



Computer Architecture

Chapter 1: Computer Abstractions and Technology

Assoc. Prof. Dinh Duc Anh Vu




Adapted from Computer Organization the Hardware/Software Interface – 5th



The industrial revolutions


1784



1st Industrial Revolution

- Mechanization, steam power, weaving looms
- Large-scale transportation with steam-powered vessels and railroads
- Replacing human and animal power with machines


1870



2nd Industrial Revolution

- Electricity, assembly line, mass production
- Internal combustion engines, automobiles
- Radio and television


1969



3rd Industrial Revolution

- Electronics, computers
- Automation
- Information technology

Now



4th Industrial Revolution

- Cyber-physical systems
- Internet of Things
- Sensor networks
- Advanced Robotics
- Big Data
- Machine Learning
- Cloud Computing
- Driverless cars
- 3D/4D printing-based manufacturing
- Blockchain transaction architecture

Chapter 1: Computer Technology

2

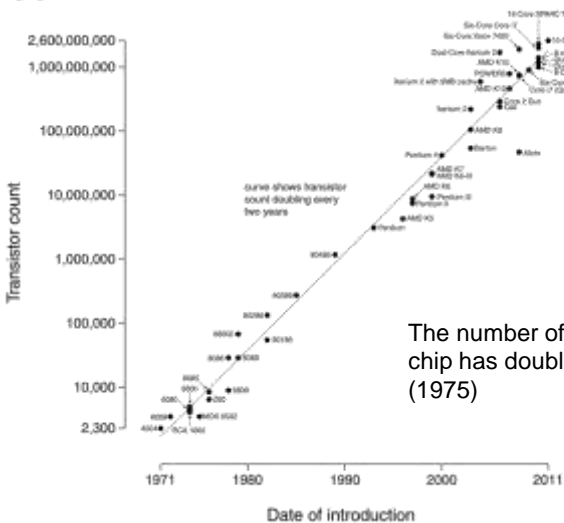


The Computer Revolution

- The third revolution along with agriculture and industry
- 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



The Moore's Law



Co-founder of Intel Corp.

The number of transistors integrated in a chip has doubled every 18-24 months (1975)



Intel Processors & Chips

- World record, in terms of the number of transistors integrated into a chip:
 - Altera FPGA device: 30+ Billions
- Intel processor
 - Core i 8th generation [Coffee Lake](#)
 - 14 nm technology
 - >1.4B transistors (6th generation – SkyLake)



Chapter 1: Computer Technology



The First “Computer”



Source: Internet

Chapter 1: Computer Technology



The First “Computer” (cont.)



The ENIAC Computer, source: US Army photo



The ENIAC Computer

- 30+ tons
- 1,500+ square feet (140 square meter)
- 18,000+ vacuum tubes
- 140+ KW power
- 5,000+ additions per second



A Brief History of Computers

- The first generation
 - Vacuum tubes
 - 1946 – 1955
- The second generation
 - Transistors
 - 1955 – 1965
- The third generation
 - 1965 – 1980
 - Integrated circuits
- The current generation
 - 1980 - ...
 - Personal computers
- *What's the next?*
 - *Quantum computers?*
 - *Memristor?*



