## 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*

1

## 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.

2

## 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 equivalent to performance?*

3

## # 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?

4

## 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?
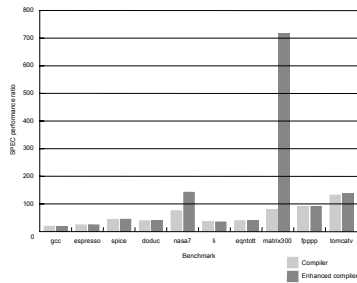
5

## 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)

6

## SPEC '89

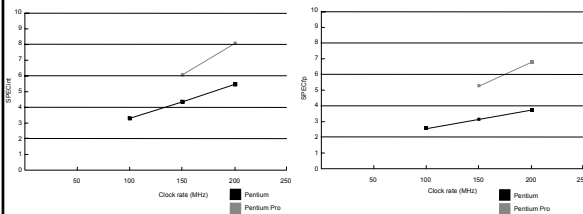- **Compiler "enhancements" and performance**

---

## SPEC '95

| Benchmark | Description |
|---|---|
| go | Artificial intelligence; plays the game of Go |
| m88ksim | Motorola 88k chip simulator; runs test program |
| gcc | The Gnu C compiler generating SPARC code |
| compress | Compresses and decompresses file in memory |
| li | Lisp interpreter |
| ijpeg | Graphic compression and decompression |
| perl | Manipulates strings and prime numbers in the special-purpose programming language Perl |
| vortex | A database program |
| tomcatv | A mesh generation program |
| swim | Shallow water model with 513 x 513 grid |
| su2cor | quantum physics; Monte Carlo simulation |
| hydro2d | Astrophysics; Hydrodynamic Naiver Stokes equations |
| mgrid | Multigrid solver in 3-D potential field |
| applu | Parabolic/elliptic partial differential equations |
| trub3d | Simulates isotropic, homogeneous turbulence in a cube |
| apsi | Solves problems regarding temperature, wind velocity, and distribution of pollutant |
| fpppp | Quantum chemistry |
| wave5 | Plasma physics; electromagnetic particle simulation |

---

## SPEC '95

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

---

## 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 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

- 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 10 seconds, what will the speedup be if half of the 10 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?
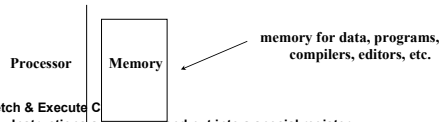
---

## 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)**

## Stored Program Concept

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

Processor    Memory

memory for data, programs, compilers, editors, etc.

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

## 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*

## Architecture Specification

- Data types:
    - bit, byte, bit field, signed/unsigned integers logical, floating point, character
- Operations:
    - data movement, arithmetic, logical, shift/rotate, conversion, input/output, control, and system calls
- # of operands:
    - 3, 2, 1, or 0 operands
- Registers:
    - integer, floating point, control
- Instruction representation as bit strings

## Characteristics of Instruction Set

- Complete
    - Can be used for a variety of application
- Efficient
    - Useful in code generation
- Regular
    - Expected instruction should exist
- Compatible
    - Programs written for previous versions of machines need it
- Primitive
    - Basic operations
- Simple
    - Easy to implement
- Smaller
    - Implementation