# CS/SE 3340
# Computer Architecture



Application

Operating System

Compiler | Firmware

Instr. Set Proc. | I/O system

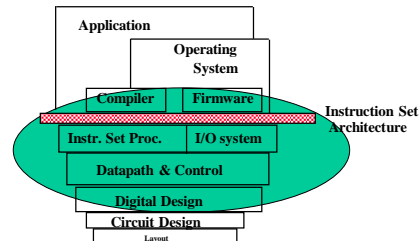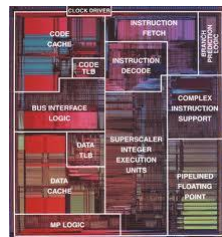Datapath & Control

Digital Design

Circuit Design

Layout
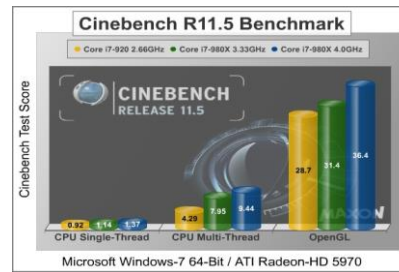
Instruction Set Architecture
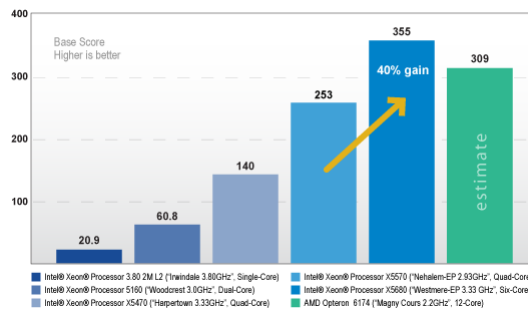
## Introduction to Computer Performance

*Adapted from "Computer Organization and Design, 4th Ed." by D. Patterson and J. Hennessy*

---

# Questions

- How to define performance?
- What is response time and what is throughput?
- What are performance metrics for computer systems?
- How to measure CPU time?
- How to reduce power consumption of a processor?
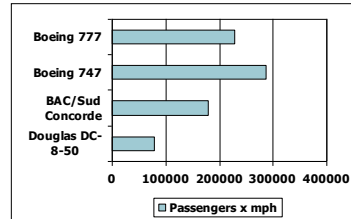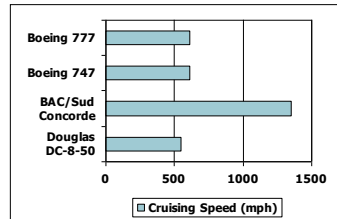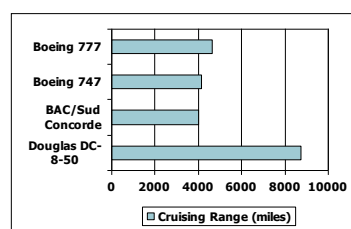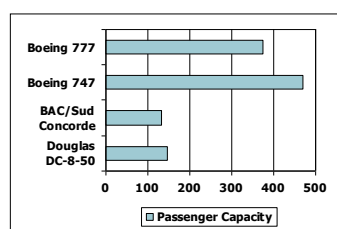- What is Amdahl's law?

2

1

3



# Defining Performance



- *Which airplane has the best performance?*

4

2

# Response Time and Throughput

- *Response time*
  - 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…

5

# 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

  *How to measure 'Execution Time'?*

6

# 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 by the CPU
    - 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

7

# CPU Clocking

- Operation of digital hardware governed by a constant-rate clock



- Clock period: duration of a clock cycle
  - e.g., 250ps = 0.25ns = $250 \times 10^{-12}$ sec
- Clock frequency (rate): cycles per second
  - e.g., 4.0GHz = 4000MHz = $4.0 \times 10^{9}$Hz

8

4

# CPU Time

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

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

9

# CPU Time Example

- Computer A: 2GHz clock, 10s CPU time
- Designing Computer B
  - Aim for 6s CPU time
  - Can do faster clock, but causes 1.2 × clock cycles
- How fast must Computer B clock be?

