



# COMPUTER ORGANIZATION AND DESIGN

## The Hardware/Software Interface

# Chapter 1

## Computer Abstractions and Technology

heeyoul choi  
[hchoi@handong.edu](mailto:hchoi@handong.edu)

School of CS and EE  
Handong Global University

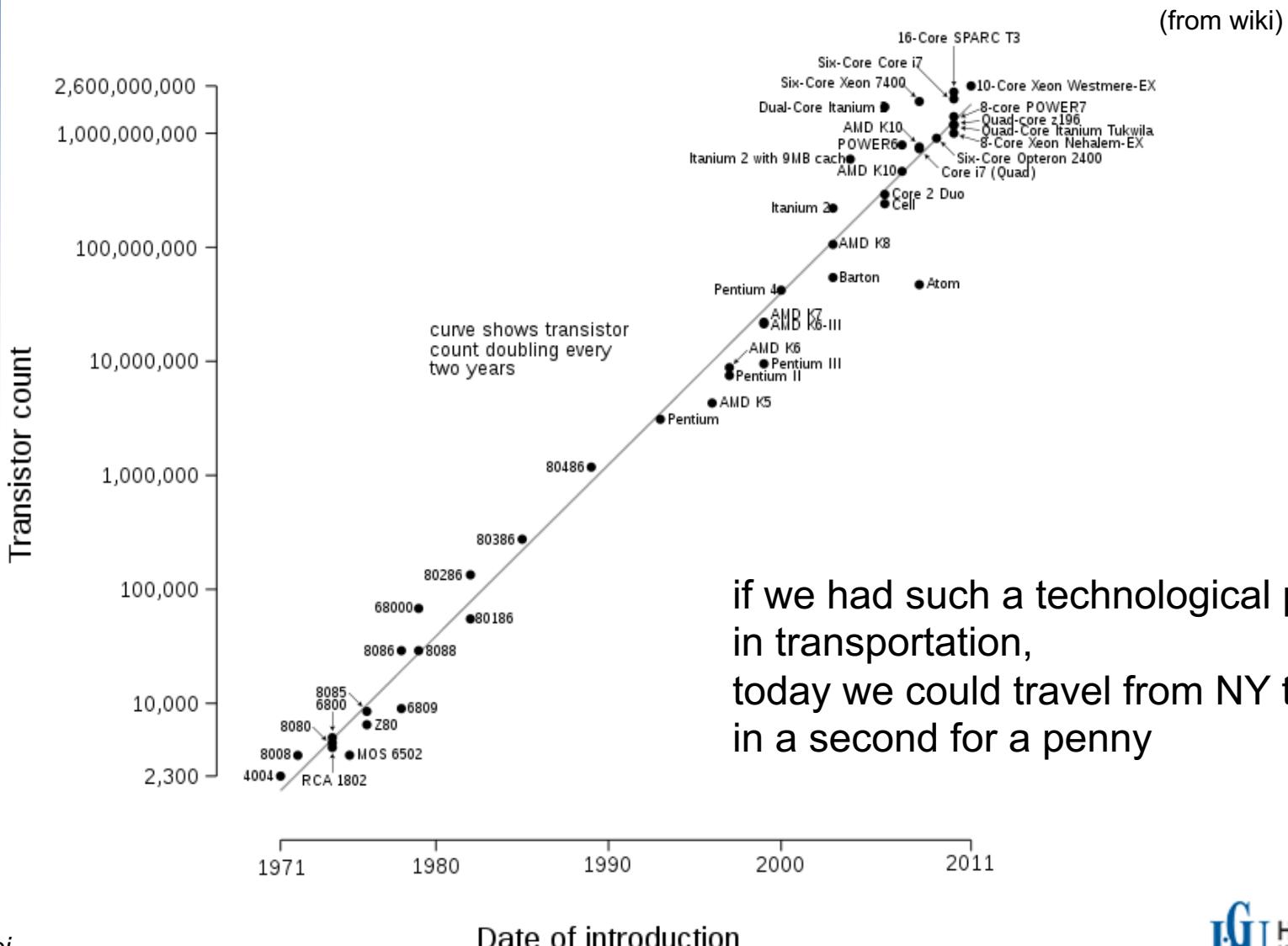
# The Computer Revolution

- 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

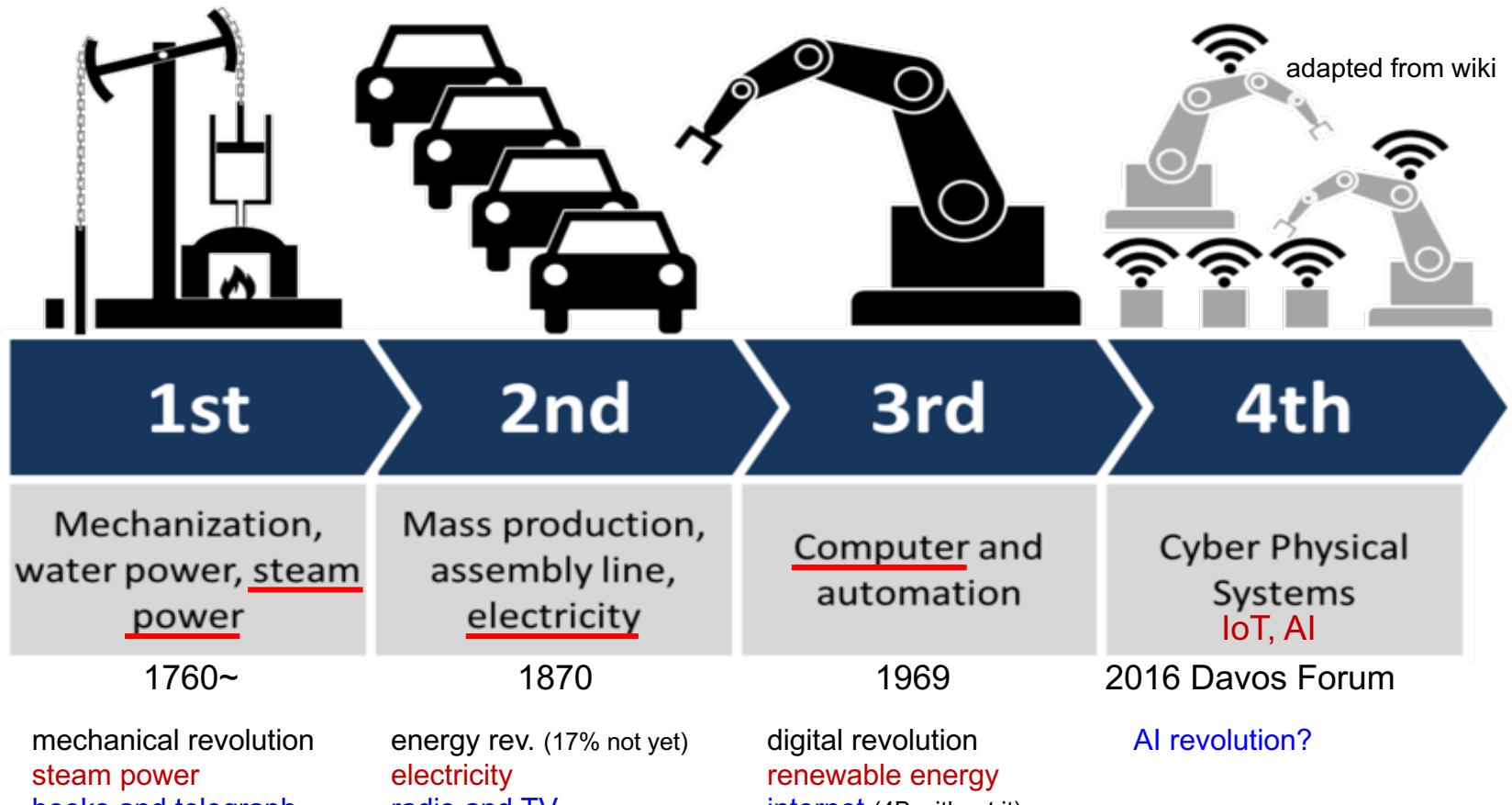
Today's SF suggests tomorrow's killer apps.

