

Single Pipeline design - "Program guided execution"

for they follow instruction order.

Page No.			
Date			

VERY LONG INSTRUCTION WORD

I

→ Beyond scalar architecture

1 Fetch multiple (but contiguous) instructions

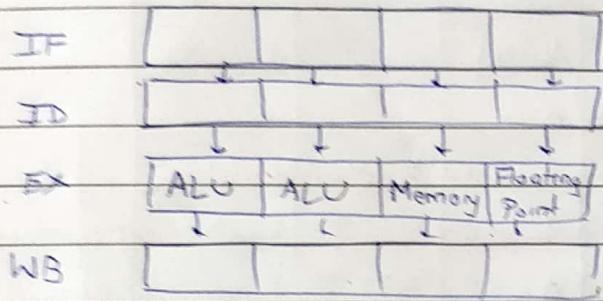
$$\text{eg} - a = b + c$$

$$d = a + e$$

$$p = q + r$$

Can execute I1 and T3 simultaneously.

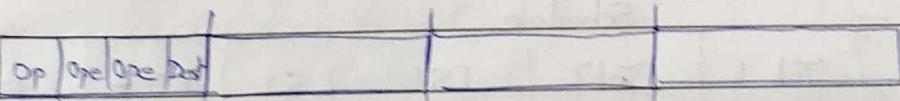
- Discover independent instructions → Software ⇒ Compiler
→ Hardware



All pipelines are specialized based on proportions of types of instructions on average

Compiler will try to improve efficiency by executing independent instructions in that specific proportion of types.

- This system uses new ISA, consisting of '4' combinations of individual instructions.



'Very Long Instruction Word' (VLIW)

- If another pipeline is added, entire ISA changes
- If specialized instruction is not found, NOP is executed in that pipeline

VLIW is tightly coupled with machine organization

- Immediate instruction dependency - No forwarding mechanism.
- Dependent instruction is executed only when previous result is written back
- Compiler is more complex in

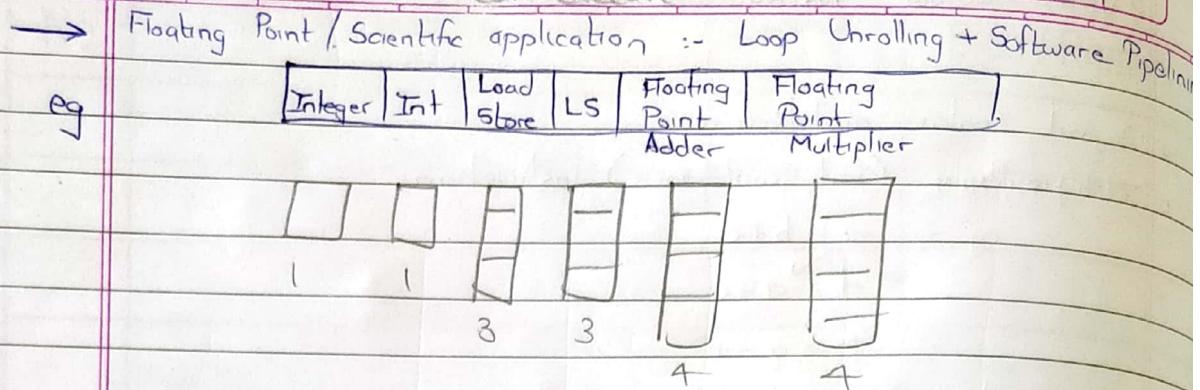
\$

Superscalar architecture is more sensitive to cache miss than scalar.

- More sensitive to Icache :- Stalling for more cycles
- Loss

D :- Even if miss, other instruction can execute

Page No.	
Date	



for ($i=0$; $i < 1000$; $i++$) { $A(i) = B(i) + C$; }

All floats

r_1 contains base $B(0)$

r_2 $A(0)$

r_3 $8000 + r_1$

f_0 (floating point register) contains C

Meth1 back: load $f_1, (r_1) 0$ I
 add $r_1, r_1, 8$ 2
 fadd f_2, f_1, f_0 3
 store $f_2, (r_2) 0$ 4
 add $r_2, r_2, 8$ 5
 bneq r_1, r_3, back 6

Schedule

Cycle	Int 1	Int 2	LS1	LS2	FPT	FP*
1	I2		I1			
2						
3						
4						
5						
6						
7						
8	I5	I6	I4			

3 cycle delay before dependency.

4 cycle delay

Blank spaces = NOPs

Number of instructions to unroll :-

Minimum = Maximum latency of any pipeline = 4

Maximum depends on resources (registers, etc.) available

Page No.

Date

∴ 8 cycles for one floating point operation (one iteration)

Terrible :-

Meth2 Unroll the for loop :-

for ($i = 0$; $i < 1050$; $i = i + 4$) {

$$A(i) = B(i) + C;$$

$$A(i+1) = B(i+1) + C.$$

2 2 ;

3 3 ;

}

Assembly	back	load $f_1, (r_1)_0$	Cycle	Total	Int1	Int2	LS1	LS2	FPT	FP
2	load $f_2, (r_1)_8$		1				I1 load 1	I2 load 2		
3	load $f_3, (r_1)_16$		2	T5			I3 load 3	I4 load 4		
4	load $f_1, (r_1)_24$		3						IG	
5	add $r_1, r_1, 32$		4						fadd 5	
6	fadd f_4, f_1, f_0		5						I7 fadd 6	
7	fadd f_5, f_2, f_0		6						I8 fadd 7	
8	fadd f_6, f_3, f_0		7						I9 fadd 8	
9	fadd f_7, f_4, f_0		8							
10	store $f_5, (r_2)_0$		9						store 5	
11	store $f_6, (r_2)_8$		10						I11 store 6	
12	store $f_7, (r_2)_16$		11		I14	I15	I13 bneq			
13	store $f_8, (r_2)_24$									
14	add $r_2, r_2, 32$									
15	bneq r_1, r_3, back									

∴ Performance = $\frac{4}{11}$ floating point operations per cycle.

∴ Higher use of cache and no. of registers. The more you unroll.

- * VLIW requires 1 cache BW than superscalar
- = ISA has ~~more~~ NOPs included (\uparrow no. of bits) in VLIW
- * VLIW does not need branch predictor

Page No.		
Date		

→ Better :- Pipeline each unrolled iteration

- Say, for eg., use LS1 for load only & LS2 for store only

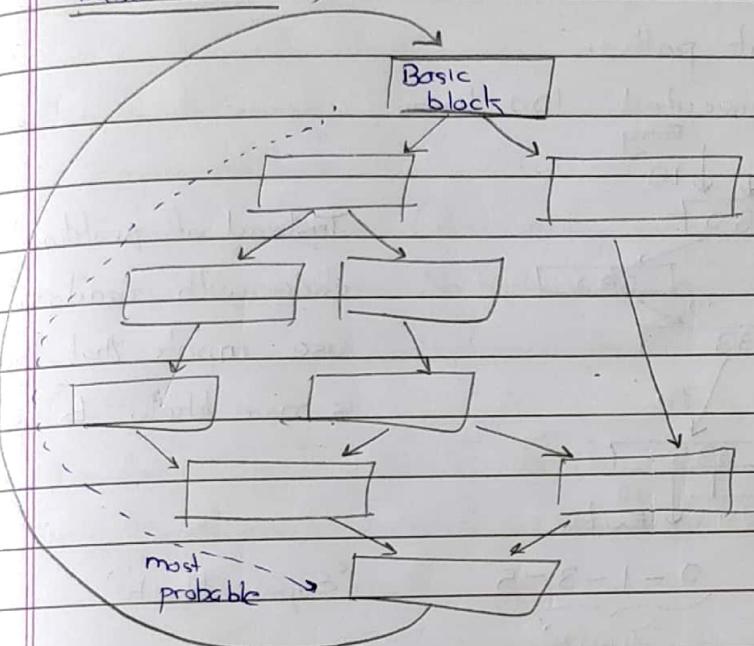
Int 1	Int 2	LS1	LS2	FP+	FP*
1		LF ₁			
2		LF ₂			
3		LF ₃			
4	addr,	LF ₄		fadd F ₅	
5		LF ₁		fadd F ₆	
6		LF ₂		fadd F ₇	
7		LF ₃		fadd F ₈	
8	addr,	LF ₄	SF ₅	fadd F ₅	This block contains
9		LF ₁	SF ₆	fadd F ₆	all instructions
10		LF ₂	SF ₇	fadd F ₇	in each unrolled
11	addr, bneq	LF ₃	SF ₈	fadd F ₈	iteration.
12	addr,	LF ₄	SF ₅	fadd F ₅	
13			SF ₆	fadd F ₆	
14			SF ₇	fadd F ₇	
15	addr, bneq		SF ₈	fadd F ₈	
16			SF ₅	:	

∴ Performance - One floating point operation per cycle

- * 'Lock-free advancement' - No checking anything ~~at~~ once instruction enters pipeline.
 - No stalling (even in branches, because we can schedule branch only when ~~instruction~~ destination is computed)
 - Stalling occurs only on cache miss.

→ Non-loop code :-

Usually, there are several branches. We can only correctly schedule instructions between consecutive branches (in one 'basic block')



While compiling, run a few sample test cases. Choose a static branch prediction, based on your findings. Then schedule instructions considering most probable path.

- If most probable path is taken, life is chill.

Most probable path = 'Super block'

- If most probable path is not taken → high penalty

- Compensation must be done - Changes of wrong prediction may have to be undone.

14/1 * Basic block \triangleq Must enter from the top, leave from the bottom

- If we try to predict just one branch ahead - no major benefit.

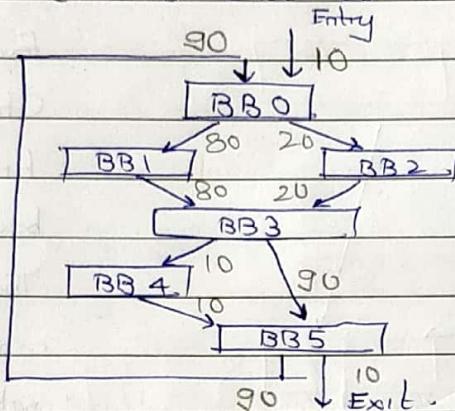
eg

$$\begin{array}{ll}
 a + b + c & (1) \\
 d \downarrow & (1) \\
 d = a + e & (2) \\
 p = q + r & (1) \quad \text{not } (1) \\
 s \downarrow & (2) \\
 s = p + t & (2) \\
 i = j + k & (3) \quad (2) \\
 m \downarrow & (3) \\
 m = i + l & (4) \quad (3)
 \end{array}$$

4 cycles 3 cycles

- To improve performance, increase no. of instructions and variety of type of instructions within a basic block (to keep the pipeline filled)
- Finding most frequent path..

