

HY425 - Embedded Systems

Metrics

- Power = $\frac{1}{2} \times \text{Capacitive load} \times \text{Voltage}^2 \times \text{Frequency}$ Watts
- Power = Current static \times Voltage static Watts
- Energy = Capacitive load dynamic \times Voltage² Joules

Reliability equations,

$$\text{MTTF} = \frac{\text{total hours of operations}}{\text{total number units fail}} = \frac{\text{number of units} \times \text{hours until fail}}{\text{how many units fail}}$$

↳ Mean Time to failure

$$\text{FIT} = \frac{10^9}{\text{MTTF}} \rightarrow \text{Failures In Time (per billion hours)}$$

$$\text{FIT}_{\text{system}} = \sum_{i=1}^n \text{FIT}_i$$

$$\text{MTTR} = \frac{\text{total maintenance time}}{\text{total number of repairs}} \rightarrow \text{Mean Time To Repair}$$

$$\text{MTBF} = \text{MTTF} + \text{MTTR} \rightarrow \text{Mean Time Before Failure}$$

$$\text{Availability} = \frac{\text{MTTF}}{\text{MTTF} + \text{MTTR}}$$

Spec Ration and performance

Design X is n times faster than design Y

$$n = \frac{\text{Execution time Y}}{\text{Execution time X}} = \frac{\frac{1}{\text{Performance Y}}}{\frac{1}{\text{Performance X}}} = \frac{\text{Performance X}}{\text{Performance Y}}$$

$$\text{Weighted arithmetic mean} = \sum_{i=1}^n \text{Weight}_i \times \text{Time}_i$$

Spec ration: Measuring against a reference computer

$$\text{Specification} = \frac{\text{Execution time reference}}{\text{Execution time A}} = \frac{\text{Performance A}}{\text{Performance reference}}$$

$$\text{Geometric Mean: } \sqrt[n]{\prod_{i=1}^n \text{Spec ration}(i)}$$

Pros and Cons Geometric Means

- Pros:
- 1) Autogeneration of program frequencies
 - 2) Dev. effort is small and so is the memory footprint ratio of functions

- Cons:
- 1) Dev. time is high due to many executions
 - 2)

Qualitative principles of design

Taking advantages of Parallelism

- Use pipeline to overlap instructions
- Use multiple execution units
- Use multiple cores
- Use multiple processors to increase throughput

Locality

- Program reuse instructions and data
- 90-10 rule: 90% of execution time spent running 10% of instructions
- Programs access data in nearby addresses (spatial)

Make the common case fast

- Trade-offs in design (e.g. performance vs power saving)
- Provide efficient design for the common case
- Amdahl's law
- Amdahl's law

$$\text{Speedup} = \frac{\text{execution time old}}{\text{execution time new}} = \frac{1}{(1 - \frac{\text{fraction unchanged}}{\text{fraction unchanged}})} + \frac{\text{fraction unchanged}}{\text{speedup unchanged}}$$

$$\text{ex. 2 cores, } 40\% \text{ parallelism} \Rightarrow \text{speedup} = \frac{1}{(1 - 0,4) + \frac{0,4}{2}}$$

2 cores, Program 1 → 25% use and 40% parallelism

Program 2 → 25% use and 60% parallelism

$$\Rightarrow \text{speedup} = 0,25 \cdot \left((1 - 0,4) + \frac{0,4}{2} \right) + 0,75 \left((1 - 0,6) + \frac{0,6}{2} \right)$$

3

• Processor performance

$$\text{CPU time} = \text{Instructions} \times \text{CPI} \times \text{Cycle time}$$

$\sim 1 \text{ GHz} = 1 \text{ ns}$
 $4 \text{ GHz} = 0,25 \text{ ns}$
 $\rightarrow 1 \text{ MHz} = 1000 \text{ ns}$

Calculate CPI

Op	Freq	CPI	F * CPI	Time
ALU	50%	1	0,5	$33\% \rightarrow \frac{0,5}{1,5}$
Load	20%	2	0,4	$27\% \rightarrow \frac{0,4}{1,5}$
Store	10%	2	0,2	$13\% \rightarrow$
Branch	20%	2	0,4	$27\% \rightarrow$
			<u>1,5</u>	

Pipeline doesn't help latency of single task, it helps throughput of entire workload

Limits of Pipelining

- Structural hazards
- Data Hazards : RAW, WAW, WAR
- Control hazards : Branches and jumps

Speedup of Pipeline

$$\textcircled{1} \quad \text{speedup} = \frac{\text{Aver. Instruction time unpipelined}}{\text{Aver. Instr. time pipelined}} = \frac{\text{CPI unpipelined}}{\text{CPI pipelined}} \times \frac{\text{Clock cycles unpiped}}{\text{Clock cycles piped}}$$

$$\textcircled{2} \quad \text{CPI pipelined} = \text{Ideal CPI} + \text{Pipeline stall clock cycles per instruction}$$

$$= 1 + \text{Pipeline stall clock cycles per instruction}$$

$$\textcircled{1} + \textcircled{2} \Rightarrow \text{speedup} = \frac{1}{1 + \text{Pipeline stalls per instruction}} \times \frac{\text{clock cycles unpipelined}}{\text{clock cycles pipelined}}$$

$$= \frac{1}{1 + \text{Pipeline stalls per instruction}} \times \text{pipeline depth}$$

→ How to avoid stalls,

- Formalizing to avoid data hazards

BUT: You can't avoid stall after load instruction if next instruction use the same register as input

lw \$1, 0(\$2)
add \$4, \$1, \$6

stall

- Software scheduling to avoid load hazards

Branch Stalls Impact

What to do with the 3 clocks stall in Branches;

- 1) Determine branch taken or not taken sooner
- 2) Compute taken branch addr. earlier.

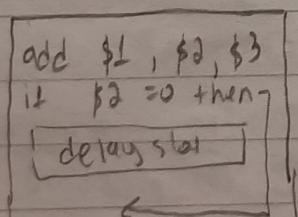
MIPS Solution: Test if register = 0 or ≠ 0

How: 1) Move zero test to Instruction Fetch / Register Fetch stage

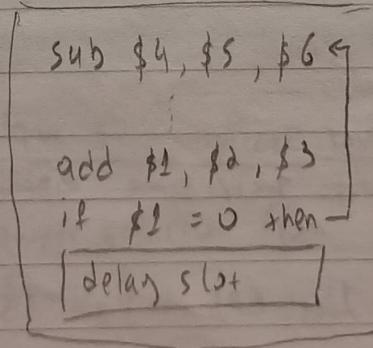
2) Add 6 to calculate new PC in Inst/Fetch / Reg. Fetch. stage

(SOSB) Always have 1 clock cycle penalty vs 3 were before

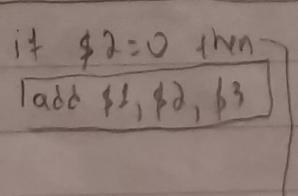
⇒ 1 delay slot after every branch



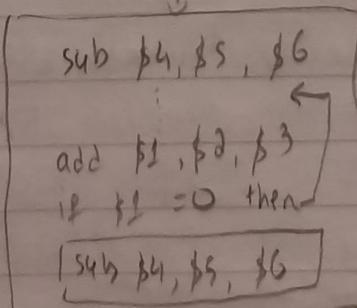
become



become



become



Exception and Interrupt

Problem: Exception appears - happens between 2 instructions (i and $i+1$)
→ The effect of all inst. up to and including i is nothing complete

a) No effect of any instruction after i can take place

The interrupt (exception) handler either abort program or restart at instruction $i+1$

HW have to remember where is the interrupt and stop the instruction at commit point.

Summary about hazards,

- Structural: need more HW resources
- Data (RAW, WAW, WAR): need forwarding, compiler scheduling
- Control (Branch, jumps): delayed Branch, prediction

Dynamic Inst.
scheduling
scalars

Instruction Level Parallelism,

Instruction level Parallelism

- Potential overlap among instructions
- Parallel or Out-of-order executions
- Requires extensions on the simple pipeline

Loop level Parallelism

- Exploit ILP between instruction from different basic blocks (e.g. iterations of a loop)

SOS Must Maintain: - Diagnose

→ Data flow: the flow of values between instruction that writes (produce) and instruction that read (consume), RAW, WAW, WAR

→ Exception behavior: changes in execution order should not change the order of exceptions or generate new ones

Scoreboard

Implications

- Solution for WAR:
 - stall registers write-back until registers have been read
 - read registers only on the Read Operand stage
- Solution for WAW:
 - detect this hazard and stall the new instruction (no issue) until the previous instruction has complete
- There is no register renaming
- Multiple instructions on execution stage
- The scoreboard keeps information about dependencies of the instructions that have been issued
- The scoreboard pipeline replace the Instruction Decode stage with 2 stages, issue and read operands

4 stages of Scoreboard

- Issue: decode instruction & record data dependences & check for structural hazards
 - Instruction issued in order (for hazard checking)
 - Don't issue if structural hazard
 - Don't issue if instruction is output dependent on any previous issue that didn't complete (no WAW hazards)
- Read Operands: wait for data hazards to be resolve, then read registers
 - All real dependencies (RAW hazards) resolved in this stage, since we wait for instructions to write back data
 - No forwarding of data

• Execution: execute the instruction

- After complete, it notice the scoreboard that it has complete

• Write result: write result on register

- Stall until there are no WAR hazards with previous instructions

Scoreboard Control

Issue: Structural hazards, WAW

ReadOper: RAW

Create: —

WriteBack: WAR

Instruction Status	Issue	Op	Read	Exec	Write Back
Fractional unit status			destination reg	source 1	source 2
				function unit pipeline F3	function unit pipeline F4
				function/unit + 2/3rd F3	function/unit + 2/3rd F4
				flag to now if F3 is ready → F3	flag to now if F4 is ready → F4
				flag to now if F3 is ready → F3	flag to now if F4 is ready → F4
			Busy	Op	Fi Fj Fr Qj Qr Rj Rk
Intiger					
Mult1					
Mult2					
Add					

1502 when we issue an inst. if the Fj or Fr is ready we mark Rj, Rk with YES, but after read oper we mark Rj, Rk with NO

Register result status	R0	R13	R30
RFU	Intiger			Mult1	

1503 when we issue an instruction we store here that Ri is ~~when~~ going to get data.

