# Defining Performance

§1.4 Performance

- Which airplane has the best performance?

# 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…

# 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

**Chapter 1 — Computer Abstractions and Technology — 29**
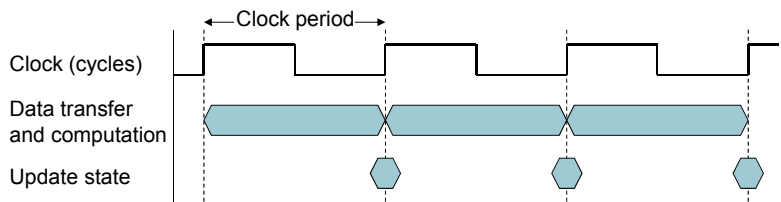
# Measuring Execution Time

- 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

**Chapter 1 — Computer Abstractions and Technology — 30**

# 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

**Chapter 1 — Computer Abstractions and Technology — 31**

# 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

**Chapter 1 — Computer Abstractions and Technology — 32**

# 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?

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

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

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

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

**Chapter 1 — Computer Abstractions and Technology — 33**

# 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

**Chapter 1 — Computer Abstractions and Technology — 34**

# 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?

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

**Chapter 1 — Computer Abstractions and Technology — 35**

# 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

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

$$\boxed{\text{Relative frequency}}$$

**Chapter 1 — Computer Abstractions and Technology — 36**

# 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

**Chapter 1 — Computer Abstractions and Technology — 37**

# Performance Summary

**The BIG Picture**

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

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

**Chapter 1 — Computer Abstractions and Technology — 38**

# Reporting Performance

- Assume 3programs and 3 systems
-       P1      P2      P3

| | P1 | P2 | P3 |
|---|---|---|---|
| A | 10 | 8 | 25 |
| B | 12 | 9 | 20 |
| C | 8 | 8 | 30 |

- Arithmetic mean
- Geometric mean

**Chapter 1 — Computer Abstractions and Technology — 39**
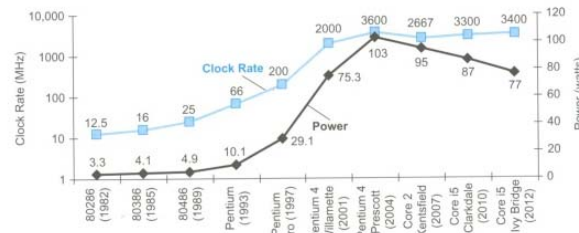
# Reporting Performance

- 2 Programs and 3 machines

| | A | B | C |
|---|---|---|---|
| P1 | 1 | 10 | 20 |
| P2 | 1000 | 100 | 20 |

**Chapter 1 — Computer Abstractions and Technology — 40**

# Power Trends

- In CMOS IC technology

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

Fanout & technology

×30

5V → 1V

×1000

**Chapter 1 — Computer Abstractions and Technology — 41**

# Reducing Power

- Suppose a new CPU has
  - 85% of capacitive load of old CPU
  - 15% voltage and 15% frequency reduction
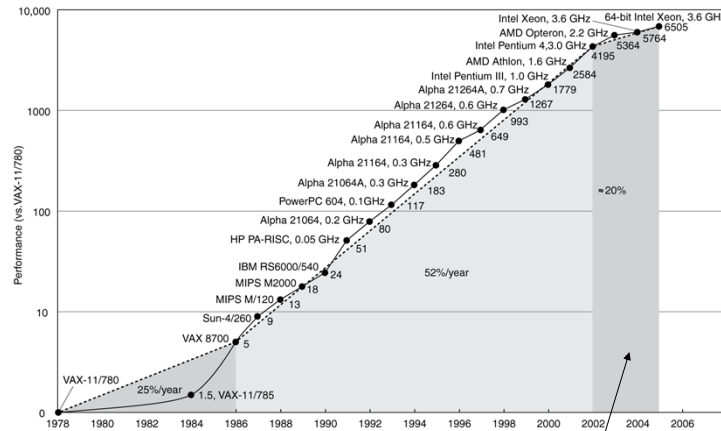
$$\frac{P_{new}}{P_{old}} = \frac{C_{old} \times 0.85 \times (V_{old} \times 0.85)^2 \times F_{old} \times 0.85}{C_{old} \times V_{old}^2 \times F_{old}} = 0.85^4 = 0.52$$