Following code is executed 100 times



Instead of 'profiling' the code with random inputs, use inputs that the code is more likely to see.

Most probable path :- $\underbrace{0 - 1 - 3 - 5}_{\text{'Trace'}}$

'Super block'

Treat super block as a basic block and schedule all of its instructions together. (better spatiality of reference) \therefore

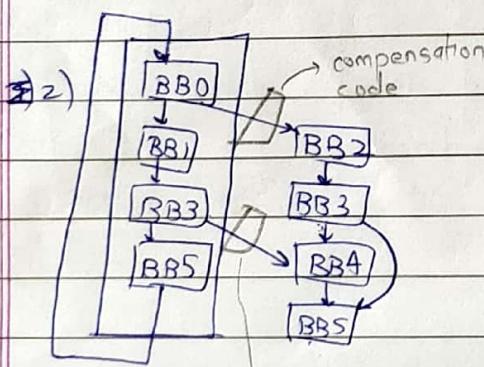
- Cache utilization will also improve, because BB2, 4 need not be brought into cache.

- Deviations from super block path.

e.g. - You want to go to BB2

- Two possible policies :-

- ~~Exit from superblock not allowed in between~~
~~Entry into superblock not allowed~~
~~(allowed only after BB5) (\uparrow compensation code)~~
~~Exit from superblock allowed~~
Exit allowed from between (BB0 or BB3)
- ~~Tree versa~~
Entry back into super block is not allowed.



Compensation code is required because incorrect instructions from further blocks have to be executed.

compensation
code

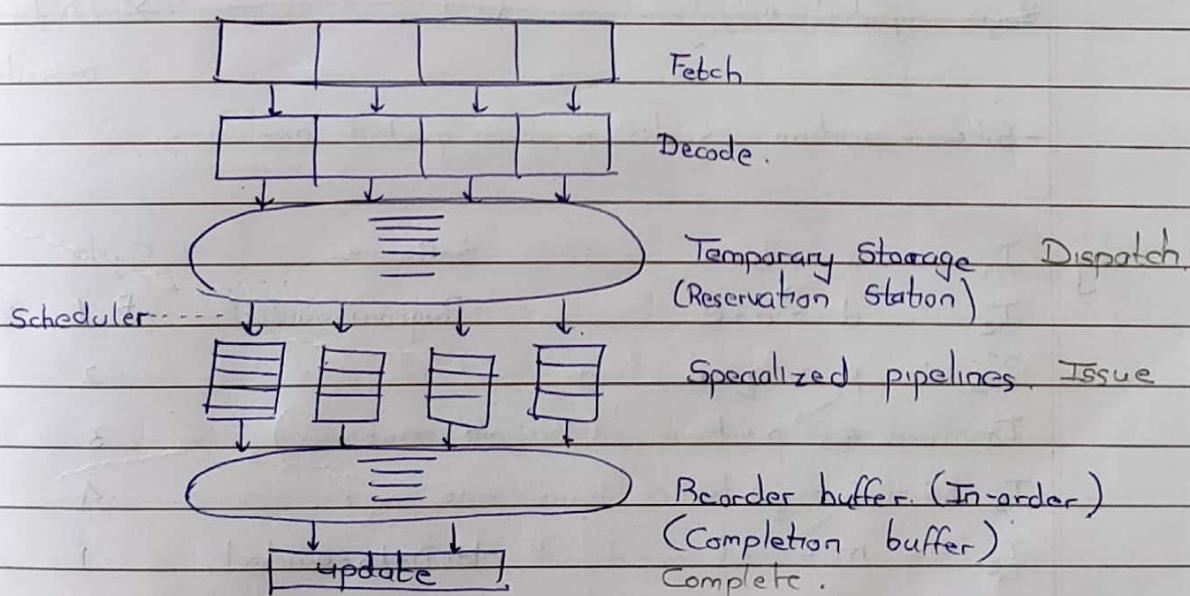
- Super-block formation
 - No loops in the middle
 - Every path's probability \geq Some predecided threshold (say 85%)
(∴ frequency of deviation is low)

→ VLIW Summary

- Simple machine
- Complex compiler (discovering independent instructions)
- ISA changed - No backward compatibility
(Code from older machines cannot be run on this)

- We don't want ISA to change.

→ Hardware, instead of compiler, should be made to discover independent instructions



Temporary storage ~ Instructions are scheduled when all their operands are available.

- It is important to maintain order of instructions while writing back (for debugging, interrupts, etc)
- * When an interrupt occurs, all further instructions have to be fetched again.
- i. Use another temporary buffer at the end of pipelines.

17/1 *

Software Pipeline

EE309

Hardware Parallelism Detection

Fetch single instruction	✓	✗
Uniform pipelines	✓	✗
Lock-step advancement	✓	✗
Program-guided execution	✓	✗ (Data-guided execution)

- Software pipeline is also program guided execution.

e.g

$$I_1 : a = b + c$$

$$I_2 : d = a + e$$

$$I_3 : p = q + r$$

$$I_4 : s = p + t$$

$$I_5 : i = j + k$$

$$I_6 : m = i + l$$

Software

Program-guided -

Cycle Instruction

1 1

2 2,3

3 4,5

4 6

⇒ Data-guided

1 1,3,5

2 2,4,6

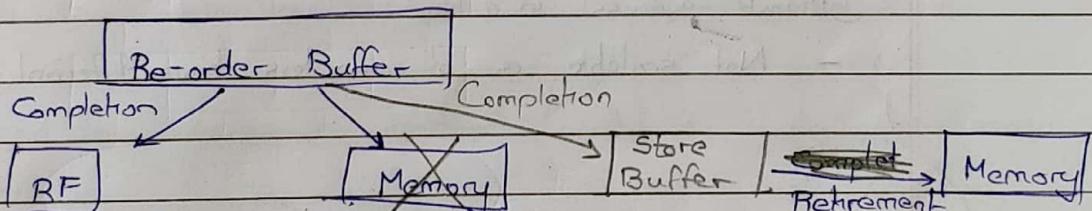


Problems:- Interrupts, Debug. → Need to write back in program order and maintain a PC as well.

Programmer perceives program-guided execution.

- How to maintain order before updating?
 - Remember (pass through every pipeline) PC of every instruction
 - X, because of branches, loops in code.
 - When you fetch an instruction, make a corresponding entry in Re-order buffer immediately
 - Instructions remain in order till decode stage
- * For stability :- No. of instructions that can be written back
 = No. of instructions fetched.

- Reorder Buffer ~ Circular Queue of instructions
 - When result of (consecutively next) instruction is available, write back and deallocate its place in reorder buffer.
- Writing back is ceremonial.
 - Any future dependent instruction can get its operand from the results stored in re-order buffer itself.



- Memory can have only one port
- Many Load instructions might be in re-order buffer
 - Load execution is always in critical path of pipeline.
 and is important for future instructions.
 If correct value is not loaded in time, it might cause a future cache miss.
 - Among Load/Store, Load must be prioritized.
 - Use a Store buffer to store results of store and write back, ~~from SB to memory only when~~
 memory is not busy

Challenges

- We want → High Bandwidth Fetch Improve prediction rate
- Cache Miss for ~~Cause~~ much instruction cache.
 - Branches - Many fetched instructions go waste.
 - When there are many branches in the code, overall probability of correct branch prediction reduces due to cascade
 - Branch target misalignment
 - Memory organization: (4 bytes each)
 - If 4 instructions are fetched at a time, memory is organized in multiples of 16 bytes. Can only fetch such 4 instructions together.
 - 0 | 1 | 2 | 3 | 4 |
 - 16 | 5 | 6 | 7 | 8 |
 - 32 | 9 | 10 | 11 | 12 |

If 4 is a branch whose destination is 11, then 9 and 10 would still have to be fetched along with 11 and 12.

- Solution 1 = $\frac{xx}{\square}$

Compiler allocates destination of branches so that they lie in the leftmost address in a block.

- Not scalable as all computers have different memory architecture.

- 2/1
- To execute 4 instructions, we also need to write back four instructions per cycle. ($\text{Inflow} = \text{Outflow}$)
 - RF can have more than one part.
 - Memory cannot, as it will increase memory size and access time.

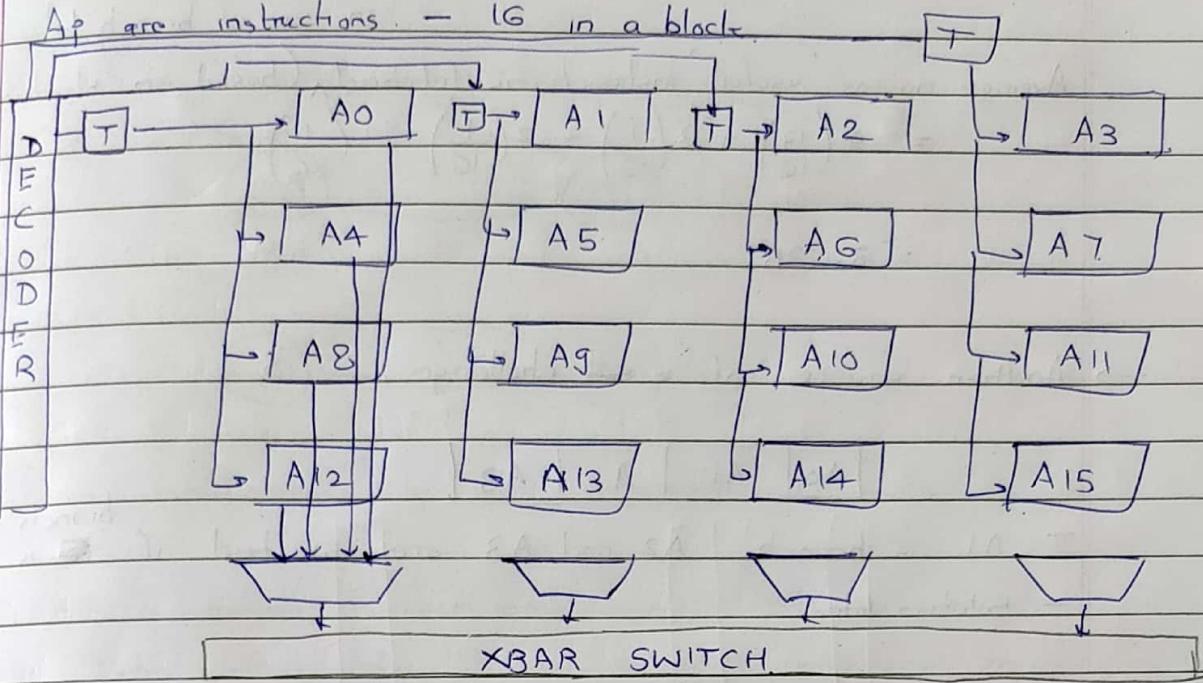
IBM's solution

One cache block has 64 bytes (16 instructions)

This block is arranged over 4 rows. (PTO)

Cache is 2-way set associative.

