

## Performance

- Measure, Report, and Summarize
- Make intelligent choices
- See through the marketing hype
- Key to understanding underlying organizational motivation

*Why is some hardware better than others for different programs?*

*What factors of system performance are hardware related?  
(e.g., Do we need a new machine, or a new operating system?)*

*How does the machine's instruction set affect performance?*

1

## Which of these airplanes has the best performance?



Airplane	Passengers	Range (mi)	Speed (mph)
Boeing 737-100	101	630	598
Boeing 747	470	4150	610
BAC/Sud Concorde	132	4000	1350
Douglas DC-8-50	146	8720	544

•How much faster is the Concorde compared to the 747?

•How much bigger is the 747 than the Douglas DC-8?

2

## Computer Performance: TIME, TIME, TIME

- Response Time (latency)
  - How long does it take for my job to run?
  - How long does it take to execute a job?
  - How long must I wait for the database query?
- Throughput
  - How many jobs can the machine run at once?
  - What is the average execution rate?
  - How much work is getting done?

• *If we upgrade a machine with a new processor what do we increase?*  
*If we add a new machine to the lab what do we increase?*

3

## 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

4

## Book's Definition of Performance

- For some program running on machine X,  

$$\text{Performance}_x = 1 / \text{Execution time}_x$$
- "X is n times faster than Y"  

$$\text{Performance}_x / \text{Performance}_y = n$$
- Problem:
  - machine A runs a program in 20 seconds
  - machine B runs the same program in 25 seconds

5

## Clock Cycles

- Instead of reporting execution time in seconds, we often use cycles

$$\frac{\text{seconds}}{\text{program}} = \frac{\text{cycles}}{\text{program}} \times \frac{\text{seconds}}{\text{cycle}}$$

- Clock "ticks" indicate when to start activities (one abstraction):



- cycle time = time between ticks = seconds per cycle
- clock rate (frequency) = cycles per second (1 Hz. = 1 cycle/sec)

A 200 Mhz. clock has a  $\frac{1}{200 \times 10^6} \times 10^9 = 5$  nanoseconds cycle time

6

## How to Improve Performance

$$\frac{\text{seconds}}{\text{program}} = \frac{\text{cycles}}{\text{program}} \times \frac{\text{seconds}}{\text{cycle}}$$

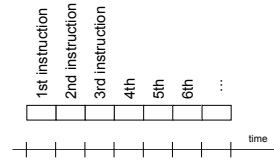
So, to improve performance (everything else being equal) you can either

- \_\_\_\_\_ the # of required cycles for a program, or
- \_\_\_\_\_ the clock cycle time or, said another way,
- \_\_\_\_\_ the clock rate.

7

## How many cycles are required for a program?

- Could assume that # of cycles = # of instructions



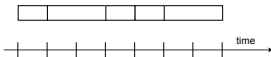
*This assumption is incorrect,*

*different instructions take different amounts of time on different machines.*

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

8

## Different numbers of cycles for different instructions



- Multiplication takes more time than addition
- Floating point operations take longer than integer ones
- Accessing memory takes more time than accessing registers
- *Important point: changing the cycle time often changes the number of cycles required for various instructions (more later)*

9

## Example

- Our favorite program runs in 10 seconds on computer A, which has a 400 Mhz. clock. We are trying to help a computer designer build a new machine B, that will run this program in 6 seconds. The designer can use new (or perhaps more expensive) technology to substantially increase the clock rate, but has informed us that this increase will affect the rest of the CPU design, causing machine B to require 1.2 times as many clock cycles as machine A for the same program. What clock rate should we tell the designer to target?"
- Don't Panic, can easily work this out from basic principles

10

## 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*

11

## 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 the variables is indicative of performance when it really isn't.

12

## CPI Example

- Suppose we have two implementations of the same instruction set architecture (ISA).

For some program,

Machine A has a clock cycle time of 10 ns. and a CPI of 2.0  
Machine B has a clock cycle time of 20 ns. and a CPI of 1.2

What machine is faster for this program, and by how much?

- If two machines have the same ISA which of our quantities (e.g., clock rate, CPI, execution time, # of instructions, MIPS) will always be identical?

13

## # of Instructions Example

- 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?

14

## MIPS example

- Two different compilers are being tested for a 100 MHz. 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 1 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?

