

COMPUTER ORGANIZATION AND DESIGN

The Hardware/Software Interface



Chapter 4

The Processor



Control & Datapath Single (long) cycle per instruction





Pipelined Datapath & Control Operation without data or control dependencies, yet

University of Crete Dept. of Computer Science CS–225 (HY–225) Computer ORganization Spring 2020 semester

Slides for §9.3 – 9.5

- §9.3 Pipelined Datapath Operation
- §9.4 Control for the Pipelined Datapath
- §9.5 Graphical representation: time-work





Data Dependences (Hazards) in Pipelines



I2 needs the new data written by I1, hence must wait for I1 to write –or at least to generate– the new data

- RAW (Read after Write) true dependence, as above
- RAR (Read after Read) not a dependence, can freely reorder the reads
- WAR (Write after Read) "antidependence": if you want to do the write (I2) early, just keep a copy of the old data and have I1 read that copy
- WAW (Write after Write) if you want to reorder them, simply abort the write of I1 (if no one reads this word between I1 and I2)

Copyright 2020 University of Crete – https://www.csd.uoc.gr/~hy225/20a/copyright.html

No Memory Data Hazards in our simple Pipeline



• Data Memory accesses are performed 'in-order' in our simple pipeline, i.e. are not reordered relative to what the program specifies, thus, no dependences of memory word accesses are ever violated

Register Accesses reordered, with Pipelining



• For each instruction that writes a destination register, if the next 2 or 3 instructions read that same register, i.e. need its result, we have to do something about it...

Actual Need–Produce Time vs. from/in–Register Time



- All we care about is actual Results 'Forwarded' from Producer to Consumer instruction
- We can 'Bypass' the 'official' loop through the Register File for immediate-use Results

ALU result to next I; Load result to next-after-next I

from ALU Instructions:



from Load Instruction:



• ALU instructions never stall the pipeline, but Load instructions will do so when immediately followed by a dependent instruction



Distance 1 dependence on LOAD: Wait!



- The instruction immediately after a Load wants to use the load'ed data
- Impossible without losing one cycle: force this instruction to wait (repeat itself on the next cycle)
- Simple, in-order Pipeline: the next instruction has to wait too!





2 extra clock cycles lost

What if the program is?: RAW dependence?



- Does the compiler know for sure if i!=j
 - (OK to reorder sd–ld) or i==j (fwd in reg.)?
- If unknown to compiler, static sch. impossible
 - => dynamic scheduling at runtime (ooo pipe)

Control Dependences (branch/jump) in Pipelines

- 'Data Dependence' = next instruction uses data (register/memory) from previous
- 'Control Dependence' = which is the next instruction depends on the previous
- Control Dependences arise from 'Control Transfer Instructions (CTI)'
- Control Transfer Instructions (CTI) are: Jump and Branch Instructions
- 'Jumps' are Unconditional CTI's: they always transfer control
- 'Branches' are Conditional CTI's: whether or not they transfer control depends on the result of a data comparison that they have to perform

Statistics (rough numbers, in a majority of programs, but NOT always so):

- Branches are about 15–16% of all ('dynamically') executed instructions in a program
 - about 2/3 of executed branches are 'taken' (successful) = $\sim 10\%$ of all instr.
 - about 1/3 of executed branches are not taken (unsuccessful) = $\sim 5\%$ of all instr.
 - most backwards branches appear in loops, and they are about 90% taken
- Jumps are about 4–5% of all executed instructions in a program
 - procedure calls are about 1%, and returns another ~1%, of all executed instr.

Copyright 2020 University of Crete – https://www.csd.uoc.gr/~hy225/20a/copyright.html (slides 1–9); Elsevier (slides 10–11)

Branch Taken example

add

fetch

► 44.