A_p are instructions - 1G in a block



- In any cycle, you can read one instruction from each ~~row~~ column
- Every column has a T-logic
- If starting address to be fetched is A2, fetch A2 and A3 from columns 3 and 4, but A4 and A5 from columns 1 and 2.
- No use if starting address is A13, A14, A15

- Need to reorder :- Instead of A₂ A₃ A₁ A₅, you seem to be reading A₁ A₅ A₂ A₃.
- Done by 'XBAR Switch'
 - This is big design. May not have 1 cycle access time.
 - Possible solution - Pipeline this fetching of instruction
 - Cycle 1 : Match from starting address
 - 2 : Fetch 4 instructions
 - 3 : Reorder
 - Bad in case of branch prediction taken or branch misprediction
 - Second solution :- Increase cycle time

In case of branch taken, and no

Average no. of 'useful' instructions fetched, (based on starting address)

$$= 4\left(\frac{1}{16}\right) + 3\left(\frac{1}{16}\right) + 2\left(\frac{1}{16}\right) + 1\left(\frac{1}{16}\right)$$

$$= 3.625$$

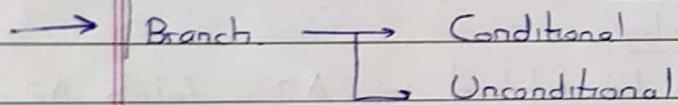
→ Another source of ~~challenge~~ challenge.

A ₀	A ₁	A ₂	A ₃
----------------	----------------	----------------	----------------

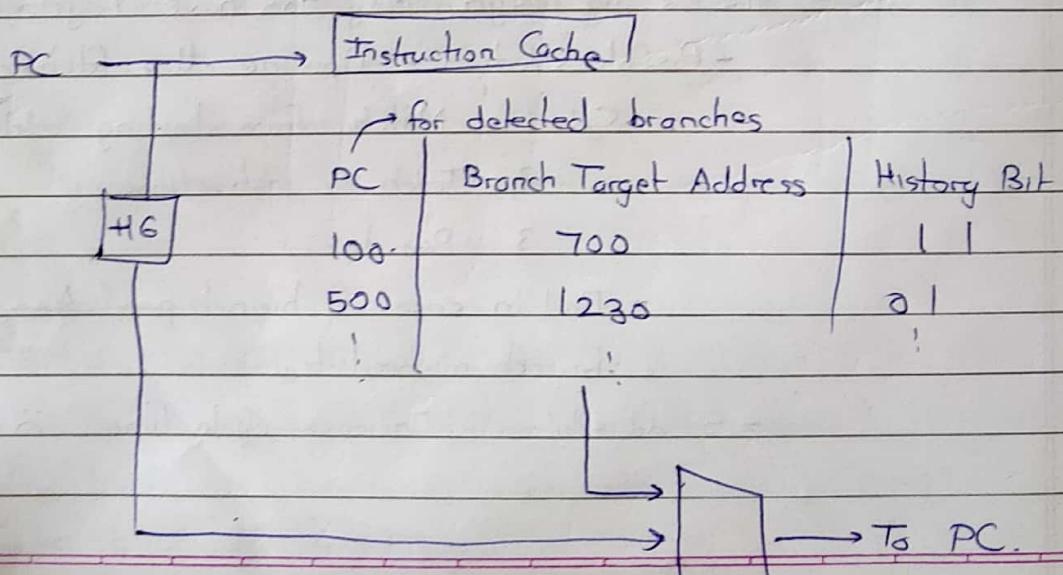
If A₁ is branch, A₂ and A₃ are wasted if ~~is~~ branch is taken.

- Solution later

BRANCH PREDICTION

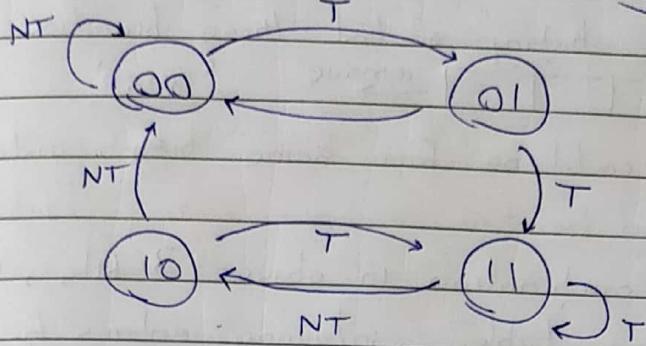


- Unconditional branch instructions can be detected and acted upon ~~by~~ in decode stage.
- For conditional branch, we need to know:-
 - 1) Direction (Taken or not taken)
 - 2) Target Address if taken.
- For direction :- Branch History Table.



cache miss

- Prediction depends on MSB of history bits. $\left\{ \begin{array}{l} O = NT \\ I = T \end{array} \right.$



'Change prediction only when mispredicted twice'.

- 85-90% better correct prediction \Rightarrow Not enough
 - Extra history bits do not increase accuracy much in

- Better :- Make prediction based on other nearby branches
(Determine correlation between branches)

eg - if ($i > 100$) { ----- }
 if ($i == 0$) { ----- }

If first branch is taken, second will never be.

- Size of History table :- ('Pattern History Table')

- Cannot be very large \Rightarrow Access time $>$ 1 cycle.
- Will have to remove some other instruction if overflow

- Lose all learning for the removed branch.

$\boxed{1}$ Separate 'Target address' and 'History' table.

$\boxed{2}$ - Instead of fully associative, make history table direct mapped

Fully associated
history table
size cannot
exceed 128

- If PC is 32 bits, only check entry corresponding

\hookrightarrow last 10 bits

- Total 1024 entries in table.

- Problematic if two branches (behaving differently)
have same last 10 bits of address.

'Aliasing'

Both will update history bits

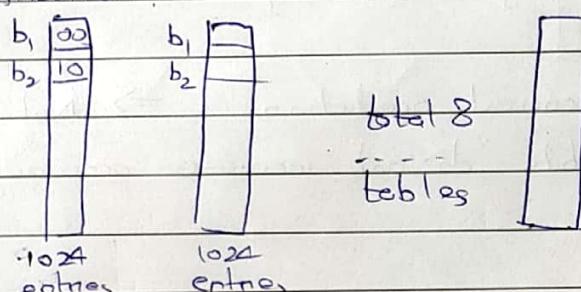
TWO LEVEL PREDICTOR 'Branch History Shift Register'

- Learning about environment :-

Maintain a separate history of last three branches
 update \rightarrow [] [] [] remove \rightarrow "BHSR"

The three entries could be from same branch instructions.

For each possible combination in above 3 bits, maintain one pattern history table, containing entries for every branch instruction.



Every combination of above 3 bits is an indication of the 'path' via which the current branch was arrived at.

- Predict your prediction based on the PHT corresponding to above 3 bits.

~~- If you consider~~

- If you consider these 8 PHTs as one long PHT, it will be indexed by 13 bits (10 from PC + 3 from BHSR)

24/1 eg

for ($j = 0$; $j < 100$; $j++$) {

 for ($i = 0$; $i < 3$; $i++$) {

(?) b₁

(?) b₂

In order

~~if~~

TTT NT T TTT NT T TTTNT

b₂ b₂ b₂ b₂ b₁ b₂ b₂ b₂

If your BHSR is 4 bits long, there will be 16 PHTs, each having one entry each for b₁ and b₂, each entry ~~containing~~ containing (two) history bits

- BHSR size :- Higher :- Checks correlation over farther branches \therefore
Exponentially high number of PHTs \therefore
- If BHSR size = 4 (say), you can either consider there to be 16 independent PHTs; ~~or~~ indexed by 10 bits of PC, or one large monolithic PHT indexed by both BHSR value and last 10 bits of PC.
- We can choose any BHSR size and no. of bits of PC such that
 \rightarrow BHSR size + no. of bits of PC = constant
= maximum possible that would allow for single cycle use of PHT.
eg BHSR size = 4, no. of bits of PC = 10.

6-SKU PREDICTOR

Better! Take BHSR size = 14, no. of bits of PC = 14.

- Index of PHT = BHSR \oplus PC bits
 - XOR operation randomizes the indices
 - High BHSR :- Farther correlation \therefore
 - Randomization + more no. of PC bits :- Lower aliasing \therefore

eg if $(i \% 2) == 0$ then b_1 (Taken if divisible)

{
if $(i \% 10) == 0$ then b_2

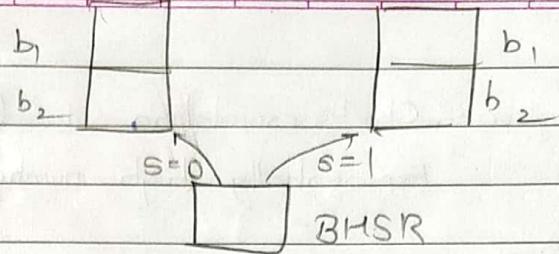
}

- BHSR size = 1 \Rightarrow 2 PHTs, with single history bit (all initially zero) (not taken))
- Inputs = {10, 15, 21, 22, 27, 37, 40}

Find prediction accuracy.

