

Welcome to CS 2323

Computer Architecture

Praveen Tammana
PhD: Univ. of Edinburgh
Postdoc: Princeton University



CA MS teams

- Connect using Microsoft Teams
 - Team name: CA-CS2323
 - Channel: General
- Will post
 - Course material, lecture videos, TA details
 - Assignments, quizzes, etc.
 - More channels to come..
- Be active in the teams group
- Help each other, ask questions



Chapter 1 — Computer Abstractions and Technology — 2

CA initial hiccups

- Overlap with other courses?
 - Please de-register and take the course next year
- Pre-requisites:
 - CS1353: Introduction to Data structure
 - ID1303: Introduction to Programming
 - If not done with both courses, please de-register
- Any trouble accessing teams account?
 - Contact Vijay (vijay.chakry@cse.iith.ac.in)



Reference books

- Mostly: Computer Organization and Design, The hardware/Software Interface 5th Edition (MIPS)
- Authors:
 - David A. Patterson
 - Joh L. Hennessy
- Optional read:
 - Computer Organization and Architecture, by Smriti Ranjan Sarangi



Chapter 1 — Computer Abstractions and Technology — 4

Course Component

- Assignments:
 - Programming assignments
- Periodic quizzes/Tutorials during class hours
- Tentative scoring policy
 - Attendance: 10% (if issues in joining, please drop a mail before the class to the TAs, they will mark it)
 - Coding assignments: 60%, Quizzes/Tutorial: 30%
 - No final exam
- Grading policy: >98 == A+, 92-98% = A, 88-92% = A-, 85-88 = B, 80—85 = B-, 75-80=C, 70-75 = C-, 50-70=D, <50=F, etc). For audit grade: 85% attendance is mandatory



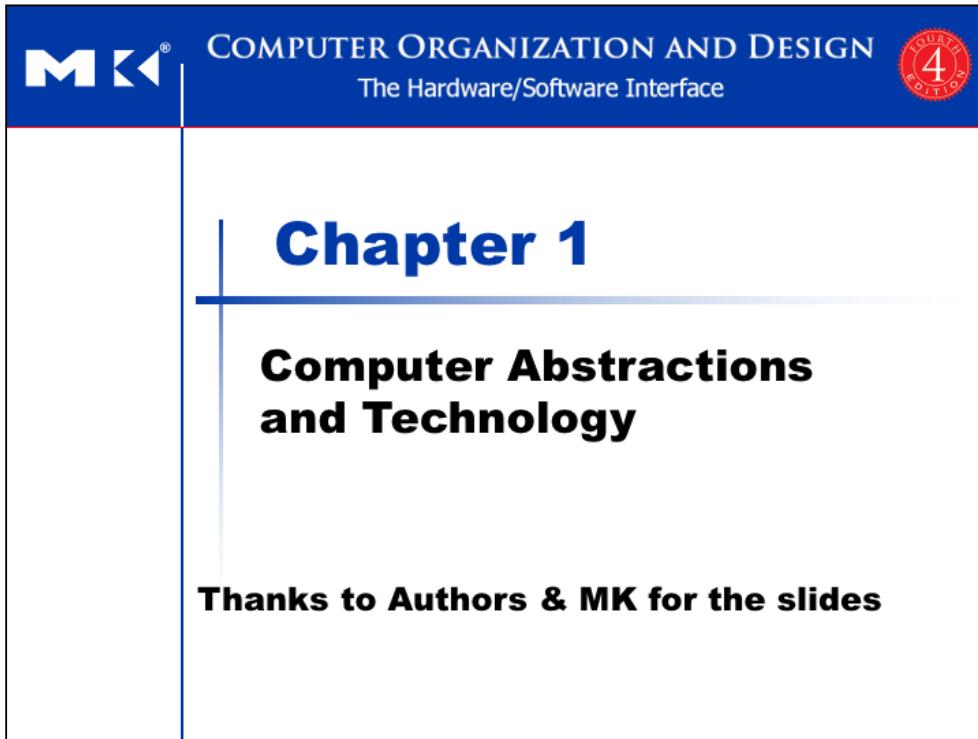
Chapter 1 — Computer Abstractions and Technology — 5

Course topics & tentative timeline

- Sep 1-7: Computer abstraction and technology
- Sep 9-25: Instructions: Language of the computer
 - Sep 17: Assignment 1 start date
- Sep 21-Oct 2: Arithmetic for Computers
- Oct 5-16: The processor
 - Oct 5: Assignment 1 end-date
 - Oct 15: Assignment 2 start-date
- Oct 19-30: Large and Fast: Exploiting Memory Hierarchy
 - Nov 2: Assignment 2 end-date
- Nov 2-11: Parallel processors from Client to Cloud



Chapter 1 — Computer Abstractions and Technology — 6



Chapter 1

**Computer Abstractions
and Technology**

Thanks to Authors & MK for the slides

- Our economy is becoming highly dependent on online services, computers, cloud, etc.
- In this course, students learn about low-level details of computer that made all application built on top of it possible.

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



Chapter 1 — Computer Abstractions and Technology — 8

- Moore's Law refers to Moore's perception that the number of transistors on a microchip doubles every two years, though the cost of computers is halved.
- Each time the cost of computing improves by order of 10, the opportunities multiply.
- Computers reduce pollution, improve fuel efficiency, air bag inflation to protect occupants in a crash
- connects ½ of the population on planet from almost every where
- The cost of computer equipment to map and analyze human DNA sequences was hundreds of millions of dollars, now it is significantly reduced and made once science fiction a reality. *DNA sequencing*, determines the order of nucleotides in a DNA strand, allows scientists to read the genetic code so they can study the normal versions of genes. After they know the order of nucleotides in normal and disease versions of a gene, they can identify which changes in the gene cause the disease.
- Web has replaced libraries and news papers
- Finding relevant information became increasingly important. It is a hardship to go without them
- Clearly, all this technology is possible because of computers we rely on every aspect of our society.

