

- Performance :-

Arithmetic / Logical	- 50%
Load	- 20%
Store	- 10%
Broadcast	- 15%
Jump	- 5%

$$\text{For multicycle, Performance} = N \times \frac{1}{\text{Cycles}} = 4.15 \times 3 \text{ns} \\ = 12.45 \\ \text{single cycle} = 12.5$$

- Depending on fractions of instruction types, either of single cycle or multicycle implementation.
 - In spite of fewer registers, single-cycle is more expensive overall in terms of hardware.
 - Single cycle :- Lowest possible CPI (1)
Multi cycle cycle time
 - Possible to strike a balance :- eg - 2 cycle implementation.

Pipelining

- Begin next instruction before previous is over; using the idle resources
 - In RISC, subtasks in every instruction are very similar - fetch instruction, fetch operands, alu operation, etc.

- It is possible to get least possible cycle time (like multicycle)
 - If we assume all instructions to be independent, then we can do pipelining and then we will get cycle time equal to $\max\{\text{sub-task time}\}$

Mealy -
Moore -

Control logic & Next-state logic depend on both input and current state
NSL is same as before
CL depends only on current state

classmate

Date _____
Page _____

Subtasks

- 1 R-type :- IF, EX, ID/OR, Instruction Decode / Operand Read, WB
- 2 LOAD :- IF, ID/OR, EX, Memory Read, Mem, WB
- 3 Store :- IF, ID/OR, EX, Mem
- 4 BEQ :- IF, ID/OR, EX
- 5 JUMP :- IF, ID, EX

circumstances set of
sub-tasks

- First 4 cycles, no instruction will be completed.
After that, one instruction will be completed every cycle.
(for independent instructions and if instructions are consecutive from top to bottom (no control flow changes in between))

$$\rightarrow \text{Performance} = N \times \underbrace{1}_{\substack{\text{CPI} \\ \text{CPI} \approx 1}} \times \underbrace{3 \text{ ns}}_{\substack{\text{Cycle time} \\ \text{Cycle time} = \max(\text{subtask time})}}$$

- Every instruction will go through all 5 stages

IF \rightarrow ID/OR \rightarrow EX \rightarrow Mem \rightarrow WB

e.g. JUMP will waste time in Mem & WB stages

- Waiting time :- Do the job in each stage, but don't write to any register (using control signals on registers)

multi cycle is example of mealy and pipelining is example of moore as in
pipelining no matter what the instruction is (ie input) it will go through all the

- Making 5 stages common for all instructions converts each stage
from a Mealy machine (which depends on instruction) to a Moore machine (which does not)

Mealy :- Control signals generated depended only on current state
Moore :- Control signals generated depended on both

- To use pipelining, all instructions must use similar type of stages, and in same order too.

	IF	ID/OP	EX	Mem	WB
Cycle 1	I_1				
2	I_2	I_1			
3	I_3	I_2	J_1		
4	I_4	I_3	I_2	I_1	
5	I_5	I_4	I_3	I_2	I_1

Latency = '5'?

not 4

→ How do you design a pipeline (determine the required subtasks)?

:- Start from multicycle :- We will have minimal resources used.

Could be an issue if resource is shared by subtasks

e.g - In multicycle, some memory is used by IF & Mem.

- ∴ Start with single cycle to get maximal resources.

→ We generate control signals through combinational decoder of ^{in a register} single cycle implementation (like always), and "remember" them. The control signals for each instruction must be propagated consecutively to each task along with the instruction itself.

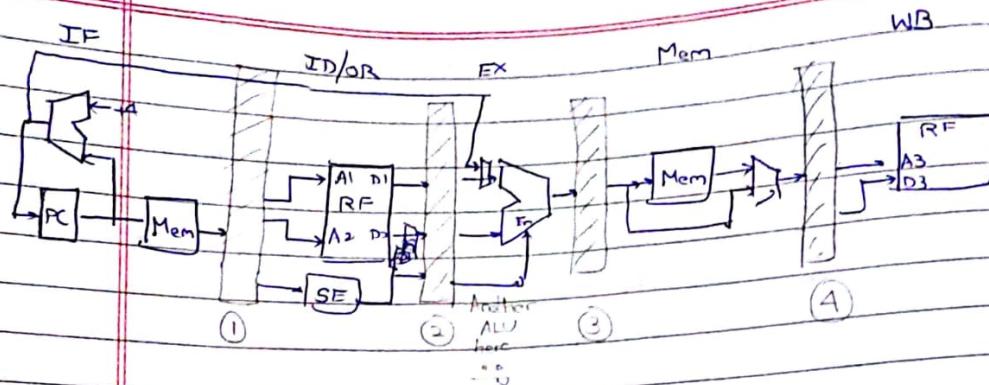
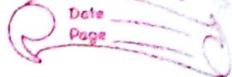
The respective subtask uses whichever control signals it needs.