PHT0

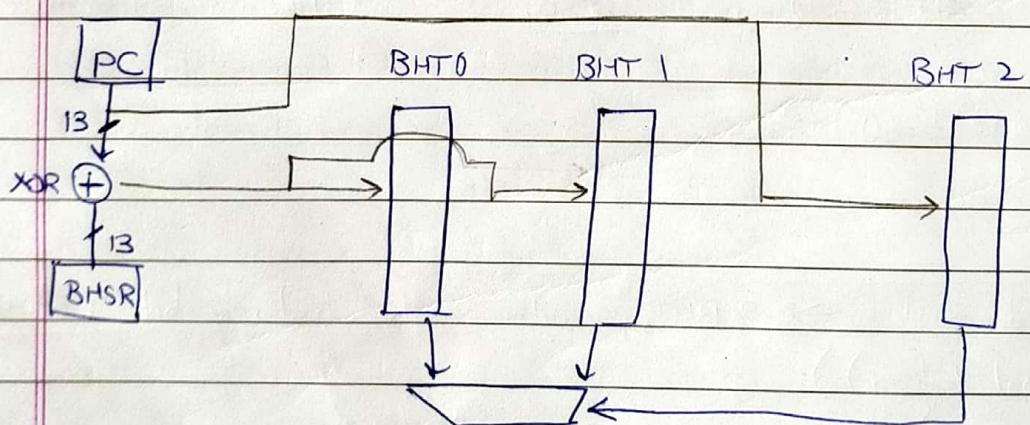
PHT1



i	10	15	21	22	27	37	40
$s=0$	b_1 predicted	0	1	0	1	0	0
	b_1 actual	1	0	1	0	0	1
	b_2 predicted	0	0	0	0	0	
	b_2 actual	0	0	0	0	0	
$s=1$	b_1 pred	0					
	b_1 actual	0					
	b_2 pred	0	1				0
	b_2 actual	1	0				1

Observation :- For $s=0$, b_2 is always correctly predicted
(if not divisible by 2 \Rightarrow not divisible by 10)

BIMODE PREDICTOR



BHT0, BHT1, BHT2 are entire PHTs, with each entry having 2 bits
BHT0, BHT1 are indexed by BHSR \oplus PC
BHT2 is indexed by PC alone.

∴ BHT0/1 and BHT2 are aliasing in different ways.

- Based on MSB in BHT2, pick BHT0 or BHT1.

Predict based on entry in picked BHT.

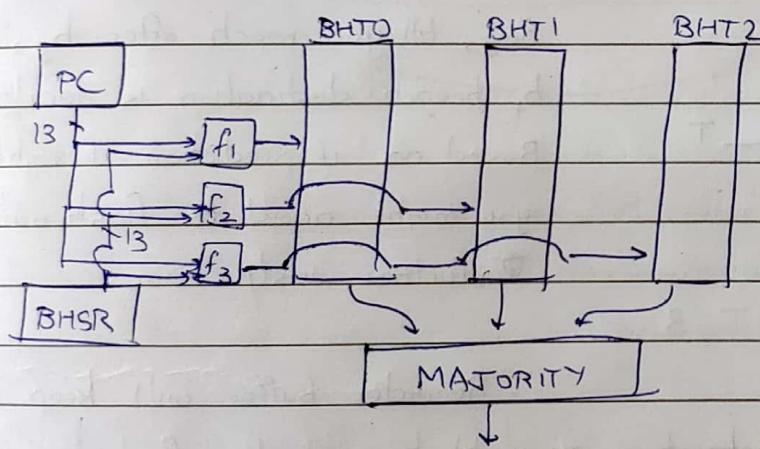
Update BHT2 and picked BHT

- ~~BHT0 will~~

- Eventually, BHT0 gets biased towards T/NT } If both aliasing branches behave
BHT1 NT/T differently.

'Dealiasing'

MAJORITY PREDICTOR

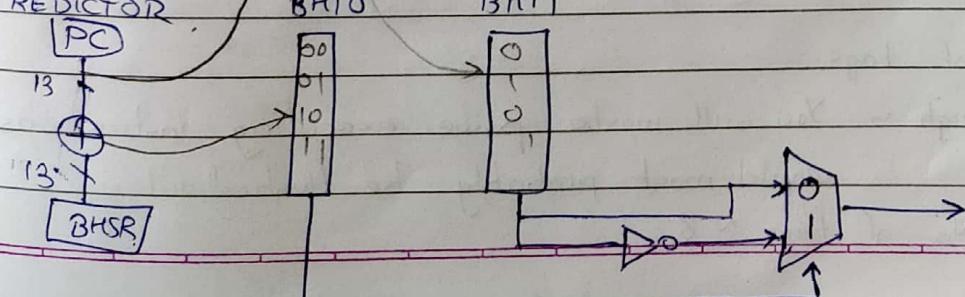


BHTs 0, 1, 2 are trained independently and indexed by different hashing functions of PC and RHSR.

∴ Aliasing is different in all 3 tables.

Final prediction is majority of all 3 predictions.

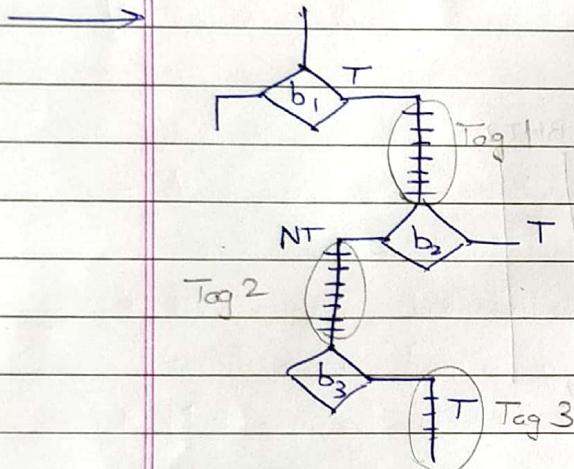
AGREE PREDICTOR



- BHT1 is filled up by compiler. This allows for better initial prediction, before hardware training is done.
- BHT0 is trained while running. Its entries tell us if prediction should agree with or differ from compiler prediction in BHT1.

11 → Strongly Agree

10 → Weakly Agree.



Say, you predicted $b_1 \rightarrow T$, $b_2 \rightarrow NT$ and $b_3 \rightarrow T$.

Say, till you reach after b_3 , your b_1 branch destination is decided.

Based on if prediction is right, you might need to flush out further instructions.

∴ Re-order buffer will keep track of whether an instruction ~~is~~ might have to be flushed out ('speculative' bit) and which Tag a speculative instruction belongs.

- If b_1 prediction is right, Tag 1 instructions are made non-speculative. Tag 1 is freed.
- If b_2 prediction was wrong, Tag 2 and Tag 3 instructions are made invalid. Tag 2 & 3 are freed.

- No. of tags :-

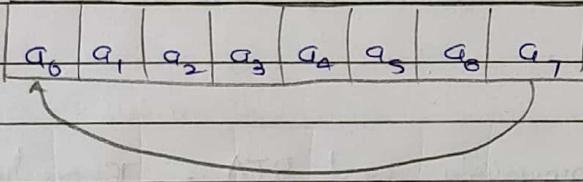
Too high :- You will waste power executing instructions that would most probably be flushed out.

∴ No. of tags ≈ 8 .

If all tags have been allocated, stop fetching new instructions till tags are freed (cause stalls)

→ "Return" Branch

- Unconditional jump, but you cannot predict destination, which depends on where that function was called from.
- store the return addresses for each ~~seq~~ function call in a circular queue on hardware
- We only know the correct address to return to in EX stage of return instruction



- If a return instruction prediction is resolved, erase its entry
- Queue will be full if there are 8 nested ~~func~~ function calls
If a ninth function is called, ~~its predicted address will overwrite~~
it has no more place on hardware queue.

Option 1: Stop predicting till a_7 is resolved (stall)

Option 2: Place a_8 in place of a_0 . While you continue executing, your returns will be resolved in reverse order (8-7-6---)

By the time you come to the 0th return, hope that its return address has been received from execution stage.

If not received, stall till received.

Predict NT - Actual T = Fine, because address will be same.
 T NT = Bad, don't have original address.
 - Every branch instruction must
 store PC.

Page No.	
Date	

28/1.

BRANCH TARGET ADDRESS

- We can make fully associative :- $PC \rightarrow BTA$
 - Size cannot exceed 128.
- Better
 - Make directly mapped :- (based on last 10 bits of PC)
 - Also store a 22-bit TAG to correctly match entire PC

TAG	BTA
(22)	(32)

To solve aliasing, you could store more tags

Multitargeted branches

Problem :- Branches that do not have a fixed BTA.

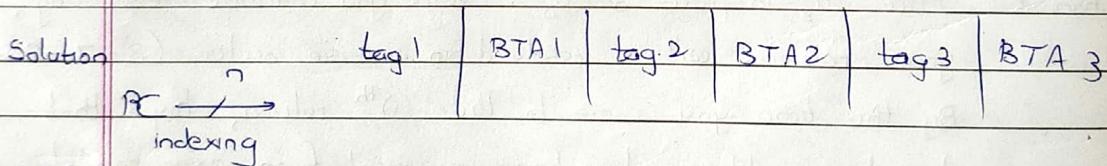
eg Switch Case

- Always store previously used BTA. If prediction is wrong, penalty.

* For example, branch could be multitargeted because its BTA could come from a register value.

* Multitargeted branches could come from object-oriented programming
 eg - Operator overloading

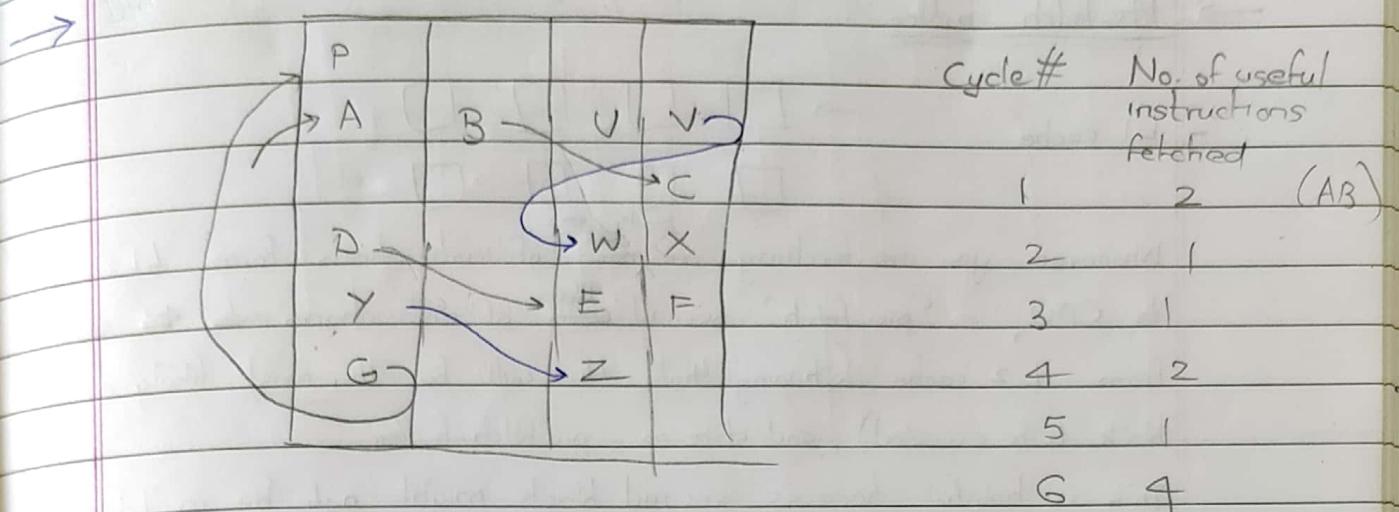
eg - Functions that behave differently based on no. of parameters



■ Multitargeted branches :- 2 of the tags will be same. BTA's will be different.

