# Chapter 1

# Computer Abstractions and Technology

# The Computer Revolution

- Progress in computer technology
  - Underpinned by Moore's Law
- Makes novel applications feasible
  - Computers in automobiles
  - Cell phones
  - Human genome project
  - World Wide Web
  - Search Engines
- Computers are pervasive

# Classes of Computers

- ## Personal computers
  - General purpose, variety of software
  - Subject to cost/performance tradeoff

- ## Server computers
  - Network based
  - High capacity, performance, reliability
  - Range from small servers to building sized

# Classes of Computers

- Supercomputers
  - High-end scientific and engineering calculations
  - Highest capability but represent a small fraction of the overall computer market

- Embedded computers
  - Hidden as components of systems
  - Stringent power/performance/cost constraints

# The PostPC Era

- Personal Mobile Device (PMD)
    - Battery operated
    - Connects to the Internet
    - Hundreds of dollars
    - Smart phones, tablets, electronic glasses
- Cloud computing
    - Warehouse Scale Computers (WSC)
    - Software as a Service (SaaS)
    - Portion of software run on a PMD and a portion run in the Cloud
    - Amazon and Google

# What You Will Learn

- How programs are translated into the machine language
  - And how the hardware executes them
- The hardware/software interface
- What determines program performance
  - And how it can be improved
- How hardware designers improve performance
- What is parallel processing

# Understanding Performance

- Algorithm
  - Determines number of operations executed
- Programming language, compiler, architecture
  - Determine number of machine instructions executed per operation
- Processor and memory system
  - Determine how fast instructions are executed
- I/O system (including OS)
  - Determines how fast I/O operations are executed

# Levels of Program Code

- ## High-level language
  - Level of abstraction closer to problem domain
  - Provides for productivity and portability

- ## Assembly language
  - Textual representation of instructions

- ## Hardware representation
  - Binary digits (bits)
  - Encoded instructions and data

High-level
language
program
(in C)

```
swap(int v[], int k)
{int temp;
    temp = v[k];
    v[k] = v[k+1];
    v[k+1] = temp;
}
```

Compiler

Assembly
language
program
(for MIPS)

```
swap:
    muli $2, $5,4
    add  $2, $4,$2
    lw   $15, 0($2)
    lw   $16, 4($2)
    sw   $16, 0($2)
    sw   $15, 4($2)
    jr   $31
```
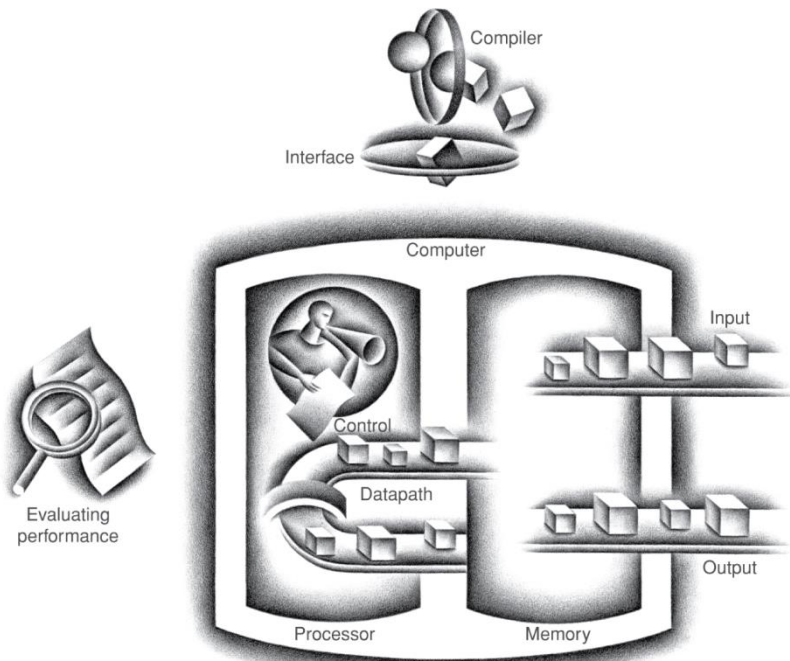
Assembler

Binary machine
language
program
(for MIPS)

```
00000000101000010000000000011000
00000000000110000001100000100001
10001100011000100000000000000000
10001100111100100000000000000100
10101100111100100000000000000000
10101100011000100000000000000100
00000011111000000000000000001000
```

# Components of a Computer

**The BIG Picture**



- Same components for all kinds of computer
  - Desktop, server, embedded
- Input/output includes
  - User-interface devices
    - Display, keyboard, mouse
  - Storage devices
    - Hard disk, CD/DVD, flash
  - Network adapters
    - For communicating with other computers

# Inside the Processor (CPU)

- Datapath: performs operations on data

- Control: sequences datapath, memory, ...

