## **Pre-requisites**

332 Advanced Computer Architecture Chapter 1

Introduction and review of Pipelines, Performance, Caches, and Virtual Memory

> January 2006 Paul H J Kelly

These lecture notes are partly based on the course text, Hennessy and Patterson's Computer Architecture, a quantitative approach (3<sup>rd</sup> ed), and on the lecture slides of David Patterson's Berkeley course (CS252)

Course materials online at <u>http://www.doc.ic.ac.uk/~phjk/AdvancedCompArchitecture.html</u>

Advanced Computer Architecture Chapter 1. p1

- This a third-level computer architecture course
- The usual path would be to take this course after following a course based on a textbook like "Computer Organization and Design" (Patterson and Hennessy, Morgan Kaufmann)
- This course is based on the more advanced book by the same authors (see next slide)
- You can take this course provided you're prepared to catch up if necessary
  - Read chapters 1 to 8 of "Computer Organization and Design" (COD) if this material is new to you
  - If you have studied computer architecture before, make sure COD Chapters 2, 6, 7 are familiar
  - See also "Appendix A Pipelining: Basic and Intermediate Concepts" of course textbook
- FAST review today of Pipelining, Performance, Caches, and Virtual Memory

Advanced Computer Architecture Chapter 1, p2

## This is a textbook-based course



- Computer Architecture: A Quantitative Approach (3rd Edition) John L. Hennessy, David A. Patterson
  - 1128 pages. Morgan Kaufmann (29 May, 2002); ISBN: 1558607242
  - Price: £ 36.99 (Amazon.co.uk, Jan 2006)
  - Publisher's companion web site: http://books.elsevier.com/companions/1558605967/
  - Web site provides seven appendices not included in print edition, including:
    - Market Processors Processors
    - Appendix H: Computer Arithmetic
    - Appendix I: Implementing Coherence Protocols

Advanced Computer Architecture Chapter 1, p3

## Who are these guys anyway and why should I read their book?

- John Hennessy:
  - Founder, MIPS Computer Systems
  - President, Stanford University
     (previous president: Condoleezza Rice)

#### David Patterson

- Leader, Berkeley RISC project (led to Sun's SPARC)
- RAID (redundant arrays of inexpensive disks)
- Professor, University of California, Berkeley



RAID-I (1989) consisted of a Sun 4/280 workstation with 128 MB of DRAM, four dualstring SCSI controllers, 28 5.25-inch SCSI disks and specialized disk striping software.



RISC-I (1982) Contains 44,420 transistors, fabbed in 5 micron NMOS, with a die area of 77 mm<sup>2</sup>, ran at 1 MHz. This chip is probably the first VLSI RISC.

## Administrivia

#### Course web site:

http://www.doc.ic.ac.uk/~phjk/AdvancedCompArchitecture.html

See web site for link to course news group:

mews:icdoc.courses.332aca.

#### Course textbook: H&P 3rd ed

Read Appendix A right away

#### Lecturer: Paul Kelly

## Course organisation

- Tutorial helper:
  - Ashley Brown PhD student working on heterogenous multicore architectures and designspace exploration
- 3 hours per week
- Nominally two hours of lectures, one hour of classroom tutorials
- We will use the time more flexibly

#### Assessment:

- E Exam
  - IN For CS M.Eng. Class, exam will take place on last Tuesday of term
  - IN For everyone else, exam will take place early in the summer term
  - In The goal of the course is to teach you how to think about computer architecture
  - In The exam usually includes some architectural ideas not presented in the lectures
- Coursework
  - IN You will be assigned a substantial, laboratory-based exercise
  - > You will learn about performance tuning for computationally-intensive kernels
  - > You will learn about using simulators, and experimentally evaluating hypotheses to understand system performance
  - In Use the machines in Studio A to get started and get help during tutorials
- Please do not use the computers for anything else during classes

Advanced Computer Architecture Chapter 1, p6

#### Ch1

Review of pipelined, in-order processor architecture and simple cache structures

#### Ch2

- Virtual memory
- Benchmarking
- Fab
- Ch3
  - Caches in more depth
  - Software techniques to improve cache performance