- 'Selection logic' - If more than one tag is same, choose one of the BTA's based on some pattern-based logic.
- Could be implemented by perceptron neural network.

We get 40% accuracy in multtarget prediction



Very few average useful instructions fetched per cycle

- Maintain a mapping in cache, based on PC of current (first) instruction 'A' in hardware

→ PC

PC of branch instruction		3 Predicted branches	Trace cache	Instructions in trace are entirely fetched and stored in trace cache
AB	AB	1 1 1 R	ABCDEF GP	
AB	AB	0 1 1	ABUVWXYZ	

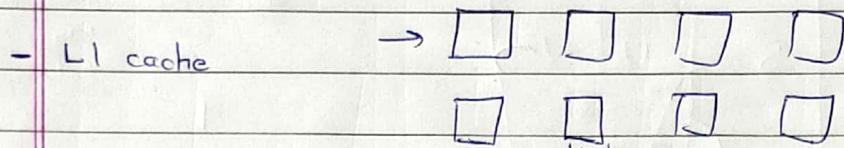
Only frequently taken paths are stored in cache, not all possible paths.

- Which traces are calculated is decided by code profiling, based on which instructions have been written back.

(Build trace cache at the time of retiring the instruction)

- Such a trace cache takes a while to be populated.

→ Pre-fetch buffer



Whenever you are fetching one ~~block~~ of instructions from L1 cache to CPU, 'pre-fetch' second ~~block~~ of instructions ~~to L1 cache~~ from L2 cache, (hoping that this will be the most likely next block to be executed), and store in pre-fetch buffer

This is helpful because second block might not be in L1 cache

→ Victim buffer

→ Instead of replacing old block in cache by new one, the old one is temporarily moved to victim buffer. Hence, if you need the victim block again you can ~~bring victim block to~~ avoid a miss by simply swapping L1 cache block and victim buffer.

- CPU can fetch one of :- L1 cache block, Trace cache, pre-fetch buffer, victim buffer, branch predictor.

- If the next required instruction is available in more than one source, preferably fetch from trace cache (because it is most likely to give most likely trace of instructions)

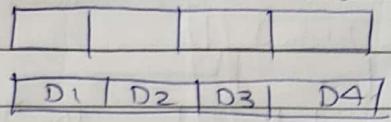
→ If there is a loop running many times, we can store all its instructions in another buffer.

→ Intel bypasses ~~Fetch and Decode~~ for further times by storing decoded instructions in trace cache.

3/1

'Intra fetch group dependence'

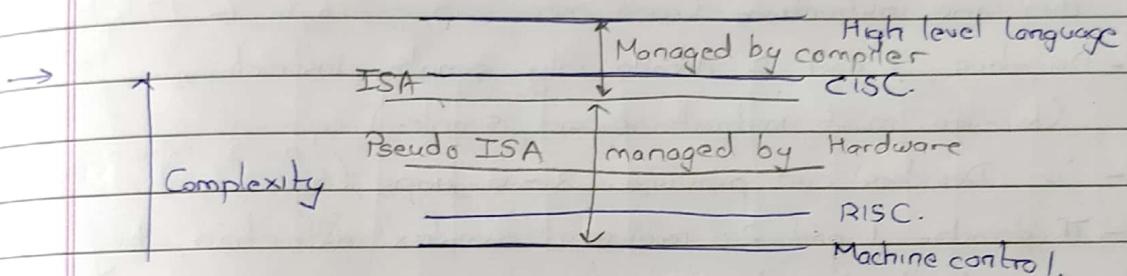
- Dependence of instructions fetched together on each other.



$2 \times C_2$ no. of dependencies possible

↳ each operand

→ This information is used during dispatch.



- Intel uses RISC-like ISA in background, but packaged for user as a CISC-like ISA.
ISA needs to be translated into pseudo-ISA by hardware.

- Every complex CISC-like instruction will be broken down into smaller RISC-like 'micro-operations'

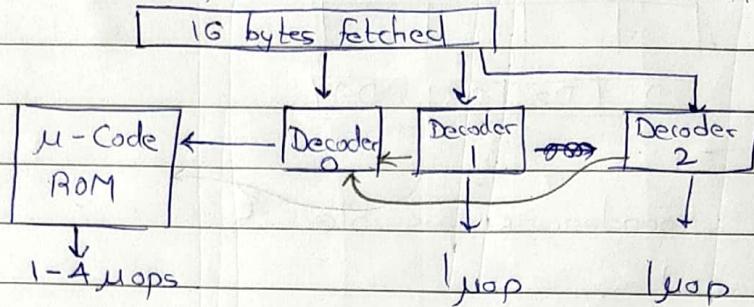
eg - add r₁, a

address	↓	}	load t ₁ , a
			add t ₂ , r ₂ , t ₁
			store t ₂ , a

- Simple instructions remain simple (eg - add r₁, r₂)

- ↑ Power consumption by hardware for translation into pseudo ISA.
↓ Efficiency because each instruction is translated one by one
- 16 bytes are fetched together = May contain variable number of CISC-like instructions

- Translation into pseudo-ISA. 'Look-up Table'



- More than one decoder because fetched bytes can contain more than one instruction (each instruction - 1 byte to 17 bytes)
- Decoders 1 and 2 decode only simple instructions.
(because most instructions are simple)
- Decoders can produce 3-6 μops. (1 or 2 instructions would have to be stored in buffer)
- If decoder 1 and 2 receive a complex instruction, they forward it to decoder 0 in next cycle.
- We aim to produce 4 μops per cycle (like before)

- Interrupts :-

All μops of one instruction will have one PC.

Interrupt may occur in the midst of μops sequence.

Maintain 'Atomic' bit :- 1, if all μops are indivisible.

Must be written back together, if at all.

- 0, if μop is simple instruction

→ Register Management - 'Write After Write' (WAW)

$$\begin{array}{l}
 \text{Definition} \\
 \begin{aligned}
 r_1 &= r_2 + r_3 && \text{Usage} \\
 r_4 &= r_1 + r_5 \\
 r_7 &= r_1 + r_6 \\
 r_1 &= r_9 + r_{10} \\
 r_5 &= r_4 + r_8
 \end{aligned}
 \end{array}$$

Both I_1 and I_4 produce

r_1 as result

IF I_2 is executed before I_1 and I_2 , I_2 might use I_4 's value for r_1 . \therefore

$I_1 - I_2$: 'True' dependency : 'Read after Write' (RAW) dependency

$I_2 - I_1$: 'False' dependency - Because of absence of lock-step advancement

Output dependency
(I_1, I_2)
'Write After Write'
(WAW) hazard

Anti-dependency
(I_2, I_1)
'Write After Read'
(WAR) Hazard

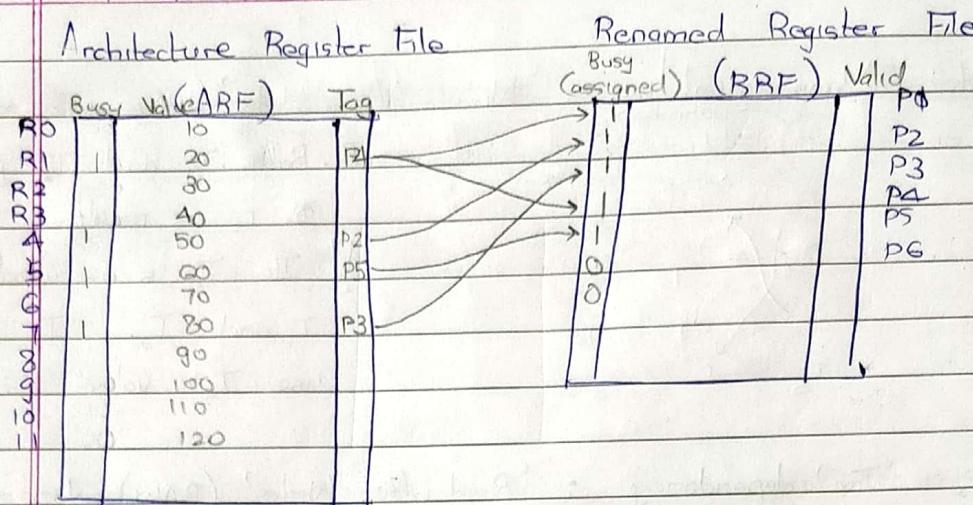
• Bad Solution: \therefore

still I_4 till I_2 is completed.

• Better:

Execute I_1 , but store its result somewhere else.

• Best :- 'Register Renaming'



Reservation Station

Inst	OP	Operand 1		Operand 2		Ready	Tag
		Value	Valid	Value	Valid		
1	ADD	30	1	40	1	1	P1
2	"	P1	0	60	1	0	P2
3	"	P11	0	70	1	0	P3
4	"	100	1	110	1	1	P4
5	"	P2	0	P4	0	0	P5

Ready = Valid 1 AND Valid 2

Re-order Buffer

Inst	PC	Dest	Tag
1	100	R1	P1
2	101	R4	P2
3	102	R7	P3
4	103	R1	P4
5	104	RS	P5

$$\begin{aligned}
 I_1 &: r_1 = r_2 + r_3 \\
 I_2 &: r_4 = r_1 + r_5 \\
 I_3 &: r_7 = r_1 + r_6 \\
 I_4 &: r_1 = r_9 + r_{10} \\
 I_5 &: r_5 = r_4 + r_1
 \end{aligned}$$

- Instructions are listed one by one in RS and RoB
- If both operands are available from registers, $\text{Valid 1} = \text{Valid 2} = \text{Ready} = 1$
- Else, we write tag (P1)

From where operand will come when P1 is valid (ie - if register is busy)

- A non-busy tag is allotted to destination register and Busy bits of both register and its tag are made 1.