Tomasulo

VS Scoreboard

- FU buffers called "Reservation station RS" and hold the register / variable values needed for execution
- The source registers are replaced with the actual values or pointers to the RSs that will produce the values. Register renumbering
 - avoids WAW, WAR hazards
 - more RSs than real registers (more optimizations)
- Results arrive to the FUs from the RSs, not via register file but over Common Data Bus (CDB), that broadcasts results to all FUs
- Load and Stores use separate FUs with RSs as well
- Integer Instructions can go past branches, allowing FP ops beyond basic block in FP queue
- Control logic & buffers are distributed with the FUs vs centralized in Scoreboard

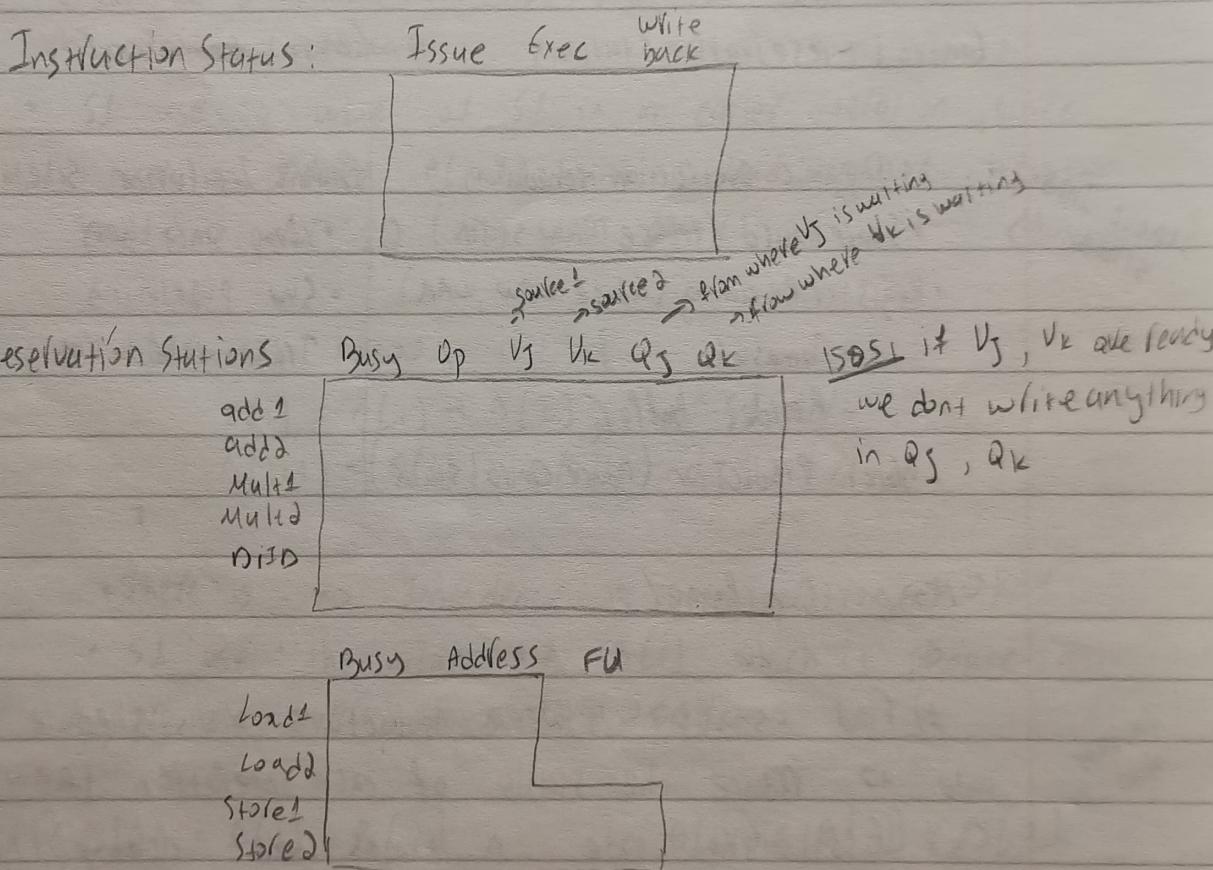
Three Stages of Tomasulo Algorithm

- Issue: Get a new instruction from the instruction queue.
 - If there is a vacant reservation station (no structural hazard) then issue instruction and send operand values or keep track of the FU/RS that will produce the operand (register renumbering)
- Execution: Execute on the FU
 - when all operands are ready in RS start execution
 - when not ready \Rightarrow watch the CDB for the result
- SOS: They can be ready multiple instruction for execution. for the FP adders it just select one of them random BUT for the Memory Unit we need 2 steps. First we calculate the base register when register is ready and then we add the offset and we store it in load or store unit. All loads and stores are execute with the right order.
- Write result (commit)
 - 1) write the result in the CDB and all waiting FUs/RSs get the result
 - 2) write to register file
 - 3) mark reservation station available

Common Data Bus Vs Typical data bus

- Typical data bus: data + destination ("go to" bus)
- Common data bus: data + source/reservation station tag ("come from" bus)
 - masters send to all slaves
 - does arbitration and the broadcast
 - 64 bits data + 4 bits of FU source address
 - write result if it matches the expected FU

Tomasulo Data Structures



Register Result Status

FU	..	F30
FU add1		Mult2

We write where the data are going to come

FU	..	F30
Reg File		

We write the data after write back

Tomasulo Vs Scoreboard

Tomasulo

- Pipelined FUs
(6 load, 3 store, 3+, 2×/÷)
- ⇒ window size ≤ 14 Instructions
- No Issue on structural hazard
- WAR: renaming avoid them
- WAW: renaming avoid them
- Broadcast result from FU
- Control: reservation stations

Scoreboard

- Multiple FUs
(2 load/store, 1+, 2×, 1÷)
- ⇒ window size ≤ 5 Instructions
- Same
- stall
- stall
- write/read registers
- central scoreboard

Static Inst.
Scheduling

Dynamic Instruction scheduling

- Scoreboard (reduce RAW stalls)
- Register Renaming (reduce WAW, WAR)
 - Tomasulo
 - Reorder Buffer
- Branch Prediction (reduce control stalls)

Static Instruction Scheduling (SW, compiled)

- Loop unrolling
- SW pipelining
- Trace scheduling

Make loop faster

- 1) Find <stalls> between instructions and schedule the loop code → reduce the stalls of each iteration (RAW)
- 2) Unroll the code n times and change the values of load & store instruction in every other iteration (± Immediate)
- 3) Find WAR, WAW and rename registers
- 4) Rewrite the code to reduce stalls > (RAW)

When is it safe to unroll?

when we don't have dependence across iterations
(loop-carried dependence)

e.g. 1) $\text{for}(i=0; i<100; i++) \{$

$$A[i+1] = A[i] + C[i] \quad /* S1 */$$

$$B[i+1] = B[i] + A[i+1] \quad /* S2 */$$

}

- S2 use the value of $A[i+1]$, compute by S1 in same iteration \Rightarrow OKAY
- S1 use a value of S1 in an earlier iteration, since i iteration produce $A[i+1]$ which is read in $i+1$ iteration
The same for S2 \Rightarrow NOT OKAY

a) $\text{for}(i=0; i<100; i++) \{$

$$A[i] = A[i] + B[i]; \quad /* S1 */$$

$$B[i+1] = C[i] + D[i]; \quad /* S2 */$$

}

- There is no dependence in same iteration \Rightarrow OKAY
- S1 use the value of $B[i]$ which is produce in previous iteration \Rightarrow NOT OKAY

1st iteration $\{ A[0] = A[0] + B[0] \}$

$$A[0] = A[0] + B[0]$$

$$\{ B[1] = C[0] + D[0] \}$$

BUT

$$B[1] = C[0] + D[0] \quad \{ 1st iteration \}$$

2nd iteration $\{ A[1] = A[1] + B[1] \}$

$$A[1] = A[1] + B[1]$$

$$\{ B[2] = C[1] + D[1] \}$$

$$B[2] = C[1] + D[1] \quad \{ 2nd iteration \}$$

3rd iteration $\{ A[2] = A[2] + B[2] \}$

$$A[2] = A[2] + B[2]$$

$$\{ B[3] = C[2] + D[2] \}$$

$$B[3] = C[2] + D[2] \quad \{ 3rd iteration \}$$

stateful
code

loop-carried dependence

Now the dependencies are in the same iteration so its OKAY to an RDO!!

```

for(i=0; i<100; i++) {
    A[i] = A[i] + B[i]; /* S1 */
    B[i+1] = C[i] + D[i]; /* S2 */
}

```



$A[0] = A[0] + B[0]$ → start-up code

```

for(i=0; i<99; i++) {

```

$B[i+1] = C[i] + D[i];$ /* S2 */

$A[i+1] = A[i+1] + B[i+1];$ /* S1 */



$B[99] = C[99] + D[99];$ → clean-up code

BSOI The longer the dependence distance, the more potential to extract parallelism

Software Pipeline

If the iterations of the loop are independent, then we can exploit ILP from different iterations of the loop

So we reduce stalls between instruction in the same iteration by executing instructions from other iterations

Iteration 0

Iteration 1

Iteration 2

Iteration 3

LD F0, 0(R1)

ADD F4, F0, F2

SD 0(R1), F4

...