- So every subtask will have a register (to store control signals) of the size equal to number of control signals flowing from previous state to that state.

Register in subtask 1 will have size equal to total number of control signals.

- PC is also propagated through the stages, as it will be required for returning control flow changes.

Instructions move in lock step manner (= FIFO, no jumping) classmate



- Signals crossing interface after TF :-
Instruction (32 bit), PC+4 (32 bit), etc

- Interface 2 :-
PC+4 (32-bit), Operand 1 (32-bit), Operand 2 (32), Sign-extended Imm(32), ALU control signals, destination register address, etc

22/10

- Interface 3 :-
ALU output (32), A3 address etc

eg

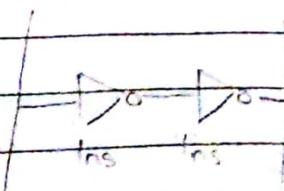
- Interface 4 :-
Destination register address, data to be written, etc

→ How many stages to divide the instruction into?

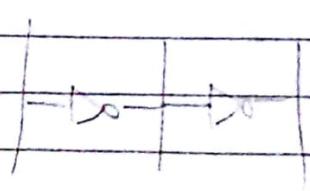
- More no. of (smaller) stages will give lower cycle time 😊
 - Becomes costlier, as you need more registers for interfaces 😭
 - Performance ↑, CPI still ≈ 1

One stage

Two stages



Cycle time = 2ns



Cycle time = 1ns

- Choose the value of no. of stages that maximize Performance / Cost
 ↗ avg cost of interface registers

$$\text{Cost} = G + (n-1)k$$

↓
 cost without pipelining

$$\text{Performance} = 1$$

$$\frac{T}{n} + q \leftarrow \begin{array}{l} \text{Same smal} \\ \text{Delay (intercon)} \end{array}$$

→ time without pipelining

$$* \frac{\partial}{\partial n} \frac{1}{(\frac{T}{n} + q)(G + (n-1)k)} = 0$$

22/10

→ Our assumptions :- (Hazards)

e.g. Dependent instructions :- $I_1 : R_1 = R_2 + R_3$

$$I_2 : R_4 = R_1 + R_5$$

In ID/OR stage, I_2 will read incorrect R_1

consecutive

• Detecting violation of independence of instructions :-

Source registers in ID/OR stage instruction = Destination register of

Ex stage instruction.

$$\text{i.e. } (R_{S1}^{\text{ID/OR}} == R_{D1}^{\text{Ex}}) \text{ || } (R_{S2}^{\text{ID/OR}} == R_{D2}^{\text{Ex}})$$

- If this violation occurs, 'freeze' I_2 in the ID/OR stage

for 3 cycles (until I_1 completes). Then resume.

- On every violation, we are losing 3 cycles

e.g. - If 10% instructions have this 'immediate' dependency.

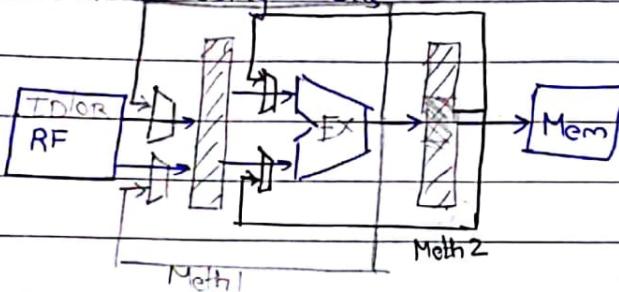
$$\text{New CPI} = 1 + 3 \times 10\% = 1.3$$



- The result of I_1 is getting written back in stage 5 (WB). But it is being produced in stage 3 (EX) itself. This result will be stored in the pipelining interface register at the boundary of EX and Mem. Let I_2 fetch the result from there.