- In this course you will learn the fundamentals of hardware/software of a computer that enable killer applications like glasses with augment reality, cashless society, and cars that can drive themselves is possible.

Classes of Computers

- Desktop computers
 - General purpose, variety of software
 - Subject to cost/performance tradeoff
- Server computers
 - Network based
 - **Super computers:** High capacity, performance, reliability
 - **Cloud computing:** Range from 10's of servers to 100K servers
(E.g., DCs built by Google, Amazon..)
- Embedded computers
 - Hidden as components of systems
 - Stringent power/performance/cost constraints



Intel i3-35K, i5-50K, i7-80K (Rs)



Super computer: >100 crores, 1000's of processors, Terabytes of memory



Chapter 1 — Computer Abstractions and Technology — 9

Although a common set of hardware technologies is used in computers ranging from smart home appliances to cell phones to the largest supercomputers, these different applications have different design requirements and employ the core hardware technologies in different ways.

Desktop computers

- emphasize delivery of good performance to single users at low cost
- usually execute third-party software.
- Intel I3-7 cost in Rs.

Servers are the modern form of what were once mainframes, minicomputers, and supercomputers, and are usually accessed only via a network.

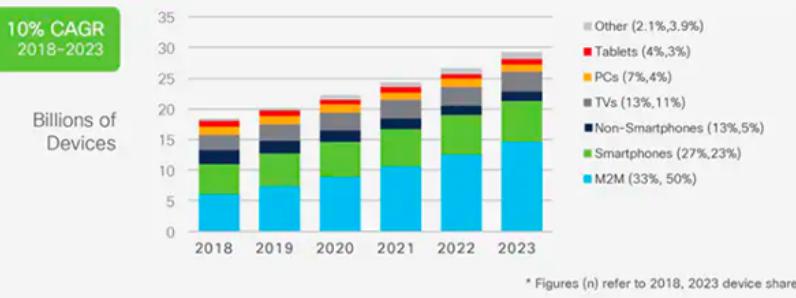
- Servers are oriented to carrying large workloads - a single complex application—usually a scientific or engineering application—or handling many small jobs, such as would occur in building a large Web server.
- In general, servers also place a greater emphasis on dependability, since a crash is usually more costly than it would be on a single-user desktop computer.
- Datacenter (DC) is a farm of servers that could be as large as 300 times cricket field.

Embedded computers: Widest range of apps and performance

- includes microprocessors in cars
- computers in television sets
- Amazon echo, smart bulbs
- processors that control modern airplane or cargoship

- + Personal mobile devices

The Processor Market



Chapter 1 — Computer Abstractions and Technology — 10

- Globally, devices and connections are growing faster (10 percent CAGR) than both the population (1.0 percent CAGR) and the Internet users (6 percent CAGR).
- A growing number of M2M applications, such as smart meters, video surveillance, healthcare monitoring, transportation, and package or asset tracking, are contributing in a major way to the growth of devices and connections.
- By 2023, M2M connections will be half or 50 percent of the total devices and connections.

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
- What is parallel processing



Chapter 1 — Computer Abstractions and Technology — 11

- How are programs written in a high-level lang, such as C or Java translated into the language of the hardware, and the hardware execute those programs?
- What is the interface between s/w and h/w? How s/w instruct hardware?
- What determines a performance & how it can be improved? program -> machine language -> h/w executes the program!
- What techniques h/w designers use to improve performance?
- Reasons for moving from sequential processing to parallel processing?
- Why one has to understand answers to these questions?
 - Without that improving program performance or evaluate what features make one computer better than another will be a complex process of trial and error!

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



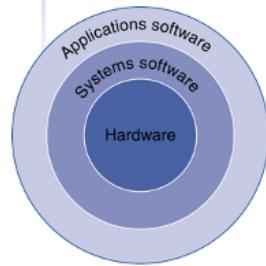
Chapter 1 — Computer Abstractions and Technology — 12

- Algorithm determines #operations
- PL, comp, arch determine #machine instructions per operation
- Processor and mem system determine how fast each instruction can be executed
- I/O system determine how fast I/O operations are executed

Below Your Program

§1.2 Below Your Program

- Application software
 - Written in high-level language
 - Java, Python, C, C++
- System software
 - Compiler: translates HLL code to machine code
 - Ex: GCC
 - Operating System: service code
 - Handling input/output
 - Managing memory and storage
 - Scheduling tasks & sharing resources
 - Ex: Linux, Windows
- Hardware
 - Processor, memory, I/O controllers
 - Ex: Intel i3-7, ARM, DRAM

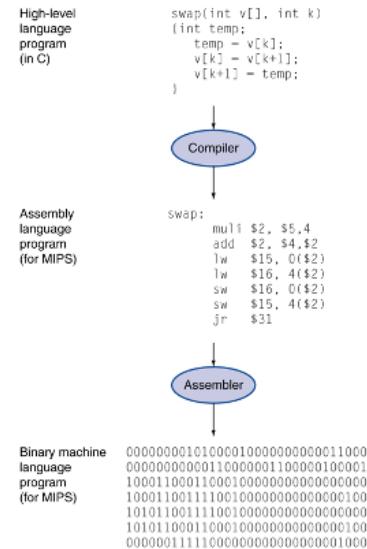


Chapter 1 — Computer Abstractions and Technology — 13

- App software (what app to banking apps to web browser)
- underlying libraries both kernel/OS and user apps - compilers, s/w that talks to h/w, memory management, process/thread scheduling code, etc.
- CPU, DRAM, I/O controllers (PCI)

