

Topics

- Computer Performance in light of
 - Power Wall
 - Complexity Wall
 - Moore's Law
 - Dennard Scaling
- Computing ideas in post-PC era
- What is meant by the Hardware / Software Interface
- Calculate Power dissipation
- Notion of Computer Performance



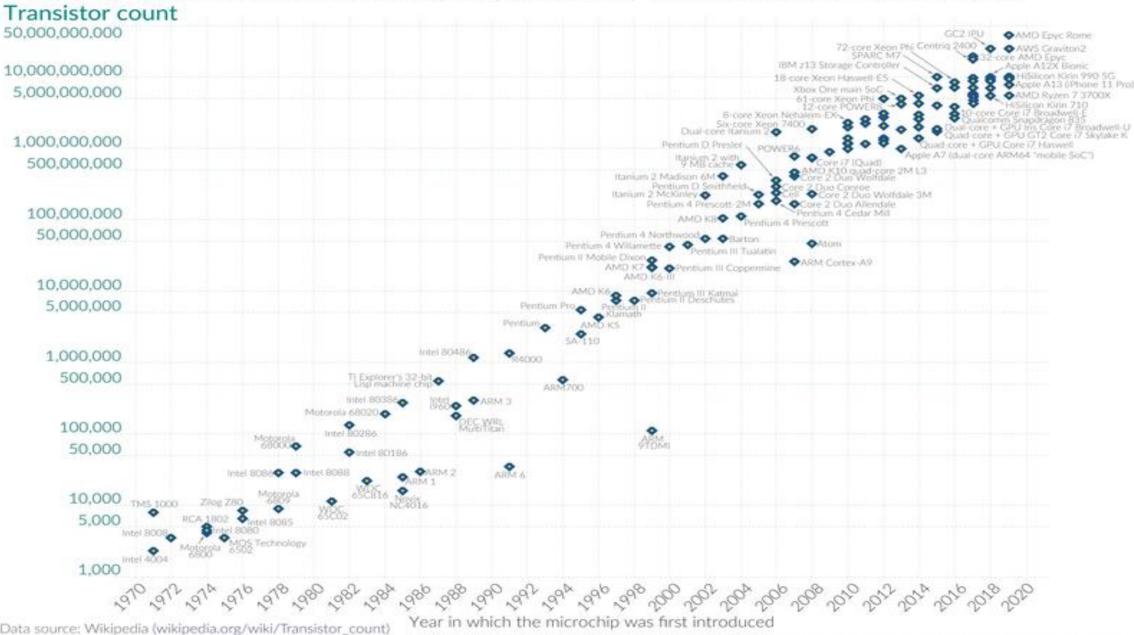
Power Wall, and Complexity Wall

Moore's Law: The number of transistors on microchips doubles every two years Our World



