## **Defining Performance**





#### CS@VT

#### **Computer Organization II**

# Computer Performance: TIME, TIME, TIME Performance 2

Response Time (latency): how long it takes to complete a task

Throughput: total # of tasks completed per unit time

For now, we will focus on response time... and define

Performance = **Execution Time** 

### **Execution Time**

### Elapsed Time

- counts everything (disk and memory accesses, I/O, etc.)
- a useful number, but often not good for comparison purposes

### CPU time

- doesn't count I/O or time spent running other programs
- can be broken up into system time, and user time

### Our focus: user CPU time

- time spent executing the lines of code that are "in" our program

## **Relative Performance**



Example: time taken to run a program

- 10s on A, 15s on B
- Execution Time<sub>B</sub> / Execution Time<sub>A</sub>
  - = 15s / 10s = 1.5
- So A is 1.5 times as fast as B (and B is 2/3 as fast as A)

wrt this program

**Computer Organization II** 

# **CPU Clocking**

Operation of digital hardware elements is governed by a constant-rate clock.



Clock period: duration of a clock cycle

- e.g.,  $250ps = 0.25ns = 250 \times 10^{-12}s$ 

Clock frequency (rate): cycles per second – e.g., 4.0GHz = 4000MHz =  $4.0 \times 10^{9}$ Hz



### **CPU** Time



Performance improved by

CS@VT

- Reducing number of clock cycles
- Increasing clock rate
- Hardware designer must often trade off clock rate against cycle count

## CPU Time Example

Computer A: 2GHz clock, 10s CPU time (to execute a particular program) Designing Computer B

- Aim for 6s CPU time (to execute the same program)
- Can increase clock rate, but that requires a 20% increase in clock cycles

How fast must Computer B clock be?

$$\operatorname{Clock} \operatorname{Rate}_{B} = \frac{\operatorname{Clock} \operatorname{Cycles}_{B}}{\operatorname{CPU} \operatorname{Time}_{B}} = \frac{1.2 \times \operatorname{Clock} \operatorname{Cycles}_{A}}{6s}$$

 $Clock Cycles_{A} = CPU Time_{A} \times Clock Rate_{A}$  $= 10s \times 2GHz = 20 \times 10^{9}$ 

Clock Rate<sub>*B*</sub> = 
$$\frac{1.2 \times 20 \times 10^9}{6s} = \frac{24 \times 10^9}{6s} = 4$$
GHz

CS@VT

**Computer Organization II** 

# How many cycles are required for a program? Performance 8

Could assume that number of cycles equals number of instructions:



This assumption is incorrect,

- different instructions may take different numbers of cycles on same machine
- same instruction may take different number of cycles on different machines

Why? Hint: remember that these are machine instructions, not lines of C code

#### CS@VT

## Now that we understand cycles...

### A given program will require

- some number of instructions (machine instructions)
- some number of cycles
- some number of seconds

We have a vocabulary that relates these quantities:

- cycle time (seconds per cycle)
- clock rate (cycles per second)
- CPI (cycles per instruction)

a floating point intensive application might have a higher CPI

- MIPS (millions of instructions per second)

this would be higher for a program using simple instructions

### Performance

Performance is determined by execution time

Do any of the other variables equal performance? # of cycles to execute program? # of instructions in program? # of cycles per second? average # of cycles per instruction? average # of instructions per second?

Common pitfall: thinking one of these variables (alone) is indicative of performance, when none really is

## **CPI** Example

Computer A: Cycle Time = 250ps, CPI = 2.0Computer B: Cycle Time = 500ps, CPI = 1.2 Same ISA For a given program, which is faster, and by how much?  $CPU Time_A = Instruction Count \times CPI_A \times Cycle Time_A$ A is faster... =I $\times$ 2.0 $\times$ 250ps =I $\times$ 500ps +CPU Time<sub>*B*</sub> = Instruction Count × CPI<sub>*B*</sub> × Cycle Time<sub>*B*</sub> =I $\times$ 1.2 $\times$ 500ps =I $\times$ 600ps  $\underline{\text{CPU Time}_B} = \frac{\text{I} \times 600 \text{ps}}{1.2} = 1.2$ ...by this much CPU Time<sub>A</sub>  $I \times 500 \text{ ps}$ 

#### CS@VT

#### **Computer Organization II**

# CPI in More Detail

If different instruction classes take different numbers of cycles

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

Weighted average CPI

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

# **CPI** Example

Alternative compiled code sequences using instructions in classes A, B, C

Class	А	В	С
CPI for class	1	2	3
IC in sequence 1	2	1	2
IC in sequence 2	4	1	1

Sequence 1: A B C A C

- IC = 5
- Clock Cycles
   = 2×1 + 1×2 + 2×3
   = 10

- Avg. CPI = 10/5 = 2.0

Sequence 2: A B A A C A

- IC = 6
- Clock Cycles
   = 4×1 + 1×2 + 1×3
   = 9
- Avg. CPI = 9/6 = 1.5

A compiler designer is trying to decide between two code sequences for a particular machine. Based on the hardware implementation, there are three different classes of instructions: Class A, Class B, and Class C, and they require one, two, and three cycles (respectively).