- Begin by executing instructions whose Ready = 1 = I_1

Store result in P1 (and make P1 Valid)

\Rightarrow Instructions 2 8 3 become Ready

Move on to I_2 , then I_3 , I_4

- When I_1 retires (in lock-step advancement),

Result transfers from P1 to R1

P1 becomes not Busy

Maintain a list of free registers

P85 | P42 | P30

Allocate new RA's from this list.

Page No.		
Date		

It checks if tag of RI is equal to the tag which is retiring. (If yes, make RI not Busy)

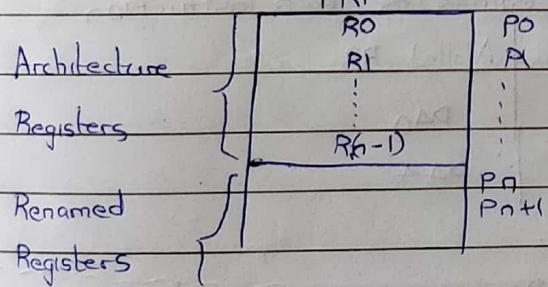
Because I+ also has destination RI, RI has tag $P_4 \neq P_1$. In this case, RI must be kept Busy.

A/2 • Value of operand I comes from a MUX with $\frac{2n+m}{2}$ inputs

- $2n$ = No. of instructions fetched per cycle.
- $2n$ = No. of operands required per cycle = No. of ports that ARF needs.
- m = No. of specialized pipelines, each of which can broadcast its result after execution.

• Instead of making separate ARF and RRF, make one monolithic RE

'PRF'



- During retirement of RR's (writing them back to AR's), significant power is used.
 - For allocation of unused RR's, looking for first non-busy register RR would require heavy, time consuming logic of Priority encoder
- Look ↑
- During retirement of RR's, don't move value in RR to corresponding AR. Rename that RR as AR.
eg - Make P40 'R0', instead of moving from P40 to R0.

- The instantaneous location in PRF of every AR is tabulated in Register Architecture Table (RAT) / Register Map Table (RMT)

RO	P40
RI	P23
:	:
:	:
R31	P54

5-bits \sim 8 bits

RAT is updated
only during retirement
of instruction.

∴ Cannot be used for
operand read.

- Speculative RMT

- Just like RAT, except that it is updated whenever execution of instructions is completed. Instruction is decoded.
- Can be used to judge where to get operands from.
- Since you are maintaining speculative RMT, you no longer need to store 'tag' for each AR.

-

- RoB contains following entries for each instruction

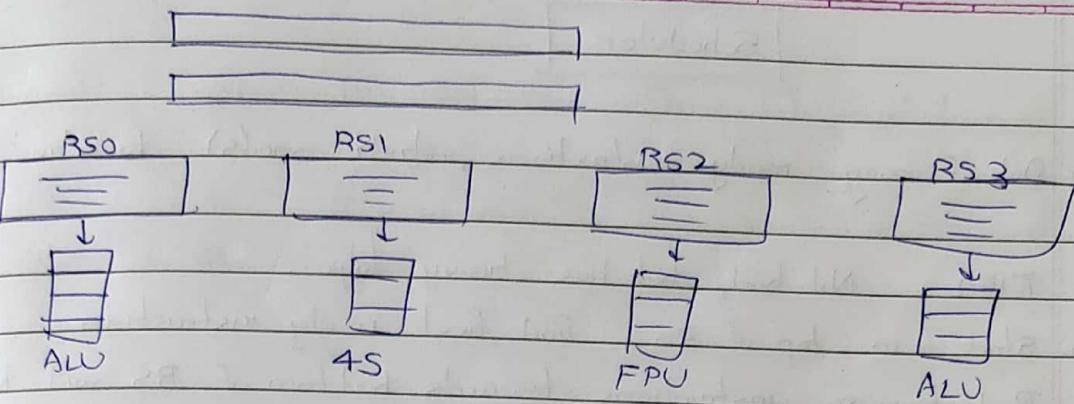
	Destination	Allotted RB
I1	RO	P40
I2	RO	P44

- ∴ Simply change RO entry in RMT to P40 instead of moving P40 value to RMT

‘Distributed Architecture’

- Instead of one Reservation Station, every specialized pipeline has its own RS, meant only for instructions that use that particular pipeline.

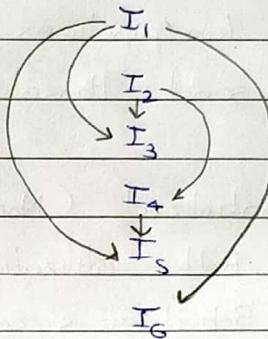
P.T.O.



- If two pipelines are identical, we need to decide which RS to send those type of instructions.
- Centralized (Non-distributed) RS :-
 - Better utilization of space
 - Scheduler is complicated : Needs to scan through many more instructions to send to different pipelines.
- Distributed RS :-
 - Scheduler is simple : Needs to schedule only one instruction (for its own pipeline)
 - Poorer utilization of space :
For some particular code, one RS might remain empty while other RS is overflowing
- Hybrid RS :-
 - Pipelines which are identical/similar are managed by same RS and scheduler.
 - Middle path

Scheduler

- Out of many ready instructions, which one(s) to schedule?
- 1. FIFO - Not bad, but has heavy logic.
- 2. Start from top of RS, find first ready instruction :-
Bad, because instructions towards bottom of RS will be starved
- 3. Maintain a dependency chain
Out of ready instructions, schedule the instruction on which most number of instructions depend.
 - 'Fanout' - Doing this will make more and more instructions ready to be scheduled.



Better! Schedule that instruction on which most ~~of~~ number of instructions having only this one dependency depend.

e.g. I_3 and I_5 will not become ready immediately after I_1 , I_6 will.

- These heuristic methods come with heavy overhead.

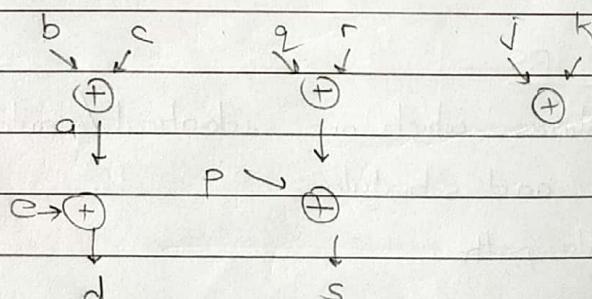
e.g. $a = b + c$

$d = a + e$

$p = q + r$

$s = p + t$

$r = j + k$



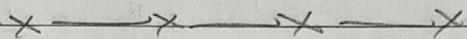
Control Flow Limit = 2 cycles . TPC = 2.5

By doing ALL we have learnt, we can achieve 'Control Flow Limit' or 'Data Flow Limit'

- Beyond Control Flow limit \rightarrow Operand prediction
 - * 40% times, we use value of register used last time.
 - * 60% times, we use one of the values used last 4 times
 - Need to verify prediction

eg Functions could return same output for same input. ~~Keep~~
Memorize the output in hardware.

* In software, it is called memorization.

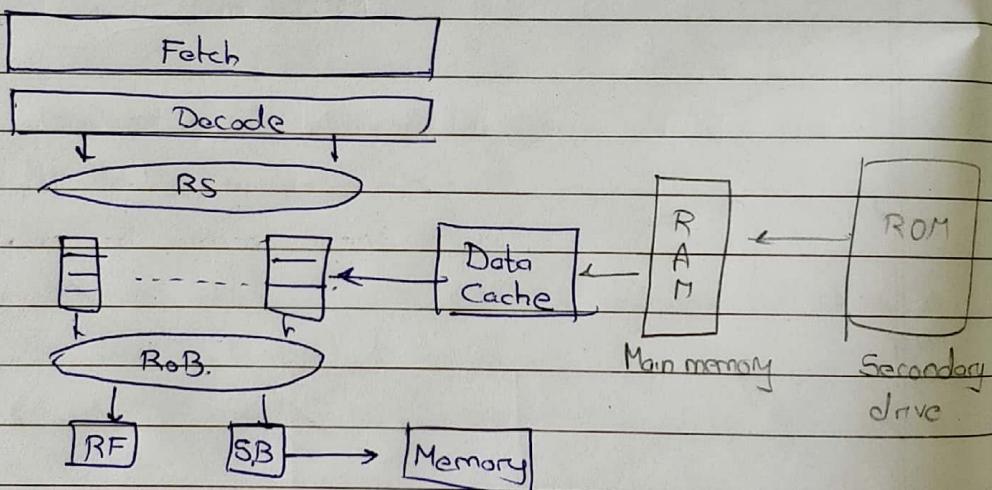


<break>

→ Memory Related Instructions.

load $r_1, (r_2)50$ // $r_1 \leftarrow M(r_2 + 50)$

store $r_5, (r_6)50$ // $M(r_6 + 5) \leftarrow r_5$



Stages of pipeline

- 1. Address computation :- $r_2 + 50$
 - "Virtual Address" - Not actual physical address.
 - VA space size = $2^{\text{No. of bits in address}}$
 - VA space divided into 'pages'
 - Pages for different programs (processes) are moved to main memory from virtual memory before executing secondary drive.
 - Page size :- High :- wastage
Low :- More memory (ROM) accesses required.
 $\sim 4\text{ kB}, 6, 16\text{ kB}$.

ROM has virtual pages

eg 40 bit secondary, 1 kB (10 bit) page size,
 4 GB (32 bit) main memory

\therefore In ROM \rightarrow [Virtual Page No.] [Offset]

30 10

In RAM \rightarrow [Physical Page No.] [Offset]

22 10

Code cannot be accessed from virtual memory. It must first be brought to main memory.

We want PPN from VPN

- Lookup Table - 'Page Table'

Virtual Page Number \rightarrow Physical Page Number
 (ROM) (RAM)

- Stored in RAM

- 2^{30} entries, each containing 22 bits \therefore (3 GB \times 22)
 \therefore those present in main memory

- Reduce size of lookup table. Use LRU policy for replacement of data.

Solution - 'Inverted Page Table': PPN (RAM) \rightarrow VPN (ROM)

- 2^{22} entries, each containing 30 bits

- Much smaller (15 MB) \therefore

- How to index? \therefore

We wanted to find physical page no. of virtual page no.

- Cannot always scan all VPN entries

Better - 'Paged Page Table' 'Translation Lookaside Buffer' (TLB)

- First time you want VPN of PPN, scan IPT
 (takes time)

- Enter this VPN \rightarrow PPN mapping into PPT, so that you don't have to scan again.

