CS 423 Computer Architecture Spring 2012

Lecture 02: Performance, MIPS Datapath, MIPS Pipeline Ozcan Ozturk

http://www.cs.bilkent.edu.tr/~ozturk/cs423/ [Adapted from *Computer Organization and Design*, Patterson & Hennessy, © 2005, UCB]

Performance Metrics

- Purchasing perspective
 - given a collection of machines, which has the
 - best performance ?
 - least cost ?
 - best cost/performance?
- Design perspective
 - faced with design options, which has the
 - best performance improvement ?
 - least cost ?
 - best cost/performance?
- Both require
 - basis for comparison
 - metric for evaluation
- Our goal is to understand what factors in the architecture contribute to overall system performance and the relative importance (and cost) of these factors

Defining (Speed) Performance

Normally interested in reducing

- Response time (aka execution time) the time between the start and the completion of a task
 - Important to individual users
- Thus, to maximize performance, need to minimize execution time

 $performance_{x} = 1 / execution_time_{x}$

If X is n times faster than Y, then

 $\frac{\text{performance}_{X}}{\text{performance}_{Y}} = \frac{\text{execution}_{time}_{Y}}{\text{execution}_{time}_{X}} = n$

- Throughput the total amount of work done in a given time
 - Important to data center managers
- Decreasing response time almost always improves throughput

Performance Factors

- Want to distinguish elapsed time and the time spent on our task
- CPU execution time (CPU time) time the CPU spends working on a task
 - Does not include time waiting for I/O or running other programs
- CPU execution time for a program = # CPU clock cycles x clock cycle time for a program

or

- CPU execution time = # CPU clock cycles for a program for a program clock rate
- Can improve performance by reducing either the length of the clock cycle or the number of clock cycles required for a program

Review: Machine Clock Rate

Clock rate (MHz, GHz) is inverse of clock cycle time (clock period)

CC = 1/CR



10 nsec clock cycle => 100 MHz clock rate

5 nsec clock cycle => 200 MHz clock rate

2 nsec clock cycle => 500 MHz clock rate

1 nsec clock cycle => 1 GHz clock rate

500 psec clock cycle => 2 GHz clock rate

250 psec clock cycle => 4 GHz clock rate

200 psec clock cycle => 5 GHz clock rate

Clock Cycles per Instruction

- Not all instructions take the same amount of time to execute
 - One way to think about execution time is that it equals the number of instructions executed multiplied by the average time per instruction
- # CPU clock cycles = # Instructions Average clock cycles for a program = for a program x Per instruction
- Clock cycles per instruction (CPI) the average number of clock cycles each instruction takes to execute
 - A way to compare two different implementations of the same ISA

	CPI for this instruction class				
	А	В	С		
CPI	1	2	3		

Effective CPI

Computing the overall effective CPI is done by looking at the different types of instructions and their individual cycle counts and averaging

Overall effective CPI =
$$\sum_{i=1}^{n} (CPI_i \times IC_i)$$

- Where IC_i is the count (percentage) of the number of instructions of class i executed
- CPI_i is the (average) number of clock cycles per instruction for that instruction class
- n is the number of instruction classes
- The overall effective CPI varies by instruction mix a measure of the dynamic frequency of instructions across one or many programs

THE Performance Equation

Our basic performance equation is then

CPU time = Instruction_count x CPI x clock_cycle

Instruction_count x CPI CPU time = -----clock_rate

or

- These equations separate the three key factors that affect performance
 - Can measure the CPU execution time by running the program
 - The clock rate is usually given
 - Can measure overall instruction count by using profilers/ simulators without knowing all of the implementation details
 - CPI varies by instruction type and ISA implementation for which we must know the implementation details

Determinates of CPU Performance

CPU time = Instruction_count x CPI x clock_cycle

	Instruction_ count	CPI	clock_cycle
Algorithm	X	X	
Programming language	X	X	
Compiler	X	X	
ISA	X	X	X
Processor organization		X	X
Technology			X

A Simple Example

Ор	Freq	CPI _i	Freq x	CPI _i			
ALU	50%	1		.5	.5	.5	.25
Load	20%	5		1.0	.4	1.0	1.0
Store	10%	3		.3	.3	.3	.3
Branch	20%	2		.4	.4	.2	.4
			$\Sigma =$	2.2	1.6	2.0	1.95

How much faster would the machine be if a better data cache reduced the average load time to 2 cycles?