$I_1 - - -$
 $I_2 I_1 - - -$
 $I_3 I_2 T_1 - -$
 $I_4 I_3 I_2 T_1 -$

Two possibilities for doing this



→ Non-consecutive dependency.

$$I_1 : \quad r_1 = r_2 + r_3$$

$$I_2 : \quad r_4 = r_5 + r_6$$

$$I_3 : \quad r_7 = r_8 + r_9$$

1 $(RS1 \text{ or } RS2)^{ID} == RD^{Mem}$:- Where to fetch operand from?

2 $(RS1 \text{ or } RS2)^{ID} == RD^{WB}$:- Obtain result from RF-D3 of WB stage

→ Multiple write conflicts.

$$r_1 = r_2 + r_3$$

$$r_4 = r_5 + r_6$$

$$r_7 = r_8 + r_9$$

- Priority order :- For conflicting dependencies, fetching result :-

Ex \rightarrow Mem $>$ WB

newest instruction \rightarrow

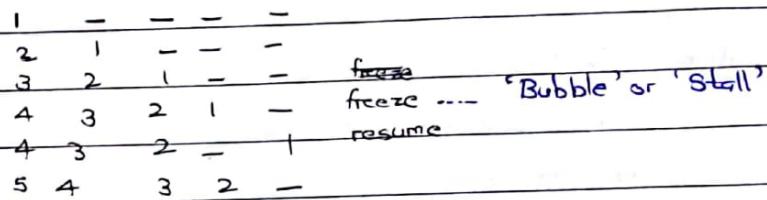
\hookleftarrow oldest instruction

→ Load instruction produces result in Mem stage.

$$I_1 : r_1 \leftarrow (R_2 + 50)$$

$$I_2 : r_3 = r_1 + r_2$$

When I_2 is in ID/OP, the required result has not been produced yet. Forwarding will not work. We have to wait (freeze) for one cycle at least in such cases.



23/10 - We are losing one cycle for every instruction of this type.

eg - $\alpha\%$ AL instructions, $\beta\%$ Load (of which $\lambda\%$ have such immediate hazard), $\gamma\%$ Store, $\delta\%$ Branch.

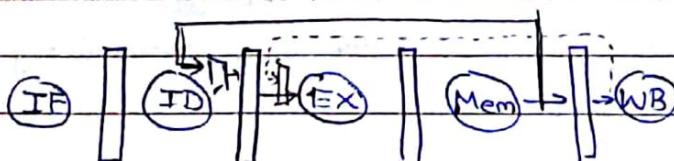
$$\text{New CPT} = 1 + \beta \lambda (1)$$

- Bubble or Stall :-

All interface registers should have an enable signal, which is made zero to stop updation.

Make enable of PC zero, too.

- For hazardous Load.



→ Minimizing the damage of Load hazard :-

eg $A = B + C + D \rightarrow \text{load } r_1, B$

$(T_1 = B + C)$
 $A = T_1 + D$

A, B, C, D, T_1
are in memory

load r_2, C
add r_3, r_2, r_1
store r_3, t
load r_4, D
add r_5, r_4, r_3
store r_5, A

Compiler optimization

load r_1, B
load r_2, C
load r_3, D
add r_3, r_2, r_1
store r_3, t
add r_4, r_3, r_1
store r_4, A

Smart compiler

IF I₂
IF I₂
I₁ b
IF I₁
instruction

→ Condition
eg I₁

→ Change in Control Flow

eg I₁ contains JUMP.

I₂ is not supposed to be executed.

IF ID/OR

- 1 - - - -
- 2 1 - - -

All * LOAD

I

- Every interface register should have an extra bit, which will decide if an instruction is to be disabled.

If I₁ is JUMP, then the disabled I₂ will go through all stages, but will not write to any registers, flags

EE Task
JMP we
Use

Meth 2 Delayed Branch Execution

After every branching/jumping instruction, compiler automatically adds a NOP instruction.

P.T.O.

If I₂ is JUMP,

- If I₂ does not depend on I₁ (unconditional jump), move I₁ to just after I₂.
- If I₂ depends on I₁ (conditional jump), add a NOP instruction to just after I₂.

