Performance



Which computer has a better performance?

- Time
- Number of tasks
- Power

- ...



Defining Performance

Which airplane has the best performance?



Chapter 1 — Computer Abstractions and Technology — 2

Response Time and Throughput

- Response time (PC user)
 - How long it takes to do a task
- Throughput (Datacenter manager)
 - 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...



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



Relative Performance

- Define Performance = 1/Execution Time
- "X is n time faster than Y"

 $Performance_{x}/Performance_{y}$

= Execution time $_{\rm Y}$ /Execution time $_{\rm X}$ = 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 (wall clock time, response 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 = 4.0×10⁹Hz

CPU Time

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

Performance improved by

A program takes 2500 clock cycles

to run on a computer with 2.5 GHz

- Reducing number of clock cycles processor. What is CPU time?
- Increasing clock rate
- Hardware designer must often trade off clock rate against cycle count



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?

$$Clock Rate_{B} = \frac{Clock Cycles_{B}}{CPU Time_{B}} = \frac{1.2 \times Clock Cycles_{A}}{6s}$$

$$Clock Cycles_{A} = CPU Time_{A} \times Clock Rate_{A}$$

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

$$Clock Rate_{B} = \frac{1.2 \times 20 \times 10^{9}}{6s} = \frac{24 \times 10^{9}}{6s} = 4GHz$$

Instruction Count and CPI

 $\begin{aligned} \text{Clock Cycles} &= \text{Instruction Count} \times \text{Cycles per Instruction} \\ \text{CPU Time} &= \text{Instruction Count} \times \text{CPI} \times \text{Clock Cycle Time} \\ &= \frac{\text{Instruction Count} \times \text{CPI}}{\text{Clock Rate}} \end{aligned}$

- 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?



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

Which code sequence executes the most instructions?

Which will be faster?

What is the CPI for each sequence?



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: 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



Performance depends on

- Algorithm: affects IC, possibly CPI (float)
- Programming language: affects IC, CPI
- Compiler: affects IC, CPI
- Instruction set architecture: affects IC, CPI, T_c

Power Trends



Power = Capacitive load \times Voltage² \times Frequency \checkmark \land \land

Chapter 1 — Computer Abstractions and Technology — 16

Reducing Power

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

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

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

Uniprocessor Performance



Chapter 1 — Computer Abstractions and Technology — 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

Manufacturing ICs



Yield: proportion of working dies per wafer



AMD Opteron X2 Wafer



- X2: 300mm wafer, 117 chips, 90nm technology
- X4: 45nm technology

Integrated Circuit Cost

```
Cost per die = \frac{\text{Cost per wafer}}{\text{Dies per wafer } \times \text{Yield}}

Dies per wafer \approx Wafer area/Die area

\text{Yield} = \frac{1}{(1+(\text{Defects per area} \times \text{Die area}/2))^2}
```

Nonlinear relation to area and defect rate

- Wafer cost and area are fixed
- Defect rate determined by manufacturing process
- Die area determined by architecture and circuit design