# The Computer Revolution

Microprocessor Transistor Counts 1971-2011 & Moore's Law



# industrial revolution



- about **energy** and **communication style**

# 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

Year	Supercomputer	Peak speed (Rmax)	Location
2018	IBM Summit	122.3 PFLOPS	Oak Ridge, U.S.
2016	Sunway TaihuLight	93.01 PFLOPS	Wuxi, China
2013	NUDT Tianhe-2	33.86 PFLOPS	Guangzhou, China
2012	Cray Titan	17.59 PFLOPS	Oak Ridge, U.S.
2012	IBM Sequoia	17.17 PFLOPS	Livermore, U.S.
2011	Fujitsu K computer	10.51 PFLOPS	Kobe, Japan
2010	Tianhe-IA	2.566 PFLOPS	Tianjin, China
2009	Cray Jaguar	1.759 PFLOPS	Oak Ridge, U.S.
2008	IBM Roadrunner	1.026 PFLOPS	Los Alamos, U.S.
		1.105 PFLOPS	

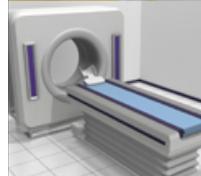
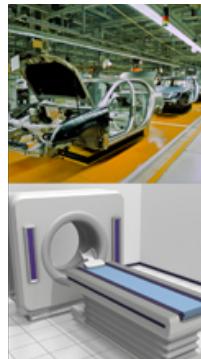
from wiki

# Classes of Computers

## ■ Embedded computers

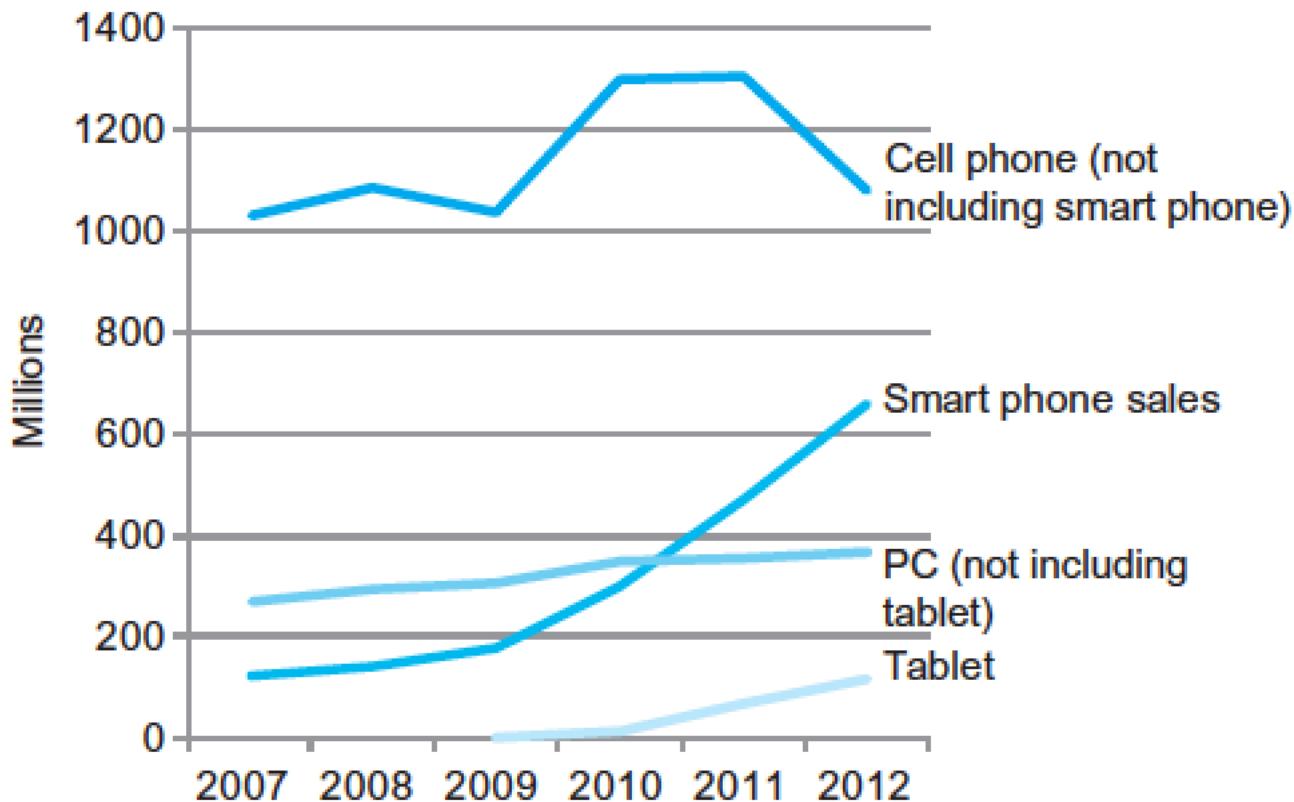
- Hidden as components of systems
  - running specific applications integrated with the HW
- Stringent power/performance/cost constraints

including microprocessors  
in cars or television sets



# The PostPC Era

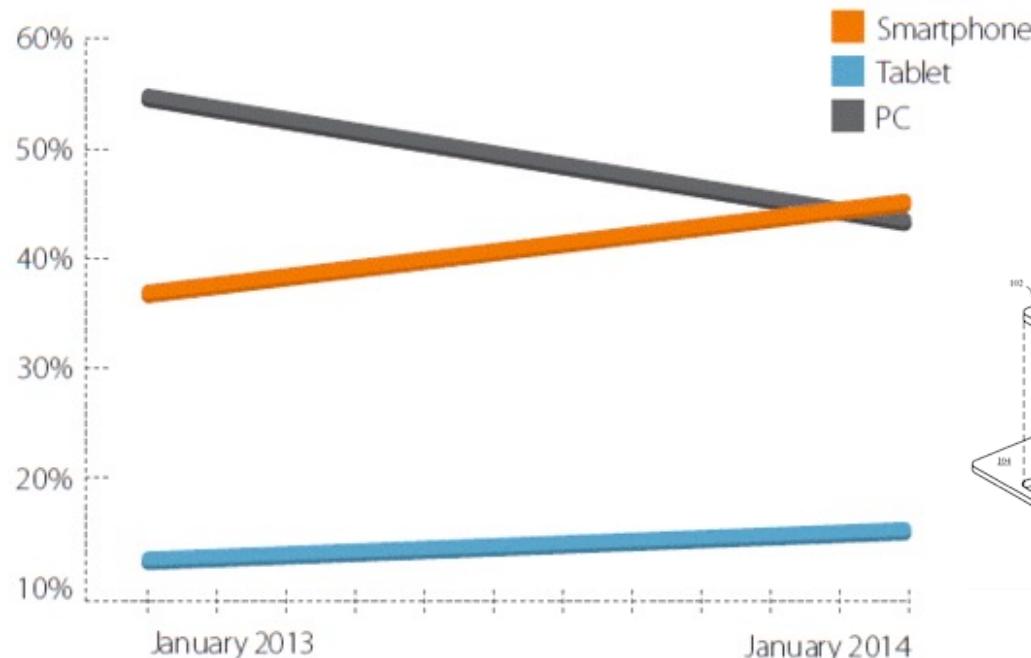
the number manufactured per year



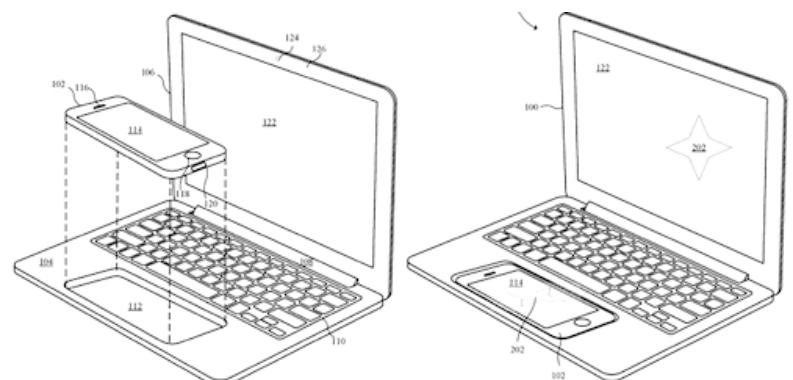
# The PostPC Era

## Share of Monthly Digital Minutes

In 2014, the amount of time spent on smartphones surpassed PCs, in terms of both total minutes spent and as a percentage of total time. Time spent on smartphones grew 7% year-over-year, while time spent on PCs decreased 8%. Tablets grew slightly at 2%.



a patent application from Apple



from blog.eogn.com

Source: comScore Multi-Platform and Mobile Metrix, U.S.; January 2013 & 2014.