#### Ch4

- Instruction-level parallelism
- Dynamic scheduling, out-of-order
- Register renaming
- Speculative execution
- Branch prediction
- Limits to ILP

#### Ch5

- Compiler techniques loop nest transformations
- Loop parallelisation, interchange, tiling/blocking, skewing

Advanced Computer Architecture Chapter 1 n5

- Uniform frameworks
- Ch6
  - Multithreading, hyperthreading, SMT
  - Static instruction scheduling
  - Software pipelining EPIC/IA-64: instruction-set support for
  - speculation and register renaming

#### 🔶 Ch7

- Shared-memory multiprocessors
- Cache coherency
- Large-scale cache-coherency; ccNUMA. COMA

#### Lab-based coursework exercise:

- Simulation study
- "challenge"
- Using performance analysis tools
- 🔹 Exam:

## Course overview

- Answer 3 questions out of 4
- Partially based on recent processor architecture article, which we will study in advance (see past papers)

- A "Typical" RISC
- ♦ 32-bit fixed format instruction (3 formats, see next slide)
- ♦ 32 32-bit general-purpose registers
- (RO contains zero, double-precision/long operands occupy a pair)
- Memory access only via load/store instructions
  - No instruction both accesses memory and does arithmetic
- All arithmetic is done on registers
- 3-address, reg-reg arithmetic instruction Subw r1,r2,r3 means r1 := r2-r3
  - registers identifiers always occupy same bits of instruction encoding
- Single addressing mode for load/store:
  - base + displacement
  - ie register contents are added to constant from instruction word, and used as address, eq "lw R2,100(r1)" means "r2 := Mem[100+r1]" no indirection
- Simple branch conditions
- Delayed branch
- SPARC, MIPS, ARM, HP PA-Risc, see: DEC Alpha, IBM PowerPC, CDC 6600, CDC 7600, Cray-1,
- Cray-2, Cray-3 Intel IA-32, IA-64 (?), Not: Motorola 68000 DEC VAX, PDP-11, IBM 360/370
- Eg: VAX matchc instruction! re Chapter 1 pi



**Register-Register** 31 26 25 21 20 16 15 11 10 65 0 Rs2 Rd Оp Rs1 Opx **Register-Immediate** 26 25 21 20 16 15 31 0 immediate Rd Op Rs1 Branch 31 26 25 21 20 16 15 0 immediate Оp Rs1 Rs2/Opx Jump / Call 31 26 25 ٥ target Оp

Q: What is the largest signed immediate operand for "subw r1,r2,X"? Q: To what range of addresses can a conditional branch jump to?

Advanced Computer Architecture Chapter 1, p9



Advanced Computer Architecture Chapter 1, p10







Advanced Computer Architecture Chapter 1, p13

**Pipelined Laundry:** 



**Pipelined Laundry:** 

- Unbalanced lengths of pipe stages reduces speedup

Advanced Computer Architecture Chapter 1, p14

## How can we apply this idea to the MIPS datapath?



## 5-stage MIPS pipeline with pipeline buffers



## Visualizing Pipelining

Figure 3.3, Page 133 , CA:AQA 2e



Advanced Computer Architecture Chapter 1, p17

## It's Not That Easy for Computers

- Limits to pipelining: Hazards prevent next instruction from executing during its designated clock cycle
  - Structural hazards: HW cannot support this combination of instructions (single person to fold and put clothes away)
  - <u>Data hazards</u>: Instruction depends on result of prior instruction still in the pipeline (missing sock)
  - <u>Control hazards</u>: Caused by delay between the fetching of instructions and decisions about changes in control flow (branches and jumps).