'Cache for VPN \rightarrow PPN mappings'

CPU works with 40 bit PC

• Steps

1 Generate address to fetch data from $(\frac{1}{2} + 50)$
virtual

2 Translate address (access TLB)

Hence obtain P main memory page number (and thus main memory address)

3 Read from main memory / data cache.

* Data and Instructions are stored in different TLBs, as they are in different pages in ROM.

If VPN entry is not present in TLB, scan IPT, find P and add entry to TLB.

- If entry is not present, TLB raises a page fault.

~~Hard~~

11/2

→ Store instructions cannot be stored in RoB.

- They have 32 bits address + 32 bits data.

- Allotting so much space for every instruction is wasteful, because not all instructions are stored.

∴ Store buffer is required.

When store enters RS, entry is made in SB instead of RoB. (ID of this entry is remembered) in its entry made in RoB.

SB			RoB		
Address	Data	Completed	Destination	RR	Valid
200	504	0	RS	P11 P41	0

Store Instruction

ID of Store Buffer

- Valid bit of RoB is set after Store is executed
- When store execution is completed, make Completed = 1
comes to top of RoB eventually
- 'Retirement' - Perform 'lazy write' back to memory (as it is not in critical path of processor)
- Only SB entries with Completed = 1 can be written back
- * SB is FIFO

- Hazard :- Load of same address following Store that hasn't been retired.
- Cannot wait for Store to be retired (as performance \downarrow a lot because of pending Load)
 - \therefore Cannot insist that all Loads and Stores are written back in program order
- Entry of Load store is made in RoB and SB when dispatched into RS.
 - Before scheduling Load, match its address against all addresses in SB.
 - If no match is found, \Rightarrow no dependency.
Allow Load to be scheduled
'Load Bypassing'
 - If match is found, there is dependency. Do not schedule Load
 - Sometimes, addresses in SB are only stored as last 10 bits of actual address.
 - If match is found, match might be false due to aliasing.
 \Rightarrow still do not allow Load Bypassing.
 - If SB addresses are entire 32 bits, and match is found, and Store in RoB is valid, Load can obtain data directly from SB (saves one memory access)

- If Store is not entirely executed and its address is not known, we have to 'speculate' matching.
- Easy speculation - No match.
- Need to detect if speculation was right or wrong and correct if ~~a~~ wrong.

- Detection

First check address of Load against 'Completed'

Store instructions.

Take no match speculation if no match is found.

- Take a 'no match' speculation. Execute Load.

- 'Load Buffer'

Execute Load as soon as address is available.

When you execute, make a note of this address in LB.

LB

Address

300

:

- When Store is 'Completed' (reached top of RoB), the address in SB is matched against all addresses in LB.
- If match(es) is found, that Load was wrongly speculated. Correction is required.

- Correction:-

- xx ① Repeat just Load and all of its dependent instructions again. We will also require all other operands of all instructions again in

- ② Treat as branch misprediction. Start fresh from that Load instruction.

- Need to store PC in LB. Or store pointer to RoB, which already has PC

- Dynamic match speculation

 - Requires PHT-like tables

 - Profile load-store behaviour.

- ① If matched before, it is likely to be matched again

- Do not execute such load-store ~~is~~ out of order.

- ② Make Load-Store pairs that are aliasing ~~is~~ (matching) ↗

LB

3/2

- * 'Total Store Order' type architecture - All stores are performed written back in programming order.

→ Load Store Aliasing

i) eg Spill Code:-

When all registers are occupied, value of some register must be replaced.

$q = p + r$ // r initially contains b . All registers are

occupied. p is in memory

$t = l + b$

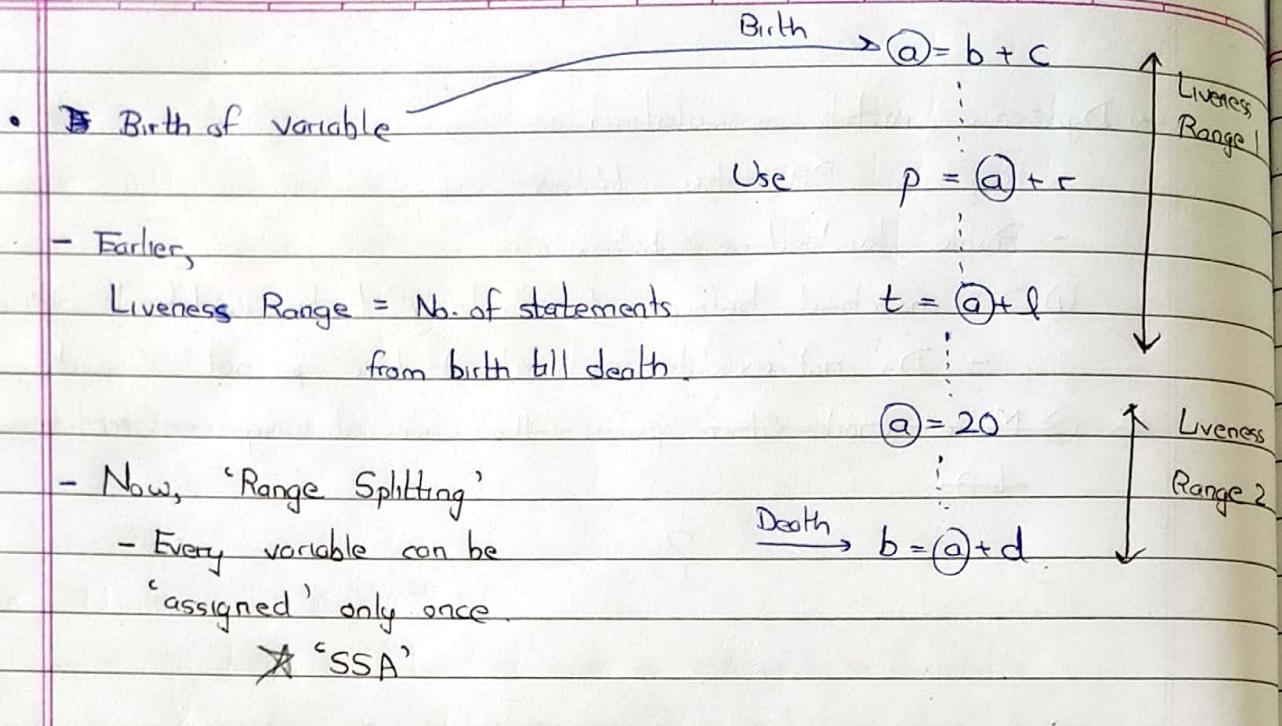
~~p~~ \rightarrow

b must be STOREd into memory/cache.

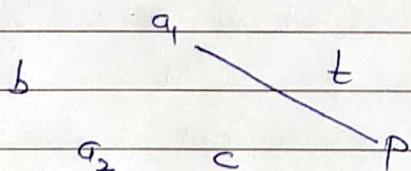
p LOADED into RF

Later, for $t = l + b$, b must be LOADED back

This generates an aliasing LOAD-STORE pair for b .

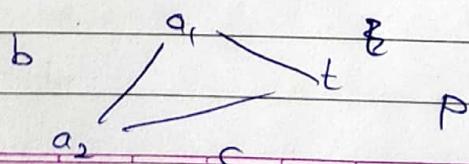


- Variables having overlapping liveness range should not use same register. Non-overlapping variables can. (eg - a_1, t, a_2)
- 'Conflict Graph' - Shows overlaps as edges



'Colouring Problem' - No. of colours required to paint the nodes, so that adjacent nodes do not have same colour

- = Minimum no. of registers required to not have to LOAD/STORE ~~to~~/from/to memory.
- If you have fewer registers than this, need to choose which register(s) to move to/from memory.
 - Alternatively, 'Anticonflict graph' - No overlaps



Lower no. of registers of ~~Y~~ Spill Code

~~Y~~ Clique \triangleq Fully connected subgraph (a_1, t, a_2)

\triangleq Can share same register.

\therefore Minimum no. of registers = No. of cliques.

- Spill code causes aliasing of load-store. (store followed by load of same address)

2) Pointers

Address is dynamically generated (only known at runtime)

- Can cause load-store aliasing.

e.g. store r_1, r_2 } if r_2 and r_5 have same value
 load r_5, r_6

• Higher Fetch Width (≥ 4)

- ~~Prefetch~~ We will fetch say 3 instructions, which are likely to contain LOAD and STORE.

In our attempt to bypass LOAD, this will cause accumulation of instructions.

\therefore We are forced to increase ports on memory.

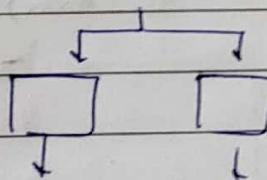
\uparrow size, cost, access time

Better - Use two identical data caches

- Must be written to simultaneously

- Can read independently

- Can perform LOADS twice as fast

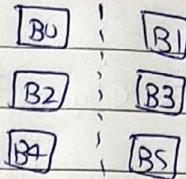


'Replication'

P.T.O.

Alternate Single Cache

- Odd-numbered and even-numbered cache blocks can be read from separately.
- Even-odd division, because we usually read from consecutive addresses.
- Cheaper, not as good as replication.



• 'Non-blocking Cache'

eg - Load X, Load Y executing simultaneously.

X is not in L1 cache, Y is.

Y should be read and not wait for X to come to L1 cache
 'Non-blocking'

- While X is being fetched, the address 'X' is stored in 'Miss Status Holding Register'

MSHR

100

300

:

- Data is brought to L1 cache in 'block' granularity.

- Multiple requests in MSHR can be serviced (possibly out of order) if they are present in same block.
- Reduces average access time.

14/2 → Memory-related Instruction Handling

- 1 Non-blocking D-cache (L1) :- MSHR
- 2 Speculative execution of loads
- 3 Multiported cache
- 4 Use of victim buffer
- 5 Prefetching.

→ Pre-fetching

--- of next likely block of data

- Can be done by software or hardware.
- Data can be prefetched into memory (L1 cache) or prefetch buffer.
 - Whenever there is a hit on block in prefetch buffer, it is swapped out moved to L1 cache before being used.
- Many ISAs provide a 'Prefetch (Block #)' instruction.
 - Prefetches into L1 D-cache.

- We may not know which block # needs to be fetched.

Prefetching
useless
instruction

- If you prefetch a block and never use it, it is bad because it replaced another block that could have been used.