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

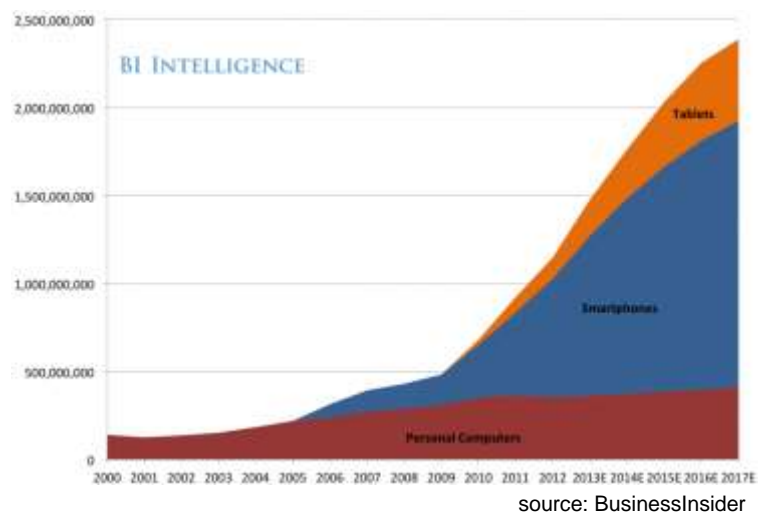


Classes of Computers

- 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





The PostPC Era

- Personal Mobile Device (PMD)
 - Battery operated
 - Connects to the Internet
 - Hundreds of dollars
 - Smart phones, tablets, electronic glasses,...
- Clouding 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



Understanding Performance

- Algorithm
 - Determines both the number of source-level statements and number of I/O operations executed
- Programming language, compiler, architecture
 - Determine number of machine instructions executed for each source-level statement
- Processor and memory system
 - Determine how fast instructions can be executed
- I/O system (including OS)
 - Determines how fast I/O operations may be executed



Quiz


1. The number of embedded processors sold every year greatly outnumbers the number of PC and even PostPC processors. Can you confirm or deny this insight based on your own experience? Try to count the number of embedded processors in your home. How does it compare with the number of conventional computers in your home?
2. As mentioned earlier, both the software and hardware affect the performance of a program. Can you think of examples where each of the following is the right place to look for a performance bottleneck?
 - The algorithm chosen
 - The programming language or compiler
 - The operating system
 - The processor
 - The I/O system and devices



Below Your Program



- Application software
 - Written in high-level language
- System software
 - Compiler: translates HLL code to machine code
 - Operating System: service code
 - Handling input/output
 - Managing memory and storage
 - Scheduling tasks & sharing resources
- Hardware
 - Processor, memory, I/O controllers



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)

```
000000010100001000000000011000
000000000001100000110000100001
100011000110001000000000000000
100011001111001000000000000100
101011001111001000000000000000
101011001100010000000000000100
00000011110000000000000001000
```


Chapter 1: Computer Technology

17



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

Chapter 1: Computer Technology

18



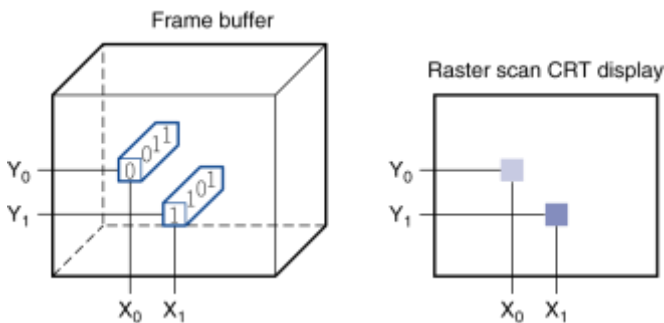
Touchscreen


- PostPC devices
- Supersedes keyboard and mouse
- Resistive and Capacitive types
 - Most tablets, smart phones use capacitive
 - Capacitive allows multiple touches simultaneously



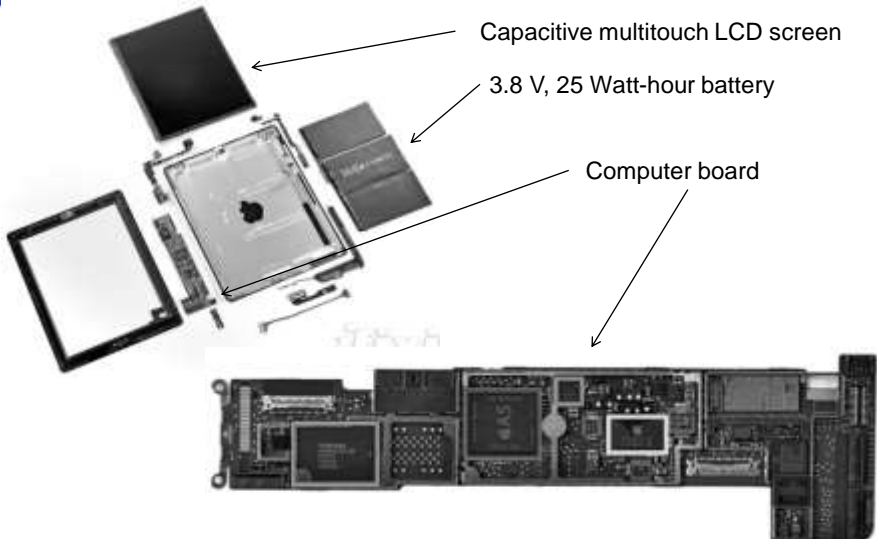
Through the Looking Glass

- LCD screen: picture elements (pixels)
 - Mirrors content of frame buffer memory