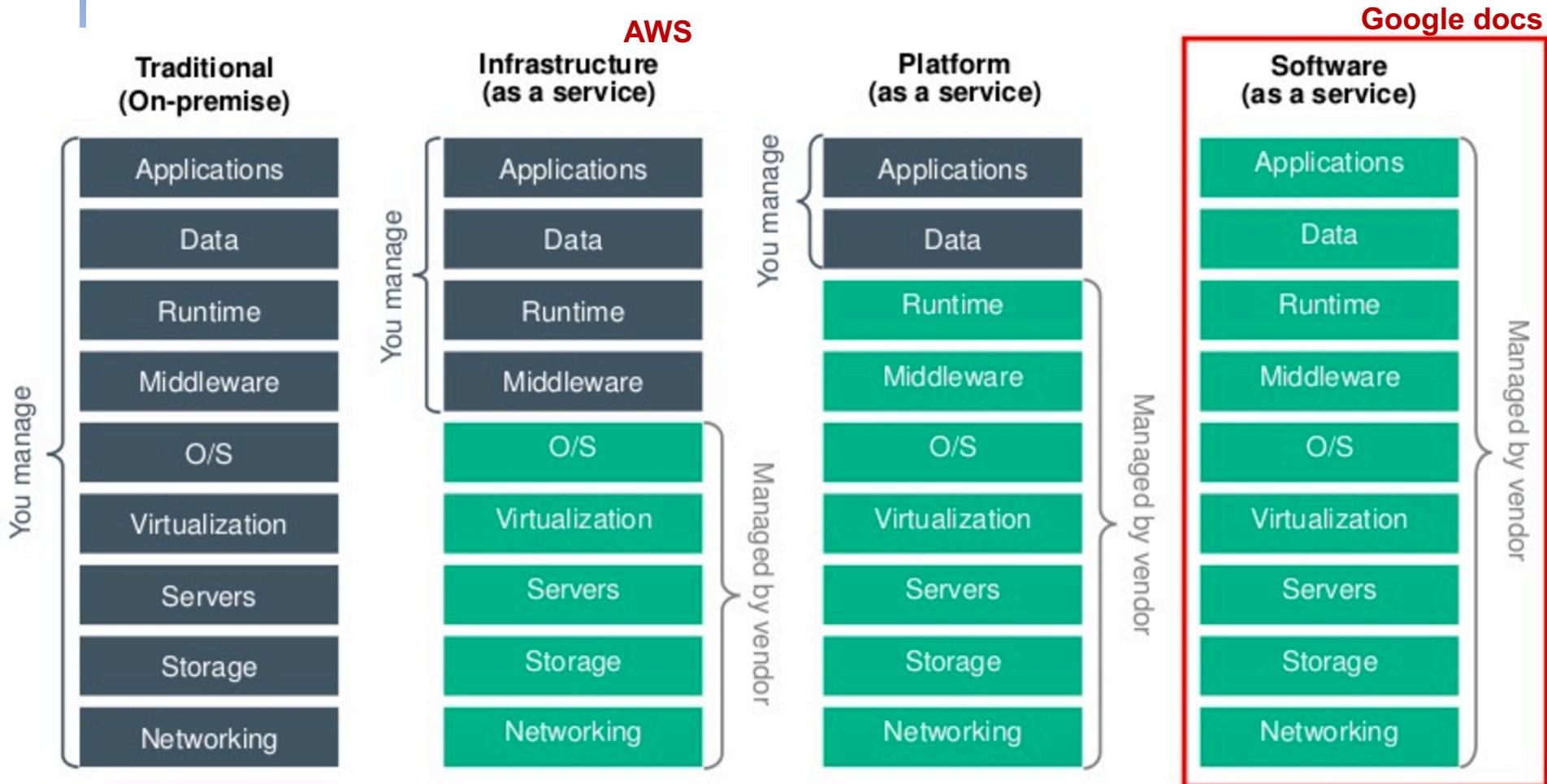
from [www.networkworld.com](http://www.networkworld.com)

# The PostPC Era

- Personal Mobile Device (PMD)
  - Battery operated
  - Connects to the Internet
  - Hundreds of dollars
  - Smart phones, tablets, electronic glasses
- Cloud computing
  - Warehouse Scale Computers (WSC) - giant datacenters
  - Software as a Service (SaaS)
  - Amazon and Google
- software runs both on a PMD and in the Cloud.
  - e.g., speech recognition in your smartphone

# Software as a Service (SaaS)

## Cloud Computing



from HP Korea Slides

# What You Will Learn

- How programs are translated into the machine language
  - And how the hardware executes them
- The hardware/software interface
- What determines program performance
  - And how it can be improved
- How hardware designers improve performance
  - and energy efficiency
- What is parallel processing

# Understanding Performance

- Algorithm
  - Determines the 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

# Eight Great Ideas

- Design for ***Moore's Law***
- Use ***abstraction*** to simplify design
- Make the ***common case fast***
- Performance *via parallelism*
- Performance *via pipelining*
- Performance *via prediction*
- ***Hierarchy*** of memories
- ***Dependability*** *via redundancy*



COMMON CASE FAST



PARALLELISM



PIPELINING



PREDICTION

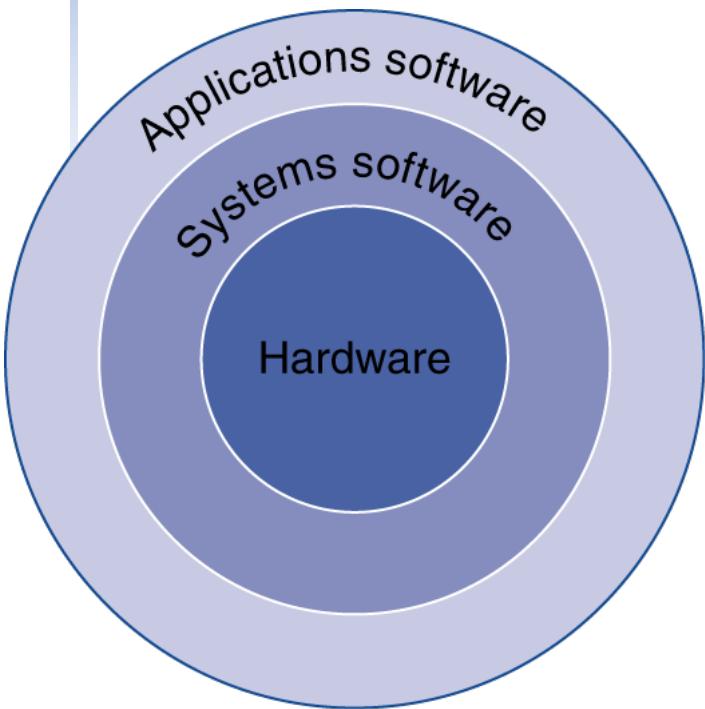


HIERARCHY



DEPENDABILITY

# Below Your Program



- Application software
  - Written in high-level language (HLL)
- 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

abstraction

## High-level language

- Level of abstraction closer to problem domain
- Provides for productivity and portability

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

- Textual representation of instructions

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

## Hardware representation

- Binary digits (bits)
- Encoded instructions and data
- machine language

Binary machine  
language  
program  
(for MIPS)

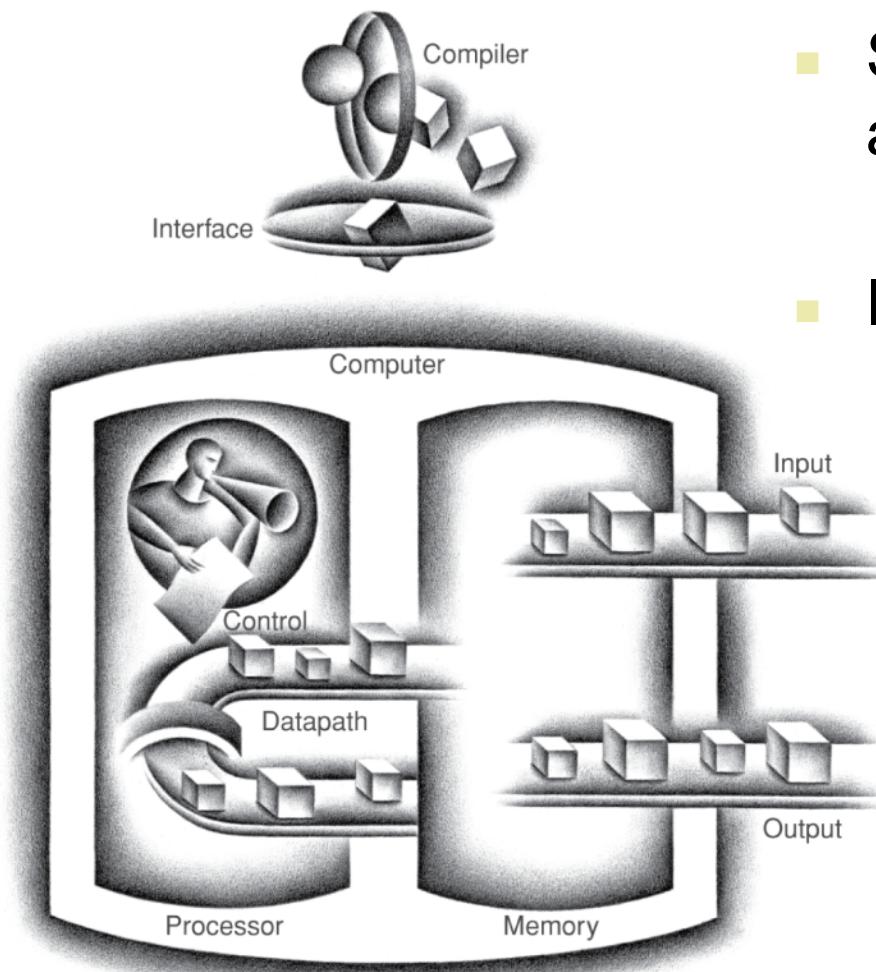
```
000000001010000100000000000011000
00000000000110000001100000100001
1000110001100010000000000000000000
1000110011110010000000000000000000
1010110011110010000000000000000000
1010110001100010000000000000000000
0000001111000000000000000000000000
```

# benefits of high-level languages

- programmers can think in a natural language and algebraic notations
  - Fortran for science computation
  - Cobol for business data processing
  - Lisp for symbol manipulation
- the productivity is improved
  - it takes less time to develop programs with fewer lines.
- programs are independent of the computer
  - compilers and assemblers translate them to the binary instructions of any computer.