CPU time new = $1.6 \times IC \times CC$ so 2.2/1.6 means 37.5% faster

How does this compare with using branch prediction to shave a cycle off the branch time?

CPU time new = $2.0 \times IC \times CC$ so 2.2/2.0 means 10% faster

What if two ALU instructions could be executed at once?

CPU time new = $1.95 \times IC \times CC$ so 2.2/1.95 means 12.8% faster

SPEC Benchmarks www.spec.org

Integer benchmarks		FP benchmarks		
gzip	compression	wupwise	Quantum chromodynamics	
vpr	FPGA place & route	swim	Shallow water model	
gcc	GNU C compiler	mgrid	Multigrid solver in 3D fields	
mcf	Combinatorial optimization	applu	Parabolic/elliptic pde	
crafty	Chess program	mesa	3D graphics library	
parser	Word processing program	galgel	Computational fluid dynamics	
eon	Computer visualization	art	Image recognition (NN)	
perlbmk	perl application	equake	Seismic wave propagation simulation	
gap	Group theory interpreter	facerec	Facial image recognition	
vortex	Object oriented database	ammp	Computational chemistry	
bzip2	compression	lucas	Primality testing	
twolf	Circuit place & route	fma3d	Crash simulation fem	
		sixtrack	Nuclear physics accel	
		apsi	Pollutant distribution	

Example SPEC Ratings



Other Performance Metrics

- Power consumption especially in the embedded market where battery life is important (and passive cooling)
 - For power-limited applications, the most important metric is energy efficiency



MIPS Basic Architecture

The Processor: Datapath & Control

- Our implementation of the MIPS is simplified
 - memory-reference instructions: lw, sw
 - arithmetic-logical instructions: add, sub, and, or, slt
 - control flow instructions: beq, j
- Generic implementation
 - use the program counter (PC) to supply the instruction address and fetch the instruction from memory (and update the PC)



- decode the instruction (and read registers)
- execute the instruction
- All instructions (except j) use the ALU after reading the registers

How? memory-reference? arithmetic? control flow?

Clocking Methodologies

- The clocking methodology defines when signals can be read and when they are written
 - An edge-triggered methodology
- Typical execution
 - read contents of state elements
 - send values through combinational logic
 - write results to one or more state elements



Assumes state elements are written on every clock cycle; if not, need explicit write control signal

 write occurs only when both the write control is asserted and the clock edge occurs

CS423 L02 Performance.18

Fetching Instructions

Fetching instructions involves

- reading the instruction from the Instruction Memory
- updating the PC to hold the address of the next instruction



- PC is updated every cycle, so it does not need an explicit write control signal
- Instruction Memory is read every cycle, so it doesn't need an explicit read control signal

Decoding Instructions

- Decoding instructions involves
 - sending the fetched instruction's opcode and function field bits to the control unit



- reading two values from the Register File
 - Register File addresses are contained in the instruction

Executing R Format Operations

□ R format operations (add, sub, slt, and, or)

	31	25	20	15	10	5 0
R-type:	ор	rs	rt	rd	sham	funct

- perform the (op and funct) operation on values in rs and rt
- store the result back into the Register File (into location rd)



 The Register File is not written every cycle (e.g. sw), so we need an explicit write control signal for the Register File

Executing Load and Store Operations

Load and store operations involves

- compute memory address by adding the base register (read from the Register File during decode) to the 16-bit signed-extended offset field in the instruction
- store value (read from the Register File during decode) written to the Data Memory
- load value, read from the Data Memory, written to the Register



Executing Branch Operations

Branch operations involves

- compare the operands read from the Register File during decode for equality (zero ALU output)
- compute the branch target address by adding the updated PC to