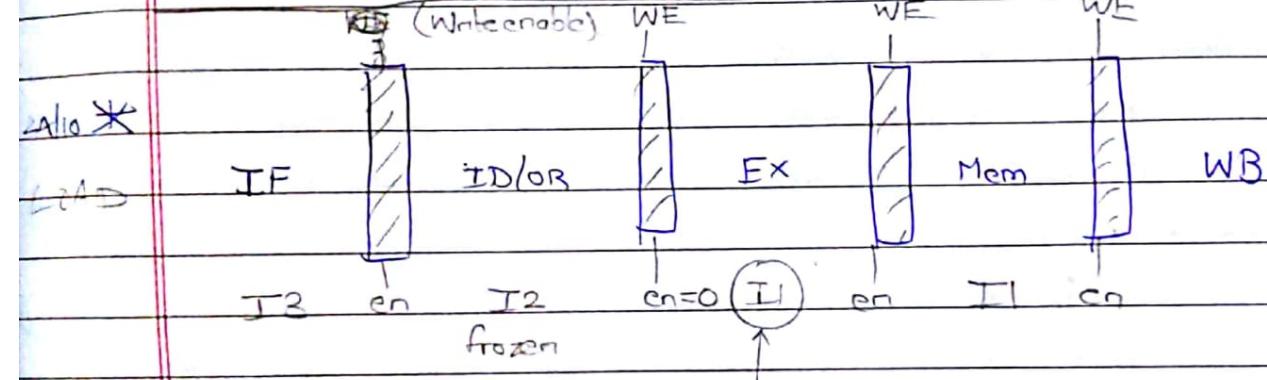
→ Conditional branching

e.g. I₁ is BEQ

1 - - - -

2 1 - - -

3 2 1 - -



Comes here because interface register 2 was not updated in previous cycle.
This extra instruction needs to be killed.

- Instead of using $(R_{S1}^{ID} == R_D^{EX}) \text{ || } (R_{S1}^{ID} == R_D^{Mem}) \text{ || } (R_{S1}^{ID} == R_D^{WB})$, we now need more checks because I₁ is duplicated.
- Use: $(R_{S1}^{ID} == R_D^{EX} \text{ & } R_{S1}^{WE} == 1) \text{ || } (R_{S1}^{ID} == R_D^{Mem} \text{ & } R_{S1}^{Mem} == 1) \text{ || } (....)$
- eg. it should be LOAD instruction
- MIM-AD = 0

Only LOAD instruction reads from data memory.

• Shift Register elaboration

If this condition is met, $IF_{EN} = ID/OR_{EN} = PC_{EN} - RF_{WR} = 0 = \text{Flags}$
 $\text{Mem-WR} = 1$

→ Unconditional JUMP.

Let ~~all~~ interface registers have a bit "valid".

If I_1 is found to be jump, make 'valid' of I_2 from (in I_1) '0'.

If valid = 0, don't let that instruction write to flags, registers, memory. (make WE of subsequent interface registers 0)

★ Assume that for JUMP instruction, we can update PC as soon as ID stage (using separate ALU.)

→ Conditional JUMP. (BEQ) (I.)

- Assume branch is not taken. Bring in I_2, I_3 serially.

When I_1 reaches EX stage, it decides whether branch is taken or not.

- If branch is to be taken, don't let I_2, I_3 update flags/registers/memory (assumption is invalidated)

- Penalty of 2 cycles

* Most ~~%~~ of the times, branches are taken

- Assume that if branch is to be taken, PC is updated at the end of EX stage.

• Our initial assumption cannot be to take the branch, as we don't know yet where to branch to.

• $\alpha \rightarrow AL, \beta \rightarrow Load, \gamma \rightarrow Store, \delta \rightarrow Branch$ (Conditional), $\epsilon \rightarrow Jump$
 ↓
 1% dependency
 ↓
 70% branches are taken

$$CPI = 1 + \beta \lambda(1) + \underset{\text{Why?}}{\epsilon}(1) + \delta(0.7)(2)$$

25/10 If branch is taken,

I1	-	-	-
I2	I1	-	-
3	2	1	-
20	-	-	1
21	20	-	1

~~I2~~

args,

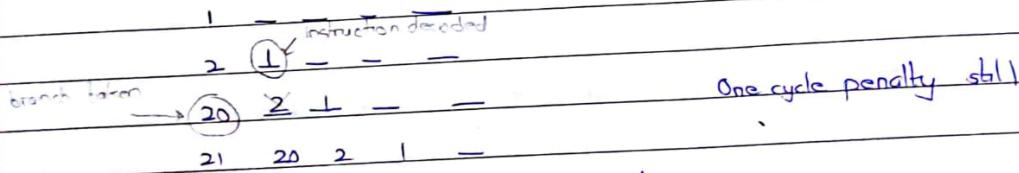
3)

- Optimization for when branch is taken:-

For first time of execution, compute proceed as before.
Compute destination of branch and store it in Look-up table.