# Components of a Computer

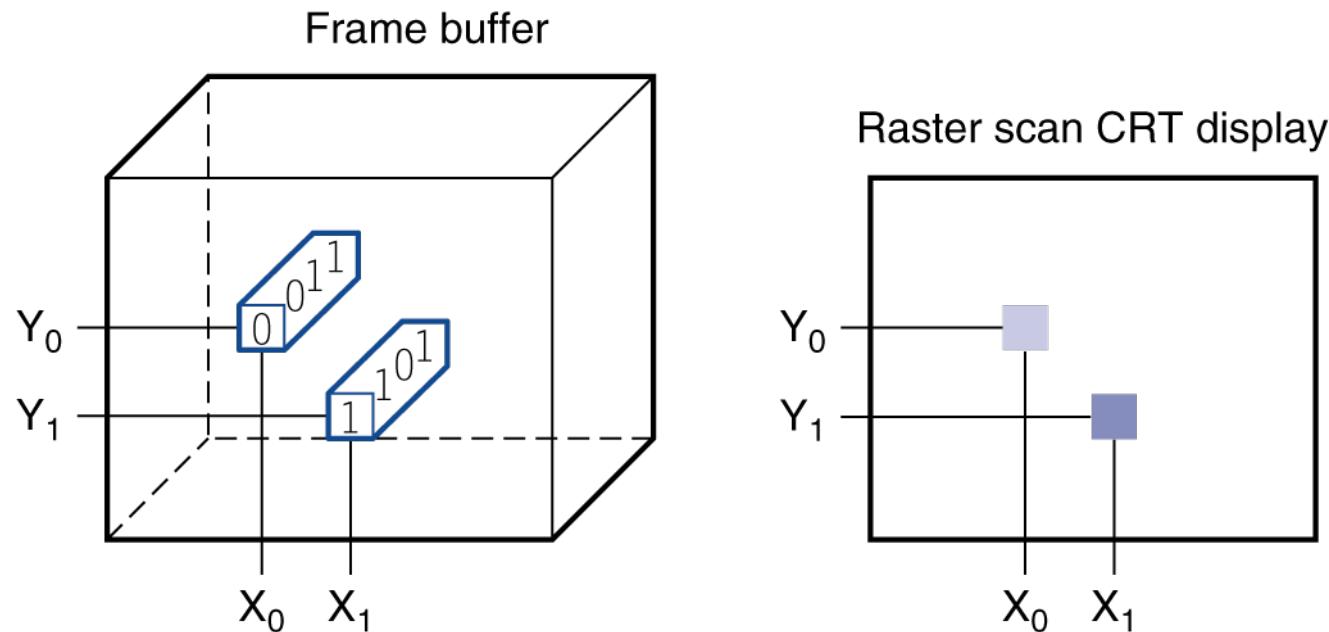
## The BIG Picture



- 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

# Through the Looking Glass

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

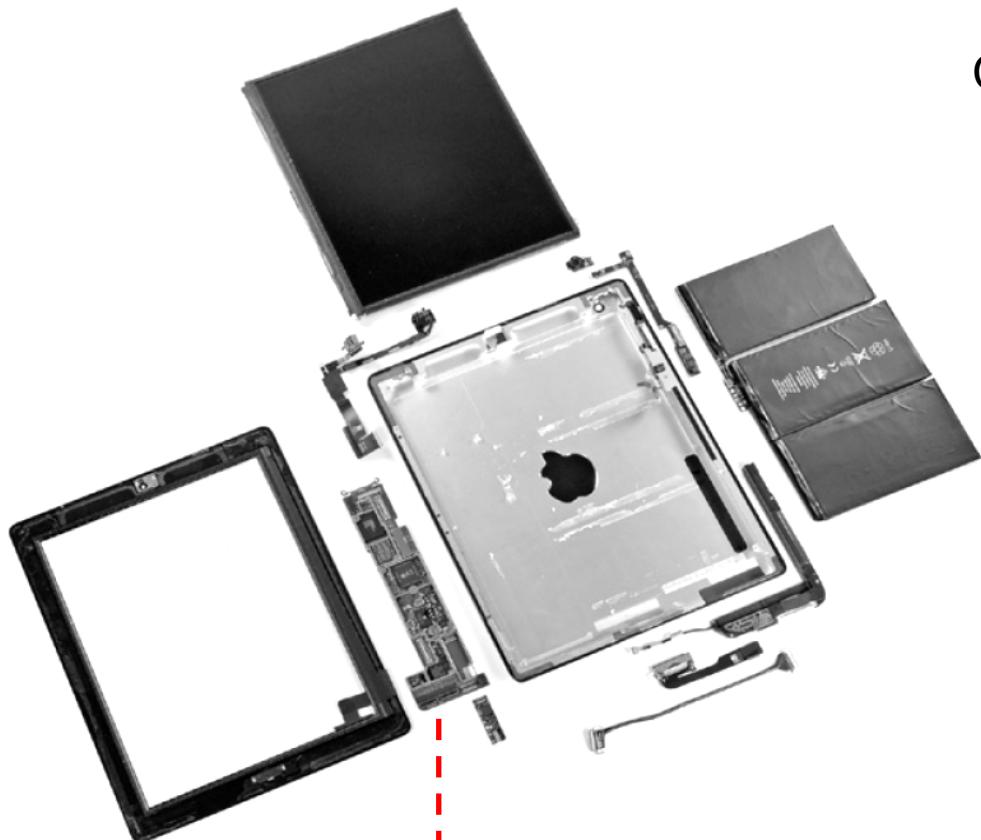


# Touchscreen

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



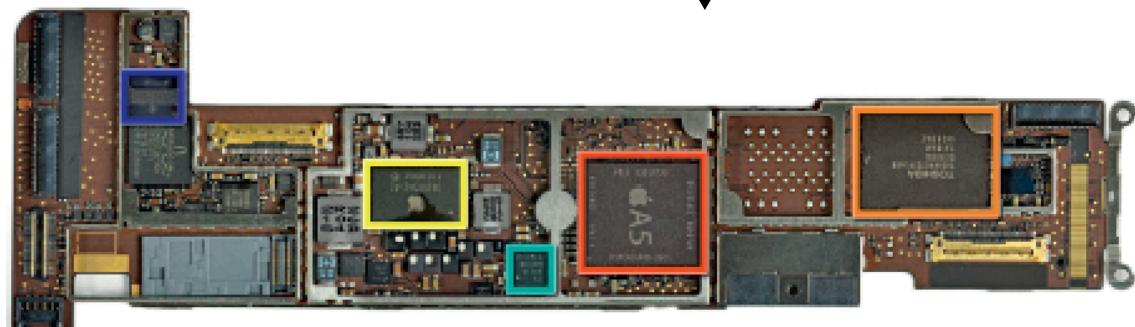
# Opening the Box



Capacitive multitouch LCD screen

3.8 V, 25 Watt-hour battery

Computer board



# Inside the Processor (CPU)

- Datapath: performs arithmetic operations
- Control: controls datapath, memory, I/O...
- Cache memory
  - Small fast SRAM memory for immediate access to data

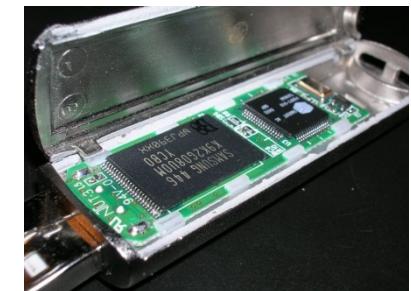
# Abstractions

## 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
  - HW that obeys the ISA abstraction
  - The details underlying and interface