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



Chapter 1 — Computer Abstractions and Technology — 14

- Computers are slaves to commands
 - However it understands only 1's and 0's – a tedious job for programmers to give instructions in 1's and 0's
 - Like alphabets
 - So we invented assembler that translates textual instructions to binary
 - Like dictionary ☺
 - Assembly language requires programmer to write one line for every instruction, forcing programmer to think like the computer.
 - Although this is better than binary, still far from the human language, say normal English
 - Today creation of HLL like c, java, python, are one of great breakthroughs
 - the creation increased productivity, and
 - make humans think in a more natural language.
 - HLL → compilers → assembly language → assembly → binary code ☺

Components of a Computer

The BIG Picture

Evaluating performance

Compiler
Interface

Computer
Processor
Memory
Control
Datapath
Input
Output

Same components for all kinds of computer

- Desktop, server, embedded

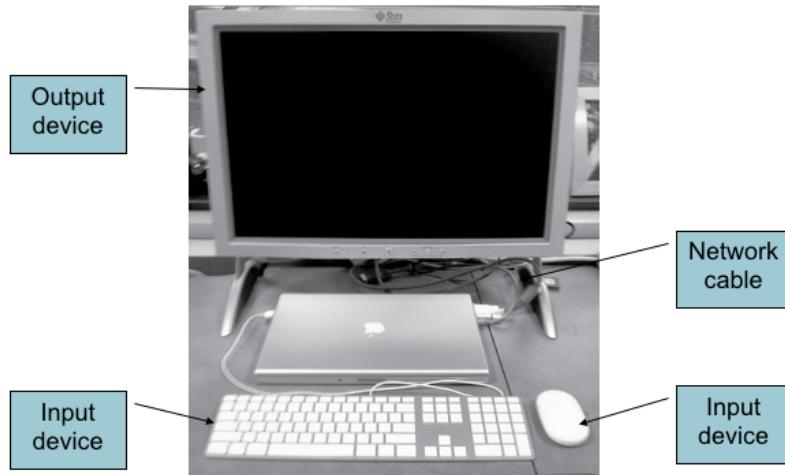
Input/output includes

- User-interface devices
 - Display, keyboard, mouse
- Storage devices
 - Hard disk, CD/DVD, flash
- Network adapters
 - For communicating with other computers

Chapter 1 — Computer Abstractions and Technology — 15

- Computer do mainly four basic operations are input, output, process and store. They are same for all kinds of computers – laptop to server to mobile to smart bulb.
- The processor (CPU) gets instructions from memory (DRAM)
- Input writes data to memory (Mouse, keyboard)
- Output reads data from memory (Harddisk, Monitor, N/w adapters)
- Control sends the signals that determine the operations of the data path, memory, input and output.

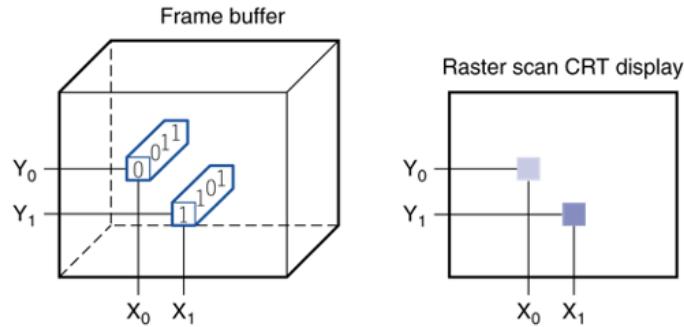
Anatomy of a Computer



Chapter 1 — Computer Abstractions and Technology — 16

Through the Looking Glass

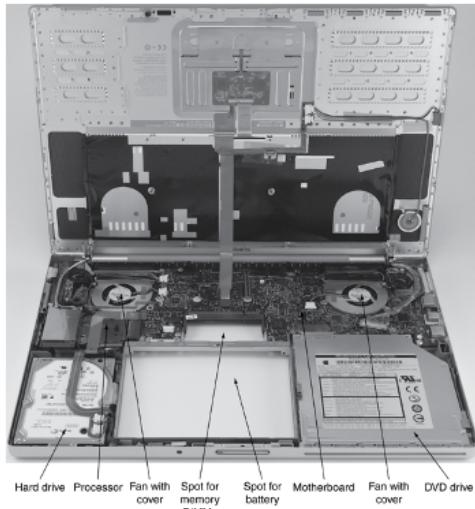
- LCD screen: picture elements (pixels)
 - Mirrors content of frame buffer memory



Chapter 1 — Computer Abstractions and Technology — 17

- 640 * 480, 2560 * 1600 matrix of bits also called pixels
- The computer hardware support for graphics consists mainly of a *frame buffer*, to store the bit map.
- The image to be represented onscreen is stored in the frame buffer, and the bit pattern per pixel is read out to the graphics display at the refresh rate.
 - shows a frame buffer with a simplified design of just 4 bits per pixel.
- Pixel x₀,y₀ contains the bit pattern 0011, which is a lighter shade on the screen than the bit pattern 1101 in pixel x₁,y₁