LD F0, -8(R1)

ADD F4, F0, F2

SD -8(R1), F4

...

→ startup code

LD F0, -16(R1)

ADD F4, F0, F2

SD -16(R1), F4

...

LD F0, -24(R1)

ADD F4, F0, F2

SD -24(R1), F4

...

RAW → stall in every iteration

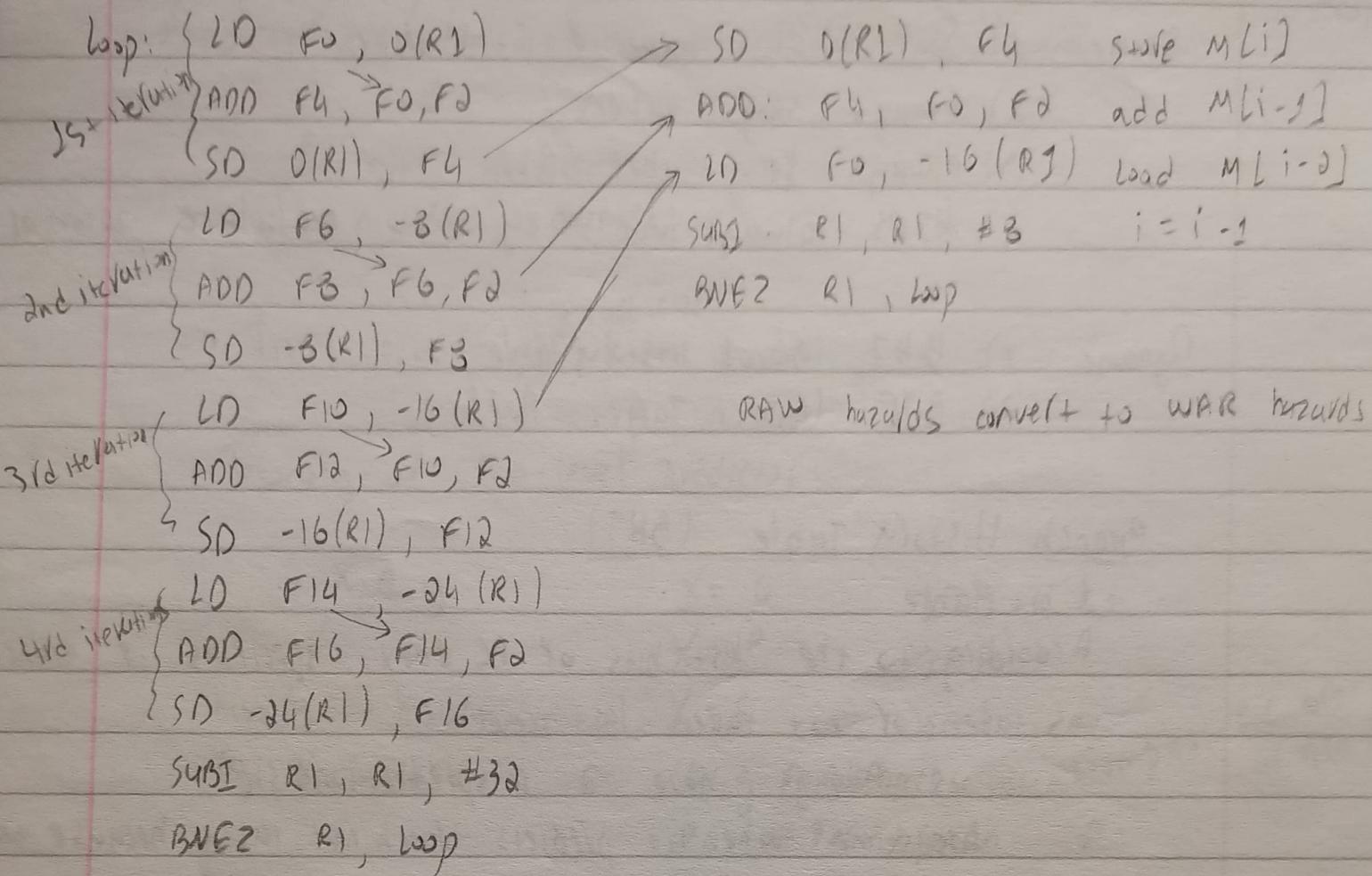
SD -24(R1), F4



finish-up code

Before: Unroll 4 times

After: Software Pipelined

Loop Unfolding vs Software pipelined

increase performance of the loop

- Loop unfolding: Pros: • Reduce cycles of each iteration takes?
- Cons: • Increase the code size
 • We need register renumbering \Rightarrow register pressure
 • The criteria for loop unfolding are very restrictive

Software pipelining: Pros: • Reduce stalls due data dependencies in iterations

- Cons: • Overhead at start and the end of loop (start-up code, end-up code)

Static Vs Dynamic Prediction

- Static:
- 1) Stall until i know if it's taken
 - 2) Always say taken
 - 3) Always say not taken
 - 4) MIPS - always use 1 cycle to calculate taken or not

- Dynamic:
- 1) BHT (Branch History Table)
 - 2) BTF (Branch Target Buffer)

Branch History Table (BHT)

- 1 bit Table

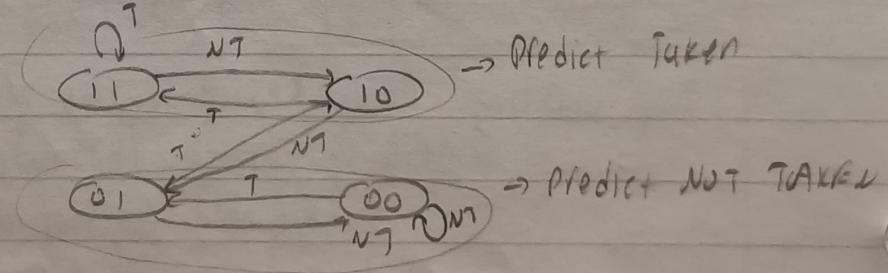
According to the last bits of PC index specifies if branch was taken or not taken last time

Limitations: Always have 2 mispredictions:

First time: When it predicts branch not taken because was 0

Last time: When it predicts branch taken because was 1

- 2 bit Table:



Add hysteresis in predictions

Correlation predictors - 2 level predictors

(13)

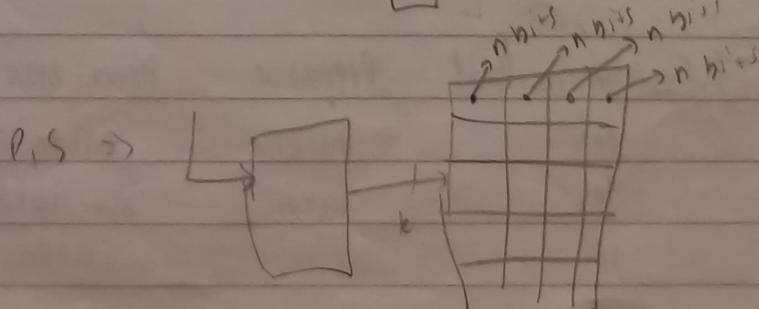
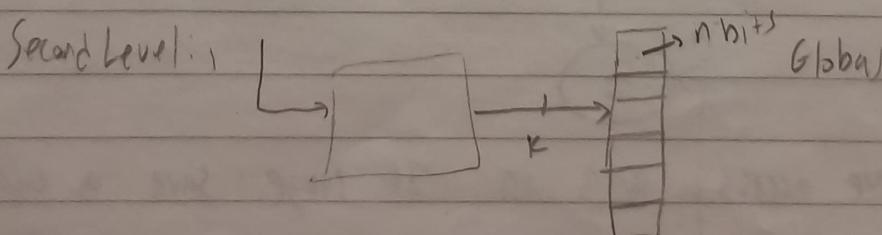
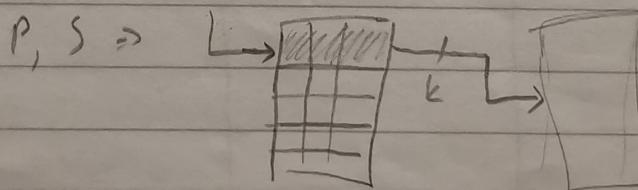
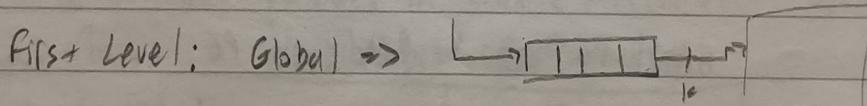
Predictor: (l_1, a) , (k, n)