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

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

$$= 10s \times 2\text{GHz} = 20 \times 10^9$$

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

10

5

# Instruction Count and CPI

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

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

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

- Instruction 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

11

# 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?

$$\text{CPU Time}_A = \text{Instruction Count} \times \text{CPI}_A \times \text{Cycle Time}_A$$
$$= I \times 2.0 \times 250\text{ps} = I \times 500\text{ps} \quad \boxed{\text{A is faster...}}$$
$$\text{CPU Time}_B = \text{Instruction Count} \times \text{CPI}_B \times \text{Cycle Time}_B$$
$$= I \times 1.2 \times 500\text{ps} = I \times 600\text{ps}$$
$$\frac{\text{CPU Time}_B}{\text{CPU Time}_A} = \frac{I \times 600\text{ps}}{I \times 500\text{ps}} = 1.2 \quad \boxed{\text{... by this much}}$$

12

6

# 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

13

# CPI Example

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

| 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
    = 2×1 + 1×2 + 2×3
    = 10
  - Avg. CPI = 10/5 = 2.0

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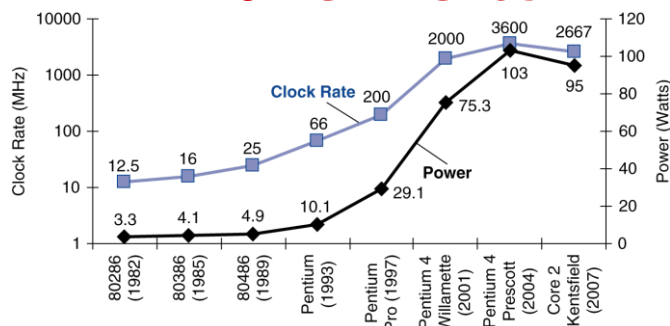    = 9
  - Avg. CPI = 9/6 = 1.5

14

# Performance Summary

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

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

15

---

# Power Trends



- In CMOS IC technology

$$Power = Capacitive\,load \times Voltage^2 \times Frequency$$

×30    5V → 1V    ×1000

16

8

# Reducing Power
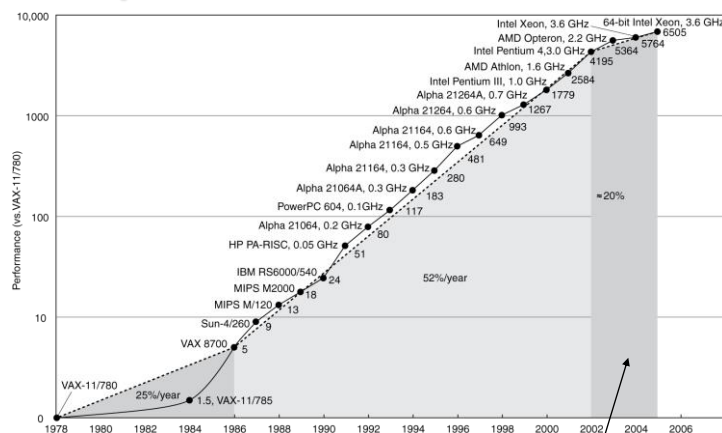
- Suppose a new CPU has
  - 85% of capacitive load of old CPU
  - 15% voltage and 15% frequency reduction

$$\frac{P_{new}}{P_{old}} = \frac{C_{old} \times 0.85 \times (V_{old} \times 0.85)^2 \times F_{old} \times 0.85}{C_{old} \times V_{old}^2 \times F_{old}} = 0.85^4 = 0.52$$

- The power wall
  - We can't reduce voltage further
  - We can't remove more heat
- How else can we improve performance?

17

# Uniprocessor Performance



Constrained by power, instruction-level parallelism, memory latency

18

# Multiprocessors

- Multicore microprocessors
  - More than one processor per chip
- Requires explicitly parallel programming
  - Compare with instruction level parallelism
    - Hardware executes multiple instructions at once
    - Hidden from the programmer
  - Hard to do
    - Programming for performance
    - Load balancing
    - Optimizing communication and synchronization

19

# Pitfall: MIPS as a Performance Metric

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

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

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

  - *CPI varies between programs on a given CPU*

20

# Pitfall: Amdahl's Law

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

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

- Example: multiply accounts for 80s/100s
    - How much improvement in multiply performance to get 5× overall?

$$20 = \frac{80}{n} + 20 \qquad \text{■ Can't be done!}$$

- Corollary: make the common case fast

21

# Summary Remarks

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

22