- Cache memory
    - Small fast SRAM memory for immediate access to data

# Response Time and Throughput

- Response time

  *decreasing response time improves throughput. Hence in case 1, both response time and throughput improve. In case 2, no one task gets done faster, so only throughput increases.*

  *If however the demand for processing in the second case was almost as large as throughput, the system might force requests to queue up. In this case, increasing throughput could improve response time, since it would decrease the waiting time in queue.*

  - How long it takes to do a task

- Throughput

  - Total work done per unit time

    - e.g., tasks/transactions/… per hour

- How are response time and throughput affected by

  - Replacing the processor with a faster version?

  - Adding more processors?

- We'll focus on response time for now…

# Relative Performance

- Define Performance = 1/Execution Time

- "X is $n$ time faster than Y"

$$\text{Performance}_X / \text{Performance}_Y$$
$$= \text{Execution time}_Y / \text{Execution time}_X = n$$

- Example: time taken to run a program

  - 10s on A, 15s on B

  - Execution Time$_B$ / Execution Time$_A$
    = 15s / 10s = 1.5

  - So A is 1.5 times faster than B

# Measuring Execution Time

- ## Elapsed time
  - Total response time, including all aspects
    - Processing, I/O, OS overhead, idle time
  - Determines system performance
- ## CPU time
  - Time spent processing a given job
    - Discounts I/O time, other jobs' shares
  - Comprises user CPU time and system CPU time
  - Different programs are affected differently by CPU and system performance

# CPU Clocking

- Operation of digital hardware 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

$$\text{CPU Time} = \text{CPU Clock Cycles} \times \text{Clock Cycle Time}$$

$$= \frac{\text{CPU Clock Cycles}}{\text{Clock Rate}}$$

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

# CPU Time Example

$$10 = \frac{X_A}{2 \times 10^9} \Rightarrow X_A = 2 \times 10^{10} \text{ clock cycles}$$

$$\Rightarrow X_B = 2.4 \times 10^{10} \text{ clock cycles}$$

- Computer A: 2GHz clock, 10s CPU time

$$\therefore 6 = \frac{2.4 \times 10^{10}}{Y_B}$$

- Designing Computer B

$$\Rightarrow Y_B = 0.4 \times 10^{10} \text{ Hz}$$
$$= 4GHz$$

  - Aim for 6s CPU time
  - Can do faster clock, but causes 1.2 × clock cycles
- How fast must Computer B clock be?

$$\text{Clock Rate}_B = \frac{\text{Clock Cycles}_B}{\text{CPU Time}_B} = \frac{1.2 \times \text{Clock Cycles}_A}{6s}$$

$$\text{Clock Cycles}_A = \text{CPU Time}_A \times \text{Clock Rate}_A$$

$$= 10s \times 2GHz = 20 \times 10^9$$

$$\text{Clock Rate}_B = \frac{1.2 \times 20 \times 10^9}{6s} = \frac{24 \times 10^9}{6s} = 4GHz$$

# Instruction Count and CPI

$$\text{Clock Cycles} = \text{Instruction } \boxed{\text{Count} \times \text{Cycles}} \text{ per Instruction}$$

$$\text{CPU Time} = \text{Instruction Count} \times \text{CPI} \times \text{Clock Cycle Time}$$

$$= \frac{\text{Instruction Count} \times \text{CPI}}{\text{Clock Rate}}$$

*→ Instruction Set Architecture*

- Instruction Count for a program
  - Determined by program, ISA and compiler
- Average cycles per instruction
  - Determined by CPU hardware
  - If different instructions have different CPI
    - Average CPI affected by instruction mix

# CPI Example

- Computer A: Cycle Time = 250ps, CPI = 2.0
- Computer B: Cycle Time = 500ps, CPI = 1.2
- Same ISA
- Which is faster, and by how much?

