $D = D_L$  bytes

switched queue

connectionless delivery mechanism to nodes

DATAGRAM DELIVERY MECHANISM

the ATM switch

connectionless delivery mechanism of ATM and other similar such protocols

connectionless delivery mechanism

connectionless delivery mechanism

connectionless delivery mechanism

connectionless delivery mechanism for connectionless delivery mechanism +

connectionless delivery mechanism for connectionless delivery mechanism

L28:- superscalar & O3 processors:-  
out-of-order

- Scalar pipeline:  $CPI \geq 1.0$ 
  - can never run more than 1 ins/cycle.

Superscalar can have  $CPI < 1$ .

- execute multiple instructions parallelly.

→ Instruction level parallelism ILP:-

1. Scalar:

say D-stage pipeline.

- Instruction overlap parallelism = D

- peak IPC = 1.0      in all stages loaded.

2. Superscalar:

say N instructions fetched at once.

- Instruction overlap parallelism = D × N

- peak IPC = N.

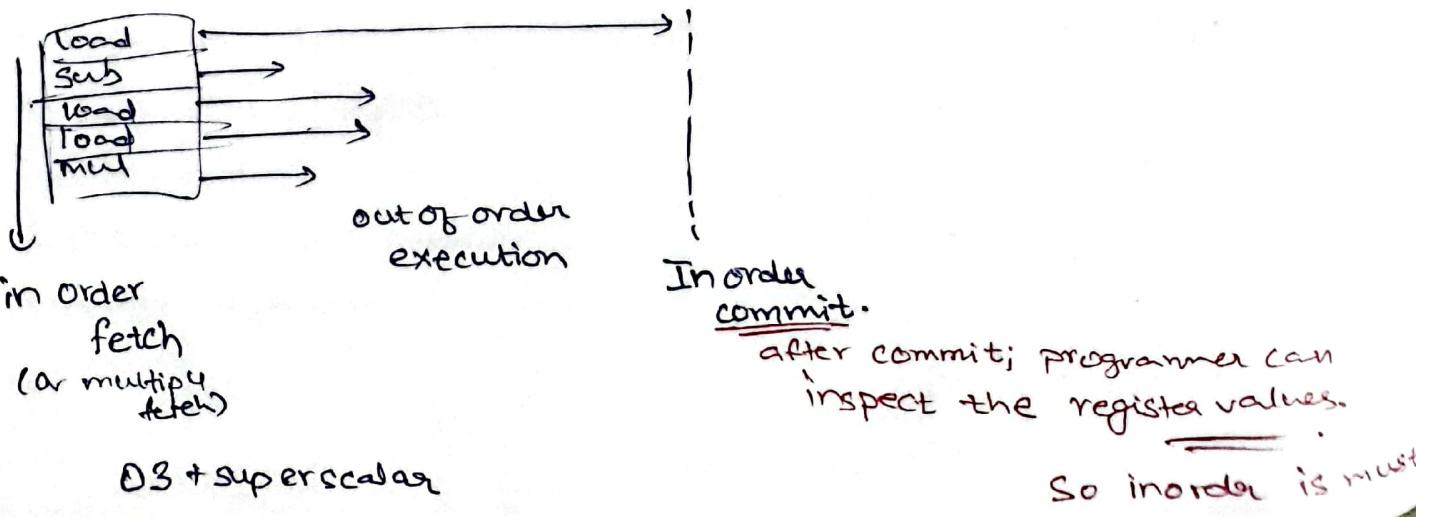
\* Hope is that we're able to fetch N independent instructions.

- complicates datapath

- complicates exceptions

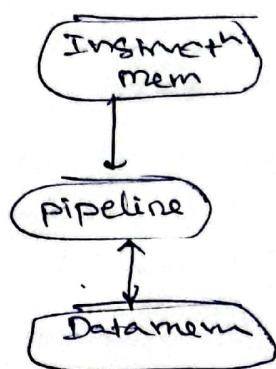
→ O3 processors:-

\* orthogonal development to superscaling/pipelining.



## L29: Memory Hierarchy:-

\* Ideally



- zero-cycle latency
- Infinite capacity
- Perfect control flow! Oh!

