# **Defining Performance**

### Which airplane has the best performance?





### Chapter 1 — Computer Abstractions and Technology — 26

# **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"

Performance<sub>x</sub>/Performance<sub>y</sub> = Execution time<sub>y</sub>/Execution time<sub>x</sub> = n

- Example: time taken to run a program
  - 10s on A, 15s on B
  - Execution Time<sub>B</sub> / Execution Time<sub>A</sub> = 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 = 250×10<sup>-12</sup>s
- Clock frequency (rate): cycles per second
  - e.g., 4.0GHz = 4000MHz = 4.0×10<sup>9</sup>Hz





CPU Time = CPU Clock Cycles × Clock Cycle Time

CPU Clock Cycles 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**

- 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

Clock Cycles = Instruction Count × Cycles per Instruction

CPU Time = Instruction Count × CPI × Clock Cycle Time

Instruction Count × CPI

**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



# **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?

 $\begin{array}{l} \mathsf{CPU Time}_{\mathsf{A}} = \mathsf{Instruction Count} \times \mathsf{CPI}_{\mathsf{A}} \times \mathsf{Cycle Time}_{\mathsf{A}} \\ = \mathsf{I} \times 2.0 \times 250 \mathsf{ps} = \mathsf{I} \times 500 \mathsf{ps} & & \mathsf{A is faster...} \\ \mathsf{CPU Time}_{\mathsf{B}} = \mathsf{Instruction Count} \times \mathsf{CPI}_{\mathsf{B}} \times \mathsf{Cycle Time}_{\mathsf{B}} \\ = \mathsf{I} \times 1.2 \times 500 \mathsf{ps} = \mathsf{I} \times 600 \mathsf{ps} \\ \hline \mathsf{CPU Time}_{\mathsf{A}} = \frac{\mathsf{I} \times 600 \mathsf{ps}}{\mathsf{I} \times 500 \mathsf{ps}} = 1.2 & & \mathsf{...by this much} \end{array}$ 



# **CPI in More Detail**

If different instruction classes take different numbers of cycles

Clock Cycles = 
$$\sum_{i=1}^{n} (CPI_i \times Instruction Count_i)$$

$$CPI = \frac{Clock Cycles}{Instruction Count} = \sum_{i=1}^{n} \left( CPI_i \times \frac{Instruction Count_i}{Instruction Count} \right)$$
Relative frequency



# **CPI Example**

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

Class	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 SummaryThe BIG Picture $CPU Time = \frac{Instructions}{Program} \times \frac{Clock cycles}{Instruction} \times \frac{Seconds}{Clock cycle}$

### Performance depends on

- Algorithm: affects IC, possibly CPI
- Programming language: affects IC, CPI
- Compiler: affects IC, CPI
- Instruction set architecture: affects IC, CPI, T<sub>c</sub>



# **Power Trends**





Chapter 1 — Computer Abstractions and Technology — 38

# **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?



# **Uniprocessor Performance**





# **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



# **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)





# SPECspeed 2017 Integer benchmarks on a 1.8 GHz Intel Xeon E5-2650L

								ć
Description	Name	Instruction Count x 10^9	CPI	Clock cycle time (seconds x 10^–9)	Execution Time (seconds)	Reference Time (seconds)	SPECratio	
Perl interpreter	perlbench	2684	0.42	0.556	627	1774	2.83	
GNU C compiler	gcc	2322	0.67	0.556	863	3976	4.61	
Route planning	mcf	1786	1.22	0.556	1215	4721	3.89	
Discrete Event simulation - computer network	omnetpp	1107	0.82	0.556	507	1630	3.21	
XML to HTML conversion via XSLT	xalancbmk	1314	0.75	0.556	549	1417	2.58	
Video compression	x264	4488	0.32	0.556	813	1763	2.17	
Artificial Intelligence: alpha-beta tree search (Chess)	deepsjeng	2216	0.57	0.556	698	1432	2.05	
Artificial Intelligence: Monte Carlo tree search (Go)	leela	2236	0.79	0.556	987	1703	1.73	
Artificial Intelligence: recursive solution generator (Sudoku)	exchange2	6683	0.46	0.556	1718	2939	1.71	
General data compression	xz	8533	1.32	0.556	6290	6182	0.98	
Geometric mean							2.36	



# **SPEC Power Benchmark**

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

$$Overall \, ssj\_ops \, per \, Watt = \left(\sum_{i=0}^{10} ssj\_ops_i\right) / \left(\sum_{i=0}^{10} power_i\right)$$



### SPECpower\_ssj2008 for Xeon E5-2650L

Target Load %	Performance (ssj_ops)	Average Power (watts)
100%	4,864,136	347
90%	4,389,196	312
80%	3,905,724	278
70%	3,418,737	241
60%	2,925,811	212
50%	2,439,017	183
40%	1,951,394	160
30%	1,461,411	141
20%	974,045	128
10%	485,973	115
0%	0	48
Overall Sum	26,815,444	2,165
Σssi ops / Σpov	ver =	12,385



# Pitfall: Amdahl's Law

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



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

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

Corollary: make the common case fast

# Fallacy: Low Power at Idle

- Look back at i7 power benchmark
  - At 100% load: 258W
  - At 50% load: 170W (66%)
  - At 10% load: 121W (47%)
- 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



## **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



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