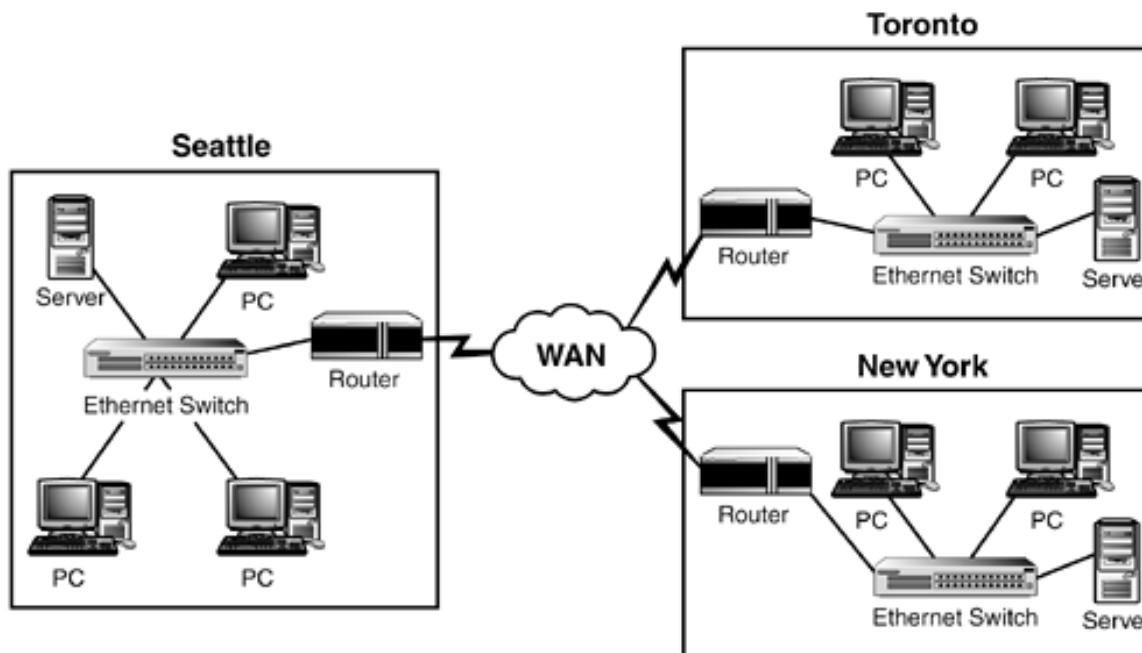
# A Safe Place for Data

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



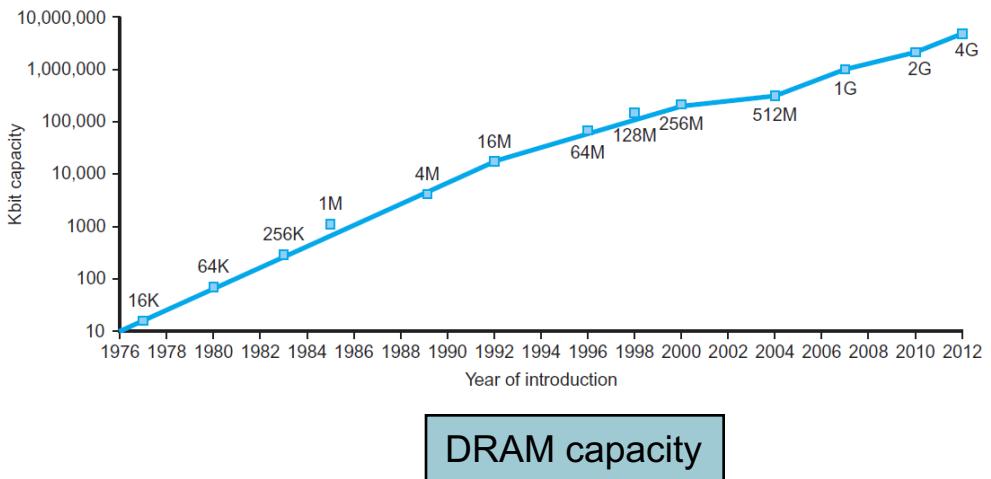
# Networks

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



# Technology Trends

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

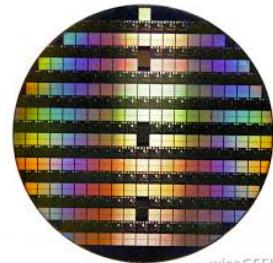
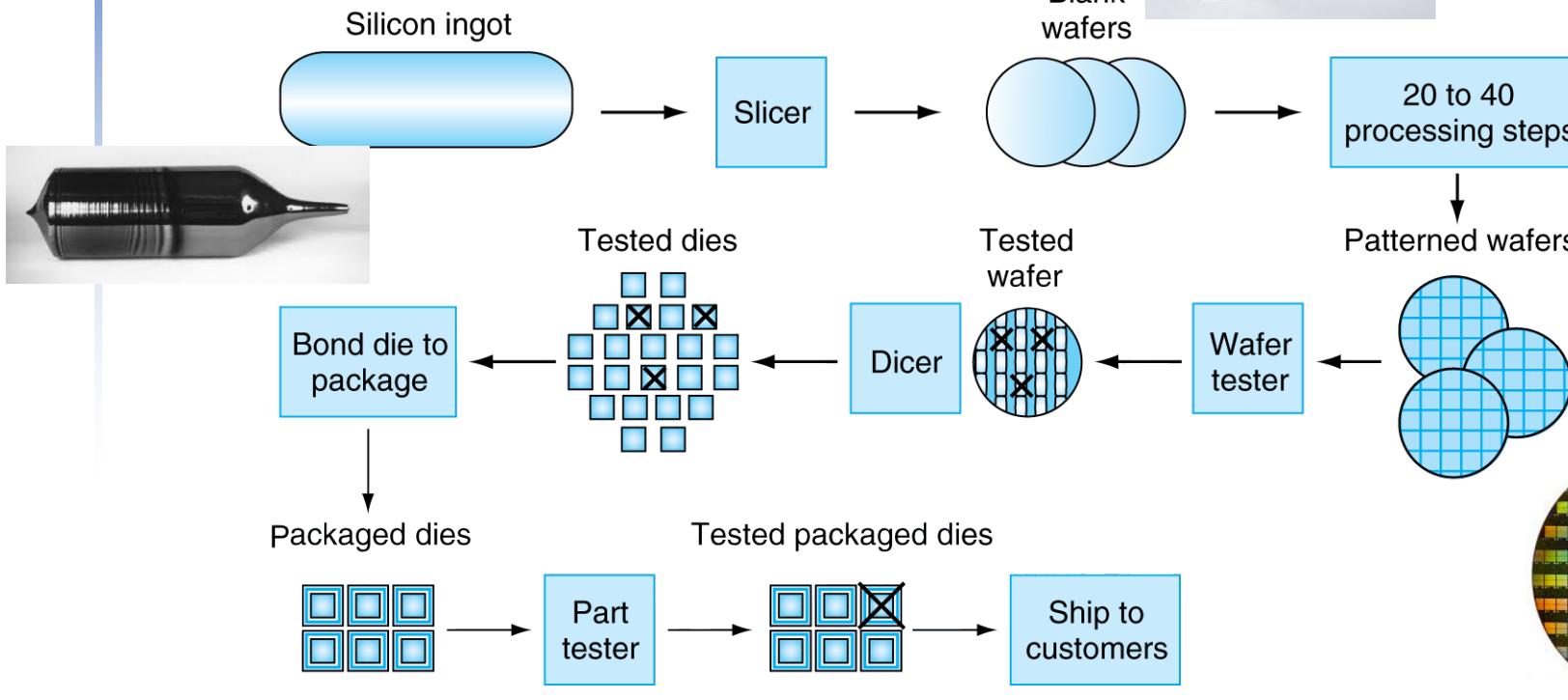


Year	Technology	Relative performance/cost
1951	Vacuum tube	1
1965	Transistor	35
1975	Integrated circuit (IC)	900
1995	Very large scale IC (VLSI)	2,400,000
2013	Ultra large scale IC (extrapolated)	250,000,000,000