Opening the Box




Capacitive multitouch LCD screen

3.8 V, 25 Watt-hour battery

Computer board

Chapter 1: Computer Technology

21



Inside the Processor (CPU)

- Datapath: performs operations on data
- Control: sequences datapath, memory, ...
- Cache memory
 - Small fast SRAM memory for immediate access to data

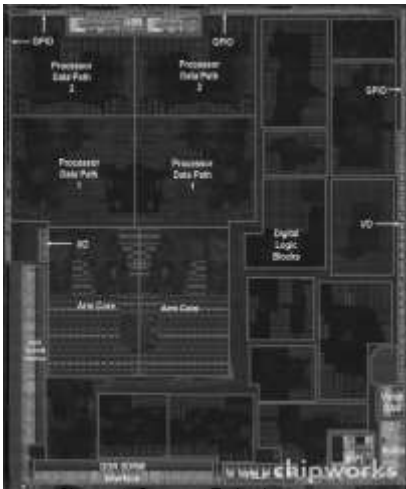
Chapter 1: Computer Technology

22



Inside the Processor

- Apple A5



Chapter 1: Computer Technology




Abstractions

The BIG Picture





- Abstraction helps us deal with complexity
 - Hide lower-level detail
- Instruction set architecture (ISA)
 - The hardware/software interface
- Application binary interface
 - The ISA plus system software interface
- Implementation
 - The details underlying and interface

Chapter 1: Computer Technology




A Safe Place for Data

- Volatile main memory
 - Loses instructions and data when power off
- Non-volatile secondary memory
 - Magnetic disk
 - Flash memory
 - Optical disk (CDROM, DVD)





Chapter 1: Computer Technology

25



Networks

- Communication, resource sharing, nonlocal access
- Local area network (LAN): Ethernet
- Wide area network (WAN): the Internet
- Wireless network: WiFi, Bluetooth



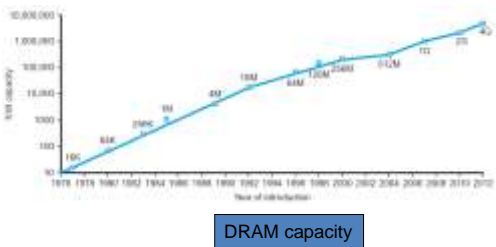
Chapter 1: Computer Technology

26



Technology Trends

- Electronics technology continues to evolve
 - Increased capacity and performance
 - Reduced cost




Year	Technology	Relative performance/cost
1951	Vacuum tube	1
1965	Transistor	35
1975	Integrated circuit (IC)	900
1995	Very large scale IC (VLSI)	2,400,000
2013	Ultra large scale IC	250,000,000,000

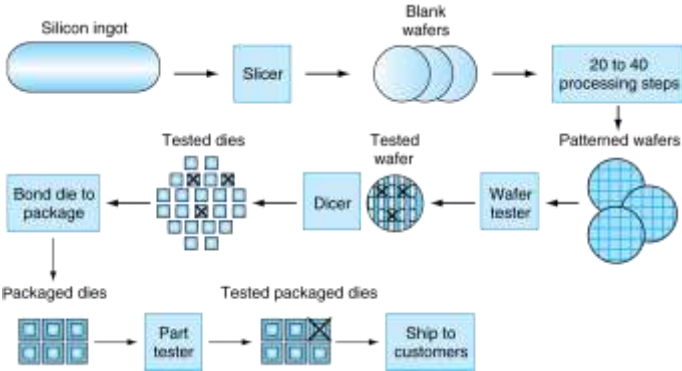


Semiconductor Technology

- Silicon: semiconductor
- Add materials to transform properties:
 - Conductors
 - Insulators
 - Switch