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

**Chapter 1 — Computer Abstractions and Technology — 42**

# Uniprocessor Performance

§1.6 The Sea Change: The Switch to Multiprocessors

Constrained by power, instruction-level parallelism, memory latency

Chapter 1 — Computer Abstractions and Technology — 43

# Multiprocessors

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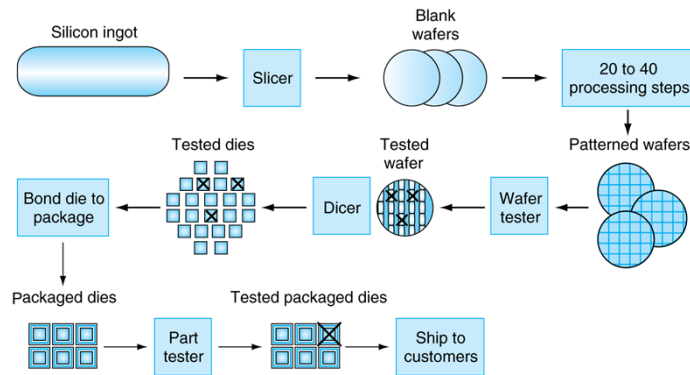  - More than one processor per chip
- Requires explicitly parallel programming
  - Compare with instruction level parallelism
    - Hardware executes multiple instructions at once
    - Hidden from the programmer
  - Hard to do
    - Programming for performance
    - Load balancing
    - Optimizing communication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

Chapter 1 — Computer Abstractions and Technology — 44
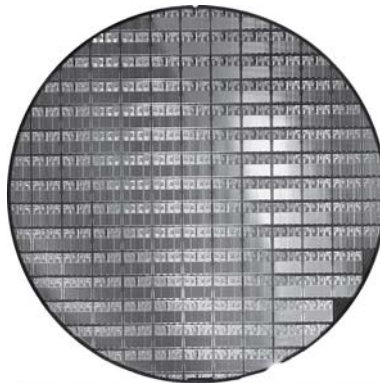
# Manufacturing ICs

- Yield: proportion of working dies per wafer

**Chapter 1 — Computer Abstractions and Technology — 45**

# AMD Opteron X2 Wafer



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

**Chapter 1 — Computer Abstractions and Technology — 46**

# Integrated Circuit Cost

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

$$\text{Dies per wafer} \approx \text{Wafer area/Die area}$$

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

- 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

**Chapter 1 — Computer Abstractions and Technology — 47**

# SPEC CPU Benchmark

- Programs used to measure performance
  - Supposedly typical of actual workload
- Standard Performance Evaluation Corp (SPEC)
  - Develops benchmarks for CPU, I/O, Web, …

- SPEC CPU2006
  - Elapsed time to execute a selection of programs
    - Negligible I/O, so focuses on CPU performance
  - Normalize relative to reference machine
  - Summarize as geometric mean of performance ratios
    - CINT2006 (integer) and CFP2006 (floating-point)

$$\sqrt[n]{\prod_{i=1}^{n} \text{Execution time ratio}_i}$$

**Chapter 1 — Computer Abstractions and Technology — 48**

# CINT2006 for Opteron X4 2356

| Name | Description | IC×10⁹ | CPI | Tc (ns) | Exec time | Ref time | SPECratio |
|---|---|---|---|---|---|---|---|
| perl | Interpreted string processing | 2,118 | 0.75 | 0.40 | 637 | 9,777 | 15.3 |
| bzip2 | Block-sorting compression | 2,389 | 0.85 | 0.40 | 817 | 9,650 | 11.8 |
| gcc | GNU C Compiler | 1,050 | 1.72 | 0.47 | 24 | 8,050 | 11.1 |
| mcf | Combinatorial optimization | 336 | 10.00 | 0.40 | 1,345 | 9,120 | 6.8 |
| go | Go game (AI) | 1,658 | 1.09 | 0.40 | 721 | 10,490 | 14.6 |
| hmmer | Search gene sequence | 2,783 | 0.80 | 0.40 | 890 | 9,330 | 10.5 |
| sjeng | Chess game (AI) | 2,176 | 0.96 | 0.48 | 37 | 12,100 | 14.5 |
| libquantum | Quantum computer simulation | 1,623 | 1.61 | 0.40 | 1,047 | 20,720 | 19.8 |
| h264avc | Video compression | 3,102 | 0.80 | 0.40 | 993 | 22,130 | 22.3 |
| omnetpp | Discrete event simulation | 587 | 2.94 | 0.40 | 690 | 6,250 | 9.1 |
| astar | Games/path finding | 1,082 | 1.79 | 0.40 | 773 | 7,020 | 9.1 |
| xalancbmk | XML parsing | 1,058 | 2.70 | 0.40 | 1,143 | 6,900 | 6.0 |
| Geometric mean | | | | | | | 11.7 |