PC	Branch Target Address
100	700
150	755
:	:

Such instructions are mostly executed many times (in a loop)
Second time onwards, change assumption to branch taken

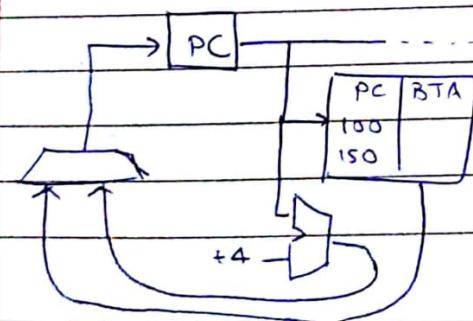


Now assumption is correct 70% of the times.

★ - New CPI = $1 + \beta \lambda(1) + \epsilon(1) + s(a, z)(\frac{1}{2})$

Verify

- Better still (for no penalty if branch is taken) detection in IF stage of T1
- Match the Lookup table with new value of PC.
If current PC is found to be present in table, fetch I20 in next consecutive cycle.



CPI = $1 + \beta \lambda(1) + \epsilon(1)$

★ Verify

1 - - -

20 1 - - -

:

- If branch is not taken (2 cycle penalty)
- We have lost original PC. Where to go ?
- We need to carry forward original $PC + 4 (\text{if } f \neq 1)$ forward
- Anyways we were carrying it forward till EX to compute
 $PC + 4 + T_{mm}$

- Further optimization :-
- Branch instructions tend to show similar behaviour as the last time they were executed.
- In the look-up table, add a bit to say if last time branch was taken or not.

Change your assumption accordingly.

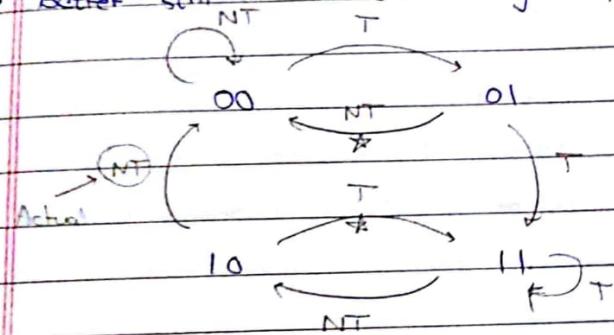
PC	BTA	History ($0 \rightarrow \text{not taken}$)

In execution stage of IF, update history bit every time.

Prediction	Actual	
0	1 Entry into loop
1	1	
1	1	{ 8 times
:	!	
1	0 Exit from loop
0	1	

We were right 80% of time. (no penalty.)

- Better still :- 2 bit history (past 2 decisions)

 $H_{t+1}H_{t+2}$

Prediction

00 → Strongly NT01 → Weakly NT10 → Weakly T11 → strongly T.

Now you will not make mistake for entry into loop (90% correct)

29/10

MEMORY SYSTEM

DRAM

- For modern memory, one memory access (not 1 as before)
 - Not counting multiple memory accesses (just one, to fetch ir), each instruction will take ~200 cycles
 - Pipelining makes no sense, as a fetched instruction will take only #stages cycles more to complete

eg

- One possible solution:-
 - Between processor and memory, add a very small, but fast SRAM memory. This has, say, only 8 bytes (worth 2 instructions) and is made from more expensive SRAM.
 - SRAM memory is accessible in single cycle (must fetch consecutively written instructions)

'Cache'

- Fetch 8 bytes (2 instructions) from memory into SRAM (~200 cycles)
 - Processor fetches I_1 (200th cycle)
 - I_2 (201st cycle) --- pipelining
 - I_3 (402nd cycle)
 - I_4 (403rd cycle) --- pipelining
- Performance almost doubles by introducing SRAM, as at least 2 instructions are being pipelined

so/10

eg Eg Program

for i in range(n):

I_1

Assume instruction

I_2

do not access memory

I_3

I_4

II SRAM is 16 bytes, 4 instructions will be fetched in 200 cycles and stored in SRAM

For every subsequent iteration, we don't need to fetch again and instructions can be pipelined as before

CPI \rightarrow 1



eg for i in range (100):

if $i \% 2 == 0$: $I_1 I_2 T_3 I_4$

else $I_5 I_6 T_7 I_8$

- With 16 bytes of SRAM, we will have to fetch and replace the four instructions again and again.

$$\text{Total no. of cycles} = (200 + 1 + 1 + 1) \times 100$$

??