# Semiconductor Technology

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

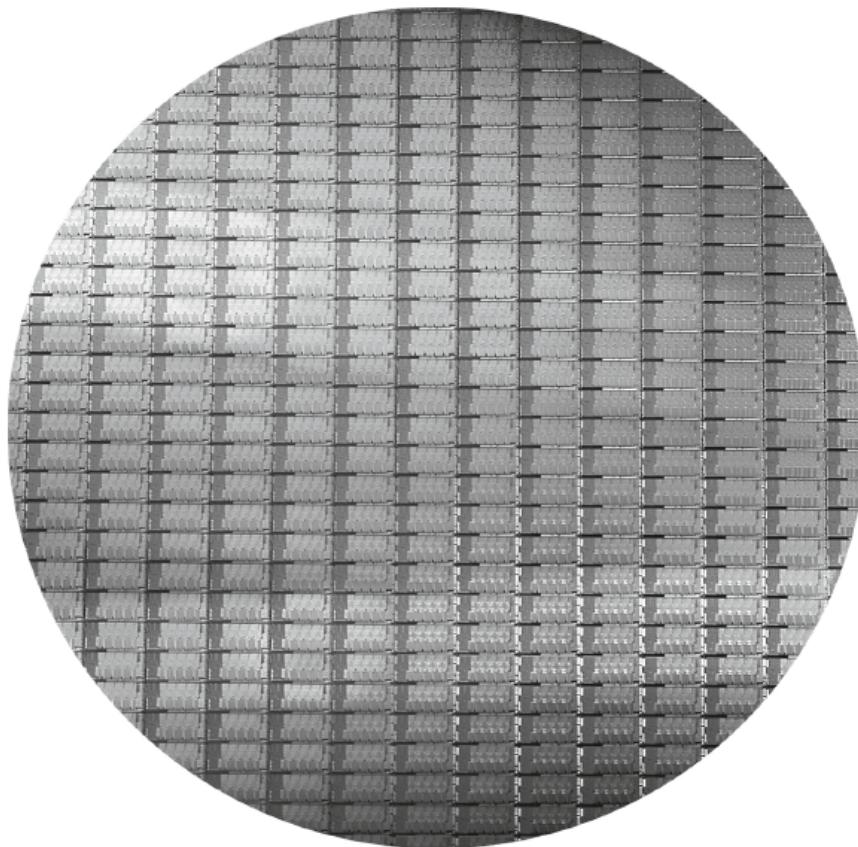
# Manufacturing ICs



wiseGEEK

- Yield: proportion of working dies per wafer

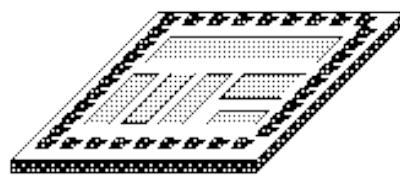
# Intel Core i7 Wafer



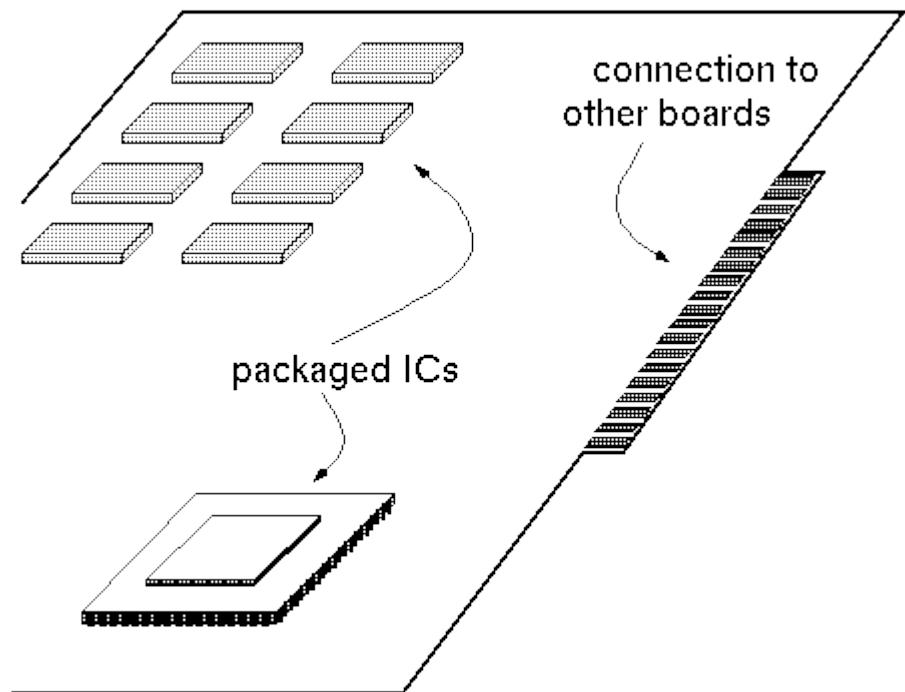
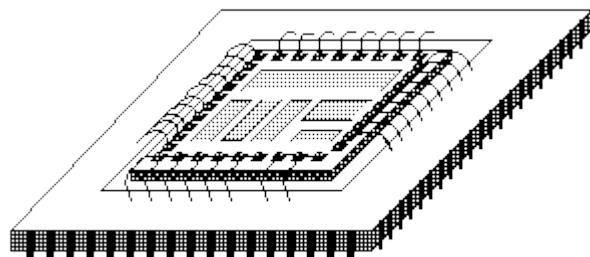
- 300mm wafer, 280 chips, 32nm technology
- Each chip is 20.7 x 10.5 mm

# Integrated Circuit

Bare Die



Chip in Package



# Integrated Circuit Cost

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

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

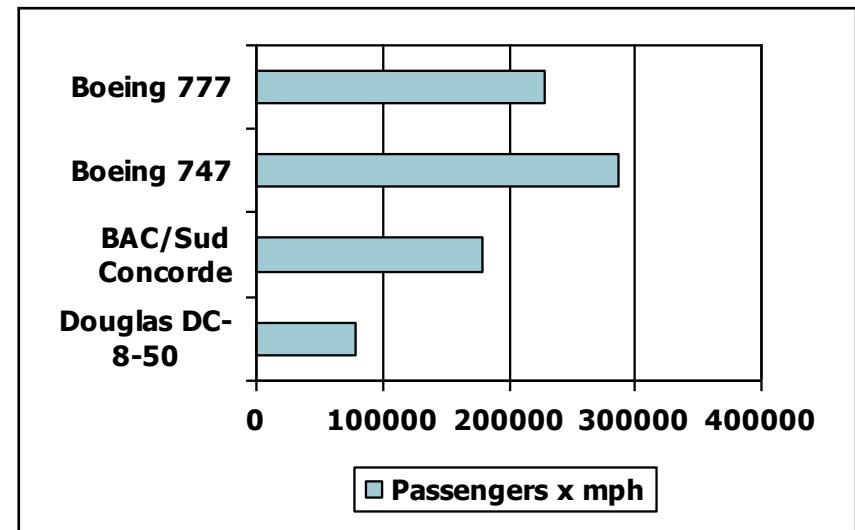
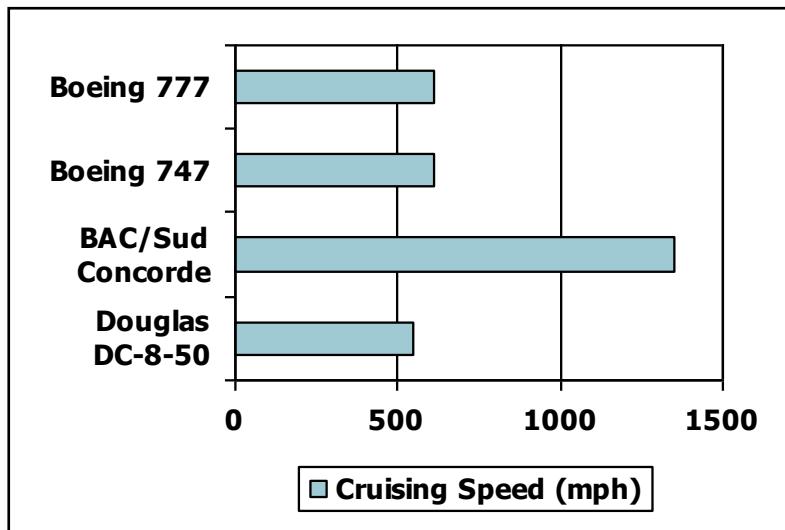
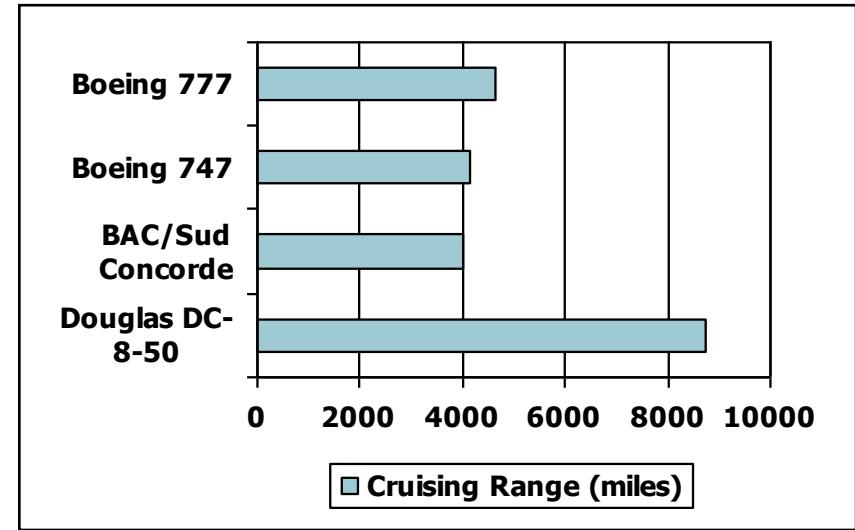
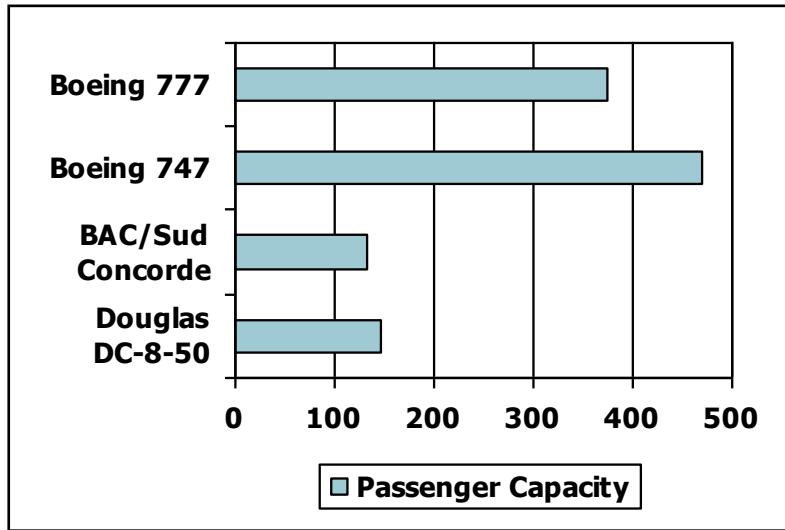
$$\text{Yield} = \frac{1}{(1 + (\text{Defects per area} \times \text{Die area}/2))^2}$$