\downarrow \downarrow
 four bits 2 bits
 of history counter

$$\Rightarrow \text{Total bits} = 2^{k+\text{add-size}} * n$$

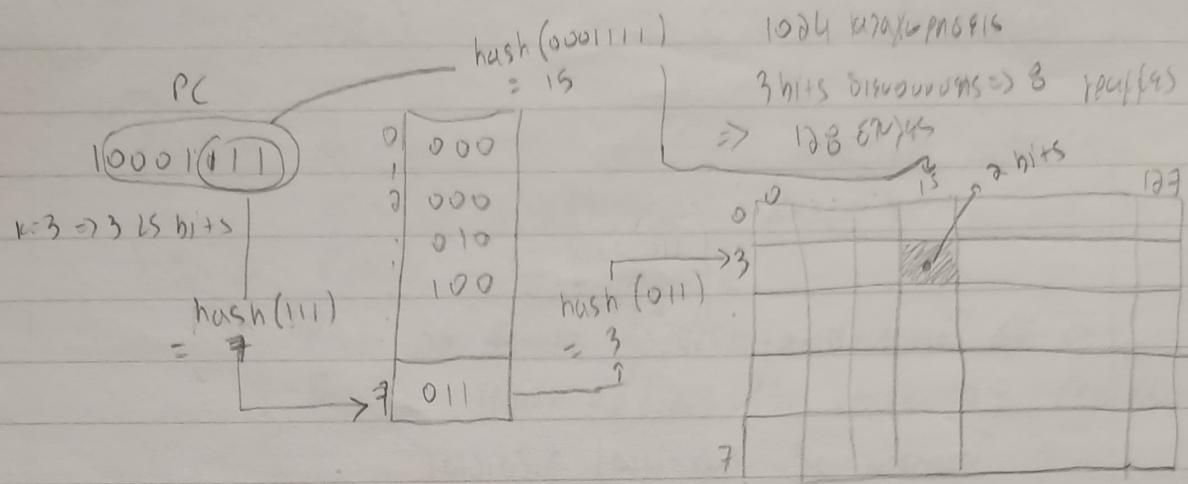
Taxonomy of two-level predictors

X	A	y	$X = G$	$y = g$	
\downarrow	\downarrow		$X = P$	$y = p$	
history	history		$X = S$	$y = s$	
register	table		$G, g \Rightarrow \text{global}$		examples of $bfa = 0$
			$P, p \Rightarrow \text{per-branch}$		or $nat \neq 0$
			$S, s \Rightarrow \text{set-associate per branch}$		



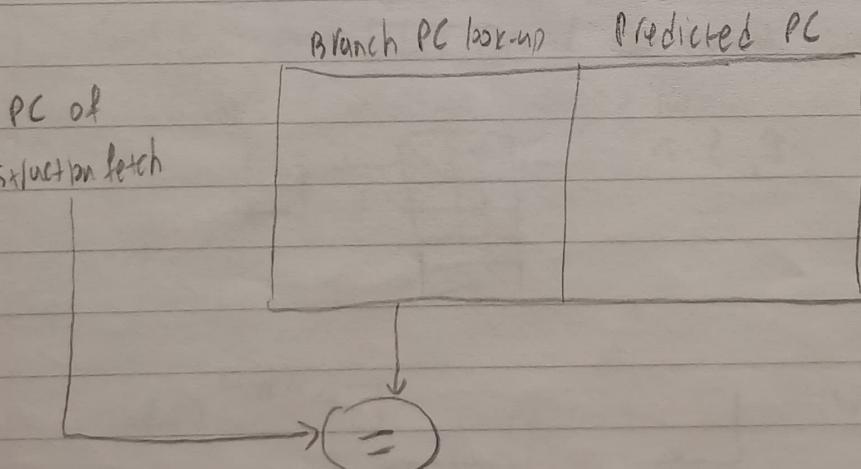
Example : (3, 2) , i.e. 1024 ways with 15 bits

PAP



Branch Target Buffer

- A look-up table with PC address of branches. When the PC matches, it provides the target address of the branch and TAKEN or NOT TAKEN
- Optimization: Stores only the taken branches



We access BTB in IF stage : save of cycle

Instruction in Buff	Prediction	Actual Branch	Penalty Cycles
Yes	taken	taken	0
Yes	taken	not taken	0
No		taken	0
No		not taken	0

Indirect Jumps

switch/case statements, goto's, function returns

Problem: Not good accuracy when procedure called from multiple sites

Solution: Return address stack

Push a return address in stack at call

Pop a return address from stack at return

Speculative Execution

Convert branches to conditionally executed instructions:

Pros:

- Reduce branch pressure - rate of missprediction

- Conditional Move (CMOV)

- May reduce code size

- May be quicker in simple arithmetic operations

Cons:

- Takes a clock even if ~~unnulled~~ "unnulled"

- Stalls if condition evaluate late

Types of Interrupts / Exceptions

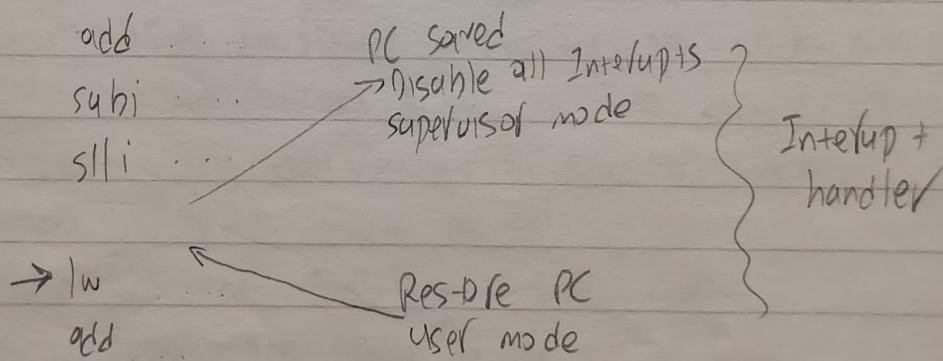
- I/O device request
- Invoking an operating system service from a user program
- Breakpoint
- Page fault
- Integer arithmetic overflow
- Memory protection violation

Precise Interrupts / Exceptions

if : 1) all instructions before this instruction has completed
 2) none of the instruction after has modified

\Rightarrow Precise interrupt

\Rightarrow We can restart the execution from the interrupt point and still get the correct results.



Imprecise Interrupts / Exceptions

if the process does not look exactly as if the instructions were executed sequentially

- Two possibilities:
- a) An instruction may completed that are later in program order
 - b) An instruction may not yet completed that are earlier in program order

Precise interrupt point requires multiple PCs when there are delayed branches

add . . .

sub

PC bne R1, there \Leftarrow Interrupt describe as $\langle PC, PC+4 \rangle$

PC+4 and . . .

To interrupt handle or not handle in branch

other >

more known how the future two stages happen
but now stuck

add

sub

PC bne R1, there

PC+4 and

\Leftarrow Interrupt describe as $\langle PC+4, \text{there} \rangle$ taken

other >

or $\langle PC+4, PC+8 \rangle$ not taken

To interrupt handle or not handle in and (no delay)

Why we need precise interrupts

slot selection in branch \Rightarrow A slot says if branch was taken or not taken \Rightarrow save both

- o Several interrupts need to be restartable
 - eg. TLB faults. fix translation and then restart faulting load/store
- o Restartability does not require precision. However, precision makes restarts much easier
- o Simplifies the Operating System (OS)
 - Less state needs to be saved away
 - Quick to restart

How to guarantee precise exceptions

- We can have exception in different stages of the pipeline
 - eg - Arithmetic exceptions occur in execution stage
 - TLB faults can occur in fetch or memory stage

⇒ How to guarantee precise exception?

- Use the pipeline
 - Each Instruction has an exception status field
 - Keep the PCs for every instruction in the pipeline
 - Check the exception when the instruction reaches the WB stage
- When an instruction reach the WB stage and has an exception
 - Save PC \Rightarrow EPC, Interrupt vector address \Rightarrow PC
 - Convert all fetched instructions to NOPs
- It works because of in-order completion WB

Precise Interrupts and Speculation

- During the Issue stage of instructions we operate as if we are predicting that all previous instructions do not generate instruction
 - Branch prediction, data prediction
 - If we speculate wrong, need to back up and restart execution to the point at which we predicted incorrectly
 - Exactly the same with precise exceptions

⇒ Use ROB

Idea behind ROB

Keep the instructions in a FIFO, with the exact order that they are issued

- Each ROB entry contains PC, dest Reg, result, exception status.