$$\text{CPU Time}_A = \text{Instruction Count} \times \text{CPI}_A \times \text{Cycle Time}_A$$

$$= I \times 2.0 \times 250ps = I \times 500ps \quad \leftarrow \boxed{\text{A is faster...}}$$

$$\text{CPU Time}_B = \text{Instruction Count} \times \text{CPI}_B \times \text{Cycle Time}_B$$

$$= I \times 1.2 \times 500ps = I \times 600ps$$

$$\frac{\text{CPU Time}_B}{\text{CPU Time}_A} = \frac{I \times 600ps}{I \times 500ps} = 1.2 \quad \leftarrow \boxed{\text{...by this much}}$$

# CPI in More Detail

- If different instruction classes take different numbers of cycles

$$\text{Clock Cycles} = \sum_{i=1}^{n} (\text{CPI}_i \times \text{Instruction Count}_i)$$

  - Weighted average CPI

$$\text{CPI} = \frac{\text{Clock Cycles}}{\text{Instruction Count}} = \sum_{i=1}^{n} \left( \text{CPI}_i \times \frac{\text{Instruction Count}_i}{\text{Instruction Count}} \right)$$

Relative frequency

# CPI Example

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

*Total Instruction Count*

*5*
*6*

| Class | A | B | C |
|---|---|---|---|
| CPI for class | 1 | 2 | 3 |
| IC in sequence 1 | 2 | 1 | 2 |
| IC in sequence 2 | 4 | 1 | 1 |

- Sequence 1: IC = 5
  - Clock Cycles $= \sum_{i=1}^{n}\left(CPI_i \times C_i\right)$
    $= 2\times1 + 1\times2 + 2\times3$
    $= 10$
  - Avg. CPI = 10/5 = 2.0

- Sequence 2: IC = 6
  - Clock Cycles
    $= 4\times1 + 1\times2 + 1\times3$
    $= 9$
  - Avg. CPI = 9/6 = 1.5

# Performance Summary

*Time = Seconds / Program*

## The BIG Picture

$$\text{CPU Time} = \frac{\text{Instructions}}{\text{Program}} \times \frac{\text{Clock cycles}}{\text{Instruction}} \times \frac{\text{Seconds}}{\text{Clock cycle}}$$

*Average CPI*

*Clock Cycle Time*

- Performance depends on
    - Algorithm: affects IC, possibly CPI
    - Programming language: affects IC, CPI
    - Compiler: affects IC, CPI
    - Instruction set architecture: affects IC, CPI, $T_c$

# Pitfall: Amdahl's Law

- Improving an aspect of a computer and expecting a proportional improvement in overall performance

$$T_{improved} = \frac{T_{affected}}{improvement\ factor} + T_{unaffected}$$

*Suppose a program runs in 100s on a computer, with multiply operations responsible for 80s/100s.*

- Example: multiply accounts for 80s/100s
  - How much improvement in multiply performance to get 5× overall? *to make the program run 5 times faster*

$$20 = \frac{80}{n} + 20$$

  - Can't be done!

- Corollary: make the common case fast

# Pitfall: MIPS as a Performance Metric

- ## MIPS: Millions of Instructions Per Second
  - ### Doesn't account for *Instruction Set Architecture*
    - Differences in ISAs between computers
    - Differences in complexity between instructions

$$MIPS = \frac{Instruction\ count}{Execution\ time \times 10^6}$$

$$= \frac{Instruction\ count}{\dfrac{Instruction\ count \times CPI}{Clock\ rate} \times 10^6} = \frac{Clock\ rate}{CPI \times 10^6}$$

  - ### CPI varies between programs on a given CPU

# Concluding Remarks

- Cost/performance is improving
  - Due to underlying technology development
- Hierarchical layers of abstraction
  - In both hardware and software
- Instruction set architecture
  - The hardware/software interface
- Execution time: the best performance measure
- Power is a limiting factor
  - Use parallelism to improve performance

Created with
Notewise

# Sections to Read from the Book

- 5<sup>th</sup> Edition sections to read
  - 1.1
  - 1.3
  - 1.4
  - 1.6
  - 1.10
  - 1.11

Created with **Notewise**