empirical  
observations

- Cost is nonlinear 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
- to reduce the cost,
  - die area should decrease
  - and manufacturing process should be improved.

# Defining Performance

- Which airplane has the best performance?



# Response Time and Throughput

- Response time (also referred to as execution time)
  - How long it takes to do a task
- Throughput (or bandwidth)
  - 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...
  - for the first few chapters

# Relative Performance

- Define Performance = 1/Execution Time
- “X is  $n$  times faster than Y”

$$\begin{aligned} \text{Performance}_X / \text{Performance}_Y \\ = \text{Execution time}_Y / \text{Execution time}_X = n \end{aligned}$$

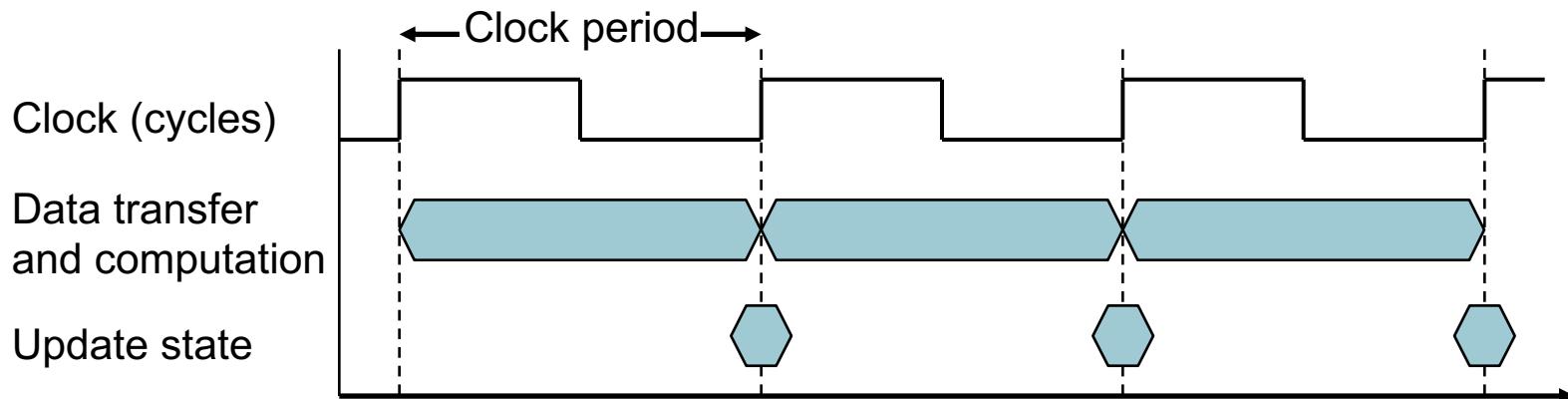
- Example: time taken to run a program
  - 10s on A, 15s on B
  - $\text{Execution Time}_B / \text{Execution Time}_A$   
 $= 15s / 10s = 1.5$
  - So, A is 1.5 times faster than B

# 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 or other jobs' shares
  - Comprises **user CPU time and system CPU time**
  - Determines CPU performance
- 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

CPU Time = CPU Clock Cycles  $\times$  Clock Cycle Time

$$= \frac{\text{CPU Clock Cycles}}{\text{Clock Rate}}$$

- Performance improved by
  - Reducing the number of clock cycles
  - Increasing the clock rate
  - Hardware designer must often trade off clock rate against cycle count for a program

# CPU Time Example

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

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

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

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

# Instruction Count and CPI

Clock Cycles = Instruction Count  $\times$  Cycles per Instruction

CPU Time = Instruction Count  $\times$  CPI  $\times$  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
- The same ISA
- Which one is faster, and by how much?

$$\begin{aligned}\text{CPU Time}_A &= \text{Instruction Count} \times \text{CPI}_A \times \text{Cycle Time}_A \\ &= I \times 2.0 \times 250\text{ps} = I \times 500\text{ps}\end{aligned}$$

A is faster...

$$\begin{aligned}\text{CPU Time}_B &= \text{Instruction Count} \times \text{CPI}_B \times \text{Cycle Time}_B \\ &= I \times 1.2 \times 500\text{ps} = I \times 600\text{ps}\end{aligned}$$

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

...by this much

# 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

# Instruction Count and CPI

- If there are 5 types of instructions:
  - Load (5 cycles)
  - Store (4 cycles)
  - R-type (4 cycles)
  - Branch (3 cycles)
  - Jump (3 cycles)
- and if a program has:
  - 50% load instructions
  - 15% store instructions
  - 25% R-type instructions
  - 8% branch instructions
  - 2% jump instructions
- then, the CPI is:

$$\text{CPI} = \frac{5 \times 50 + 4 \times 15 + 4 \times 25 + 3 \times 8 + 3 \times 2}{100} = 4.4$$

# CPI Example

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

Instruction 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  
 $= 2 \times 1 + 1 \times 2 + 2 \times 3$   
 $= 10$
  - Avg. CPI =  $10/5 = 2.0$
- Sequence 2: IC = 6
  - Clock Cycles  
 $= 4 \times 1 + 1 \times 2 + 1 \times 3$   
 $= 9$
  - Avg. CPI =  $9/6 = 1.5$