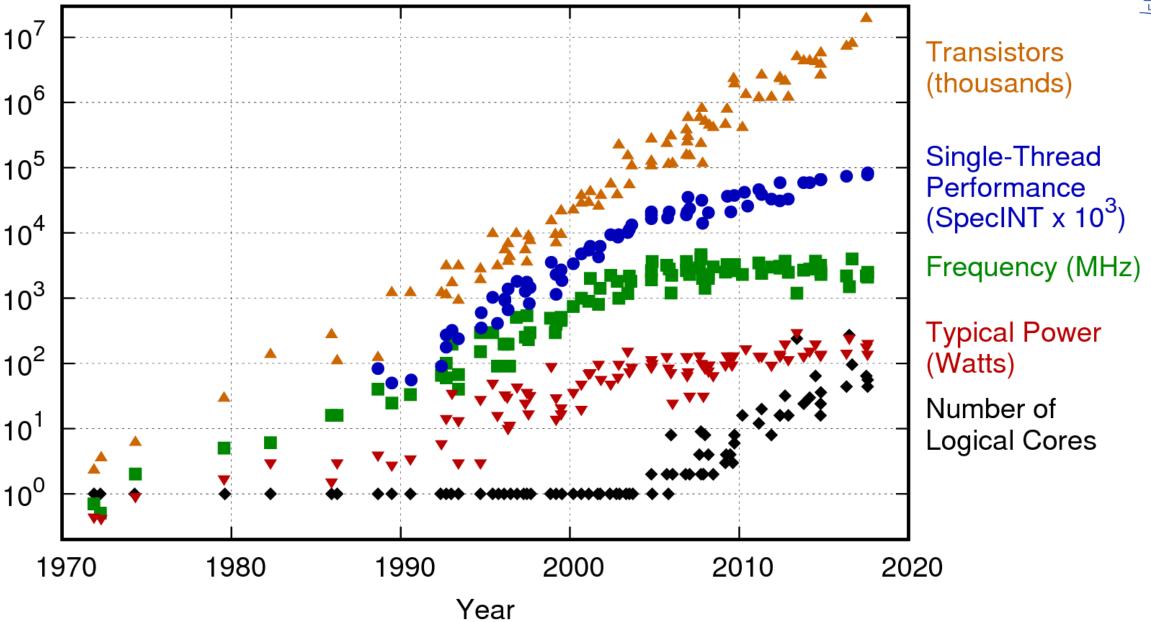
Moore's law describes the empirical regularity that the number of transistors on integrated circuits doubles approximately every two years.

This advancement is important for other aspects of technological progress in computing – such as processing speed or the price of computers.

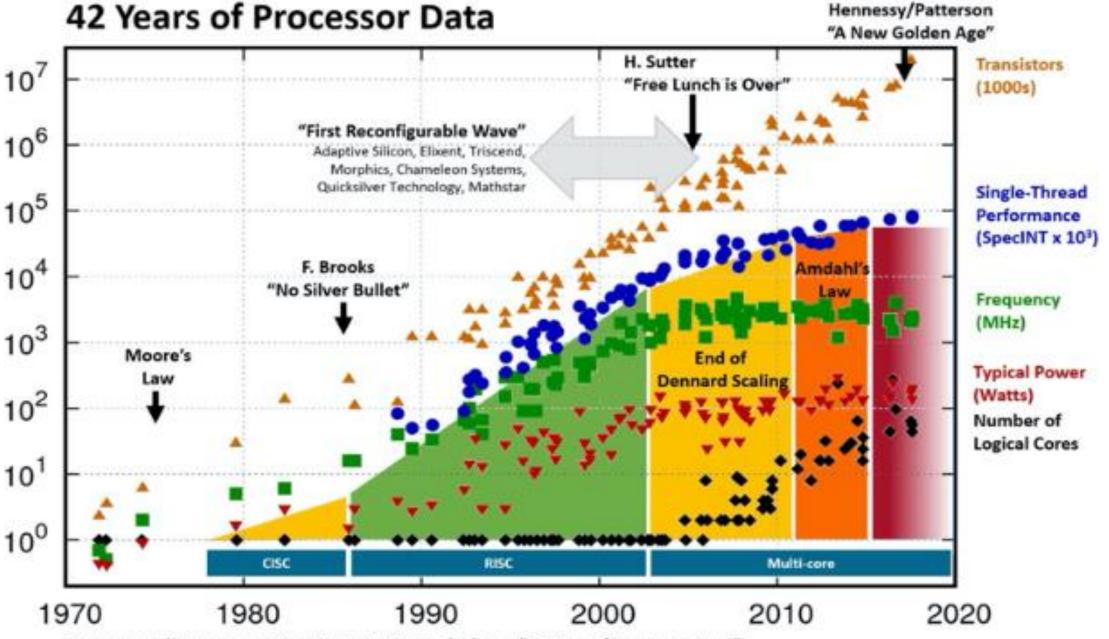


42 Years of Microprocessor Trend Data





Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten New plot and data collected for 2010-2017 by Mr. Original data col