Executing Jump Operations

- Jump operation involves
 - replace the lower 28 bits of the PC with the lower 26 bits of the fetched instruction shifted left by 2 bits



Creating a Single Datapath from the Parts

- Assemble the datapath segments and add control lines and multiplexors as needed
- Single cycle design fetch, decode and execute each instructions in one clock cycle
 - no datapath resource can be used more than once per instruction, so some must be duplicated (e.g., separate Instruction Memory and Data Memory, several adders)
 - multiplexors needed at the input of shared elements with control lines to do the selection
 - write signals to control writing to the Register File and Data Memory

Cycle time is determined by length of the longest path

Fetch, R, and Memory Access Portions



Adding the Control

- Selecting the operations to perform (ALU, Register File and Memory read/write)
- Controlling the flow of data (multiplexor inputs)



- register
- addr. of register to be written is in one of two places in rt (bits 20-16) for lw; in rd (bits 15-11) for R-type instructions
- offset for beq, lw, and sw always in bits 15-0

Single Cycle Datapath with Control Unit



R-type Instruction Data/Control Flow



Load Word Instruction Data/Control Flow



Load Word Instruction Data/Control Flow



Branch Instruction Data/Control Flow



Branch Instruction Data/Control Flow



Adding the Jump Operation



Single Cycle Disadvantages & Advantages

- Uses the clock cycle inefficiently the clock cycle must be timed to accommodate the slowest instruction
 - especially problematic for more complex instructions like floating point multiply



May be wasteful of area since some functional units (e.g., adders) must be duplicated since they can not be shared during a clock cycle

but

Is simple and easy to understand

Multicycle Datapath Approach

- Let an instruction take more than 1 clock cycle to complete
 - Break up instructions into steps where each step takes a cycle while trying to
 - balance the amount of work to be done in each step
 - restrict each cycle to use only one major functional unit
 - Not every instruction takes the *same* number of clock cycles
- In addition to faster clock rates, multicycle allows functional units that can be used more than once per instruction as long as they are used on *different* clock cycles, as a result
 - only need one memory but only one memory access per cycle
 - need only one ALU/adder but only one ALU operation per cycle

Multicycle Datapath Approach, con't

□ At the end of a cycle

• Store values needed in a later cycle by the current instruction in an internal register (not visible to the programmer). All (except IR) hold data only between a pair of adjacent clock cycles (no write control signal needed)



IR – Instruction RegisterA, B – regfile read data registers

MDR - Memory Data Register
ALUout - ALU output register

 Data used by subsequent instructions are stored in programmer visible registers (i.e., register file, PC, or memory)

The Multicycle Datapath with Control Signals



Multicycle Control Unit

- Multicycle datapath control signals are not determined solely by the bits in the instruction
 - e.g., op code bits tell what operation the ALU should be doing, but not what instruction cycle is to be done next
- Must use a finite state machine (FSM) for control
 - a set of states (current state stored in State Register)
 - next state function (determined by current state and the input)
 - output function (determined by current state and the input)



The Five Steps of the Load Instruction



□ IFetch: Instruction Fetch and Update PC

- Dec: Instruction Decode, Register Read, Sign Extend Offset
- Exec: Execute R-type; Calculate Memory Address; Branch Comparison; Branch and Jump Completion
- Mem: Memory Read; Memory Write Completion; Rtype Completion (RegFile write)
- WB: Memory Read Completion (RegFile write) INSTRUCTIONS TAKE FROM 3 - 5 CYCLES!

Multicycle Advantages & Disadvantages

Uses the clock cycle efficiently – the clock cycle is timed to accommodate the slowest instruction step

Multicycle implementations allow functional units to be used more than once per instruction as long as they are used on different clock cycles

but

Requires additional internal state registers, more muxes, and more complicated (FSM) control

Single Cycle vs. Multiple Cycle Timing

Single Cycle Implementation:



How Can We Make It Even Faster?