Branches ~ 20% of code print.

- zero cycle latency
- Infinite capacity
- Infinite bandwidth.

### SRAM:

fast

costly

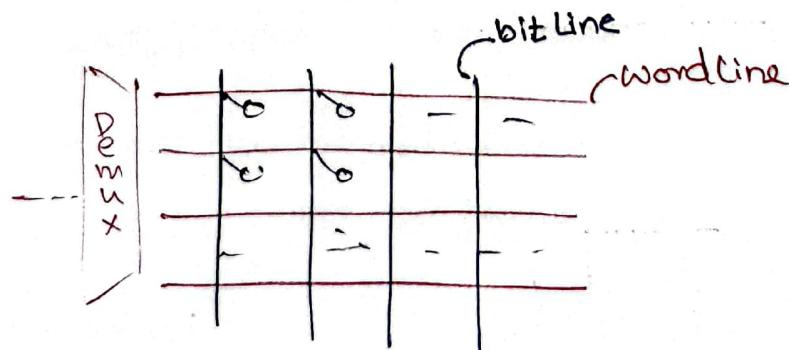
no refreshing

- 2 inverters
- 6 transistors per bit

### DRAM:

- refreshing needed; as capacitor leaks
- capacitor charge
- 1 transistor + 1 capacitor per bit

### \* Dram organisation

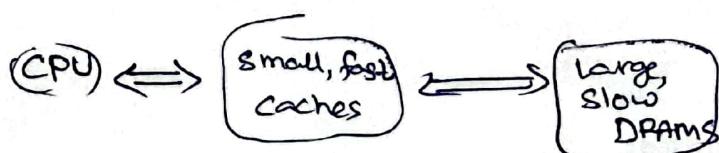


\* Bigger is slower  
Faster is costly.

\* But we want fast and big!

Hence

- memory hierarchy,
- Caches...



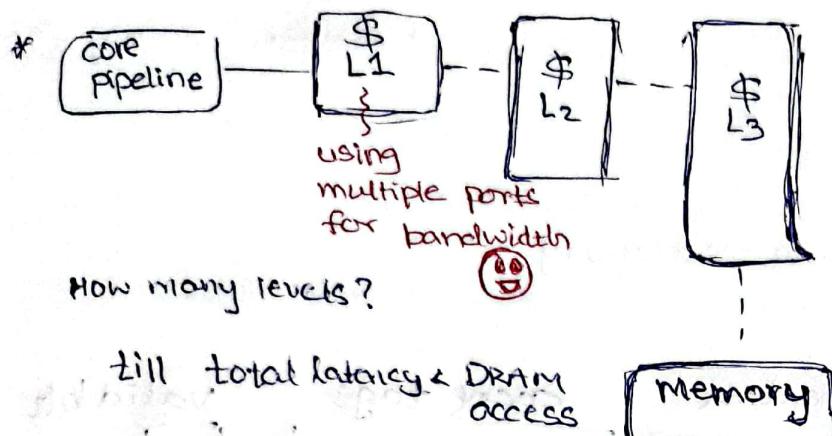
L-30: Caches \$ / coz, costly jobs.

\* DRAM access takes like 200-300 clockcycles.

Need to improve this.

→ Introduce caches.

need to have less latency,  
more bandwidth



\* Our utmost priority is getting better latency.

L1: latency + bandwidth

L2: latency

L3: Capacity... used by multicores.

Personal doubt?

Isn't bandwidth =  $\frac{1}{\text{latency}}$ ?

Nah! Bcoz; even if speed of transmission is same; length of message matters!

Low BW:	3	2	1
High BW:	7	6	5 4 3 2 1

or add parallel channels!