Opening the Box



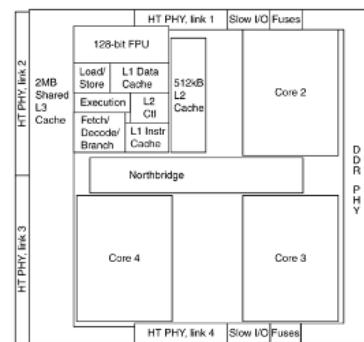
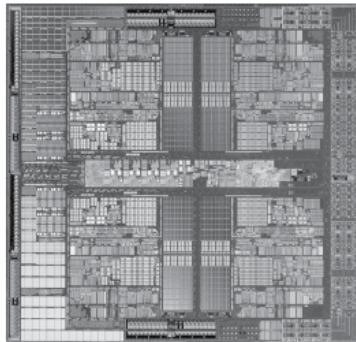
Chapter 1 — Computer Abstractions and Technology — 18

- The shiny box with the white label on the lower left is a 100 GB SATA hard disk drive,
- and the shiny metal box on the lower right side is the DVD drive.
- The hole between them is where the laptop battery would be located.
- The small hole above the battery hole is for memory DIMMs. Right figure is a closeup of the DIMMs, which are inserted from the bottom in this laptop.
- Above the battery hole and DVD drive is a printed circuit board (PC board), called the *motherboard*, which contains most of the electronics of the computer. The board is composed of three pieces: the piece connecting to the I/O devices mentioned earlier, the memory, and the processor.

- The two shiny circles in the upper half of the picture are two fans with covers.
- The processor is the large raised rectangle just below the left fan.

Inside the Processor

■ AMD Barcelona: 4 processor cores

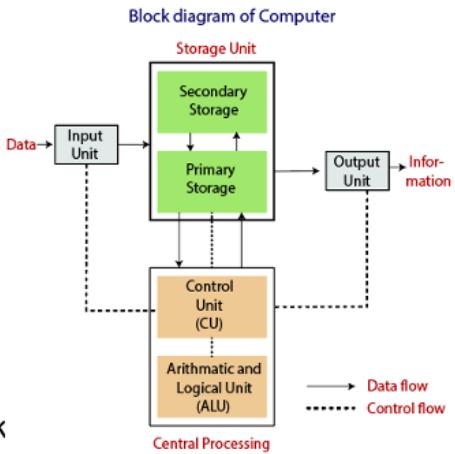


Chapter 1 — Computer Abstractions and Technology — 19

- The lefthand side is a microphotograph of the AMD Barcelona processor chip, and the righthand side shows the major blocks in the processor.
- This chip has four processors or “cores”. The microprocessor in the laptop has two cores per chip, called an Intel Core 2 Duo.

Block diagram of computer

- Datapath: performs operations on data
- Control: sequences datapath, memory, ...
- Cache memory
 - Small fast SRAM memory for immediate access to data
- Primary storage – DRAM
 - Large
- Secondary storage – Hard disk



Chapter 1 — Computer Abstractions and Technology — 20

- The **datapath** performs the arithmetic operations, and **control** tells the datapath, memory, and I/O devices what to do according to the wishes of the instructions of the HLL program.
- **Cache memory** consists of a small, fast memory that acts as a buffer for the DRAM memory.
- Cache is built using a different memory technology, **static random access memory (SRAM)**. SRAM is faster but less dense, and hence more expensive, than DRAM (more in next lectures)

A Safe Place for Data

- Volatile main/primary memory
 - Loses instructions and data when power off
 - Cache (SRAM), DRAM
- Non-volatile secondary memory
 - Magnetic disk (also hard disk)
 - Flash memory (also SSD)
 - Optical disk (CDROM, DVD)



Chapter 1 — Computer Abstractions and Technology — 21

- volatile memory/main memory used to hold data and programs while they are running and
- this nonvolatile memory/ secondary memory used to store data and programs between runs
- **Flash memory**, a nonvolatile semiconductor tor memory, is used instead of disks in mobile devices such as cell phones and is increasingly replacing disks in music players and even laptops.
- Access times for magnetic disks are much slower than for DRAMs: disks typically take 5–20 milliseconds, while DRAMs take 50–70 nanoseconds—making DRAMs about 100,000 times faster.
- Thus, there are three primary differences between magnetic disks and main memory: disks are nonvolatile because they are magnetic; they have a slower access time because they are mechanical devices; and they are cheaper per gigabyte because they have very high storage capacity at a modest cost.
- Many have tried to invent a technology cheaper than DRAM but faster than disk to fill that gap, but many have failed.
- Flash memory, however, is a serious challenger. This semiconductor memory is nonvolatile like disks and has about the same bandwidth, but latency is 100 to 1000 times faster than disk.

Networks

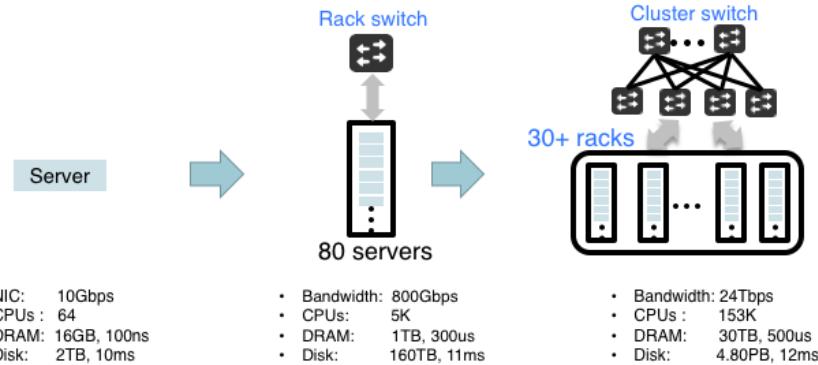
- Communication and resource sharing
- Local area network (LAN): Ethernet
 - Within a building
- Wide area network (WAN: the Internet)
- Wireless network: WiFi, Bluetooth



Chapter 1 — Computer Abstractions and Technology — 22

Networks enable compute and storage at massive scales

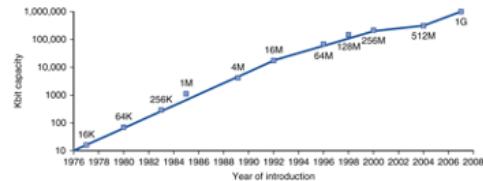
The network has thousands of switches and millions of links
(Datacenters built by Google, facebook, Amazon, MS..)