Manufacturing ICs



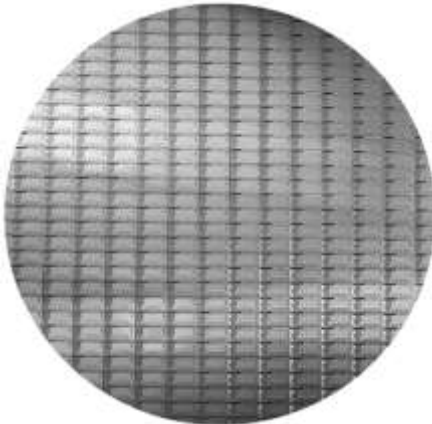
- Yield: proportion of working dies per wafer

Chapter 1: Computer Technology

29



Intel Core i7 Wafer



- 300mm wafer, 280 chips, 32nm technology
- Each chip is 20.7 x 10.5 mm

Chapter 1: Computer Technology

30



Integrated Circuit Cost

$$\text{Cost per die} = \frac{\text{Cost per wafer}}{\text{Dies per wafer} \times \text{Yield}}$$
$$\text{Dies per wafer} \approx \text{Wafer area} / \text{Die area}$$
$$\text{Yield} = \frac{1}{(1 + (\text{Defects per area} \times \text{Die area} / \alpha))^\alpha}$$

- 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



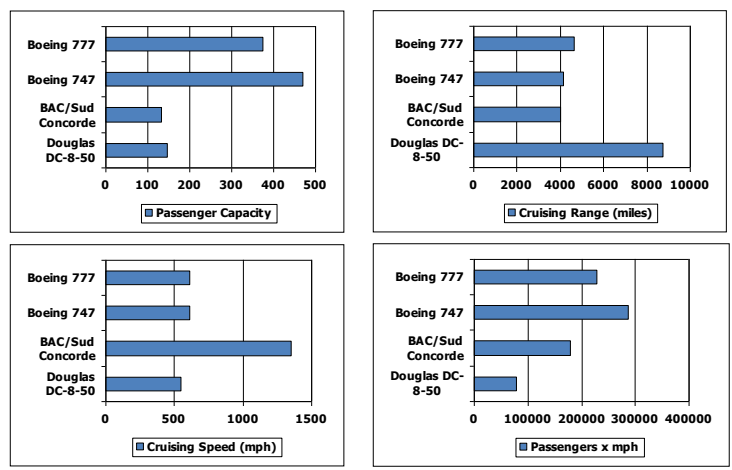
Quiz

- A key factor in determining the cost of an integrated circuit is volume. Which of the following are reasons why a chip made in high volume should cost less?
 1. With high volumes, the manufacturing process can be tuned to a particular design, increasing the yield.
 2. It is less work to design a high-volume part than a low-volume part.
 3. The masks used to make the chip are expensive, so the cost per chip is lower for higher volumes.
 4. Engineering development costs are high and largely independent of volume; thus, the development cost per die is lower with high-volume parts.
 5. High-volume parts usually have smaller die sizes than low-volume parts and therefore have higher yield per wafer



Defining Performance

- Which airplane has the best performance?



Response Time and Throughput

- Response time (execution time)
 - How long it takes to do a task
 - The total time required for the computer to complete a task, including disk accesses, memory accesses, I/O activities, operating system overhead, CPU execution time, and so on
- 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”

$$\begin{aligned} &\text{Performance}_x / \text{Performance}_y \\ &= \text{Execution time}_y / \text{Execution time}_x = n \end{aligned}$$

- Example: time taken to run a program
 - 10s on A, 15s on B
 - $\text{Execution Time}_B / \text{Execution Time}_A$
 $= 15\text{s} / 10\text{s} = 1.5$
 - 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