L3:

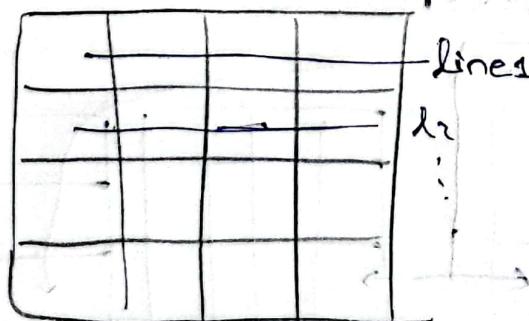
\* We always see cache \$ in terms of lines:

whenever we write

into  
cache;

we write a line.

core



Typical line

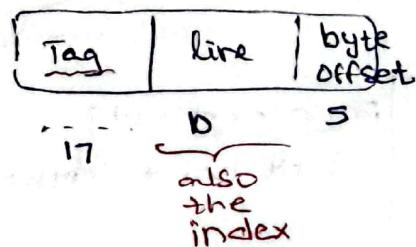
64 to 128 bytes

Eg: 1024 lines with 32 bytes each

$2^{10}$

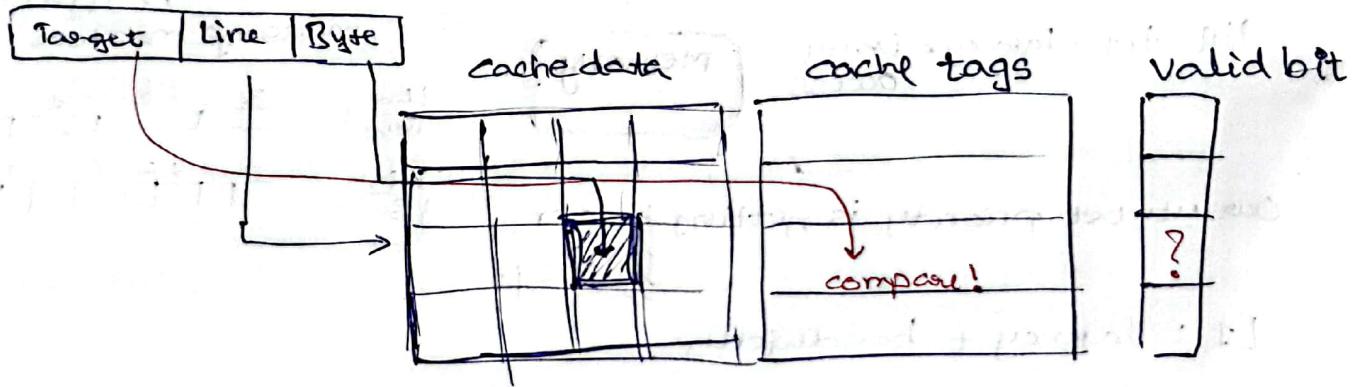
$2^5$

for Address of 82-bit



Hence;  
we need to compare tags  
finally.

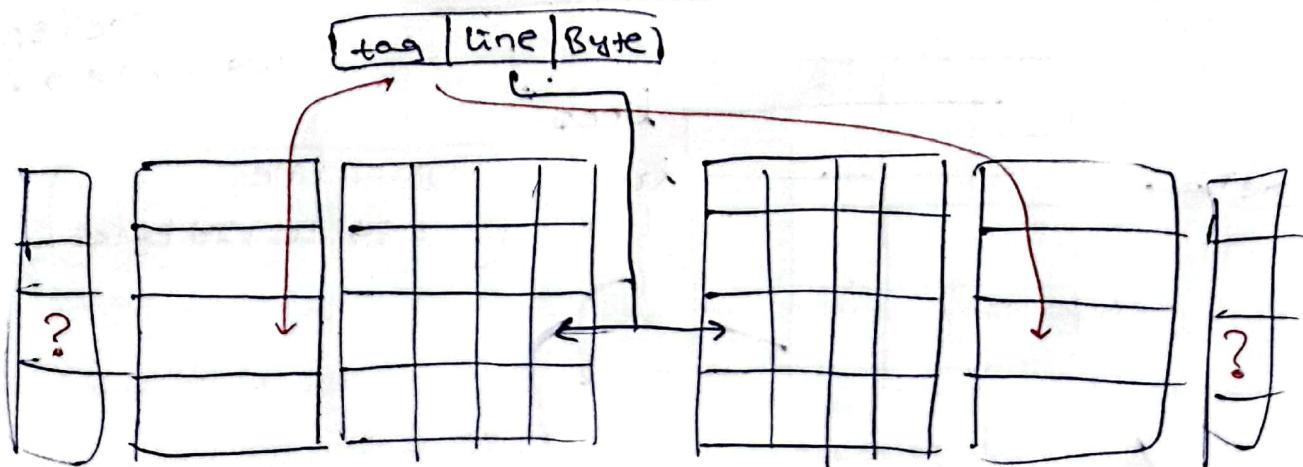
1) Directly Mapped Cache: 1-way associativity