Reference: Jeff Dean. Designs, Lessons and Advice from
Building Large Distributed Systems.

Technology Trends

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



DRAM capacity

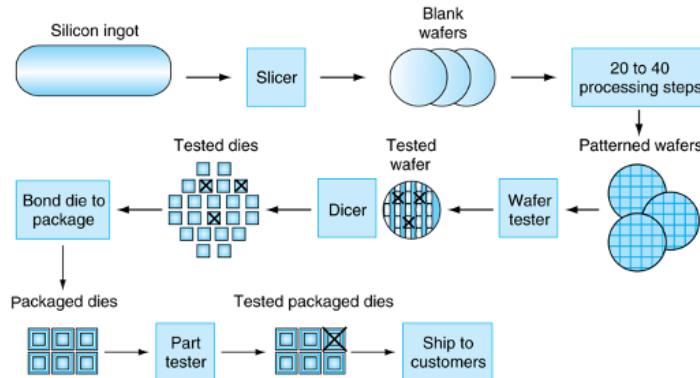
Year	Technology	Relative performance/cost
1951	Vacuum tube	1
1965	Transistor	35
1975	Integrated circuit (IC)	900
1995	Very large scale IC (VLSI)	2,400,000
2005	Ultra large scale IC	6,200,000,000



Chapter 1 — Computer Abstractions and Technology — 24

- Moore's law -- that the number of transistors on a microchip doubles every two years, though the cost of computers is halved.
- This table estimates relative performance per unit cost for each technology
 - A transistor is simply an on/off switch controlled by electricity
 - The IC combine dozens to hundreds of transistors into a single chip
 - VLSI : hundreds to millions
- The DRAM shows the growth in DRAM capacity since 1977 - quadrupled capacity every 3 year

Manufacturing ICs



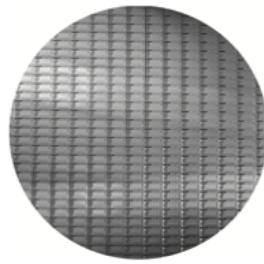
- Yield: proportion of working dies per wafer



- Let's understand how ICs are manufactured
- Silicon is a substance found in sand, because it does not conduct electricity well, it is called a semiconductor
- With special chemical process and adding some materials to silicon, tiny areas transform into one of three devices
 - Excellent conductors of electricity (using either microscopic copper or aluminum wire)
 - Excellent insulators from electricity (like plastic sheathing or glass)
 - Areas that can conduct or insulate under special conditions (as a switch)
- Transistors fall in the last category. VLSI is just billions of combinations of conductors, insulators, and switches manufactured in a small packaging
- Processing:
 - Silicon Ingots: 8-12 inches in diameter and 12-24 long
 - Sliced into wafers
 - Each wafer goes through a series of chemical processing steps, during which chemical patterns create transistor, insulator, and conductors
 - A single flaw in the wafer or in one of the dozens processing steps can result in that area of the wafer failing

- It is impossible to make process perfect, so a best way to work around imperfection is to place many components on a single wafer
- The partitioned wafer is diced into components called dies or informally called chips
- Dicing enables to discard dies, instead of whole wafer
- yield of a process: % of good dies from total number of dies
- Once you've found good dies, they are connected to the input/output pins of a package, using a process called *bonding*.
- These packaged parts are tested a final time, since mistakes can occur in packaging, and then they are shipped to customers.

Intel core i7



- 300mm (12 inch) wafer
- 100% yields = 280 chips, 32nm technology



Chapter 1 — Computer Abstractions and Technology — 26

- The number of dies on this 300 mm (12 inch) wafer at 100% yield is 280 chips/dies, each 20.7 by 10.5 mm.
- This die uses a 32-nanometer technology, which means that the smallest features are approximately 32 nm in size (e.g., transistor size)

Integrated Circuit Cost

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

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

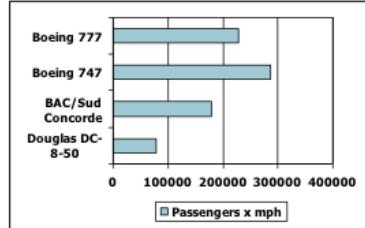
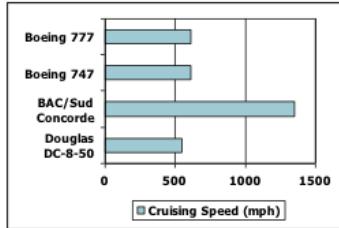
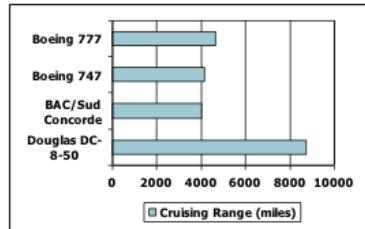
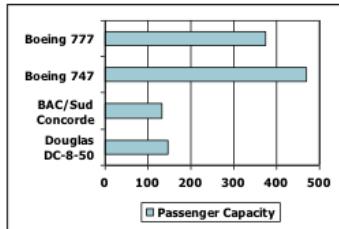


Chapter 1 — Computer Abstractions and Technology — 27

- The cost of an integrated circuit rises quickly as the die size increases, due both to the lower yield and the smaller number of dies that fit on a wafer.
- The first equation is straightforward to derive.
- The second is an approximation, since it does not subtract the area near the border of the round wafer that cannot accommodate the rectangular dies
- The final equation is based on empirical observations of yields at integrated circuit factories, with the exponent related to the number of critical processing steps.