- With 32 bytes of SRAM, we can get CPI $\rightarrow 1$, but only if all 8 instructions are consecutively written. Otherwise same as before.

$$(200 + 1 + 1 + 1 + 1 + 1 + 1 + 1) \times 59 + (1 + 1 + 1 + 1 + 1 + 1 + 1 + 1) \times 49$$

- With 2 SRAM chunks of 16 bytes each, we can get CPI $\rightarrow 1$. The instructions in each individual chunk must be consecutive.

$$(200 + 1 + 1 + 1 + 1 + 200 + 1 + 1 + 1) + (1 + 1 + 1 + 1 + 1 + 1 + 1 + 1) \times 48$$

- The SRAM chunks are hidden from programmer and known as 'Cache'.

so/10

→ Bandwidth of Memory :-

$$\frac{\text{Data}}{\text{sec}} = \frac{\text{Data}}{\text{Instruction}} \times \frac{\text{Instruction}}{\text{Cycle}} \times \frac{\text{Cycle}}{\text{Sec}}$$

$$= \left(\frac{\# \text{Bytes}}{\text{Instruction}} + \frac{\# \text{Bytes of data}}{\text{Instruction}} \right) \times \text{IPC} \times \text{freq}$$

$$= (4B + 4B) \times 1 \times 3 \times 10^9$$

$\approx 24 \text{ GBps}$.

max IPC = 1

This is for peak demand from memory (each instruction requires 4B data to be fetched too)

- Final CPI depends on no. of 'stalls' required by the program
- Instructions may stall, for reading from or writing to memory, for 200 cycles
 - Stalling also happens to fetch instructions not present in cache

- Instructions can be 'hit' or 'miss'

present in cache
No penalty

not present
200 cycle penalty

$$\bullet \text{ No. of stalls} = \text{No. of misses} \times \text{penalty}$$

$$= \# \text{Instruction} \times \frac{\# \text{miss}}{\text{Instruction}} \times \text{penalty}$$

$$= \# \text{Inst} \times \frac{\# \text{Memory references}}{\text{Inst}} \times \frac{\# \text{miss}}{\# \text{Mem ref}} \times \text{penalty}$$

$$= \frac{\text{Inst Count}}{(\text{IC})} \times \frac{\# \text{Mem ref}}{\text{Inst}} \times \underbrace{\frac{\text{Miss rate}}{\text{Inst}} \times \text{penalty}}$$

To the hands of memory designers

→ Miss Rate reduction

1 Increase no. of blocks in SRAM.

2 Replacing contents of cache :-

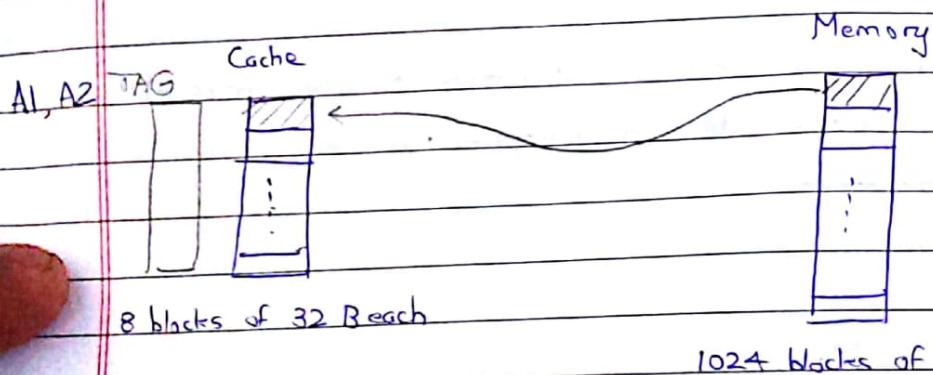
Preferentially, do not remove instructions that are conservatively likely to follow current instruction (spatial likeliness) or that are likely to be executed again at some later point in time (temporal likeliness)

'Locality of reference'

- Increasing the number of blocks in SRAM beyond a certain point will increase access time from 1 cycle to 2 cycles. :-

→ Questions to be answered

- 1 In which block of cache to place?
- 2 How to identify where required instruction is placed in cache?
- 3 Which cache instruction to replace when overwriting?
- 4 When to write data (obtained from processor) from cache to memory?