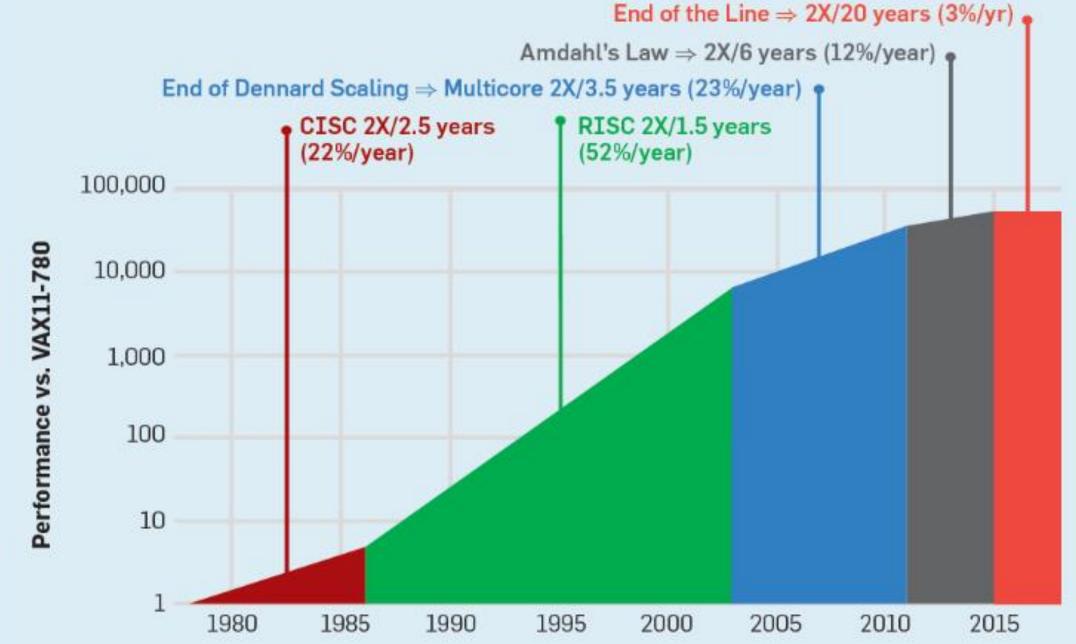


Hennessy/Patterson

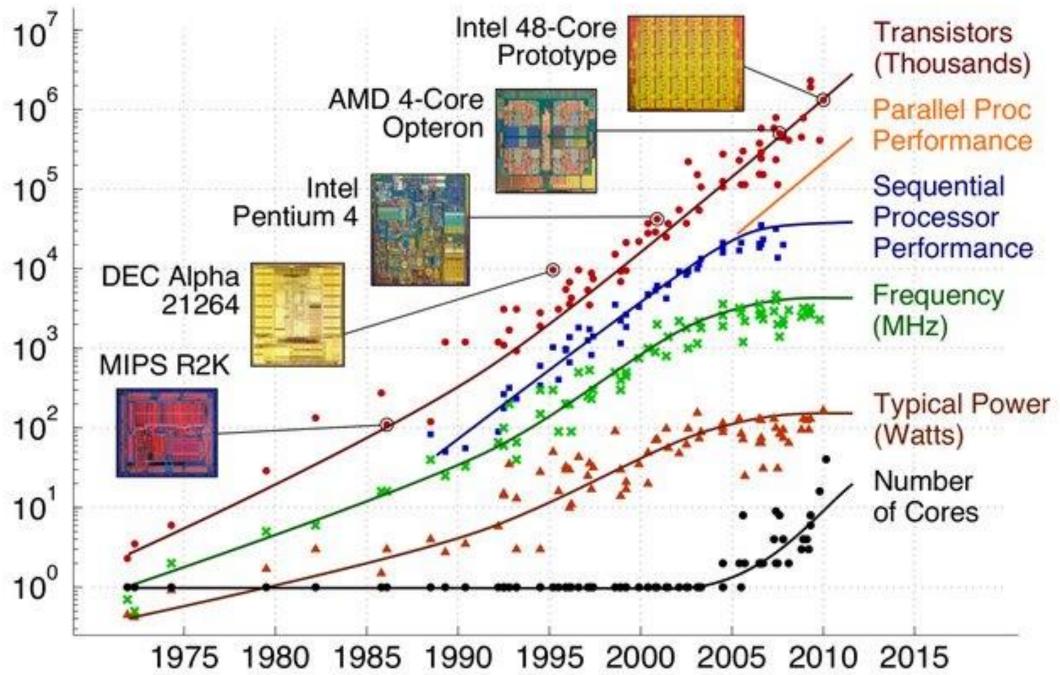
Hennessy and Patterson, Turing Lecture 2018, overlaid over "42 Years of Processors Data" https://www.karlrupp.net/2018/02/42-years-of-microprocessor-trend-data/; "First Wave" added by Les Wilson, Frank Schirrmeister Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batter New plot and data collected for 201 Computer Organization and Assembly Language Lecture 2 Spring 2025