2) 2-way associative caches:

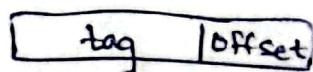
- for a single line index, we got multiple data+tag blocks

Set  $\equiv$  bunches of lines, with a index.



Similarly; 4-way associative, n-way associative...

\* Fully associative Caches:-



Single line.

size of offset

\* Knobs: Line Size, associativity, cache size.

Tradeoffs: latency, complexity, energy

L32:

- but line has many blocks

→ cache misses cuz they are costly:-

line = block in n-way associativity.

cold miss: misses that would occur, even with  $\infty$  cache

- first reference to a line.

capacity

\* capacity miss: misses that would occur, even with perfect (Belady's) replacement policy.

conflict

- cache size insufficient.

\* conflict miss:

misses that wouldn't occur, with ideal full associativity.

- many mapped to same line

→ on miss; replace a block:

\* Each block has priorities (dynamic ...).

Ideal policy: Belady's OPT policy

\* Replace the one, which comes in the farthest future.

**impractical!**

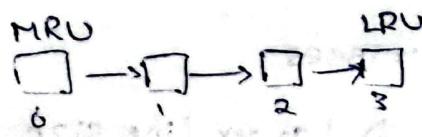
\* 3 decisions

Insertion :-

Promotion :-

Eviction :-

1) Simple LRU:-



Insertion: at MRU

eviction: at LRU

Promotion: bring to MRU

\* we need priority bits per block in a line.

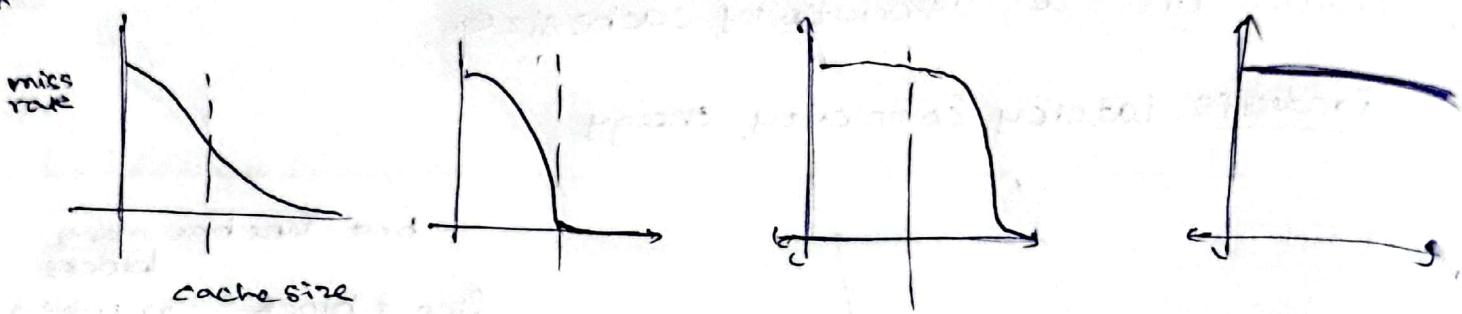
\* 16-way cache;

each block → 4 bits

Expensive! 😠

\* LRU causes thrashing when working set > cache size.

at some duration of time...



- \*) a) cache friendly workload    b) cache fitting workload    c) cache thrashing workload    d) streaming workload.

\* ALSO; LRU not effective for shared caches

L3 cache.

