ECE 552 / CPS 550 Advanced Computer Architecture I

Lecture 8 Instruction-Level Parallelism – Part 1

Benjamin Lee Electrical and Computer Engineering Duke University

www.duke.edu/~bcl15 www.duke.edu/~bcl15/class/class_ece252fall12.html



27 September – Homework #2 Due

- Use blackboard forum for questions
- Attend office hours with questions
- Email for separate meetings

2 October – Class Discussion

Roughly one reading per class. Do not wait until the day before!

- 1. Srinivasan et al. "Optimizing pipelines for power and performance"
- 2. Mahlke et al. "A comparison of full and partial predicated execution support for ILP processors"
- 3. Palacharla et al. "Complexity-effective superscalar processors"
- 4. Yeh et al. "Two-level adaptive training branch prediction"



Pipelining becomes complex when we want high performance in the presence of...

- Long latency or partially pipelined floating-point units
- Memory systems with variable access time
- Multiple arithmetic and memory units

MIPS Floating Point

- Interaction between floating-point (FP), integer datapath defined by ISA
- Architect separate register files for floating point (FPR) and integer (GPR)
- Define separate load/store instructions for FPR, GPR
- Define move instructions between register files
- Define FP branches in terms of FP-specific condition codes



FPU requires much more hardware than integer unit

Single-cycle FPU a bad idea

- Why?
- It is common to have several, different types of FPUs (Fadd, Fmul, etc.)
- FPU may be pipelined, partially pipelined, or not pipelined

Floating-point Register File (FPR)

- To operate several FPUs concurrently, FPR requires several read/write ports





Functional units have internal pipeline registers

- Inputs to a functional unit (e.g., register file) can change during a long latency operation
- Operands are latched when an instruction enters the functional unit



Latency of main memory access usually greater than one cycle and often unpredictable

- Solving this problem is a central issue in computer architecture

Improving memory performance

- Separate instruction and data memory ports, no self-modifying code
- Caches -- size L1 cache for single-cycle access
- Caches -- L1 miss stalls pipeline
- Memory interleaving memory allows multiple simultaneous access
- Memory bank conflicts stall the pipeline









Implications of multi-cycle instructions

- FPU or memory unit requires more than one cycle
- Structural conflict in execution stage, if FPU or memory unit is not pipelined

Different functional unit latencies

- Structural conflict in writeback stage due to different latencies
- Out-of-order write conflicts due to variable latencies

How to handle exceptions?









Superscalar In-Order Pipeline



• ECE 552 / CPS 550



Consider executing a sequence of instructions of the form: $Rk \leftarrow (Ri) op (Rj)$

Data Dependence

R3 ← (R1) op (R4) R5 ← (R3) op (R4) # RAW hazard (R3)

Anti-dependence

R3 ← (R1) op (R2) R1 ← (R4) op (R5)

WAR hazard (R1)

Output-dependence

R3 ← (R1) op (R2) R3 ← (R6) op (R7)

WAW hazard (R3)



Range and Domain of Instruction (j)

R(j) = registers (or other storage) modified by instruction j D(j) = registers (or other storage) read by instruction j

Suppose instruction k follows instruction j in program order. Executing instruction k before the effect of instruction j has occurred can cause...

RAW hazard if $R(j) \cap D(k) \neq \emptyset$ # j modifies a register read by kWAR hazard if $D(j) \cap R(k) \neq \emptyset$ # j reads a register modified by kWAW hazard if $R(j) \cap R(k) \neq \emptyset$ # j, k modify the same register



Data hazards due to register operands can be determined at decode stage

Data hazards due to memory operands can be determined <u>only after computing effective address</u> in execute stage

store	$M[R1 + disp1] \leftarrow R2$
load	$R3 \leftarrow M[R4 + disp2]$

(R1 + disp1) == (R4 + disp2)?





RAW Hazards WAR Hazards WAW Hazards





Valid Instruction Orderings in-order 6 13 5 I_4 I_2 1 out-of-order $|_4$ I_5 6 l₂ 1 I_3 out-of-order $|_1$ $|_2$ $|_4$ **1**₅ I₆ l₃



	I ₁					
I ₁	DIVD	f6,	f6,	f4	<u>Latency</u> 4	I ₂
I_2	LD	f2,	45(r3)		1	I ₃
l ₃	MULTD	fO,	f2,	f4	3	I
I_4	DIVD	f8,	f6,	f2	4	
I_5	SUBD	f10,	fO,	f6	1	I ₅
I ₆	ADDD	f6,	f8,	f2	1	

Let k indicate when instruction k is issued. Let \underline{k} denote when instruction k is completed. **1**6

	Out-of-					
I ₁	DIVD	f6,	f6,	f4	<u>Latency</u> 4	I ₂
I_2	LD	f2,	45(r3)		1	I_3
l ₃	MULTD	fO,	f2,	f4	3	
I_4	DIVD	f8,	f6,	f2	4	
I_5	SUBD	f10,	fO,	f6	1	I ₅
۱ ₆	ADDD	f6,	f8,	f2	1	I ₆
	in-order comp	1 2		<u>1</u> <u>2</u> 3 4	<u>3</u> 5 <u>4</u>	6 <u>5</u> <u>6</u>
	out-of-order comp	12	<u>2</u> 3	<u>1</u> 4 <u>3</u> 5	<u>5</u> <u>4</u> 6 <u>6</u>	