LUMS









CPU Architecture Today Heat becoming un unmanageable problem



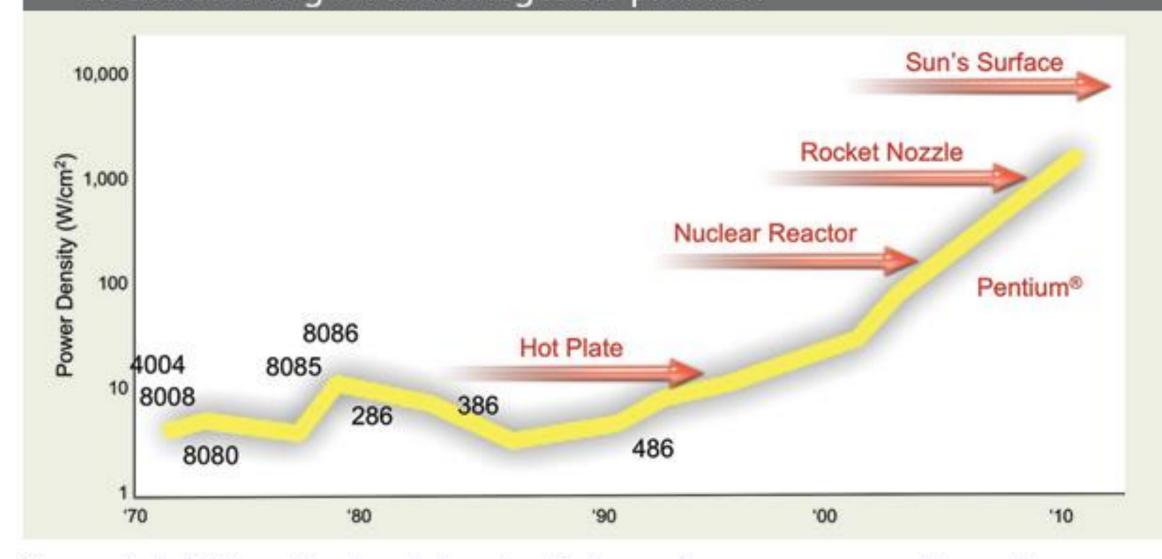


Figure 1. In CPU architecture today, heat is becoming an unmanageable problem. (Courtesy of Pat Gelsinger, Intel Developer Forum, Spring 2004)

Computer Organization and Assembly Language Lecture 2 Spring 2025

CPU: From GHz to multi-core



Moore's Law:

- the number of transistors on an IC doubles every two years.
 - Less space, more complexity.
 - Shorter gates, higher clock rate.