- When an instruction completes execution then the results are placed in the allocated entry in the ROB
 - Supplies operands to other instruction between execution complete & commit \Rightarrow more registers like RS
 - Tag results with ROB number instead of RS
- The instructions change the machine state at commit stage not on the WB \Rightarrow in-order commit \Rightarrow values at head of ROB are placed in registers.
 \Rightarrow We can cancel/squash speculatively executed instructions during mispredicted branches or exception

Four stages of Tomasulo with ROB

- 1) Issue: Get Instruction from Op Queue
 - if there are free reservation stations & ROB slot, issue instruction & send operands & ROB no. for destination (dispatch)
- 2) Execution: Execute Instr. in the Execution Unit
 - When the values of the 2 source regs are ready then execute instructions, otherwise watch CDB for result. When both in RS execute. (check RAW problem)
- 3) Write Result: End of execution (WB)
 - Write on CDB to all awaiting FUs & ROB. Mark RS available
- 4) Commit: Update the dest. Reg with the value from the ROB
 - When instr. at head of ROB & result present, update register file with the result (or store to memory) and remove instr. from ROB. Misspredicted branches or exception, flushes ROB

Memory Disambiguation

WAW/WAR hazards:

like hazards in Register file, we must avoid hazards through memory

- WAR/WAW are eliminated because the actual updating in memory occurs in order. When a store is at the head of the ROB, no earlier loads or stores can still be pending

RAW hazards:

- St $V(R2), R5$
- Id $R6, V(R3)$

What if $V(R2) == V(R3) \Rightarrow$ RAW \Rightarrow stall

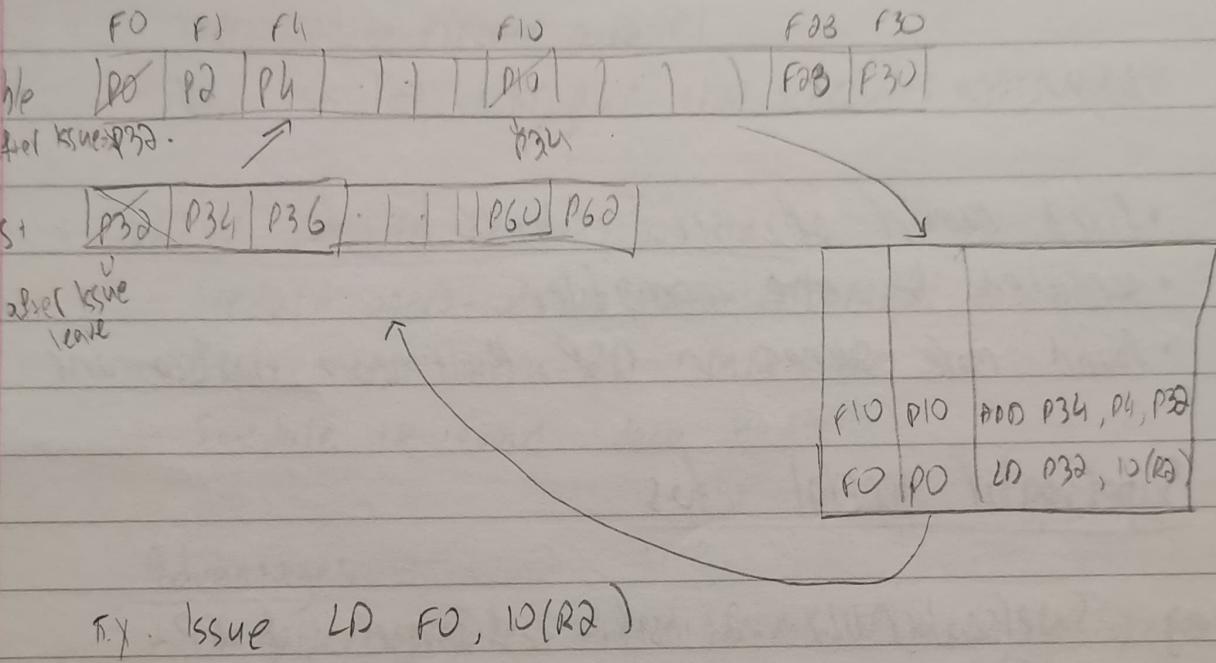
Solution: We might guess at whatever or not they are dependent (dependence speculation) and use ROB to fixup if we are wrong

\Rightarrow HW support for memory disambiguation

- Store buffer: keeps all pending stores in-order
 - Keep track of address and value (when they are available)
 - FIFO ordering will retire stores in program order
- When issuing a load
 - If any store prior to load is waiting for its address \Rightarrow stall load
 - If load addr matches earlier store addr (associative lookup)
 - then we have memory induced, RAW
 - if store value is available \Rightarrow return value
 - if store value is not available \Rightarrow return ROB number of source
 - Otherwise, send out request to memory
- Stores commit in-order \Rightarrow NO WAR, NO WAW hazards

Explicit register renaming

Upon issue, every instruction that writes a register allocates a new physical register from the freelist



BUT: what happens with Branches

There is HW feature to save the curr. Map Table and the free list. So we have a checkpoint, if we make wrong speculation to go back in previous state

Superscalar

- Varying multiple instr/cycle, 1-way, 2-way ... 8-way
- Schedule by compiler (statically scheduled) or by HW (dynamically scheduled)

VLIW

- fixed number of instr. (4-16)
- scheduled by the compiler
- found more success in DSP, Multimedia Applications

Superscalar general infos

e.g. Superscalar DLX: 2 instr., 1 Integer, 1 FP
 ⇒ fetch 64-bits: Integer + FP
 ⇒ Can only issue FP if there is Integer Instr.
 ⇒ More ports to FP registers to do FP load & FP op in parallel

Advantage: Increase the number of instr. which can be in the pipeline at one time and hence the level of parallelism

Disadvantage: The larger number of instr. in pipeline, increase the potential for data & control dependencies ⇒ stalls

Advantage Vs VLIW: The ability to extract some parallelism from less structured code, and easily cache all forms of data

Challenges: 1) While Inte/FP split is simple for the HW but get CPI 0.5 only if: - 50% FP operation
 - No hazards

2) If more instr. issued at the same time ⇒ greater difficulty of decode & issue

VLIW general Infos

Tweakoff instruction space for simple decoding:

- The long instruction word has room for many operations
- All operations the compiler puts in the long instruction word are independent \rightarrow execute parallel
- Need compiling technique that schedules across several branches

VLIW compiler Responsibilities:

- Schedule operations to maximize parallel execution
- Guarantees intra-inst. parallelism
- Schedule to avoid data hazards

Advantages

- Compiler prepares fixed packets of multiple operations that give the full plan of execution
 - Dependencies are determined by compiler and used to schedule according to function unit latencies
 - Function units are assigned by compiler and correspond to the position within the instr. packet
 - HW doesn't have to "rediscover" dependencies

Disadvantages

- Object-code compatibility: recompile all code for every machine
- Object-code size:
 - Instr. padding wastes int. memory
 - loop unrolling / software predicates code
- Scheduling variable latency memory operations
- Knowing branch probabilities: profiling requires a significant extra step in build process
- Scheduling for statically unpredictable branches

IA-64 Instruction format

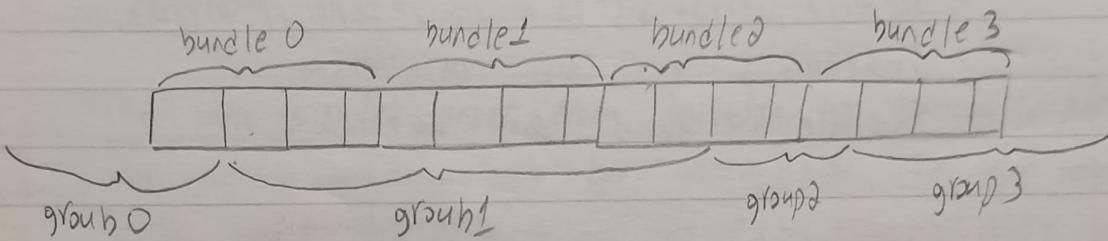
Instr. 2	Instr. 1	Instr. 0	Template
----------	----------	----------	----------

128-bit instruction bundle ($41 \times 3 + 5$)

Template bits

3 instr \Rightarrow 41 bit each

- Template bits: describe the types of instructions and the grouping of these instructions with other in adjacent bundles



- Each group contains instructions that can execute in parallel and the boundary of a group (stop) is indicated by the template

IA-64 Predicated Execution

- Problem: Misspredicted branches
limit JLP

b0 Instr 1

Instr 2

hr a==b, b2 if

- Solution: Eliminate hard to predict branches

with predicated execution.

b1 Instr 3 else

Instr 4

jmb b3

Instr 1

Instr 2

pl = (a != b), pd = (a == b)

(pl) Instr 3 || (pd) Instr 5

(pd) Instr 4 || (pd) Instr 6

Instr 7

Instr 8

1 basic block

Predicated

h2

Instr 5 then

Instr 6

h3 Instr 7

Instr 8

in basic blocks

Branch Predication

- Is an aggressive compilation technique to generate code with higher degree of instruction level parallelism
- It lets operations from both branches of a conditional branch to be executed in parallel, to increase amount of parallel operation
- In this way branches are eliminated and replaced by conditional execution

The idea is: let instr. from both branches go on in parallel, before the branch condition has been evaluated. The HW takes care that only those corresponding to the right branch will be finally committed