High cache miss rates

Chapter 1 — Computer Abstractions and Technology — 49

# CINT2006 for 2.66 GHz i7 920

| Name | Description | IC×10⁹ | CPI | Tc (ns) | Exec time | Ref time | SPECratio |
|---|---|---|---|---|---|---|---|
| perl | Interpreted string processing | 2,252 | 0.60 | 0.376 | 508 | 9,770 | 19.2 |
| bzip2 | Block-sorting compression | 2,390 | 0.70 | 0.376 | 629 | 9,650 | 15.4 |
| gcc | GNU C Compiler | 794 | 1.20 | 0.376 | 358 | 8,050 | 22.5 |
| mcf | Combinatorial optimization | 221 | 2.66 | 0.376 | 221 | 9,120 | 41.2 |
| go | Go game (AI) | 1,274 | 1.10 | 0.376 | 527 | 10,490 | 19.9 |
| Hmmer | Search gene sequence | 2,616 | 0.60 | 0.376 | 590 | 9,330 | 15.8 |
| sjeng | Chess game (AI) | 1,948 | 0.80 | 0.376 | 586 | 12,100 | 20.7 |
| libquantum | Quantum computer simulation | 659 | 0.44 | 0.376 | 109 | 20,720 | 190.0 |
| h264avc | Video compression | 3,793 | 0.50 | 0.376 | 713 | 22,130 | 31.0 |
| omnetpp | Discrete event simulation | 367 | 2.10 | 0.376 | 290 | 6,250 | 21.5 |
| astar | Games/path finding | 1,250 | 1.00 | 0.376 | 470 | 7,020 | 14.9 |
| xalancbmk | XML parsing | 1,045 | 0.70 | 0.376 | 275 | 6,900 | 25.1 |
| Geometric mean | | | | | | | 25.7 |

Chapter 1 — Computer Abstractions and Technology — 50

# SPEC Power Benchmark

- Power consumption of server at different workload levels
  - Performance: ssj_ops/sec
  - Power: Watts (Joules/sec)

$$\text{Overall ssj\_ops per Watt} = \left( \sum_{i=0}^{10} \text{ssj\_ops}_i \right) \Big/ \left( \sum_{i=0}^{10} \text{power}_i \right)$$

Chapter 1 — Computer Abstractions and Technology — 51

# SPECpower_ssj2008 for X4

| Target Load % | Performance (ssj_ops/sec) | Average Power (Watts) |
|---|---|---|
| 100% | 231,867 | 295 |
| 90% | 211,282 | 286 |
| 80% | 185,803 | 275 |
| 70% | 163,427 | 265 |
| 60% | 140,160 | 256 |
| 50% | 118,324 | 246 |
| 40% | 920,35 | 233 |
| 30% | 70,500 | 222 |
| 20% | 47,126 | 206 |
| 10% | 23,066 | 180 |
| 0% | 0 | 141 |
| Overall sum | 1,283,590 | 2,605 |
| ∑ssj_ops/ ∑power | | 493 |

Chapter 1 — Computer Abstractions and Technology — 52

# Pitfall: Amdahl's Law

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

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

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

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

- Corollary: make the common case fast

**Chapter 1 — Computer Abstractions and Technology — 53**

---

# Fallacy: Low Power at Idle

- Look back at X4 power benchmark
  - At 100% load: 295W
  - At 50% load: 246W (83%)
  - At 10% load: 180W (61%)
- Google data center
  - Mostly operates at 10% – 50% load
  - At 100% load less than 1% of the time
- Consider designing processors to make power proportional to load

**Chapter 1 — Computer Abstractions and Technology — 54**

## 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

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

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

  - CPI varies between programs on a given CPU

Chapter 1 — Computer Abstractions and Technology — 55

## Concluding Remarks

§1.9 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

Chapter 1 — Computer Abstractions and Technology — 56