The first code sequence has 5 instructions: 2 of A, 1 of B, and 2 of C The second sequence has 6 instructions: 4 of A, 1 of B, and 1 of C.

Which sequence will be faster? How much?

What is the CPI for each sequence?

### **MIPS** example

Two different compilers are being tested for a 4 GHz. machine with three different classes of instructions: Class A, Class B, and Class C, which require one, two, and three cycles (respectively). Both compilers are used to produce code for a large piece of software.

The first compiler's code uses 5 million Class A instructions, 1 million Class B instructions, and 2 million Class C instructions.

The second compiler's code uses 10 million Class A instructions, 1 million Class B instructions, and 1 million Class C instructions.

Which sequence will be faster according to MIPS?

Which sequence will be faster according to execution time?

### Amdahl's Law

Execution Time After Improvement =

Execution Time Unaffected + (Execution Time Affected / Amount of Improvement)

Example:

Suppose a program runs in 100 seconds on a machine, with multiply instructions responsible for 80 seconds of this time.

How much do we have to improve the speed of multiplication if we want the program to run 4 times faster?

How about making it 5 times faster?

Principle: Make the common case fast

## Example 1

Suppose we enhance a machine making all floating-point instructions run five times faster.

If the execution time of some benchmark program before the floating-point enhancement is 10 seconds, what will the speedup be if half of the 10 seconds is spent executing floating-point instructions?

### Example 2

We are looking for a benchmark to show off the new floating-point unit described above, and want the overall benchmark to show a speedup of 3 (i.e., take 1/3 as long to run).

One benchmark we are considering runs for 90 seconds with the old floating-point hardware.

How much of the execution time would floating-point instructions have to account for in this program in order to yield our desired speedup on this benchmark?

### Remember

Performance is specific to a particular program/s

- Total execution time is a consistent summary of performance

For a given architecture performance increases come from:

- increases in clock rate (if that yields no adverse CPI affects)
- improvements in processor organization that lower CPI
- compiler enhancements that lower CPI (different distribution of instructions) and/or instruction count
- algorithm/language choices that affect instruction count and/or instruction distribution

Pitfall: expecting improvement in one aspect of a machine's performance to affect the total performance



Performance depends on

- Algorithm: affects IC, possibly CPI
- Programming language: affects IC, CPI
- Compiler: affects IC, CPI
- Instruction set architecture: affects IC, CPI, T<sub>c</sub>

### Benchmarks

### benchmark

A program or collection of programs selected for use in comparing computer performance

Kinds of benchmarks: Kernels (e.g. matrix multiply) Toy programs (e.g. sorting) Synthetic benchmarks (e.g. Dhrystone) Benchmark suites (e.g. SPEC06fp, TPC-C)

A benchmark can be used to isolate different aspects of hardware performance.

## **SPEC Integer Benchmarks**

### Performance 22

Description	Name	Instruction Count x 10 <sup>9</sup>	CPI	Clock cycle time (seconds x 10 <sup>-9</sup> )	Execution Time (seconds)	Reference Time (seconds)	SPECratio
Interpreted string processing	perl	2252	0.60	0.376	508	9770	19.2
Block-sorting compression	bzip2	2390	0.70	0.376	629	9650	15.4
GNU C compiler	gcc	794	1.20	0.376	358	8050	22.5
Combinatorial optimization	mcf	221	2.66	0.376	221	9120	41.2
Go game (AI)	go	1274	1.10	0.376	527	10490	19.9
Search gene sequence	hmmer	2616	0.60	0.376	590	9330	15.8
Chess game (AI)	sjeng	1948	0.80	0.376	586	12100	20.7
Quantum computer simulation	libquantum	659	0.44	0.376	109	20720	190.0
Video compression	h264avc	3793	0.50	0.376	713	22130	31.0
Discrete event simulation library	omnetpp	367	2.10	0.376	290	6250	21,5
Games/path finding	astar	1250	1.00	0.376	470	7020	14.9
XML parsing	xalancbmk	1045	0.70	0.376	275	6900	25.1
Geometric mean	-	वेत	-	0 <b>.</b>	-		25.7

FIGURE 1.18 SPECINTC2006 benchmarks running on a 2.66 GHz Intel Core i7 920. As the equation on page 35 explains, execution time is the product of the three factors in this table: instruction count in billions, clocks per instruction (CPI), and clock cycle time in nanoseconds.

SPECratio is simply the reference time, which is supplied by SPEC, divided by the measured execution time. The single number quoted as SPECINTC2006 is the geometric mean of the SPECratios.

### Single-processor Evolution



CS@VT

#### **Computer Organization II**

### Why? The Power Wall



Intel VP Patrick Gelsinger (ISSCC 2001):

If scaling continues at present pace, by 2005, high speed processors would have power density of nuclear reactor, by 2010, a rocket nozzle, and by 2015, surface of sun.