Prefetch
buffer
is less
bad

- We cannot know which block # to be fetched in unstructured program.

eg - Can be useful for dense matrix computation

eg - ~~Can~~ Prefetched data is useless for sparse matrix multiplication

most elements zero.

eg - Mostly used in loops.

- "Cache pollution" - Bringing in useless blocks into cache.

- Can occur because of both hardware and software prefetching

- When to prefetch? (Total time to fetch \approx 150 cycles)

- Too early - Can be replaced before use. (mostly useless)

- Timely - Fetch only as early as is needed to make data available when needed.

- Late - Data does not reach on time.

Cache miss is registered.

But still saves some cycles as opposed to actual cache miss

- Where to prefetch?

- Software can only prefetch to software L1 D-cache
- Hardware can prefetch to cache or prefetch buffer

Software

e.g. 100 x 100 matrix multiplication

$$\begin{bmatrix} \dots & \dots \end{bmatrix} \begin{bmatrix} \vdots \\ \vdots \end{bmatrix}$$

One computation :- 100 multiplication followed by 99 addition

Say :- Operation of row 1 on all columns takes 100 cycles

Say :- Fetching one row from memory takes 200 cycles

- ∴ Fetching next row : Too early
- ∴ Every row & prefetches next-to-next row : Timely
(Non-prefetched rows will cause cache miss)
- Cache misses are served preferentially to prefetching stuff

Hardware

e.g. You are fetching from an array : A(1), A(5), A(9) ...

'Stride' - Detect pattern :- 4 elements apart

Prefetcher - All these fetches are coming from same load.

$$\text{load } r_1, (r_2) r_3$$

$$r_3 = r_3 + 16$$

- Hardware maintains load PC v/s Address v/s Difference

PC	Address	'Stride' = Difference	Confidence
500	132	16	1
	100	-	0

When stride remains same for some predecided no. of executions make Confidence = 1, and start speculating which address will be needed next.

- Per-bit fetch - (\uparrow energy, BW, performance)

Whenever there is a hit in cache, prefetch the block that contains the earliest speculated data that is not already in the current block of cache.

- Can create traffic, because you are checking if next required block is fetched at every hit (hit rate is high)

- Per-miss fetch:-

Whenever there is a miss in cache, prefetch the *next* block

- Must have one miss, unlike per-bit prefetch

- For timely fetch, we should not always fetch the very next block

'Prefetch depth' :- How many blocks away are you fetching your block from?

'Prefetch degree' :- How many blocks to prefetch at a time.

- stride prefetcher is certainly useful for arrays

- Also useful for linked list.

- Different patterns :- $A(1), A(5), A(9), A(29), A(35), A(3S)$

$\underbrace{A_1}_4, \underbrace{A_2}_4, \underbrace{A_3}_4, \underbrace{A_4}_2, \underbrace{A_5}_2, \underbrace{A_6}_2, \dots$

18/2

DATA PREDICTION

Going beyond data flow limit.

- Can advance future instructions to current cycle.
 - Thus, resource utilization can be increased to beyond 30%.
- Need to validate predicted value
 - 'No-validation'
 - We are sure about prediction. We have checked its history and the value will not change
‘Instruction Reuse’
 - Most values do not change for a long time
 - eg - Array operation ($B + D$) :- B is constant.
 - eg - Most of the comparison of values is against zero
 - Sometimes, there is a pattern in the values over iterations
 - eg : $i, i+4, i+8 \dots$
 - eg - $a = b * 16$
If b does not change, neither will a .
 - eg - Factorial (n)

eg

As you iteratively calculate, store in a table

PC	Argument	Result
100	3	6
100	5	120
:	:	:

★ This should not be done when operating on global variables

- Not advantageous as they can have too many values.

* Local variables are the ones stored in stack and accessed using SP.

- Such tabulation of function results need not be validated.
- Tabulation can be done for functions or even single (critical) instructions.

- Predictions must always be validated

LAST - Table for results produced

LAST VALUE PREDICTION	PC	Value	Last Value	
			Confidence Counter	Last Value
	100	786	0001	
	:	:	0001	

- Instructions dependent on this I_{100} can fetch the result from here

- Alternatively, First column could be ~~value~~ PC of dependent instruction and it will map against its operand value used to be used.

- When same result is obtained same predicated no. of times, make Confidence Counter ± 1 .

- If $CC \geq 2$, that value can be used

- Used value will be flagged as 'speculative'.

- If some other value is obtained, make $CC = 0$.

- New result can be matched with table value as soon as execution is complete

STRIDE
BASED
PREDICTOR

- Sometimes, values are in AP. Modify table.

PC	Value	Stride	CC
100	786	16	01
:	:	:	:

- If same stride occurs again, $CC + = 1$.

Eg:- Instruction :- $r_1 = r_1 + 6$.

- If values remain constant, stride remains 0 \Rightarrow

Eg 30, 40, 50, n n m m , 30, 40, 50, m m m

Stride predictor fails to recognize pattern

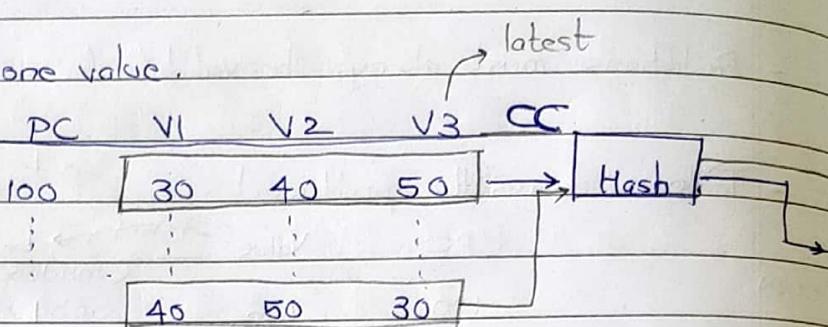
P.T.O.

FINITE

CONTEXT • Profile more than one value.

METHOD
(FCM)

Value Table :-



- Output of hash function gives index within prediction table
- If predicted value is correct, CC for that pattern $t=1$
incorrect $= 0$

* eg - Fetching 8 instructions per cycle :- TPC < 3

→ Poor utilization of resources

(~70% resources idle at all times)

- Poor utilization because of :-
- 1) Dependence of instructions
 - 2) Constraint to write back in program order
 - 3) Cache misses

* 'Software Memorization'

- Even software maintains past values to predict future values for functions

* To increase performance, we can increase independence of instructions.

Simultaneously fetch 8 instructions from more than one program, which will be inherently independent.

$$8 = 4 + 4 \quad \text{or} \quad 8 = 1 + 1 + 1 + 1 + 1 + 1 + 1 + 1$$

P.T.O.

Single Pipeline.

Prediction Table.	Cycle	IF	ID	EX	Mem	WB
	1	I ₁ ¹				
	2	I ₂ ²	I ₁ ¹			
	3	I ₃ ³	I ₂ ²	I ₁ ¹		
	4					
	5					
	6					

- Each program perceives multicycle execution.
- No need for branch predictors, forwarding mechanism, etc.
- Higher throughput.
- But lower performance of single program.
- Superscalar can have more such pipelines.
- eg - Fetching 8 instructions per cycle :- IPC \approx 6.

★ Register conflict

- How to improve performance of even a single program?

eg :-
$$\begin{bmatrix} \text{---} \\ A \end{bmatrix} \times \begin{bmatrix} \text{---} \\ B \end{bmatrix} = \begin{bmatrix} \text{---} \\ C \end{bmatrix}$$

- Computation of every element in C is independent.
- "Threads" \triangleq Such independent streams within a single program.
- Can be executed in parallel.

- 7/3 • If we execute threads of one program parallelly, we can get ↓ execution time along with ↑ throughput

- Executing a process requires the correct PC, stack, architecture registers, code segment, data segment, file descriptors, state

'Memory Image'

- State = { Active, Ready, Wait }

- Active → Ready transition :- (time expired)

Store PC, registers, stacks, etc.

- When there is a long-duration wait event,
Active → Wait

(eg - Page Fault)

- When wait event is over

Wait → Ready

(eg - Page Fault is serviced)

- Forking

- Making new memory image takes time in
: Make threads

- Threads of same process share code segment, data segment,

- Different PC, RF, stack.

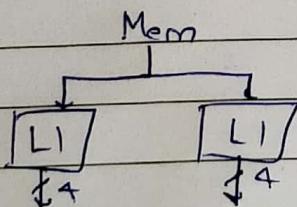
- Switching from one thread to another

eg - When long-duration wait occurs, like L2 or L3 cache miss.

- Store CPU context, retrieve new CPU context, execute.

↓
PC, RF, Stack.

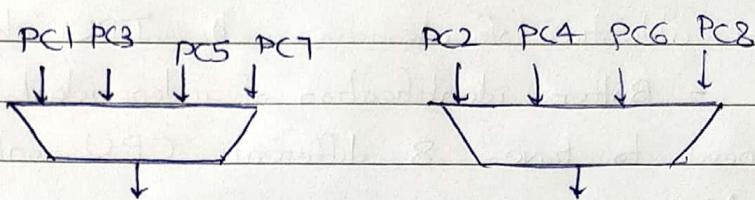
- When same thread runs for a long-time, \Rightarrow poor resource utilization. \therefore
 - * Switching threads only for long-duration wait \triangleq Coarse grain switching
- Fetch one instruction from each of 8 threads at a time
 - Better utilization of resources, b. IPC close to 8 \therefore
 - Better identification of independent instructions
 - We need to have 8 different CPU contexts (8 PCs, 8 ARFs,
 - We are fetching instructions from different locations
 - \Rightarrow Memory must have 8 parts (although of $\frac{1}{8}$ th size)
 - * Increasing no. of parts also increases static power consumption
 - 'Banking' - Say memory has 2 parts
 - Memory can have 'n' partitions and 'n' parts, each part can only give from a corresponding partition.
 - Threads need not be running at different partitions.
 - \therefore Choose to have only 2 parts, each of which can fetch 4 contiguous instructions from 2 different locations



- 'Round Robin'

Fetch 4 instructions each of T₁, T₂, then T₃, T₄, ... per cycle.

- This has poorer identification of independent instructions, but is still fine.



ii - Because each thread has PC at different locations that may map to same location in L1/L2/L3 cache, we have a huge problem of cache misses.

* All threads are sharing the same cache hierarchy.

∴ L1 hit rate ↓.

∴ Need to increase associativity of L1 cache

- Because increasing associativity can increase access time of cache, we need to 'pipeline' L1 cache

- Stage 1:- Tag Matching

- Stage 2:- Reading