In memory, every byte can be addressed as
cache

Block No.	Offset
10 bits (mem)	5 bits
3 bits (cache)	

- We will only be transferring blocks as a whole from memory to cache. After transferring, offset part of each byte remains same. Block no. will be changed accordingly
 - We maintain a mapping of which block no. of memory is transferred to which block no. of cache.
 - Simplest Mapping = 'Direct Mapping'

$$(\text{Mem Block No.}) \% 8 = (\text{Cache Block No.})$$

e.g. - Mem Block 45 is transferred to Cache block 5.

P.T.O.

i.e - Mem Block

(Cache block no.) = Last 3 bits of (Mem block no.)

∴ (Cache address) = Last 8 bits of (Mem address)

- The first seven bits of Mem address (which were lost in transferring to cache) will be stored in an adjacent array to cache. (Tag)
 - TAG is an 8-element array of 7-bit elements each (one element for each block)

1/ii • Checking if required memory data/instruction is present in cache can be done with a comparator in TAG.
→ bits will have to check

→ Fully Associative Mapping

- Any memory block number can be placed in any cache block number.
- Better utilization. :-)
- Now, TAG needs to have all 10 bits of memory block number element
- Now, comparators have to check TAGs at every cache block
 - May not happen in 1 cycle :-)

→ Take a middle path = 'Set Associative'

- Divide cache into multiple sets

e.g - 4 sets of $\frac{2}{2}$ blocks each

'~~2~~ 2-way set associative'

- Mapping of memory to each 'set' in cache is fixed.
But, within that set, any block can be used.

$$\text{(Cache set no.)} = \text{Last 2 bits of (Mem block no.)}$$

$\therefore \text{Tag element size} = 10 - 2 = 8 \text{ bits.}$

A1 Write-Through Policy

Write back into memory whenever you are writing from processor
to cache

2 Write-Back Policy

Write to memory only when that block needs to be rewritten.

A3 Replacement

- In direct mapping, you don't have a choice. Need to rewrite in some block no.
- In fully associative mapping or set associative mapping,
 - Add in empty block if present,
 - If not block is empty, we want to replace the block that is least useful.
 - Possibility - FIFO

But first-in block need not be least useful.

- 'Least recently used'

eg - In 2-way set associative

Use an LRU bit per set of 2 blocks

$LRU = 0 \therefore \text{Replace } 0 \cdot \text{ Make } LRU = 1$ 0 $\left. \begin{array}{l} \{ \text{Keep toggling} \\ \text{LRU bit} \end{array} \right.$

D₁

D₂

classmate
Date _____
Page _____

If you are using Write Back policy, maintain a 'dirty' bit for each block.

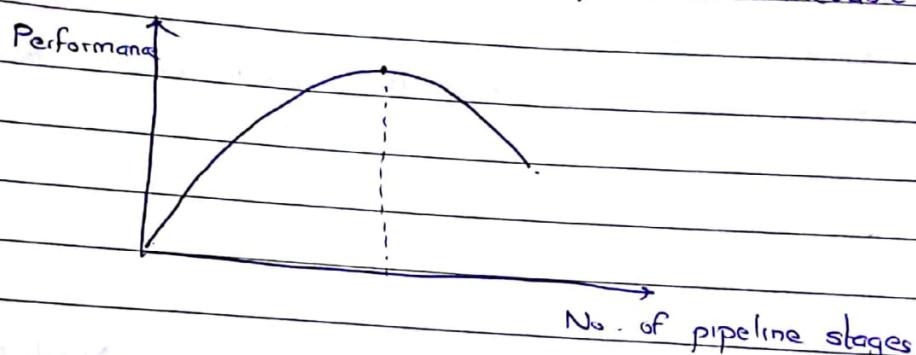
Make dirty = 1 if block was updated.

If dirty = 1, write back to memory before replacement.



IBM's experiment for pipelined design.

- Given the typical programs we run, least achievable CPI = 1.15



- What we have studied:- "Scalar pipeline"
 - 1) Fetch one instruction at a time.
 - 2) All instructions pass through the same 'Unified' pipeline
 - 3) Instructions advance in lock step.

- Lift the restriction of fetching one instruction at a time.
Then ideal CPI ≈ 1

- Fetch multiple instructions at a time.