Advanced Computer Architecture Chapter 1, p18

#### One Memory Port/Structural Hazards Figure 3.6, Page 142, CA:AQA 2e



#### One Memory Port/Structural Hazards Figure 3.7, Page 143, CA:AQA 2e







Advanced Computer Architecture Chapter 1, p22

## Three Generic Data Hazards

 Write After Read (WAR) Instr<sub>J</sub> writes operand <u>before</u> Instr<sub>I</sub> reads it

> I: sub r4,r1,r3 J: add r1,r2,r3 K: mul r6,r1,r7

- Called an "anti-dependence" by compiler writers. This results from reuse of the name "r1".
- Can't happen in MIPS 5 stage pipeline because:
  - All instructions take 5 stages, and
  - Reads are always in stage 2, and
  - Writes are always in stage 5

Advanced Computer Architecture Chapter 1, p24

## Three Generic Data Hazards

#### Read After Write (RAW)

Instr<sub>J</sub> tries to read operand before  $Instr_J$  writes it

I: add r1,r2,r3 J: sub r4,r1,r3

 Caused by a "Dependence" (in compiler nomenclature). This hazard results from an actual need for communication.

## Three Generic Data Hazards





- Called an "output dependence" by compiler writers This also results from the reuse of name "r1".
- Can't happen in MIPS 5 stage pipeline because:
  - All instructions take 5 stages, and
  - Writes are always in stage 5
- Will see WAR and WAW in later more complicated pipes

Advanced Computer Architecture Chapter 1. p25

HW Change for Forwarding

Figure 3.20, Page 161, CA:AQA 2e



Advanced Computer Architecture Chapter 1. p26









## Software Scheduling to Avoid Load Hazards



Advanced Computer Architecture Chapter 1. p30

## Software Scheduling to Avoid Load Hazards







## Example: Branch Stall Impact

- If 30% branch, Stall 3 cycles significant
- Two part solution:
   Determine branch taken or not sooner, AND
   Compute taken branch address earlier
- MIPS branch tests if register = 0 or  $\neq$  0
- MIPS Solution:
  - Move Zero test to ID/RF stage
     Adder to calculate new PC in ID/RF stage
  - 1 clock cycle penalty for branch versus 3



Advanced Computer Architecture Chapter 1. p33

Advanced Computer Architecture Chapter 1, p35



- #1: Stall until branch direction is clear
- #2: Predict Branch Not Taken
  - Execute successor instructions in sequence
  - Squash" instructions in pipeline if branch actually taken
  - Advantage of late pipeline state update
  - 47% MIPS branches not taken on average
  - PC+4 already calculated, so use it to get next instruction

#3: Predict Branch Taken

- 53% MIPS branches taken on average
- But haven't calculated branch target address in MIPS
  - MIPS still incurs 1 cycle branch penalty
  - In Other machines: branch target known before outcome

## Four Branch Hazard Alternatives

#### #4: Delayed Branch

Define branch to take place AFTER a following instruction

branch instruction sequential successor<sub>1</sub> sequential successor<sub>2</sub> ...... sequential successor<sub>n</sub> branch target if taken

 $\geq$  Branch delay of length  $\it n$ 

I slot delay allows proper decision and branch target address in 5 stage pipeline

5 stage pipeline MIPS uses this BEQZ R1, L1 SW R3, X SW R4, X L1: LW R5, X

## **Delayed Branch**

1.1: targe

## Four Branch Hazard Alternatives

#### #4: Delayed Branch

Define branch to take place **AFTER** a following instruction

branch instruction sequential successor <sub>1</sub> sequential successor <sub>2</sub>	
sequential successor_n branch target if taken	Branch delay of length n

#### ■ 1 slot delay allows proper decision and branch target address in

	MIPS uses this LW MIPS uses this BE SV L1: LW	V R3, #100 V R4, #200 QZ R1, L1 V R3, X V R4, X V R5,X	implements	If (R1==0) X=100 Else X=200 R5 = X	
--	---	---	------------	--	--

# Now, review basic performance issues in processor design

## Which is faster?

Plane	Washington DC to Paris	Speed	Passengers	Throughput (pmph)	First flew in February, 1 <u>9</u> 69
Boeing 747	6.5 hours	610 mph	470	286,700	
BAC/Sud Concorde	3 hours	1350 mph	132	178,200	First Play in Becember 1967