15

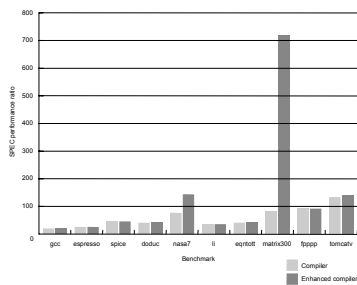
## Benchmarks

- Performance best determined by running a real application
  - Use programs typical of expected workload
  - Or, typical of expected class of applications
    - e.g., compilers/editors, scientific applications, graphics, etc.
- Small benchmarks
  - nice for architects and designers
  - easy to standardize
  - can be abused
- SPEC (System Performance Evaluation Cooperative)
  - companies have agreed on a set of real program and inputs
  - can still be abused (Intel's "other" bug)
  - valuable indicator of performance (and compiler technology)

16

## SPEC '89

- Compiler "enhancements" and performance



17

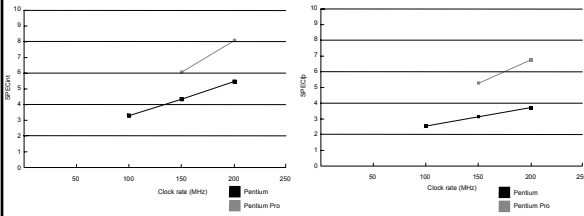
## SPEC '95

Benchmark	Description
gcc	Artificial intelligence: plays the game of Go
int88mm	MicroVax 88k chip simulator: runs test program
gcc	The GNU C compiler: generates SPARC code
compress	Compresses and decompresses file in memory
li	Li: test interpreter
liop	Graphic compression and decompression
perl	Manipulates strings and prime numbers in the special purpose programming language Perl
corlex	A database program
lispobj	A math generation program
swim	Shallow water model with 513 x 513 grid
autot	Quantum physics: Monte Carlo simulation
hydro2d	Aerodynamics: hydrodynamic Navier Stokes equations
mpid	Multi-grid solver in 3-D potential field
spmv	Parabolic/elliptic partial differential equations
fluid	Simulates isotropic, homogeneous turbulence in a cube
spci	Solves problems regarding temperature, wind velocity, and distribution of pollutant
fpop	Quantum chemistry
levent	Plasma physics: electromagnetic particle simulation

18

## SPEC '95

Does doubling the clock rate double the performance?  
Can a machine with a slower clock rate have better performance?



19

## Amdahl's Law

Execution Time After Improvement =

$$\text{Execution Time Unaffected} + (\text{Execution Time Affected} / \text{Amount of Improvement})$$

• **Example:**

"Suppose a program runs in 100 seconds on a machine, with multiply 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

20

## Example

- Suppose we enhance a machine making all floating-point instructions run five times faster. If the execution time of some benchmark before the floating-point enhancement is 100 seconds, what will the speedup be if half of the 100 seconds is spent executing floating-point instructions?
- 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. One benchmark we are considering runs for 100 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?

21

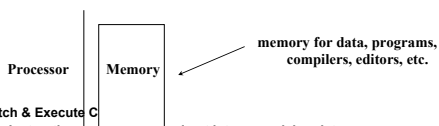
## Remember

- 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 (without adverse CPI affects)
  - improvements in processor organization that lower CPI
  - compiler enhancements that lower CPI and/or instruction count
- Pitfall: expecting improvement in one aspect of a machine's performance to affect the total performance
- You should not always believe everything you read! Read carefully! (see newspaper articles, e.g., Exercise 2.37)

22

## Stored Program Concept

- Instructions are bits
- Programs are stored in memory
  - to be read or written just like data



- Fetch & Execute Cycle
  - Instructions are fetched and put into a special register
  - Bits in the register "control" the subsequent actions
  - Fetch the "next" instruction and continue

23

## Instructions:

- Language of the Machine
- More primitive than higher level languages
  - e.g., no sophisticated control flow
- Very restrictive
  - e.g., MIPS Arithmetic Instructions
- We'll be working with the MIPS instruction set architecture
  - similar to other architectures developed since the 1980's
  - used by NEC, Nintendo, Silicon Graphics, Sony

*Design goals: maximize performance and minimize cost, reduce design time*

24