Up until now, we assumed user or compiler <u>statically</u> examines instructions, detecting hazards and scheduling instructions

Scoreboard is a hardware data structure to <u>dynamically</u> detect hazards







Seymour Cray, 1963

- Fast, pipelined machine with 60-bit words
- 128 Kword main memory capacity, 32-banks
- Ten functional units (parallel, unpipelined)
- Floating-point: adder, 2 multipliers, divider
- Integer: adder, 2 incrementers
- Dynamic instruction scheduling with scoreboard
- Ten peripheral processors for I/O

-More than 400K transistors, 750 sq-ft, 5 tons, 150kW with novel Freon-based cooling

- Very fast clock, 10MHz (FP add in 4 clocks)
- Fastest machine in world for 5 years
- Over 100 sold (\$7-10M each)



Thomas Watson Jr., IBM CEO, August 1963

"Last week, Control Data....announced the 6600 system. I understand that in the laboratory developing the system there are only 34 people including the janitor. Of these, 14 are engineers and 4 are programmers...Contrasting this modest effort with our vast development activities, I fail to understand why we have lost our industry leadership by letting someone else offer the world's most powerful computer."

To which Cray replied...

"It seems like Mr. Watson has answered his own question."





• ECE 552 / CPS 550



When is it safe to issue an instruction?

- Suppose a data structure tracks all instructions in all functional units

Before issuing instruction, issue logic must check:

- Is the required functional unit available? Check for structural hazard.
- Is the input data available? Check for RAW hazard.
- Is it safe to write the destination? Check for WAR, WAW hazard
- Is there a structural hazard at the write back stage?



In issue stage, instruction j consults the table

- Functional unit available? Check the busy column

- RAW? Search the Dest column for j's sources - WAR? Search the Src1/2 columns for j's destination
- WAW? Search the Dest column for j's destination

If no hazard, add entry and issue instruction

Upon instruction write back, remove entry

Name	Busy	Ор	Dest	Src1	Src2	
Int						
Mem						
Add1						
Add2						
Add3						
Mult1						
Mult2						



Assume instructions issue in-order

Assume issue logic does not dispatch instruction if it detects RAW hazard or busy functional unit

Assume functional unit latches operands when the instruction is issued



Can the dispatched instruction cause WAR hazard?

- No, because instructions issue in order and operands are read at issue

No WAR Hazards

- No need to track source-1 and source-2

Can the dispatched instruction cause WAW hazard?

- Yes, because instructions may complete out-of-order

Do not issue instruction in case of WAW hazard

- In scoreboard, a register name occurs at most once in 'dest' column



Busy[FU#]: a bit-vector to indicate functional unit availability (FU = Int, Add, Mutl, Div)

WP[#regs]: a bit-vector to record the registers to which writes are pending

- Bits are set to true by issue logic
- Bits are set to false by writeback stage

- Each functional unit's pipeline registers must carry 'dest' field and a flag to indicate if it's valid: "the (we, ws) pair"

Issue logic checks instruction (opcode, dest, src1, src2) against scoreboard (busy, wp) to dispatch

- FU available? Busy - RAW? WP[s
 - WAR?
 - MAMʻS

Busy[FU#] WP[src1] or WP[src2]

Cannot arise WP[dest]

	Busy-Functional Units Status										Writes Pending (WP)
	Int(1)	Add(1)	I M	ult	(3)		D۱	/(4)	WB	
t0	I_1					f6					f6
t1	<i>I</i> ₂ f2						f6				<mark>f6</mark> , f2
t2								f6		f2	f6, f2 <u>I</u> 2
t3	I_3		f0						f6		<mark>f6</mark> , f0
t4				f0						f6	f6, f0 <u>I</u> 1
t5	I_4				f0	f8					f0, f8
t6							f8			f0	f0, f8 <u>I</u> 3
t7	I_5	f10						f 8			f8, f10
t8									f8	f10	f8, f10 <u>I</u> 5
t9										f8	f8 <u>I</u> 4
t10	I_6	f6									f6
t11										f6	f6 <u>I</u> ₆

f2

DIVD f6, f6, $egin{array}{c} I_1 \ I_2 \ I_3 \end{array}$ f4 f2, 45(r3) LD f2, MULTD f0, f4 I_4 **f**2 f8, f6, DIVD I_5 f0, f6 **SUBD** f10, I_6 f6, f8, ADDD

Instruction Issue Logic FU available? Busy[FU#] RAW? WP[src1] or WP[src2] WAR? Cannot arise WAW? WP[dest]



Detect hazards dynamically

Issue instructions in-order Complete instructions out-of-order

Increases instruction-level-parallelism by

- More effectively exploiting multiple functional units
- Reducing the number of pipeline stalls due to hazards



These slides contain material developed and copyright by

- Arvind (MIT)
- Krste Asanovic (MIT/UCB)
- Joel Emer (Intel/MIT)
- James Hoe (CMU)
- John Kubiatowicz (UCB)
- Alvin Lebeck (Duke)
- David Patterson (UCB)
- Daniel Sorin (Duke)