Defining Performance

- Which airplane has the best performance?



Chapter 1 — Computer Abstractions and Technology — 28

How do you define performance?

- speed? highest cruising speed → concorde
- range? longest range → DC-8
- Capacity? Largest capacity → Boeing 747
- Which is faster?:
 - For one passenger? → concorde is the winner
 - For many passengers? → boeing 747 is the winner (#passenger x mph)
- Running a same program on two diff desktop comp? which is faster?
 - Desktop computer that gets job done first
- Running a datacenter with several serveres running jobs submitted by webusers
 - fast datacenter is the one which complets most jobs during a day

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...



Chapter 1 — Computer Abstractions and Technology — 29

- As an individual computer user, you are interested in reducing **response time**—the time between the start and completion of a task—also referred to as execution time
- Datacenter managers are often interested in increasing **throughput** or **bandwidth**—the total amount of work done in a given time.
- Hence, in most cases, we will need different performance metrics to benchmark personal mobile devices, which are more focused on response time, versus servers, which are more focused on throughput.
- How are time and throughput affected by
 - replacing faster processor
 - adding more processor
- Ans: Decreasing response time almost always improves throughput. Hence, in case 1, both response time and throughput are improved. In case 2, no one task gets work done faster, so only throughput increases.
- when it comes In comp perf, we primarily focus on response time..

Relative Performance

- Define Performance = $1/\text{Execution Time}$
- “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
 - $\text{Execution Time}_B / \text{Execution Time}_A$
 $= 15s / 10s = 1.5$
 - So A is 1.5 times faster than B



- To maximize performance, we want to minimize response time or execution time for some task.
- So performance is inverse of execution time
- To compare performance of two computers, e.g., how many times X is fast compared to Y, to get this we take the ratio: perf-x/perf-y or exec-time-y/exec-time-x
- See example...

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

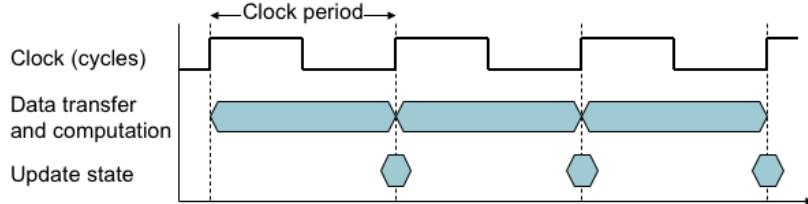


Chapter 1 — Computer Abstractions and Technology — 31

- Program execution time is measured in seconds per program.
- Response time or elapsed time includes all aspects – CPU processing, disk accesses, memory accessss, I/O delays, OS overhead
- We want to distinguish between the elapsed time and the time over which the processor is working on our behalf, called CPU time (eg., no waiting time for I/O or running other programs)
- CPU time can be further divided into CPU time spent on program, called user CPU time and the CPU time spend in the OS performing tasks on behalf of the program, called system CPU time
- CPU performance: refer to user CPU time
- System performance: refer to elapsed time on an unloaded system
- We will focus on CPU performance in this chapter

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., $4.0\text{GHz} = 4000\text{MHz} = 4.0 \times 10^9\text{Hz}$



Chapter 1 — Computer Abstractions and Technology — 32

- Although time is what we care about, computer designers use a measure that relates to how fast the hardware can perform basic functions.
- All computers are constructed using a clock that determines when events take place in the hardware
- These discrete time intervals are called clock cycles or clock periods or clock ticks- the time for one clock period, usually of the processor clock, which runs at a constant time
- Clock period: length of each clock cycle
- Clock frequency (rate): #cycles per second – inverse of clock period

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



- Users and designers often examine performance using different metrics. If we could relate these different metrics, we could determine the effect of a design change on the performance as experienced by the user.
- CPU performance = CPU execution time → relates to basic metrics – clock cycles and cycle time
- A hardware designer can improve performance either by reducing the #clock cycles required for a program or length of the clock cycle (cycle time)
- Later chapters discuss the tradeoff between cycle rate vs 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 $1.2 \times$ clock cycles
- How fast must Computer B clock be?

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

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

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



Chapter 1 — Computer Abstractions and Technology — 34

Our favorite program runs in 10 seconds on computer A, which has a 2 GHz clock. We are trying to help a computer designer build a computer, B, which will run this program in 6 seconds. The designer has determined that a substantial increase in the clock rate is possible, but this increase will affect the rest of the CPU design, causing computer B to require 1.2 times as many clock cycles as computer A for this program. What clock rate should we tell the designer to target?

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



Chapter 1 — Computer Abstractions and Technology — 35

- We haven't discussed how many instructions needed for the program
- Compiler generates instructions to execute, and the computer had to execute the instructions (C-program \rightarrow Assembly instructions \rightarrow Instructions in byte code)
- So the execution time depends on the #instructions in a program
- CPU clock cycles = instruction count \times average clock cycles per instruction (CPI)
- Each instruction may take different number of clock cycles -- Same instruction take different time on different hardware implementations, but the instruction count could be the same.
- therefore, CPU time = instruction count \times CPI \times clock cycle time or $(\text{instruction count} \times \text{CPI}) / \text{clock rate}$
- These formulas are particularly useful because they separate the three key factors that affect performance. We can use these formulas to compare two different implementations or to evaluate a design alternative if we know its impact on these three parameters.