Multithreading Cost

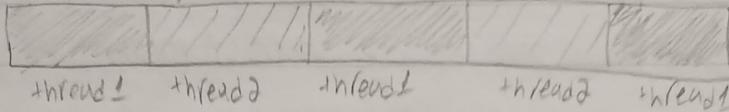
- Each thread requires its own user state. Many CPU resources are split or shared
 - PC
 - Register buffer
 - Physical registers
 - Load/Store buffer
 - Prefetch & Instr. buffers
 - Issue buffers
- Also needs its own system state
 - virtual memory page table base register
 - exception handling registers

HW multithreading alternatives

• Fine-Grain Multithreading

- Switches processor context every thread cycle
- Context belong to same address space

Cycle-1 Cycle-2 Cycle-3 Cycle-4 Cycle-5



Switch every thread cycle:

- Need fast HW switch between contexts
 - Multiple PCs are register files

- Implemented with round-robin scheduling, skipping stalled threads

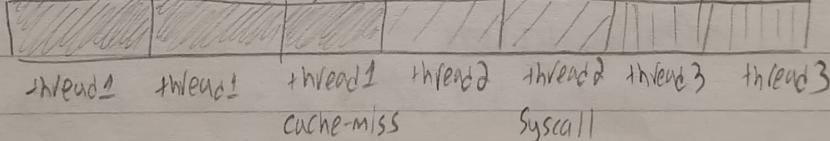
- Hide both short and long stalls

- Delays all threads, even if they have no stalls

• Coarse-Grain Multithreading

- Switches processor context upon long-latency event
- Context may belong to different address space

Cycle-1 Cycle-2 Cycle-3 Cycle-4 Cycle-5 Cycle-6 Cycle-7



Switches upon long-latency events

- Can afford slower context switch than fine-grain multithread
- Threads are not slowed down
 - Other threads run when one have stall

- Pipeline start-up cost upon thread switching

- Processor issues instructions from one thread

• Simultaneous Multithreading (SMT)

- Techniques presented so far have all been "vertical" multithreading where each pipeline stage works on one thread at a time
- SMT uses fine-grain control but allow instr. from multiple threads to enter execution on same clock time.

Simultaneous Multithreading SMT

(29)

- Add multiple contexts and fetch engines and allow instr. fetched from different threads to issue simultaneously (hardware)
- Utilize wide out-of-order superscalar processor issue queue to find instr. to issue from multiple threads
- OOO instr. window already has most of the circuitry required to schedule from multiple threads
- Any single thread can utilize whole machine
- Shared HW mechanism
 - Large set of virtual reg. can hold register sets of independent thread
 - Renaming provides unique register identifiers to different threads
 - Out-of-order completion of instructions from different threads allowed
 - No cross-thread RAW, WAW, WAR hazards
 - Separate re-order buffer per thread

+ Page 21 Single Multithreading

+ Page 26-33

SIMD: single instr. multiple data

Vector Processors

- Vector is one dimension array
- A vector processor is one whose instr. operate on vectors rather than scalar (single data) values
- Basic Requirements:
 - Need vector registers
 - Need vector length register (VLEN)
 - Need vector stride register (VSTR)
 - ↓
distance in memory between two elements in vector

- A vector instr. perform an operation on each element in consecutive cycles
 - Vector FUs are pipelined
 - Each pipeline stage operates on a different data element
- Vector instr. allow deeper pipeline
 - No HW interlock needed within vector
 - No control flow within a vector
 - Known stride allows easy address calculation for all vector elements
 - Enables prefetching of vectors into reg/cache/memory

Vector processor limitations

- Memory bandwidth can easily become a bottleneck, especially if
 - Compute/memory operation balance is not maintained
 - data is not mapped appropriately to memory banks

Vector Registers

- Each vector data register holds N values of M -bits
- Vector control registers: VLEN, VSTR, VMASK
- VLEN is N
- VSTR is the distance in memory between two elements in vector
- VMASK indicates which elements of vector to operate on
 - eg. $VMASK[i] = (V_k[i] == 0)$

Loading/Storing Vectors from/to Memory

- Requires loading/storing multiple elements
- Elements separated from each other by a constant distance (stride)
- Elements can be loaded consecutive cycles if we can start the load of one element per cycle
 - ↳ We achieve this with Memory Banks

Vector Memory System

- Memory is divided into banks that can be accessed independently. Banks share address and data buses
- Next address = previous address + stride
- if (stride == 1)
 - && (consecutive elements interleaved across banks)
 - && (number of banks > bank latency)
 - ⇒ we can sustain 1 element/cycle throughout

Vector Stripping

- What if data elements > elements in vector register
 - Eg. 527 data elements, 64-element VREG
 - ⇒ 8 iteration where VLEN = 64
 - + 1 iteration where VLEN = 15
 - ⇒ Called vector stripmining

Gather / Scatter Operations

- What if vector data is not stored in a strided fashion in memory (irregular memory access to a vector)
 - Idea: Use index vector to combine/pack elements into VREG
 - called: gather/scatter operations
- Vector loads/stores use an index vector which is added to the base register to generate the addresses

Eg. Index Vector Data Vector (to store) Striped Vector (in Memory)

0	3,14	Base + 0	3,14
2	6,50	Base + 1	X
6	71,20	base+2	6,50
7	2,5	base+3	X
		base+5	X
		base+6	71,20
		base+7	2,5

Conditional Operations in loop VMASK

- What if some operations should not be executed on a vector

eg. $\text{for } (j=0; j < n; j++)$

$\text{if } (a[j] \neq 0) \text{ then } b[j] = a[j] * b[j]$,

- Idea: VMASK

- VMASK is a bit mask determining which data element should not be acted upon

VLD $V0 = A$

VLD $V1 = B$

VMASK = ($V0 \neq 0$)

VMUL $V1 = V0 * V1$

VST $B = V1$

- Another example

$\text{for } (i=0; i < 64; i++)$

$\text{if } (a[i] \geq b[i])$

$c[i] = a[i]$

else

$c[i] = b[i]$

A	B	VMASK
---	---	-------

1	2	0
---	---	---

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

3	2	1
---	---	---

4	10	0
---	----	---

5	-4	0
---	----	---

0	-3	1
---	----	---

→ Steps to execute this loop:

- 1) Compute A, B to get VMASK

- 2) Masked store of A into C

- 3) Complement VMASK ($\text{sub}(\text{not}(a))$)

- 4) Masked store of B into C

Some Issues

• Stride & Banking

- As long as they relative prime to each other and there are enough banks to cover bank access latency, we can sustain 1 element/cycle

• Storage of Matrics

- Row Major or Column Major
- We need different stride to access a row vs column

⇒ How to optimize

- More banks
- Better data layout to match the access patterns
- Better mapping of address to bank

Vector Summary

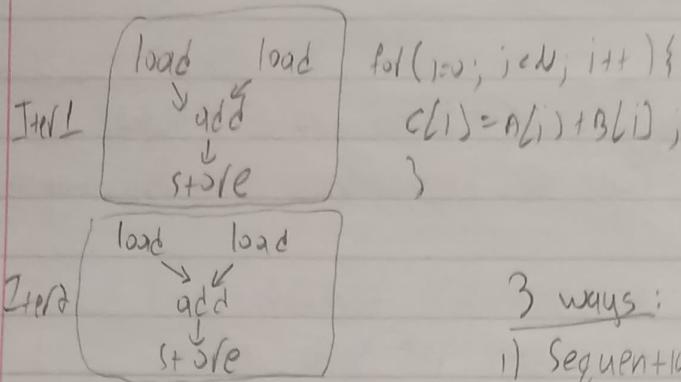
- Good to exploit regular data-level parallelism
 - Same operation performed on many data elements
 - Improve performance, simplify design (no intra-vector dependencies)
- Performance improvement limited by vectorizability of code

Masked Vector Insts

- Simple Implementation
execute all N operations, turn off result writeback according to Mask
- Density-Time Implementation

Scan mask vector and only execute elements with non-zero Mask
 \Rightarrow O(n) operations because n independent operations, upon $\Sigma_{i=1}^n$ go to vector
 you want to Mask or equal 1 and then have branch, after which go to 1
 until $\Sigma_{i=1}^n$ is 0. Revisa o Segundo appunto italiano
 about the number of operations this implementation does for repeated
 vector operations

GPlus

How can you exhibit parallelism here?

3 ways:

- 1) Sequential (SISD)
- 2) Data-parallel (SIMD)
- 3) Multithreaded (MIMD/SPMD)

1) Sequential (SISD)

- Pipelined processor
- Out-of-order execution
- Superscalar processor

2) Data parallel (SIMD)

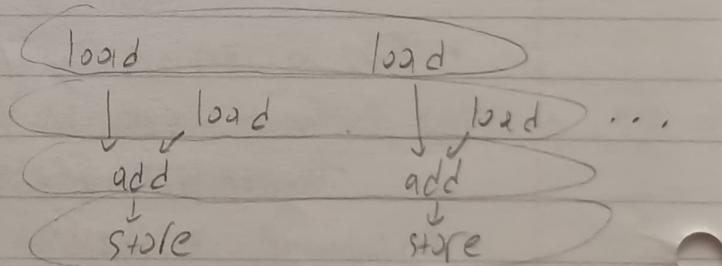
- Vector processor

- VLD $A \rightarrow V_1$

- VLD $B \rightarrow V_2$

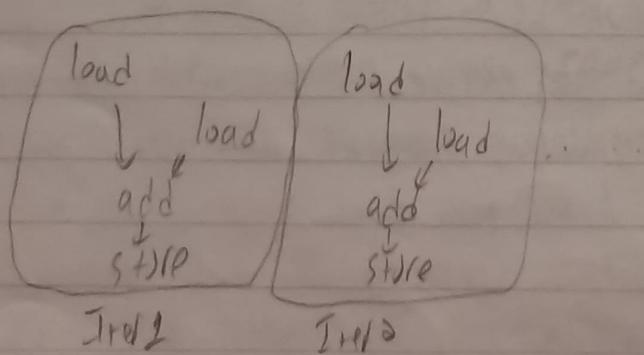
- VADD $V_1 + V_2 \rightarrow V_3$

- VST $V_3 \rightarrow C$



Programmer or compiler generates a SIMD instr. to execute the same instr. from 41) iterations across different data

3) Multithreaded



Programmer or compiler generates a thread to execute each iteration. Each thread does the same thing (but on different data.)