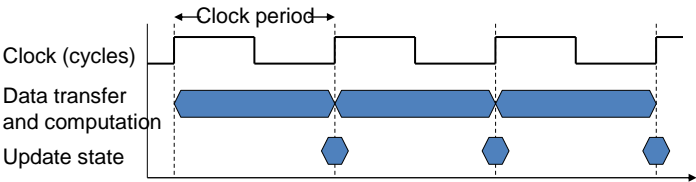
Measuring Execution Time

- 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., $250\text{ps} = 0.25\text{ns} = 250 \times 10^{-12}\text{s}$
- Clock frequency (rate): cycles per second
 - e.g., $4.0\text{GHz} = 4000\text{MHz} = 4.0 \times 10^9\text{Hz}$



CPU Time

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

- 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?

$$\begin{aligned} \text{Clock Rate}_B &= \frac{\text{Clock Cycles}_B}{\text{CPU Time}_B} = \frac{1.2 \times \text{Clock Cycles}_A}{6\text{s}} \\ \text{Clock Cycles}_A &= \text{CPU Time}_A \times \text{Clock Rate}_A \\ &= 10\text{s} \times 2\text{GHz} = 20 \times 10^9 \\ \text{Clock Rate}_B &= \frac{1.2 \times 20 \times 10^9}{6\text{s}} = \frac{24 \times 10^9}{6\text{s}} = 4\text{GHz} \end{aligned}$$



Instruction Count and CPI

Clock Cycles = Instruction Count × Cycles per Instruction

CPU Time = Instruction Count × CPI × 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



CPI Example

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

CPU Time_A = Instruction Count × CPI_A × Cycle Time_A

= I × 2.0 × 250ps = I × 500ps

CPU Time_B = Instruction Count × CPI_B × Cycle Time_B

= I × 1.2 × 500ps = I × 600ps

$\frac{\text{CPU Time}_B}{\text{CPU Time}_A} = \frac{I \times 600\text{ps}}{I \times 500\text{ps}} = 1.2$



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

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



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



Exercise

- A program consists of 1000 instructions in which:
 - 30% load/store instructions, CPI = 2.5
 - 10% jump instructions, CPI = 1
 - 20% branch instructions, CPI = 1.5
 - The rest are arithmetic instructions, CPI = 2.0
- The program is executed on a 2 GHz CPU
 - a) What is execution time (CPU time) of the program?
 - b) What is the weight average CPI of the program?
 - c) If load/store instructions are improved so that their execution time is reduced by a factor of 2, what is the speed-up of the system?



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

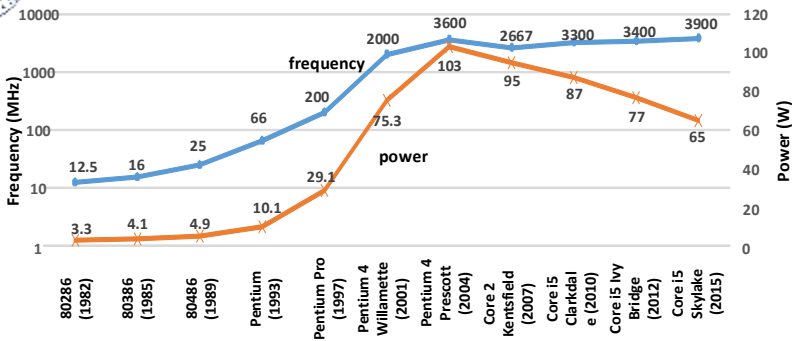


Quiz

- A given application written in Java runs 15 seconds on a desktop processor. A new Java compiler is released that requires only 0.6 as many instructions as the old compiler. Unfortunately, it increases the CPI by 1.1. How fast can we expect the application to run using this new compiler?