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?

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

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

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



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
 $= 2 \times 1 + 1 \times 2 + 2 \times 3$
 $= 10$
- 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

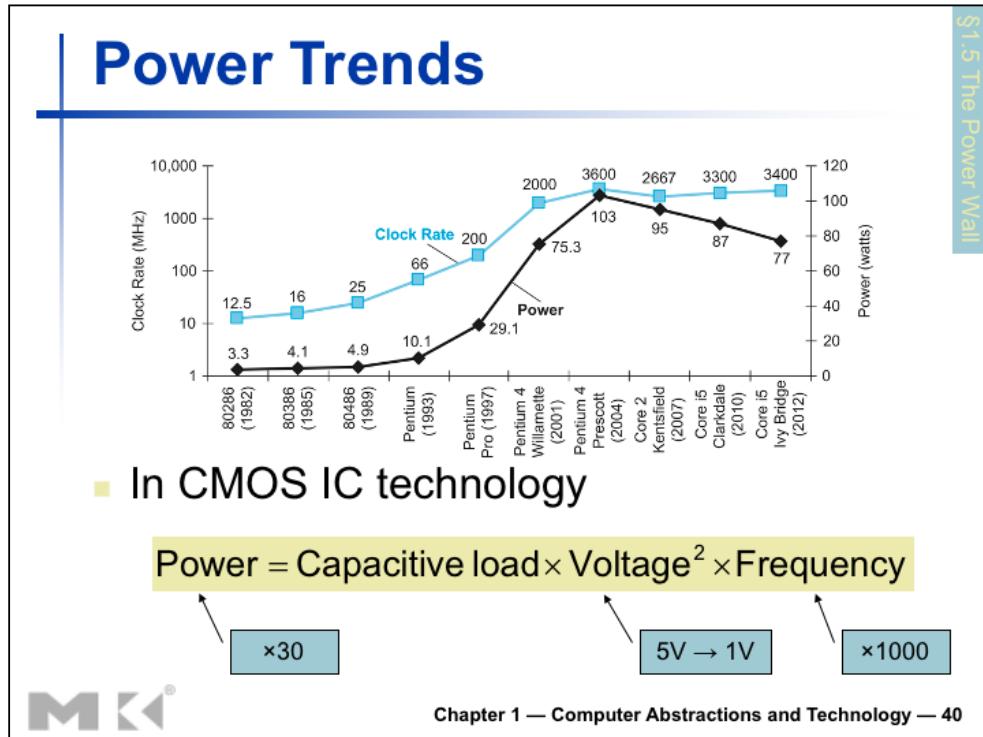
$$\text{CPU Time} = \frac{\text{Instructions}}{\text{Program}} \times \frac{\text{Clock cycles}}{\text{Instruction}} \times \frac{\text{Seconds}}{\text{Clock cycle}}$$

- 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



Chapter 1 — Computer Abstractions and Technology — 39

- The algorithm determines the number of source program instructions executed and hence the number of processor instructions executed. The algorithm may also affect the CPI, by favoring slower or faster instructions. For example, if the algorithm uses more divides, it will tend to have a higher CPI.
- The programming language certainly affects the instruction count, since statements in the language are translated to processor instructions, which determine instruction count. The language may also affect the CPI because of its features; for example, a language with heavy support for data abstraction (e.g., Java) will require indirect calls, which will use higher CPI instructions.
- The efficiency of the compiler affects both the instruction count and average cycles per instruction, since the compiler determines the translation of the source language instructions into computer instructions. The compiler's role can be very complex and affect the CPI in complex ways.
- The instruction set architecture affects all three aspects of CPU performance, since it affects the instructions needed for a function, the cost in cycles of each instruction, and the overall clock rate of the processor.



- Shows increase in clock rate and power of eight generations of intel processors
 - The Pentium 4 made a dramatic jump in clock rate and power but less so in performance. The Prescott thermal problems led to the abandonment of the Pentium 4 line. The Core 2 line reverts to a simpler pipeline with lower clock rates and multiple processors per chip. The Core i5 pipelines follow in its footsteps.
- Power and clock rate are correlated
- The reason for slowing down is we ran into practical power limit for cooling commodity processors
 - Battery life can trump performance in the personal mobile devices (laptop, mobiles)
 - Architects of servers, try to reduce costs of powering and cooling 100k servers as the costs are high at this scale
- Like CPU time for measuring performance, joules (joules/sec) is a better measure than a power rate like watts
- CMOS is the dominant technology for Integrated circuits in processor
- For CMOS, primary source of energy consumption is so-called dynamic energy – which is energy consumed when transistors switch states from 0 to 1 and viceversa

- Energy is proportional to capacitive load x voltage[^]
- Power required per transistor is: power is proportional to capacitive load x voltage^{^2} x frequency switched
 - Frequency switched is a function of clock rate
 - Capacitive load per transistor is a function of both #transistors connected to output (fanout) and the technology
- How clock rate grow by a factor of 1000 while power grew by only a factor of 30?
 - Energy and thus power can be reduced by lowering the voltage, possible with new generation of technology
 - Since power is a function of the voltage squared, it worked
 - Voltage was reduced about 15% per generation, In 20 years, voltages have gone from 5V to 1V, power increase is only 30 times