- Time to run the task (ExTime) - Execution time, response time, latency
- Tasks per day, hour, week, sec, ns ... (Performance) - Throughput, bandwidth

(Background material not covered in lectures)

Advanced Computer Architecture Chapter 1. p39

- Where to get instructions to fill branch delay slot?
  - Before branch instruction
  - From the target address: only valuable when branch taken
     From fall through: only valuable when branch not taken
- Compiler effectiveness for single branch delay slot:
  - Fills about 60% of branch delay slots
  - About 80% of instructions executed in branch delay slots useful in computation
  - About 50% (60% × 80%) of slots usefully filled
- Delayed Branch downside: 7-8 stage pipelines, multiple instructions issued per clock (superscalar)

#### Canceling branches

- Branch delay slot instruction is executed but write-back is disabled if it is not supposed to be executed
- Two variants: branch "likely taken", branch "likely not-taken"
- allows more slots to be filled

Advanced Computer Architecture Chapter 1, p38

before Blt R1,L1 fallthru

## Definitions

- Performance is in units of things per sec
   bigger is better
- •If we are primarily concerned with response time

<pre>performance(x) =</pre>	1
	execution_time(x)

" X is n times faster than Y" means

n =  $\frac{\text{Performance(X)}}{\text{Performance(Y)}} = \frac{\text{Execution_time(Y)}}{\text{Execution_time(X)}}$ 

(Background material not covered in lectures)

## Aspects of CPU Performance (CPU Law)

CPU time	= Seconds	= Instructions x	Cycles x	Seconds
	Program	Program	Instruction	Cycle

	Inst Count	CPI	Clock Rate
Program	X		
Compiler	X	(X)	
Inst. Set.	X	X	
Organization		X	X
Technology			X

Freq Cycles CPI(i) (% Time)

.5

.4

.2

.4

1.5

#### (Background material not covered in lectures)

Base Machine (Reg / Reg)

50% 1

20% 2

10% 2

20%

Typical Mix of

instruction types in program

2

Ор

ALU

Load

Store

Branch

Advanced Computer Architecture Chapter 1. p41

Example: Calculating CPI

(33%)

(27%)

(13%)

(27%)

## Cycles Per Instruction (Throughput)

## "Average Cycles per Instruction"

### "Instruction Frequency"

$$CPI = \sum_{j=1}^{n} CPI_{j} \times F_{j} \quad \text{where } F_{j} = \frac{I_{j}}{\text{Instructio n Count}}$$

(Background material not covered in lectures)

Advanced Computer Architecture Chapter 1, p42

## **Example: Branch Stall Impact**

<ul> <li>Assume CPI = 1.0 ignoring branches</li> <li>Assume solution was stalling for 3 cycles</li> <li>If 30% branch, Stall 3 cycles</li> </ul>					
<ul><li>Op</li><li>Other</li><li>Branch</li></ul>	Freq	Cycles	CPI(i)	(% Time)	
	70%	1	.7	(37%)	
	30%	4	1.2	(63%)	

> new CPI = 1.9, or almost 2 times slower

(Background material not covered in lectures)

Advanced Computer Architecture Chapter 1, p43

(Background material not covered in lectures)

## Example 2: Speed Up Equation for Pipelining

Example 3: Evaluating Branch Alternatives (for 1 program)

CPI<sub>bibelined</sub> = Ideal CPI + Average Stall cycles per Inst

 $Speedup = \frac{Ideal \ CPI \times Pipeline \ depth}{Ideal \ CPI + Pipeline \ stall \ CPI} \times \frac{Cycle \ Time_{unpipelined}}{Cycle \ Time_{pipelined}}$ 

For simple RISC pipeline, Ideal CPI = 1:

Speedup =  $\frac{\text{Pipeline depth}}{1 + \text{Pipeline stall CPI}} \times \frac{\text{Cycle Time}_{\text{unpipelined}}}{\text{Cycle Time}_{\text{ninelined}}}$ 

(Background material not covered in lectures)