- Split the multiple instruction cycle into smaller and smaller steps
 - There is a point of diminishing returns where as much time is spent loading the state registers as doing the work
- Start fetching and executing the next instruction before the current one has completed
 - Pipelining modern processors are pipelined for performance
 - Remember the performance equation:
 CDU time CDU* (

CPU time = CPI * CC * IC

□ Fetch (and execute) more than one instruction at a time

• Superscalar processing – stay tuned

A Pipelined MIPS Processor

- Start the next instruction before the current one has completed
 - improves throughput total amount of work done in a given time
 - instruction latency (execution time, delay time, response time time from the start of an instruction to its completion) is *not* reduced



- clock cycle (pipeline stage time) is limited by the slowest stage
- for some instructions, some stages are wasted cycles

Single Cycle, Multiple Cycle, vs. Pipeline



MIPS Pipeline Datapath Modifications

- □ What do we need to add/modify in our MIPS datapath?
 - State registers between each pipeline stage to isolate them



Pipelining the MIPS ISA

- What makes it easy
 - all instructions are the same length (32 bits)
 - can fetch in the 1st stage and decode in the 2nd stage
 - few instruction formats (three) with symmetry across formats
 - can begin reading register file in 2nd stage
 - memory operations can occur only in loads and stores
 - can use the execute stage to calculate memory addresses
 - each MIPS instruction writes at most one result (i.e., changes the machine state) and does so near the end of the pipeline (MEM and WB)

What makes it hard

- structural hazards: what if we had only one memory?
- control hazards: what about branches?
- data hazards: what if an instruction's input operands depend on the output of a previous instruction?

Graphically Representing MIPS Pipeline



Can help with answering questions like:

- How many cycles does it take to execute this code?
- What is the ALU doing during cycle 4?
- Is there a hazard, why does it occur, and how can it be fixed?

Why Pipeline? For Performance!

Time (clock cycles)



Can Pipelining Get Us Into Trouble?

□ Yes: Pipeline Hazards

- structural hazards: attempt to use the same resource by two different instructions at the same time
- data hazards: attempt to use data before it is ready
 - An instruction's source operand(s) are produced by a prior instruction still in the pipeline
- control hazards: attempt to make a decision about program control flow before the condition has been evaluated and the new PC target address calculated
 - branch instructions

Can always resolve hazards by waiting

- pipeline control must detect the hazard
- and take action to resolve hazards

A Single Memory Would Be a Structural Hazard

Time (clock cycles)



Fix with separate instr and data memories (I\$ and D\$)

How About Register File Access?

Time (clock cycles)



CS423 L02 Performance.53

Spring, 2012

Register Usage Can Cause Data Hazards

Dependencies backward in time cause hazards



Read before write data hazard

Loads Can Cause Data Hazards

Dependencies backward in time cause hazards



□ Load-use data hazard

One Way to "Fix" a Data Hazard



Another Way to "Fix" a Data Hazard



Forwarding with Load-use Data Hazards



Will still need one stall cycle even with forwarding

Branch Instructions Cause Control Hazards

Dependencies backward in time cause hazards



One Way to "Fix" a Control Hazard



Corrected Datapath to Save RegWrite Addr

Need to preserve the destination register address in the pipeline state registers



MIPS Pipeline Control Path Modifications

- □ All control signals can be determined during Decode
 - and held in the state registers between pipeline stages



Other Pipeline Structures Are Possible

□ What about the (slow) multiply operation?

- Make the clock twice as slow or ...
- let it take two cycles (since it doesn't use the DM stage)



- What if the data memory access is twice as slow as the instruction memory?
 - make the clock twice as slow or ...
 - let data memory access take two cycles (and keep the same clock rate)



Sample Pipeline Alternatives



IM access

DM write

reg write

exception

reg 2 access

Summary

- All modern day processors use pipelining
- Pipelining doesn't help latency of single task, it helps throughput of entire workload
- Potential speedup: a CPI of 1 and fast a CC
- Pipeline rate limited by slowest pipeline stage
 - Unbalanced pipe stages makes for inefficiencies
 - The time to "fill" pipeline and time to "drain" it can impact speedup for deep pipelines and short code runs
- Must detect and resolve hazards
 - Stalling negatively affects CPI (makes CPI less than the ideal of 1)