- If all instructions are independent,

$$a = b + c$$

$$d = e + f$$

$$g = h + i$$

$$j = k + l$$

Then, fetch, say, 2 instructions per cycle.
Duplicate the entire datapaths to make two pipelines.
(RF, memory needs to have more I/O ports)

(Relax 1)
(Relax 2)

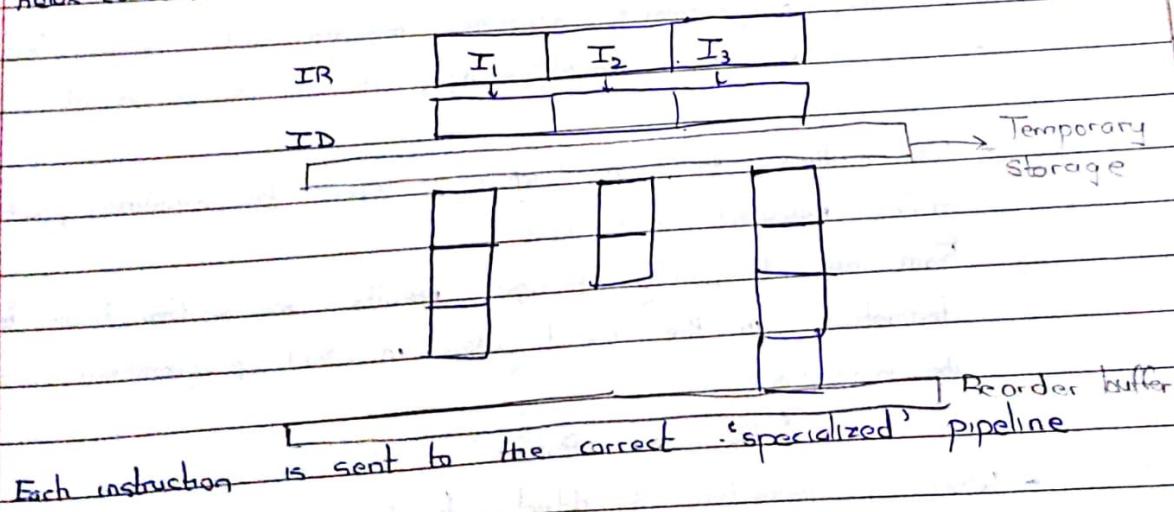
$$IPC > 2, CPI = 0.5$$

- If there is dependency between instructions,

$$\begin{array}{ll} I_1 & a = b + c \\ I_2 & d = a + e \\ I_3 & p = q + r \\ I_4 & s = p + t \\ I_5 & i = j + k \end{array}$$

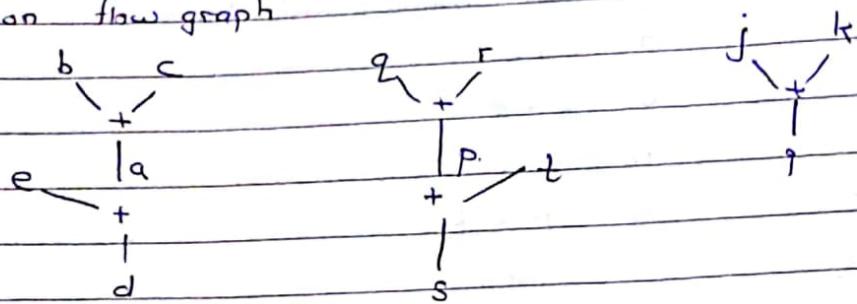
- Relax condition 3 :- Remove lock step advancement
- By-pass I_2 and I_4 for the time being.
- Execute I_1 and I_3 simultaneously first, I_2 and I_4 simultaneously next.

- Relax condition 2 :- Different pipelines have different number of stages.



- Between decoder and pipelines, make a temporary storage to store, say, 100 decoded instructions
 - In this storage, which instructions are independent is determined
 - Independent instructions are sent for execution simultaneously first.

- Instruction flow graph



Can be done in minimum 2 cycles $\Rightarrow IPC_{max} = 2.5, CPI = 0.5$

- Create the following mechanism:-

- Each instruction in temporary storage will check if its operands have been produced and are available.
Once they are available, it is ready to be sent to pipeline.

- But now, programmer's intended program order is violated
 - Will cause issues in debugging, interrupt execution, etc.

- Use another, temporary storage after the pipelines, which stores executed instructions.

From this temporary storage, results are written back to destination in the exact order in which programmer wrote the instructions.

'Re-order Buffer'

- When any instruction is added to temporary storage 1, a note is made in re-order buffer, so that instruction order is remembered