Reducing Power

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

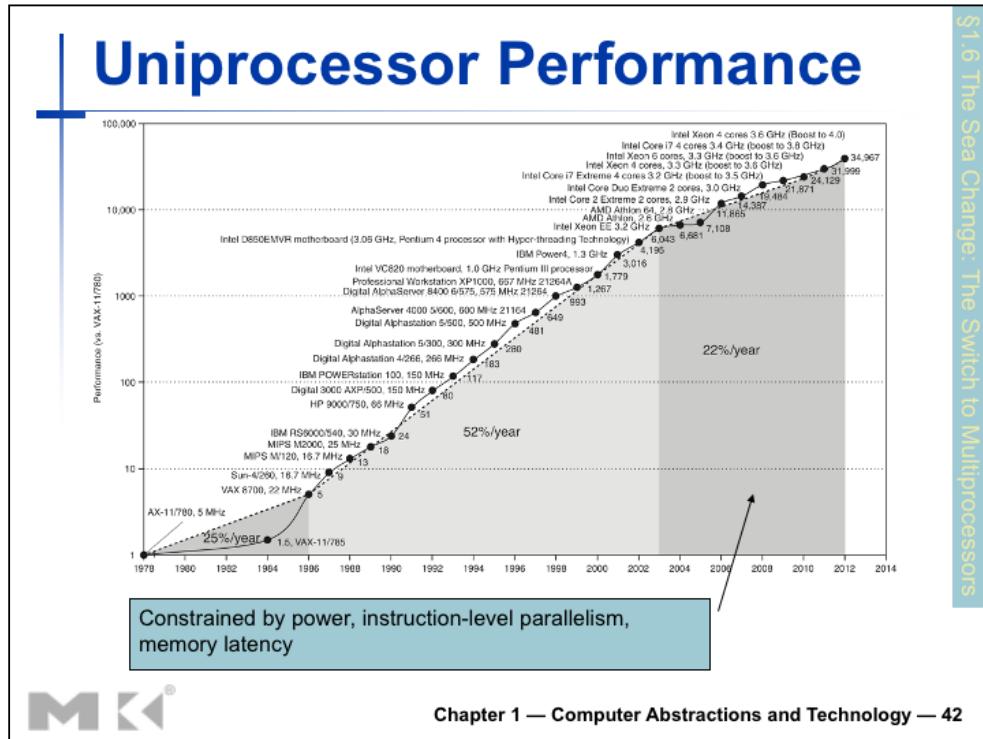
$$\frac{P_{\text{new}}}{P_{\text{old}}} = \frac{C_{\text{old}} \times 0.85 \times (V_{\text{old}} \times 0.85)^2 \times F_{\text{old}} \times 0.85}{C_{\text{old}} \times V_{\text{old}}^2 \times F_{\text{old}}} = 0.85^4 = 0.52$$

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



Chapter 1 — Computer Abstractions and Technology — 41

- Suppose we developed a new, simpler processor that has 85% of the capacitive load of the more complex older processor. Further, assume that it has adjustable voltage so that it can reduce voltage 15% compared to processor B, which results in a 15% shrink in frequency. What is the impact on dynamic power?
- Hence, the new processor uses about half the power of the old processor



- Prior to the mid-1980s, processor performance growth was largely technology-driven and averaged about 25% per year.
- The increase in growth to about 52% since then is attributable to more advanced architectural and organizational ideas. The higher annual performance improvement of 52% since the mid-1980s meant performance was about a factor of seven higher in 2002 than it would have been had it stayed at 25%.
- Since 2002, the limits of power, available instruction-level parallelism, and long memory latency have slowed uniprocessor performance recently, to about 22% per year.

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 — 43

- - In the past, programmers could rely on innovations in hardware, architecture, and compilers to double performance of their programs every 18 months without having to change a line of code. Today, for programmers to get significant improvement in response time, they need to rewrite their programs to take advantage of multiple processors.
- a “quadcore” microprocessor is a chip that contains four processors or four cores.
- Two types of parallelism:
 - Instruction-level parallelism already used heavily in uniprocessors (no changes are required to programs, hidden from the programmer)
 - Multicore/multiprocessor parallelism (changes required to programs to improve performance, programmers should rewrite their programs for multiprocessors)
- Parallel programming, hard because
 - increases difficulty of programming, divide application so each processor gets same amount of work to do, overhead of scheduling and coordination
-

Intel core i7 2.66GHz performance

Description	Name	Instruction Count $\times 10^9$	CPI	Clock cycle time (seconds $\times 10^{-9}$)	Execution Time (seconds)	Reference Time (seconds)	SPECratio
Interpreted string processing	perl	2252	0.60	0.376	508	9770	19.2
Block-sorting compression	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



- Benchmark – A program selected for use in comparing computer performance
- SPEC: System performance evaluation cooperative – created standard sets of benchmarks for modern computer systems
- Table: Execution time is the product of the three factors in this table: instruction count in billions, clocks per instruction (CPI), and clock cycle time in nanoseconds. SPECratio is simply the reference time, which is supplied by SPEC, divided by the measured execution time. The single number quoted as SPECINTC2006 is the geometric mean of the SPECratios.

Concluding Remarks

- Cost/performance is improving
 - Due to underlying technology development
- Execution time: the best performance measure ($IC \times CPI \times Clock\ period$)
- Power is a limiting factor
 - Use parallelism to improve performance



Chapter 1 — Computer Abstractions and Technology — 45

- it's a safe bet that cost/performance will be much better than they are today.
- using the execution time of real programs as the metric is a reliable method of determining and reporting performance
- Conserving power while trying to increase performance has forced the hardware industry to switch to multicore microprocessors, thereby forcing the software industry to switch to programming parallel hardware. **Parallelism** is now required for performance.