Strategy of the 80s and 90's:

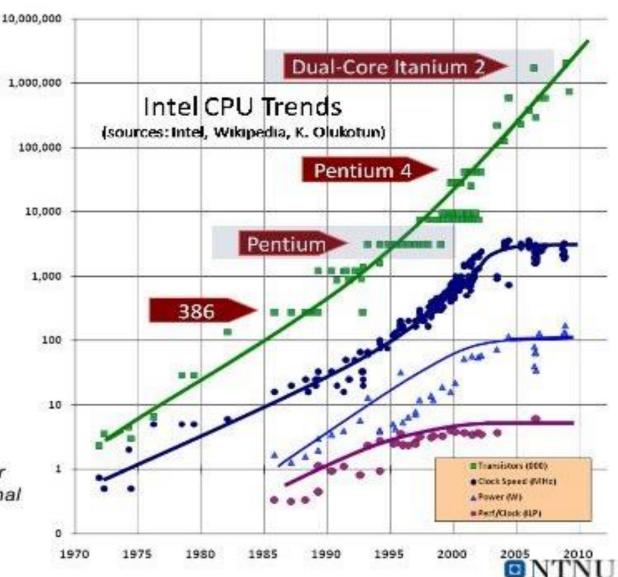
- · Add more complexity!
- Increase the clock rate!

Pollack's Rule:

 The performance increase is ~ square root of the increased complexity. [Borkar 2007]

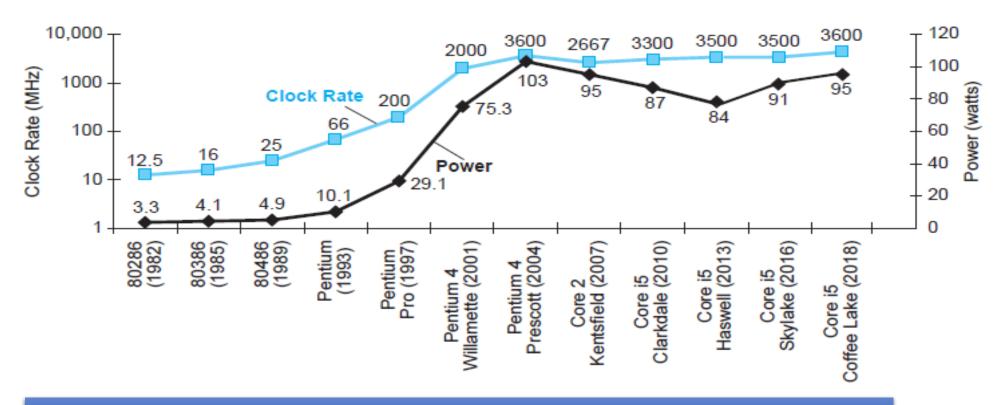
The Power Wall:

 Increasing clock rate and transistor current leakage lead to excess power consumption, while RC delays in signal transmission grow as feature sizes shrink. [Borkar et al. 2005]



CPU Power Trends





In CMOS IC technology

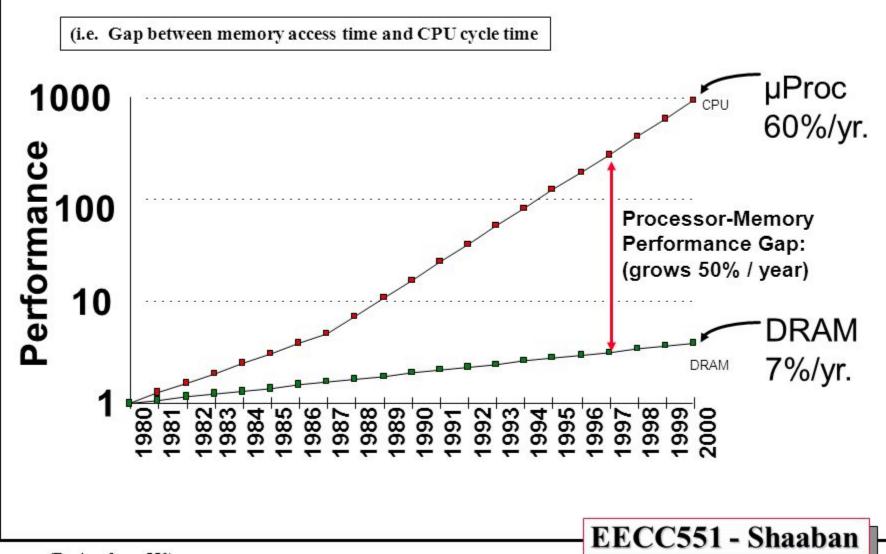
Power = Capacitive load × Voltage ² × Frequency

x30

x1000







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?