Advanced Computer Architecture Chapter 1, p45

Assuming Conditional & Unconditional branches make up 14% of the total instruction count, and 65% of them change the PC

1.31

(Background material not covered in lectures)

Delayed branch 0.5 1.07

Advanced Computer Architecture Chapter 1. p46

## Example 4: Dual-port vs. Single-port

- Machine A: Dual ported memory ("Harvard Architecture")
- Machine B: Single ported memory, but its pipelined implementation has a 1.05 times faster clock rate

```
Ideal CPI = 1 for both
```

Loads are 40% of instructions executed

```
SpeedUp<sub>A</sub> = Pipeline Depth/(1 + 0) x (clock<sub>unpipe</sub>/clock<sub>pipe</sub>)

= Pipeline Depth

SpeedUp<sub>B</sub> = Pipeline Depth/(1 + 0.4 x 1) x (clock<sub>unpipe</sub>/(clock<sub>unpipe</sub>/ 1.05)

= (Pipeline Depth/1.4) x 1.05

= 0.75 x Pipeline Depth

SpeedUp<sub>A</sub> / SpeedUp<sub>B</sub> = Pipeline Depth/(0.75 x Pipeline Depth) = 1.33
```

Machine A is 1.33 times faster

(Background material not covered in lectures)

Advanced Computer Architecture Chapter 1, p47

Now, Review of Memory Hierarchy



#### Recap: Who Cares About the Memory Hierarchy?



The Principle of Locality

The Principle of Locality:
 Programs access a relatively small portion of the address space at any instant of time.



## Memory Hierarchy: Terminology

- Hit: data appears in some block in the upper level (example: Block X)
  - Hit Rate: the fraction of memory access found in the upper level
  - Hit Time: Time to access the upper level which consists of RAM access time + Time to determine hit/miss
- Miss: data needs to be retrieved from a block in the lower level (Block Y)
  - Miss Rate = 1 (Hit Rate)
  - Miss Penalty: Time to replace a block in the upper level + Time to deliver the block the processor
- Hit Time << Miss Penalty (500 instructions on 21264!)</p>



## 1 KB Direct Mapped Cache, 32B blocks



Advanced Computer Architecture Chapter 1 p5

Advanced Computer Architecture Chapter 1 p54

Byte 992 31

Byte Select

•

••

Ex: 0x00

## 1 KB Direct Mapped Cache, 32B blocks

#### For a 2 \*\* N byte cache:





Direct-mapped cache - storage

## Direct-mapped Cache - structure

- Capacity: C bytes (eg 1KB)
- Blocksize: B bytes (eg 32)
- Byte select bits: 0..log(B)-1 (eg 0..4)
- Number of blocks: C/B (eg 32)
- Address size: A (eg 32 bits)
- Cache index size: I=log(C/B) (eg log(32)=5)
- Tag size: A-I-log(B) (eg 32-5-5=22)



## Two-way Set Associative Cache

- N-way set associative: N entries for each Cache Index
  - N direct mapped caches operated in parallel (N typically 2 to 4)
- Example: Two-way set associative cache
- Cache Index selects a "set" from the cache
- The two tags in the set are compared in parallel
- Data is selected based on the tag result



## Disadvantage of Set Associative Cache

- N-way Set Associative Cache v. Direct Mapped Cache:
  - N comparators vs. 1
  - Extra MUX delay for the data
  - Data comes AFTER Hit/Miss
- In a direct mapped cache, Cache Block is available BEFORE Hit/Miss:
   Possible to assume a hit and continue. Recover later if miss.