+4

36:

fetch

► 40:

- In modern processors, branch latency is quite long
- In our simple pipeline, branch latency is 2 cycles (read registers; compare) (with MIPS-style comparisons (beq/bne only) it could even be 1 cycle)

DM

branch! (2 or more cycles)

sd

• Example here with 3–cycle branch latency

ALU

fetch



► 48: and fetch • need to abort +4 noop noop noop noop ► 52: fetch speculative execution +4 - 56 before it causes fetch **ALU** DM WB ld permanent damage: before DM and WB stages 76: fetch ALU DM xor

WB

- In this example, each taken branch causes the loss of 3 extra clock cycles
- About 2/3 of all executed branches are taken, so this is a heavy loss

Branch Target Buffer (BTB)

 A small table – a cache, like a hash table – containing pairs of (instruction) addresses for which there is statistical evidence that their next–PC is something other than PC+4

PC of a jump or branch–likely instruction;

Target PC to which this instruction usually went, in the past.



| | ' |
|-----|-----|
| | |
| 260 | 200 |
| 40 | 72 |
| 88 | 120 |
| 180 | 160 |
| | |

- A 'best approximation' not necessarily correct information
- Branches that are believed not-taken are NOT entered into the BTB
- Like IM –the Instruction Cache– this will oftentimes 'overflow': old pairs are removed to make room for more recent ones
- May be complemented with a small hardware stack:
 - on every call (jal ra,...), push the return address;
 - on every return (jr ra), pop an address and predict jumpin to that one
- In parallel with each Fetch, search the fetched instruction's PC value in the BTB

When the BTB prediction is Correct

• When a matching BTB entry is found, use its Prediction; else, fetch from PC+4

- - -





• When Prediction is Correct, NO extra clock cycles are lost!

When the BTB prediction is Wrong

• Prediction says: After fetching from 40, fetch from 72

ALU

fetch

BTB

^76:

36:

fetch

BTB

260

40

88

180

(BTB)

40:

_ _ _

200

72

120

160

add

fetch

≁72:

• But this time, the branch ends up going the other way: to 44

DM

ld

fetch

BTB



• When Mispredicted, branches cost 3 extra clock cycles in this pipeline 9

BTB

Relative Performance

- Define Performance = 1/Execution Time
- "X is n time faster than Y"

 $Performance_{x}/Performance_{y}$

= Execution time $_{\rm Y}$ / Execution time $_{\rm X}$ = n

- Example: time taken to run a program
 - 10s on A, 15s on B
 - Execution Time_B / Execution Time_A = 15s / 10s = 1.5
 - So A is 1.5 times faster than B



Performance SummaryThe BIG Picture $CPU Time = \frac{Instructions}{Program} \times \frac{Clock cycles}{Instruction} \times \frac{Seconds}{Clock cycle}$

Performance depends on

- Algorithm: affects IC, possibly CPI
- Programming language: affects IC, CPI
- Compiler: affects IC, CPI
- Instruction set architecture: affects IC, CPI, T_c



CPI in More Detail

If different instruction classes take different numbers of cycles

Clock Cycles =
$$\sum_{i=1}^{n} (CPI_i \times Instruction Count_i)$$

$$CPI = \frac{Clock Cycles}{Instruction Count} = \sum_{i=1}^{n} \left(CPI_i \times \frac{Instruction Count_i}{Instruction Count} \right)$$
Relative frequency



Average Access Time - Caches

- Hit time is also important for performance
- Average memory access time (AMAT)
 - AMAT = Hit time + Miss rate × Miss penalty
- Example
 - CPU with 1ns clock, hit time = 1 cycle, miss penalty = 20 cycles, I-cache miss rate = 5%
 - AMAT = 1 + 0.05 × 20 = 2ns
 - 2 cycles per instruction



Measuring Cache Performance

- Components of CPU time
 - Program execution cycles
 - Includes cache hit time
 - Memory stall cycles
 - Mainly from cache misses
- With simplifying assumptions:

Memory stall cycles

= Memory accesses Program × Miss rate × Miss penalty

 $= \frac{\text{Instructions}}{\text{Program}} \times \frac{\text{Misses}}{\text{Instruction}} \times \text{Miss penalty}$