Warp

- A set of threads using the same instruction are dynamically grouped into a warp (e.g. execute at the same PC but on different data)
- Fine-Grained Multithreading
 - One instr. per thread in pipeline at a time
 - Interleave warp execution to hide latencies
- Register values of all threads stay in register file
- FGM/T enables long latency tolerance (overxn)

Warp-based SIMD vs Traditional SIMD

- Traditional SIMD contains a single thread
 - Sequential instr. execution (lock-step operations in a SIMD instr.)
 - Programming model is SIMD (no extra threads) → SW needs to know VLEN
 - ISA contains vector/SIMD instr.
- Warp-based SIMD consists of multiple scalar threads executing in a SIMD manner (same instr. executed by all threads)
 - Does not have to be lock step
 - Each thread can be treated individually (placed on diff. warp)
 - ⇒ Programming model not SIMD
 - SW does not need to know vector length
 - Enables multithreading and flexible dynamic grouping of threads
 - ISA is scalar → SIMD operations can be formed dynamically
 - Essentially, is IS SIMD programming model implemented on SIMD hardware

SPMD

multiple instr. streams execute the same program.

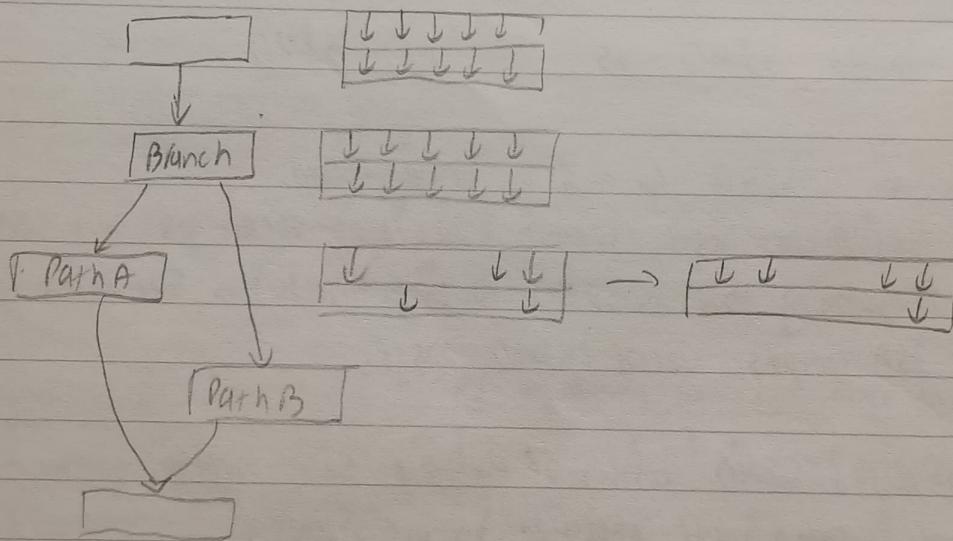
- Each program stream works on different data
 - a) can execute a different control flow-path at run-time

Two major advantages of SIMD

- 1) Can treat each thread separately (can execute each thread independently on any type of scalar pipeline \rightarrow SIMD processing)
- 2) Can group threads into warps flexibly (can group threads that are supposed to truly execute the same instr. \rightarrow dynamically obtain and maximize benefits of SIMD processing)

Dynamic Warp formation / merging

Idea: Dynamically merge threads executing the same instr. (after branch)



Large Warps & Two-level warp scheduling

- Two reasons for GPU resources be underutilized (utilization)
- Branch divergence $\xrightarrow{\text{fix}}$ Large warps
- long latency operations $\xrightarrow{\text{fix}}$ two-level warp scheduling

Cache Measures

$$AMAT = \text{Hit time} + \text{Miss Rate} * \text{Miss Penalty}$$

$$\text{CPU time} = (\text{CPU execution clock cycles} + \text{Memory stall cycles}) * \text{clock cycle time}$$

$$\Rightarrow \text{CPU time} = IC * \left(CPI + \frac{\text{memory stall cycles}}{\text{instruction}} \right) * \text{clock cycle time}$$

$$\Rightarrow \text{CPU time} = IC * \left(CPI + \text{Miss Rate} * \frac{\text{memory accesses}}{\text{instruction}} * \text{Miss Penalty} \right) * \text{clock cycle time}$$

example: - CPI = 1

- Data accesses are only load & store (explain 50% of instr)
- Miss Penalty: 25 CC
- Miss rate: 2%

$$\begin{aligned} \text{Memory Stall} &= IC * \text{Miss Rate} * \frac{\text{memory accesses}}{\text{instr.}} * \text{miss penalty} \\ &= IC * 0,02 * (1 + 0,5) * 25 \\ &= IC * 0,75 \end{aligned}$$

$$\begin{aligned} \text{CPU time} &= (IC * (CPI + \text{Memory Stall})) * \text{clock cycles} \\ &= (IC * (1 + 0,75)) * \text{clock cycles} \\ &= IC * 1,75 * \text{clock cycles} \end{aligned}$$

Direct Mapped

- Pros:
- Simple, low complexity, low power consumption
 - Fast hit time
 - Data available before cache determines hit or miss
 - Hit/Miss check done in parallel with data retrieval

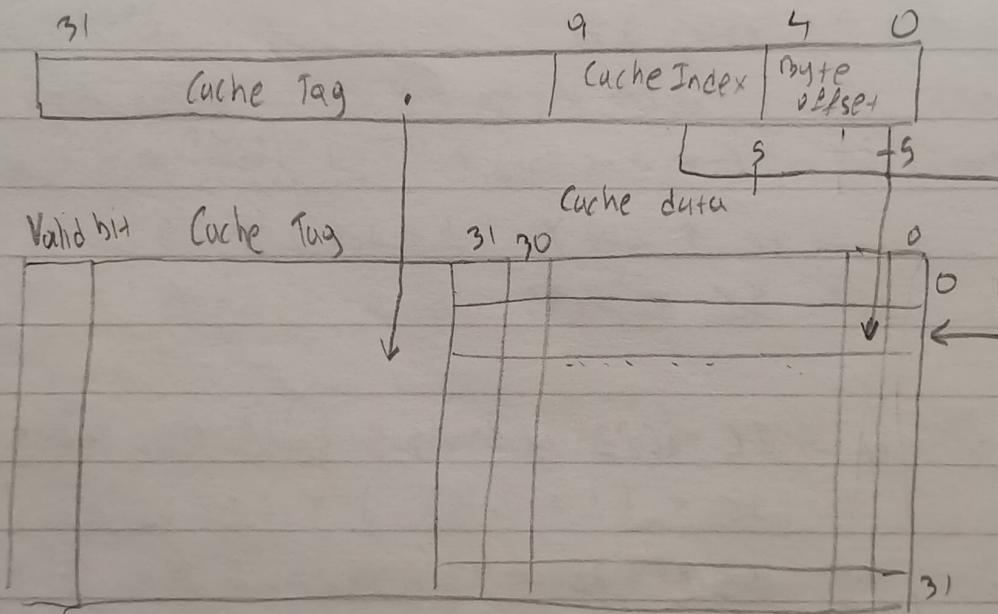
- Cons:
- Conflicts between blocks mapped to same block in cache

Two-way associative cache

- Pros:
- Choice of mapping memory block to different cache blocks in a set
 - Reduction of conflicts

- Cons:
- Increase complexity - comparators, multiplexor, parallel tag comparison
 - Increase power consumption
 - Increase hit time, due to comparators and multiplexor
 - Data available after cache determines hit or miss

How is block found in the cache



which Block is replaced on a miss

- Direct-Mapped \Rightarrow only 1 block to choose
- Set-associative cache
 - \rightarrow Random
 - \rightarrow LRU
 - \rightarrow FIFO

What happens on a write

	Write-through	Write-back
Police	Data word written to cache block, is also written to next lower level memory	Write new data word only to 1 cache block. Update lower level just before a written block leaves cache, so not lose true value
Debugging	Easier	Harder
Can read misses force writes	No	Yes (used to slow some reads, now write-buffer)
Do repeat writes toach lower level	Yes, memory busier	No

Write buffers - hold address & data

- An entry may contain multiple words
- If there is an empty entry, the data and address are written to the buffer, and the CPU is done with the write
- if buffer contains other modified blocks, check to see if new address matches one already in buffer - if so combine the new data with that entry
- if buffer full and no address match, cache and CPU wait for an empty entry to appear (meaning some entry has been written to main memory)
- Merging improves memory efficiency, since multi-word writes usually postes than one word at a time