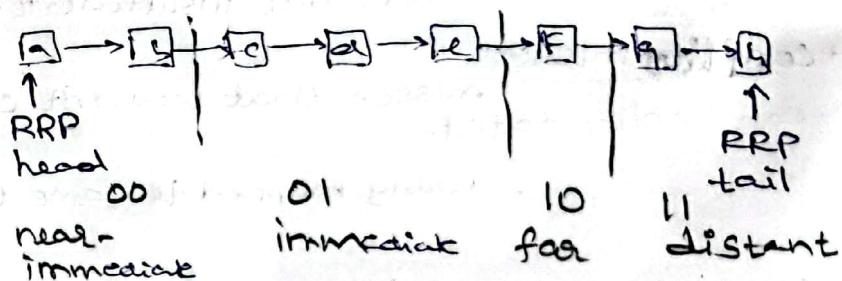
\* LRU is not effective

RRIP:

re-reference interval prediction policy

\* RRPV

- pred. values



We'll group into 4 groups.  
∴ need only 2 bits/Block

Evict: value of block = 11

Promotion: make block value = 0

Belief:

New cache block will not be re-referenced soon.

Insertion: insert with RRPV = 2

not 0 like LRU

→ Knobs & misses:

1) Larger cache size

- reduces capacity misses, conflict misses
- Hit time increases latency

2) Higher associativity

- reduces conflict misses
- May increase hit time  
(more comparators)

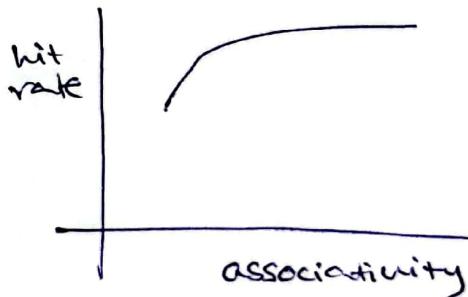
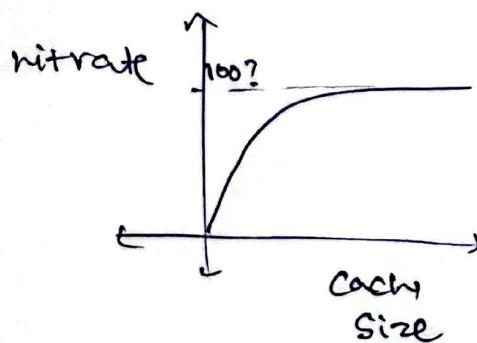
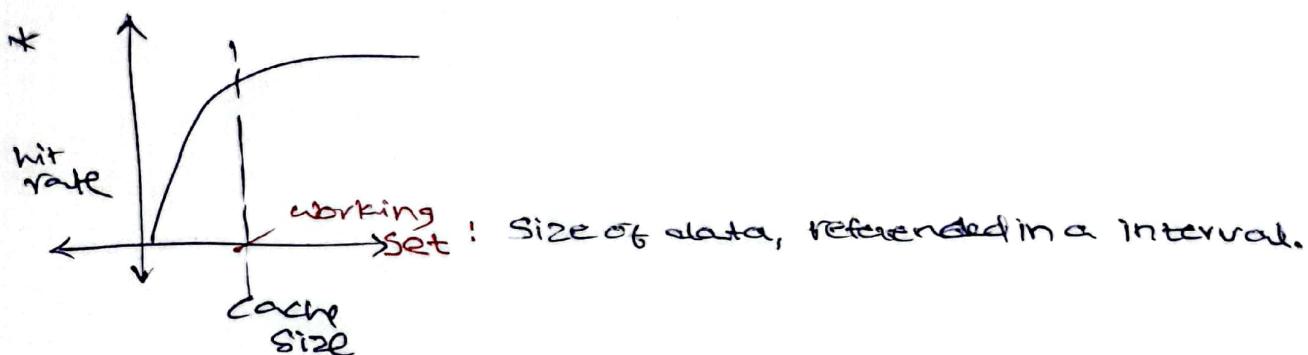
3) Larger Line Size

- reduces cold misses  
Since we load 1 line at a time.
- increases miss penalty  
(more to be loaded)
- increases conflict misses

not sure

\* too few block size: spatial locality not exploited

too large block size: waste info copied, more energy.



But hit time also increases.

L1: low associativity

L3: high associativity

\* Misses per kilo instructions

(MPKI)

per 1000 instructions is a normalizing factor.