Computing Ideas in Post-PC Era



Current High-Performance Computing Architectures Post-Moore's-Law Era

Model Based Abstract Software Development for Complex Systems

Parallel Processing through Multicores

Instruction Level Parallelism through Pipelining

Speculation and Prediction in RISC Architecture

Performance Evaluation to make the 'Common Case Fast'

Use variety of memory hierarchies to improve performance

Cloud based Computing and Storage

Role of Network in Computer Performance



Upcoming Computing Architectures In Post-PC Era

- Mobile and Embedded Computing
- Domain Specific Architectures E.g. Google Tensor Processor
- Domain Specific Languages and Compilers E.g. Tensorflow
- Open Instruction Set Architectures RISC V
- Agile Computing through Reconfigurable Hardware and Customized Instructions
- Security in CPU Hardware

The Scale In Computing





Decimal term	Abbreviation	Value	Binary term	Abbreviation	Value	% Larger
COTTIL	Abbioviation	value				
kilobyte	KB	10 ³	kibibyte	KiB	210	2%
megabyte	MB	10 ⁶	mebibyte	MiB	220	5%
gigabyte	GB	10 ⁹	gibibyte	GiB	230	7%
terabyte	TB	1012	tebibyte	TiB	240	10%
petabyte	PB	1015	pebibyte	PiB	250	13%
exabyte	EB	1018	exbibyte	EiB	2 ⁶⁰	15%
zettabyte	ZB	1021	zebibyte	ZiB	270	18%
yottabyte	YB	1024	yobibyte	YiB	280	21%

FIGURE 1.1 The 2^x vs. 10^y bytes ambiguity was resolved by adding a binary notation for all the common size terms. In the last column we note how much larger the binary term is than its corresponding decimal term, which is compounded as we head down the chart. These prefixes work for bits as well as bytes, so *gigabit* (Gb) is 10⁹ bits while *gibibits* (Gib) is 2³⁰ bits.