Reduce Miss Penalty1) Multi-level caches

• Motivation

→ Bigger caches bridge gap between CPU & DRAM

Performance Analysis:

$$\text{AMAT} = \text{Hit Time L1} + \text{Miss Rate L1} * \text{Miss Penalty L1}$$

$$\text{Miss Penalty L1} = \text{Hit Time L2} + \text{Miss Rate L2} * \text{Miss Penalty L2}$$

• Local VS Global Miss Rate

→ local miss rate: Number of misses in the cache

divided by the number of access to this cache

→ global miss rate: Number of misses in the cache divided by
the number of accesses issued from the CPU

Performance Analysis:

$$\text{Average Memory Stall per Instruction} = \text{Miss per instruction L1} * \\ \text{hitTime L2} + \text{Miss per Inst., L2} * \text{Miss Penalty L2}$$

- example:
- Write-back first level cache
 - 40 L1 misses per 1000 Inst.
 - 20 L2 misses per 1000 Inst.
 - L2 cache miss penalty = 100 cycles
 - L1 hit time = 1 cycle
 - L2 cache hit time = 10 cycles
 - 1,5 memory references per instr.

$$\text{Miss Rate L1} = \frac{40}{1000} = 0,04$$

$$\text{Miss Rate L2} = \frac{20}{40} = 0,5$$

$$\text{AMAT} = \text{HitTime L1} + \text{Miss Rate L1} * (\text{HitTime L2} + \text{Miss Rate L2} * \text{Miss Penalty L2}) \\ = 1 + 0,04 (10 + 0,5 * 100) = 3,4 \text{ cycles}$$

$$\text{Average memory stall per instr} = \text{Miss per instr L1} * \text{hitTime L2} + \text{Miss per Inst., L2} * \text{Miss Penalty L2} \\ = \frac{40}{1000} * 1,5 * 10 + \frac{20}{40} * 1,5 * 100 = 3,6 \text{ cycles}$$

Inclusion Property

- Multi-level inclusion \Rightarrow L1 data always present in L2, L2 in L3
 - \rightarrow Simplifies cache consistency check in multiprocessors
 - \rightarrow Complicates the use of different block size for L1 and L2
 - Second-level cache miss must invalidate all first-level blocks
- Multi-level exclusion \Rightarrow L1 data never present in L2, L2 in L3 ...
 - \rightarrow Reasonable choice for systems with small L2 cache relative to L1 cache
 - \rightarrow Effective expansion of total caching space with a slower cache

2) Critical word first & Early restart

- Critical word first
 - \rightarrow Cache block size tends to increase to exploit spatial locality
 - \rightarrow Any given reference needs only one word from a multi-word block
 - \rightarrow CWF fetches requested word first and sends it to processor
 - \rightarrow Processor continues execution while rest of the block is fetched
- Early restart
 - \rightarrow fetches words in the order stored in the block
 - \rightarrow As soon as critical word arrives, sends to processor and processor restarts

3) Giving priority to read misses over writes

- Write-through caches:
 - \rightarrow On read miss check contents of write buffer
 - \rightarrow If no conflict then let missing read bypass pending write
 - All desktop & servers give reads priority over writes
- Write-back caches:
 - \rightarrow slow path: write dirty block to memory, then fetch a new block from memory to cache
 - \rightarrow fast path: write dirty block to write buffer, then fetch a new block from memory to cache, then write back dirty block

4) Write Buffer & Victim Cache

• Write buffer organization

- Processor blocks on write if write buffer full
- Processor checks write address with address in write buffer
- Processor merges writes to same address if address is present in write buffer

Ex. Assume WB with 4 entries, 64-bit word each

WB with no merging

Write address	V	V	V	V
100	1 Mem [100]	0	0	0
108	1 Mem [108]	0	0	0
116	1 Mem [116]	0	0	0
124	1 Mem [124]	0	0	0

WB with merging

Write addresses	V	V	V	V
100	1 Mem [100]	2 Mem [103]	2 Mem [116]	1 Mem [124]
108	0	0	0	0
116	0	0	0	0
124	0	0	0	0

• Victim cache

Tiny caches holds evicted cache blocks

- Small (eg 4 entry) full associative buffer for evicted blocks
- Propose to reduce impact of conflicts on direct-mapped caches

Victim cache + direct-mapped cache \approx associative cache

On sequential access the direct-mapped cache tags and victim cache find the required value sequentially

5) Non-blocking or Lookup free caches

Idea: → Allow for hits while serving a miss (hit-under-miss)

→ Allow more than one outstanding miss (miss-under-miss)

Make sense: → When the processor can handle > 1 pending load/store
superscalar processor

→ When the cache serves > 1 processor or other cache

→ When the lower level allows for multiple pending accesses

Difficulties: → Handling multiple misses at the same time

→ Handling loads to pending misses

→ Handling stores to pending misses

- Miss status handling register

- keeps track of

- Outstanding cache misses

- Pending load & stores that refer to that cache block

- Fields of MSHR

- Valid bit

- Cache block address

- Issued bit (1 if already request issued to memory)

- for each pending load or store

- Valid bit

- Type (load/store) & Format (byte/halfword...)

- Block offset

- Destination reg for load or store buffer entry

1	27	1	1	3	5	5	load/store 0
Valid	Block address	Issued	Valid	Type	Block offset	Destination	load/store 1

• Non-blocking caches: operations

• On cache-miss:

- Search MSHRs for pending access to same cache block
 - if yes, just allocate new load/store entry
- (if no) Allocate free MSHR
 - update block address and first load/store entry
- if no MSHR or load/store entry free, stall

• When one word of cache line become available

- check which load/store are waiting for it
 - forward data to LSU
 - mark load/store as invalid

→ write word in the cache

• When last word for cache line is available

→ mark MSHR as invalid

⑥ Multi-parted Caches

Idea: → Allow for multiple accesses in parallel

Implemented with: → True multiplexing
2 ways → multiple banks

Difficulties: → Interaction with parallel accesses (especially stores)

• True multiplexing: → Use 2-port per tag stage

→ Problem: → large area increase

→ hit time increase

• Multiple banks: → Partition address space into multiple banks

→ Benefits: → accesses can go in parallel if no conflicts

→ Challenges: → conflicts

→ distribution network

→ bank utilization

Reduce Miss Rate

3 C's model

Characterization of cache misses:

- Compulsory miss: Miss that happens due to the first access to a block since program began. Cold-start miss
- Capacity miss: Miss that happens because a block that has been fetched in the cache needed to be replaced due to limited capacity. Block has been fetched, replaced and re-fetched
- Conflict miss: Miss that happens because address of block maps to some location in the cache with other block in memory. Block has been fetched, replaced, re-fetched and cache has invalid locations that could hold the block if a different address mapping scheme were used

i) Increase the block size

• Spatial locality

- Larger block size usually reduce compulsory misses
- Increase miss penalty, since processor need to fetch more data
- May increase conflict misses, if spatial locality is poor
(most words in fetched not used)
- May increase capacity misses, if spatial locality is poor
(most words in fetched not used)

ii) Larger caches

Pros: → Reduction of capacity misses

Cons: → longer hit time

→ increase area, power and cost

3) Increase Associativity

- Higher associativity increase hit time
- increase cycle time

4) Prefetching

Idea: fetch data to cache before processor requests them

→ can address cold misses

→ can be done by programmer, compiler, hardware

Characteristics of ideal prefetching

- You only prefetch data that are truly needed
- You issue prefetch requests early enough
- You don't issue prefetch requests too early

• SW prefetching

Issues: → Take up issue slots

→ Take up system bandwidth

→ Must have non-blocking caches

→ Prefetch distance depends on specific system implementation

→ Not easy to use for pointer base structures

• HW prefetching

Same goals as SW prefetching

Major Questions: → where to place a prefetch engine
 → what to prefetch
 → when to prefetch
 → where to place prefetched data

• Simple sequential prefetching

- On cache miss, fetch N sequential memory blocks
 - Exploit spatial locality in both instr & data
 - Exploit high bandwidth for sequential accesses
- (called "Spatial Prefetch")
- Extend to N sequential memory blocks
- Problem: slow down once every N cache lines

• Stream prefetching

→ Is a continuous version of prefetching

- Stream buffer can fit N cache lines
- On miss, start fetching N sequential cache lines
- On hit, move cache line to cache, start fetching $N+1$
- In other words tries to stay N cache lines ahead

• Stride prefetching

- Idea: Detect & prefetch strided accesses

- Stride detection using a PC-based table
 - for each PC remember stride
 - stride detection:
 - remember the last address used for this PC
 - compare to currently used address at this PC
 - track confidence using a two bit saturating counter

5) Compiler Optimizations

- Code transformation to improve:

- Spatial locality, through higher utilization of fetched cache blocks
- Temporal locality, through reduction of the reuse

distance of cache blocks

→ Examples: loop interchange, loop blocking, loop fusion

• Data layout and data structure transformations to improve:

→ spatial locality, through high utilization of fetched cache blocks

→ Examples array merging, structure/object class member reordering in memory

Reduce Hit time

- small & simple caches
- virtually addresses caches
- pre-defined caches
- trace caches