→ Average Memory Access Time:

$$AMAT = \underline{\text{Hit time}} + \text{Miss Rate} (\text{miss Penalty})$$

even

misses need this time to check... 😊

cool!

ALSO

$$= (\text{Hit time})_1 + \underset{1}{\text{miss rate}} (\text{miss Penalty}_1)$$

$$+ (\text{Hit time})_2 + \underset{2}{\text{miss rate}} (\text{miss Penalty}_2)$$

③

L33:

\* L1 data cache:

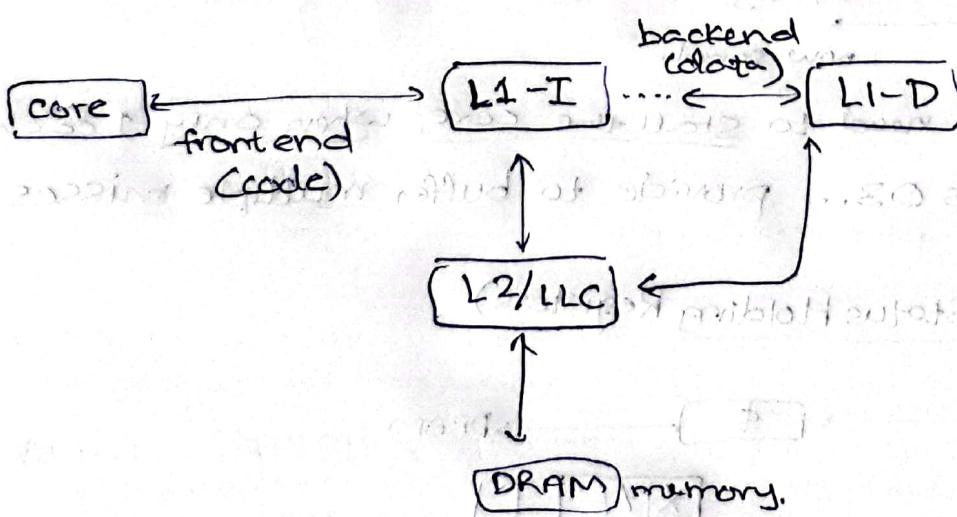
- \* High bandwidth
- low associativity to fast hit time

Victim Cache:

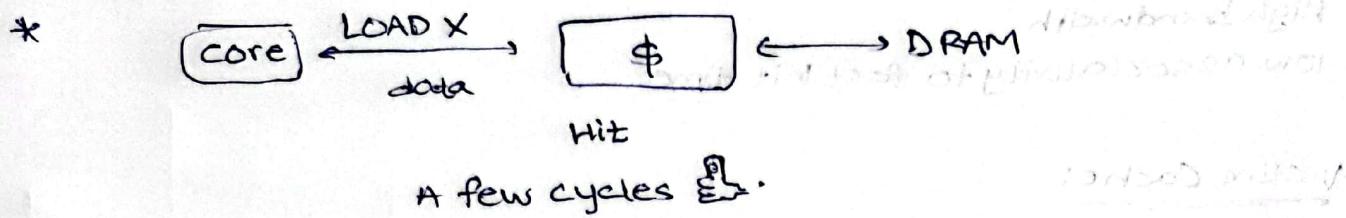
Store some evicted lines of L1 cache  
for better hit rate.

\* L1 instruction cache

- special cache that packs multiple non-contiguous blocks into one contiguous trace cache line
- Single fetch brings in multiple basic blocks
- Indexed by start address & next n branch predictions.



## → core cache, DRAM interaction:-



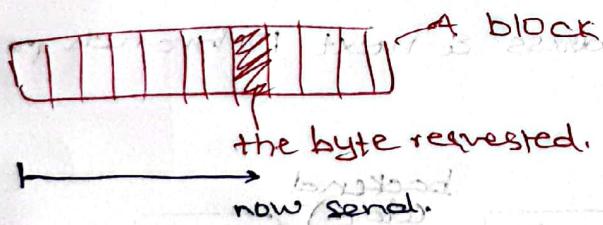
\* if we miss;

a) Critical word first:

\* respond with the word/byte requested to the core,  
 & then continue fetching rest of the block.

b) Early restart:

\* fetch the words/bytes in normal order, but  
as soon as requested word arrives, send to core.