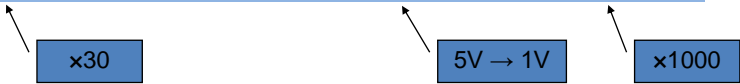


Power Trends



- In CMOS IC technology

Power = Capacitive load × Voltage² × Frequency





Reducing Power

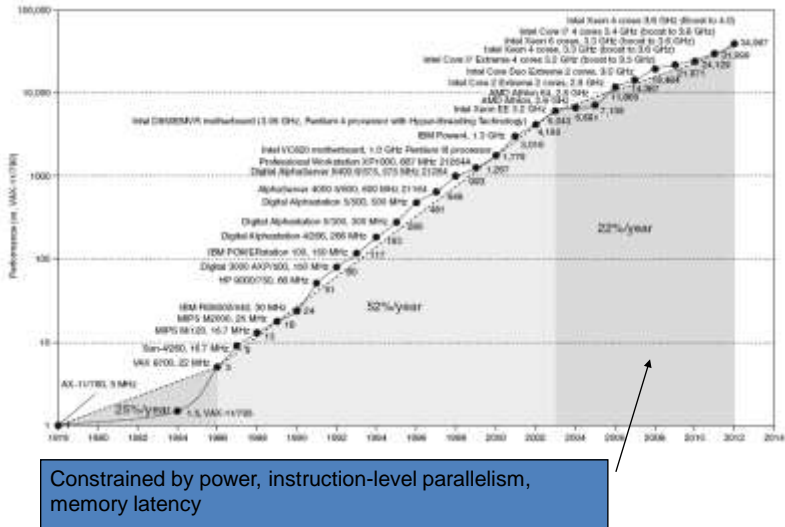
- Suppose a new CPU has
 - 85% of capacitive load of old CPU
 - 15% voltage and 15% frequency reduction

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

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



Uniprocessor Performance





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)

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

Chapter 1: Computer Technology



CINT2006 for Intel Core i7 920

Description	Name	Instruction Count x 10 ⁶	CPI	Clock cycle time (seconds x 10 ⁻⁹)	Execution Time (seconds)	Reference Time (seconds)	SPECratio
Interpreted string processing	perl	2252	0.60	0.376	508	9770	19.2
Block-sorting compression	ltp2	2390	0.70	0.376	629	9650	15.4
GNU C compiler	gcc	794	1.20	0.376	358	8050	22.5
Combinatorial optimization	ncf	221	2.66	0.376	221	9120	41.2
Go game (AI)	go	1274	1.10	0.376	527	10490	19.9
Search gene sequence	hmmer	2616	0.60	0.376	590	9330	15.8
Chess game (AI)	sjeng	1948	0.80	0.376	586	12100	20.7
Quantum computer simulation	ibquantum	659	0.44	0.376	109	20720	190.0
Video compression	h264enc	3793	0.50	0.376	713	22130	31.0
Discrete event simulation library	smnetpp	367	2.10	0.376	290	6250	21.5
Games/path finding	astar	1250	1.00	0.376	470	7020	14.9
XML parsing	svalenchmk	1045	0.70	0.376	275	6900	25.1
Geometric mean	-	-	-	-	-	-	25.7



SPEC Power Benchmark

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

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



SPECpower_ssj2008 for Xeon X5650

Target Load %	Performance (ssj_ops)	Average Power (Watts)
100%	865,618	258
90%	786,688	242
80%	698,051	224
70%	607,826	204
60%	521,391	185
50%	436,757	170
40%	345,919	157
30%	262,071	146
20%	176,061	135
10%	86,784	121
0%	0	80
Overall Sum	4,787,166	1,922
Σssj_ops/Σpower =		2,490



Pitfall: Amdahl's Law

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

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

- Amdahl's Law: A rule stating that the performance enhancement possible with a given improvement is limited by the amount that the improved feature is used.
 - It is a quantitative version of the law of diminishing returns
- Example: multiply accounts for 80s/100s
 - How much improvement in multiply performance to get 5x overall?
 - Can't be done: $20 = \frac{80}{n} + 20$
- 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

$$\begin{aligned} \text{MIPS} &= \frac{\text{Instruction count}}{\text{Execution time} \times 10^6} \\ &= \frac{\text{Instruction count}}{\frac{\text{Instruction count} \times \text{CPI}}{\text{Clock rate}} \times 10^6} = \frac{\text{Clock rate}}{\text{CPI} \times 10^6} \end{aligned}$$

- CPI varies between programs on a given CPU



Quiz

Consider the following performance measurements for a program:

Measurement	Computer A	Computer B
Instruction count	10 billion	8 billion
Clock rate	4 GHz	4 GHz
CPI	1.0	1.1

- a. Which computer has the higher MIPS rating?
- b. Which computer is faster?



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