## Basic cache terminology

Example: Intel Pentium 4 Level-1 cache (pre-Prescott)

- Capacity: 8K bytes (total amount of data cache can store)
- Block: 64 bytes (so there are 8K/64=128 blocks in the cache)
- Sets: 4 (addresses with same index bits can be placed in one of 4 ways)
- Ways: 32 (=128/4, that is each RAM array holds 32 blocks)
- Index: 5 bits (since 25=32 and we need index to select one of the 32 ways)
- Tag: 21 bits (=32 minus 5 for index, minus 6 to address byte within block)
- Access time: 2 cycles, (.6ns at 3GHz; pipelined, dual-ported [load+store])



## 4 Questions for Memory Hierarchy

- Q1: Where can a block be placed in the upper level? (Block placement)
- Q2: How is a block found if it is in the upper level? (Block identification)
- Q3: Which block should be replaced on a miss? (Block replacement)
- Q4: What happens on a write? (Write strategy)



Advanced Computer Architecture Chapter 1. p61

## Q2: How is a block found if it is in the upper level?



Tag on each block

No need to check index or block offset



Increasing associativity shrinks index, expands tag

Advanced Computer Architecture Chapter 1. p63

## Q3: Which block should be replaced on a miss?

- Easy for Direct Mapped
- Set Associative or Fully Associative:
  - Random
  - LRU (Least Recently Used)

Assoc:	2-w	ay	4-wa	у	8-wa	ау
Size	LRU	Ran	LRU	Ran	LRU	Ran
16 KB	5.2%	5.7%	4.7%	5.3%	4.4%	5.0%
64 KB	1.9%	2.0%	1.5%	1.7%	1.4%	1.5%
256 KB	1.15%	1.17%	1.13%	1.13%	1.12%	1.12%

## Q4: What happens on a write?

- <u>Write through</u>—The information is written to both the block in the cache and to the block in the lower-level memory.
- <u>Write back</u>—The information is written only to the block in the cache. The modified cache block is written to main memory only when it is replaced.
  - is block clean or dirty?
- Pros and Cons of each?
  - WT: read misses cannot result in writes
  - WB: no repeated writes to same location
- WT always combined with write buffers so that don't wait for lower level memory

Advanced Computer Architecture Chapter 1, p65

## Write Buffer for Write Through



- A Write Buffer is needed between the Cache and Memory
  - Processor: writes data into the cache and the write buffer
     Memory controller: write contents of the buffer to memory
- Write buffer is just a FIFO:
  - Typical number of entries: 4
     Works fine if: Store frequency (w.r.t. time) << 1 / DRAM write cycle</li>
- Memory system designer's nightmare:
  - Store frequency (w.r.t. time) -> 1 / DRAM write cycle
     Write buffer saturation

Advanced Computer Architecture Chapter 1, p66

## A Modern Memory Hierarchy

- By taking advantage of the principle of locality:
  - Present the user with as much memory as is available in the cheapest technology.
  - Provide access at the speed offered by the fastest technology.



anced Computer Architecture Chapter 1. p67

## Summary #1/4: Pipelining & Performance

Just overlap tasks; easy if tasks are independent

◆ Speed Up ≤ Pipeline Depth; if ideal CPI is 1, then:

Speedup = 
$$\frac{\text{Pipeline depth}}{1 + \text{Pipeline stall CPI}} \times \frac{\text{Cycle Time}_{unpipelined}}{\text{Cycle Time}_{pipelined}}$$

- Hazards limit performance on computers:
  - Structural: need more HW resources
  - Data (RAW, WAR, WAW): need forwarding, compiler scheduling
     Control: delayed branch, prediction
- Time is measure of performance: latency or throughput
- CPI Law:

CPU time	= Seconds	= Instructions >	Cycles x	Seconds
	Program	Program	Instruction	Cycle

## Summary #2/4: Caches

- The Principle of Locality:
  - Program access a relatively small portion of the address space at any instant of time.
    - Temporal Locality: Locality in Time
    - In Spatial Locality: Locality in Space
- Three Major Categories of Cache Misses:
  - **Compulsory Misses:** sad facts of life. Example: cold start misses.
  - Capacity Misses: increase cache size
  - **<u>Conflict Misses</u>**: increase cache size and/or associativity.
- Write Policy:
  - Write Through: needs a write buffer.
  - Write Back: control can be complex
- Today CPU time is often dominated by memory access time, not just computational work. What does this mean to Compilers, Data structures, Algorithms?