# 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,  $T_c$  (clock cycle time)

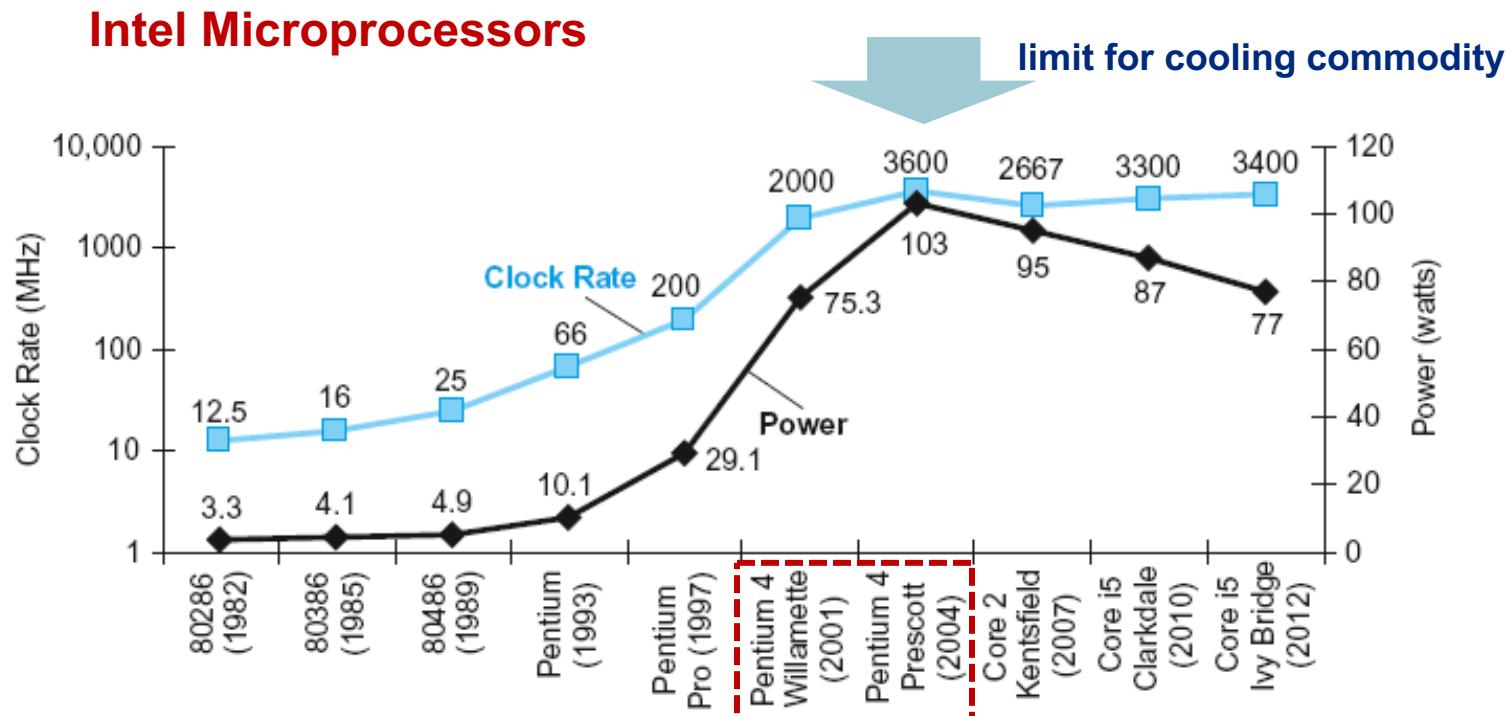
# Performance Summary

from textbook p. 40

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? Pick the right answer from the three choices below:

- a.  $\frac{15 \times 0.6}{1.1} = 8.2 \text{ sec}$
- b.  $15 \times 0.6 \times 1.1 = 9.9 \text{ sec}$
- c.  $\frac{15 \times 1.1}{0.6} = 27.5 \text{ sec}$

# Power Trends



- In CMOS IC technology

$$\text{Power} = \text{Capacitive load} \times \text{Voltage}^2 \times \text{Frequency}$$

$$w_c(t) = \frac{1}{2} C v_c^2(t)$$

×30

5V → 1V

×1000

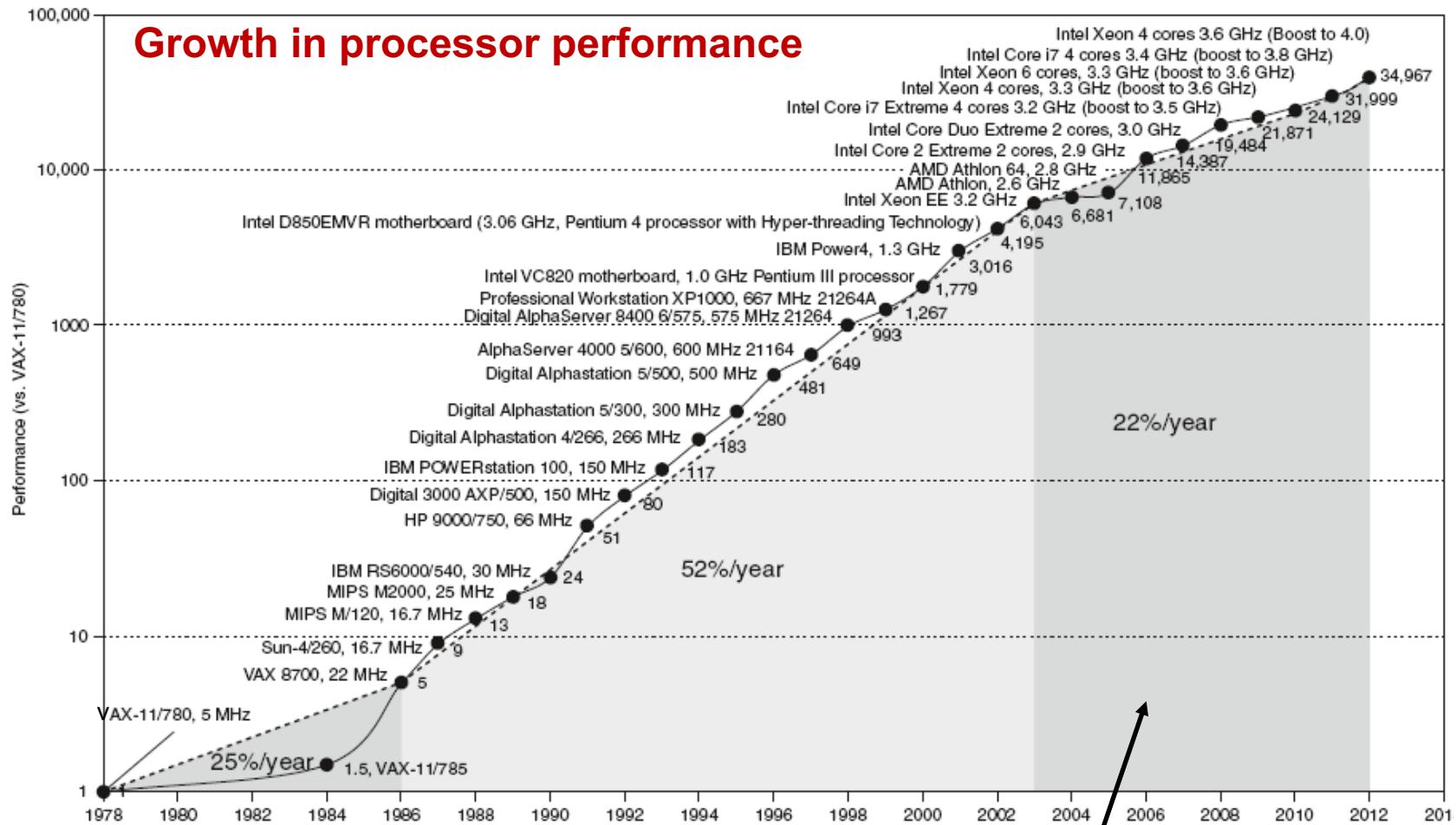
# Reducing Power

- Suppose a new CPU has
  - 85% of capacitive load of old CPU
  - 15% voltage and 15% frequency reductionwhat is the impact on dynamic power?

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

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

# Uniprocessor Performance



Constrained by power, instruction-level parallelism,  
memory latency

# Multiprocessors

- Multicore microprocessors
  - More than one processor (or core) per chip
- Requires explicitly parallel programming
  - Compare with instruction level parallelism (pipelining)
    - Hardware executes multiple instructions at once
    - Hidden from the programmer
  - Hard to do
    - Programming for performance
    - load balancing and synchronization
    - optimizing communication (overhead)

# SPEC CPU Benchmark

- Programs used to measure performance
  - Supposedly typical of actual workload
- Standard Performance Evaluation Cooperative (SPEC)
  - Develops benchmarks for CPU, I/O, Web, ...
- SPEC CPU2006
  - CPU time to execute a selection of programs
    - Negligible I/O, so focuses on CPU performance
  - Normalize relative to reference machine
  - Summarize as the geometric mean of performance ratios
    - CINT2006 (integer) and CFP2006 (floating-point)

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

# CINT2006 for Intel Core i7 920

Description	Name	Instruction Count x 10 <sup>9</sup>	CPI	Clock cycle time (seconds x 10 <sup>-9</sup> )	Execution Time (seconds)	Reference Time (seconds)	SPECratio
Interpreted string processing	perl	2252	0.60	0.376	508	9770	19.2
Block-sorting compression	bzip2	2390	0.70	0.376	629	9650	15.4
GNU C compiler	gcc	794	1.20	0.376	358	8050	22.5
Combinatorial optimization	mcf	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	libquantum	659	0.44	0.376	109	20720	190.0
Video compression	h264avc	3793	0.50	0.376	713	22130	31.0
Discrete event simulation library	omnetpp	367	2.10	0.376	290	6250	21.5
Games/path finding	astar	1250	1.00	0.376	470	7020	14.9
XML parsing	xalancbmk	1045	0.70	0.376	275	6900	25.1
Geometric mean	-	-	-	-	-	-	25.7

25.7x faster

# SPEC Power Benchmark

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

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

# 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
$\Sigma ssj\_ops / \Sigma power =$		2,490

# Pitfall: Amdahl's Law

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

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

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

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

- 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
  - 33% of the peak power at 10% of the load (2012)
- 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

# Pitfall: MIPS as a Performance Metric

from textbook p. 51

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

- Which computer has the higher MIPS rating?
- 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

“Science must begin with myths,  
and the criticism of myths”

- Sir Karl Popper

from textbook p. 49