\* So, do we really need to stall the core when only 1 cache miss?  
 the core is O3... provide to buffer multiple misses.

## → MSHRs (Miss Status Holding Registers)



K-entry MSHRs allow K-outstanding misses;

Memory Level parallelism

MLP

Are all fetched parallelly?

## → What about write?

- \* On hit, update cache and set dirty bit to 1.

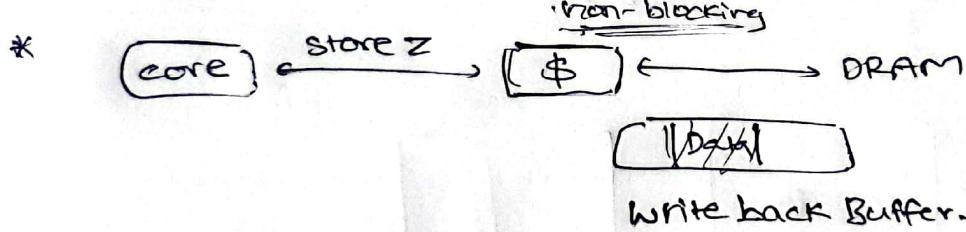
other bits are valid bit.  
replace priority bits.

### a) Write-through cache :-

on a hit; write into a cache & also into next-level in hierarchy.  
immediate

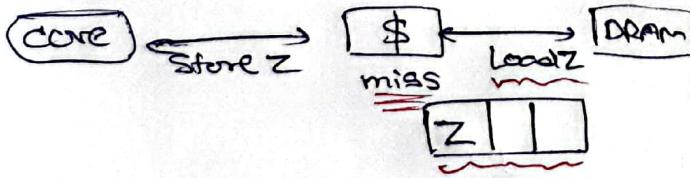
### b) Write-back cache :-

on hit, write into cache only. During eviction; write into next level.  
later, during replacement



In general, stores are not critical for performance...

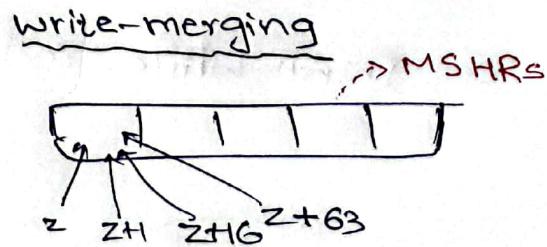
## → Write-miss:



- \* store gets converted into load;  
& data is allocated into cache

write-allocate policy

used, along with write-back policy.



All these store requests are merged into one.

Meh! even load merge then.

\* CPU time = CPU execution cycles + Memory-stalls

= readstalls  
+ writestalls

= # Read/writes × Read/write miss rate × miss penalty.

L34)

Capacitor network analysis

matrix mult:

K vary

$$\begin{bmatrix} \text{---} \\ \text{---} \end{bmatrix} \cdot \begin{bmatrix} \parallel \\ \parallel \end{bmatrix} = \begin{bmatrix} \text{---} \\ \text{---} \end{bmatrix}$$

missate      migrate      migrate      1.25  
0.25                  1                  0

i vary

$$\begin{bmatrix} \text{---} \\ 0 \end{bmatrix} \times \begin{bmatrix} \parallel \\ \parallel \end{bmatrix} = \begin{bmatrix} \text{---} \\ \text{---} \end{bmatrix}$$

missate      migrate      0.50  
0                  0.25                  0.25

i vary

$$\begin{bmatrix} \parallel \\ 1 \end{bmatrix} \times \begin{bmatrix} \text{---} \\ 0 \end{bmatrix} = \begin{bmatrix} \text{---} \\ 1 \end{bmatrix}$$

missate      migrate      2.0  
1                  0                  1

think while cooling!



parallel branch with 10Ω

one side bypassed

new loop from 10Ω

total voltage drop across 10Ω

had either one p or no loops

20.02 ohm + 10 ohm = 30.02 ohm  
10 ohm = 10 ohm  
30.02 ohm = 30.02 ohm

resistor point is offboard x resistor is