The Hardware / Software Interface

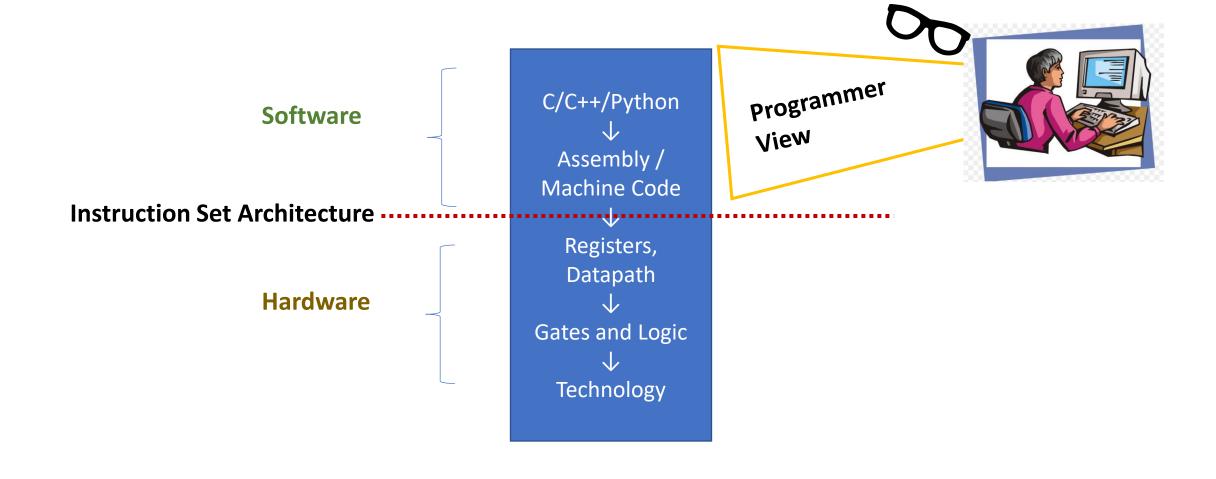
Interface of Hardware and Software



Hardware or software component	How this component affects performance		
Algorithm	Determines both the number of source-level statements and the number of I/O operations executed		
Programming language, compiler, and architecture	Determines the number of computer instructions for each source-level statement		
Processor and memory system	Determines how fast instructions can be executed		
I/O system (hardware and operating system)	Determines how fast I/O operations may be executed		

Hardware / Software Interface





Why do we need Abstraction?

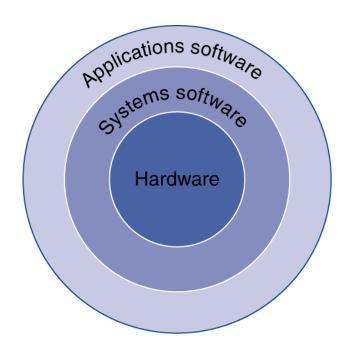


Delving into the depths reveals more information

 An abstraction omits unneeded detail, helps us cope with complexity

Underneath a Computer Program

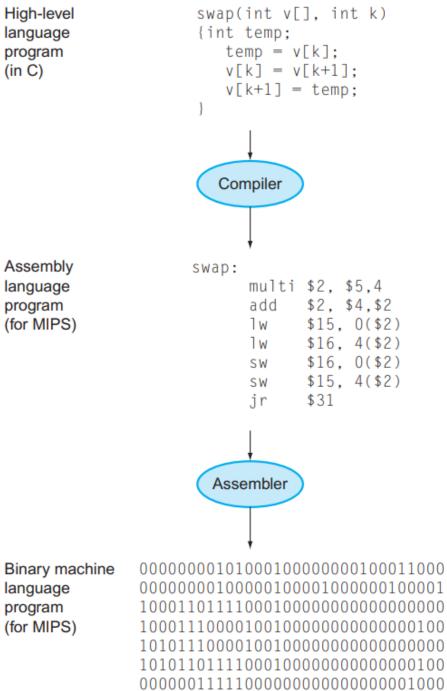




- Application software
 - Written in high-level language
- System software
 - Compiler: translates HLL code to machine code
 - Operating System: service code
 - Handling input/output
 - Managing memory and storage
 - Scheduling tasks & sharing resources
- Hardware
 - Processor, memory, I/O controllers

Compiling a Program

Example of Software Abstraction





What is an Instruction Set Architecture?

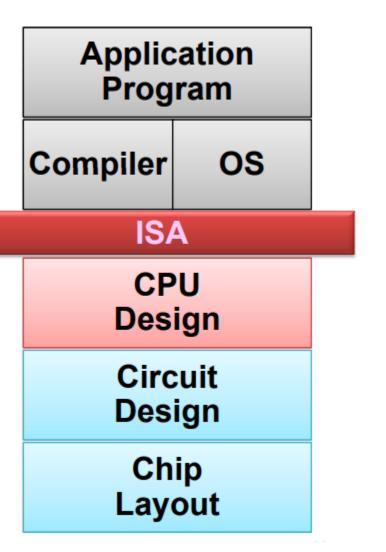


- Lowest level of Computer Architecture that is visible to a programmer
- ISA comprises instructions in Assembly Language and Corresponding Binary Machine Code
- It is the language of the CPU machine
- Primitive syntax compared to a high-level language
- It is easily interpreted and understood by the CPU
- Designed to maximize performance
- Designed to minimize cost and design time

Role of ISA within a Computer



- Assembly Language View
 - Processor state (RF, mem)
 - Instruction set and encoding
- Layer of Abstraction
 - Above: how to program machine HLL, OS
 - Below: what needs to be built
 - tricks to make it run fast

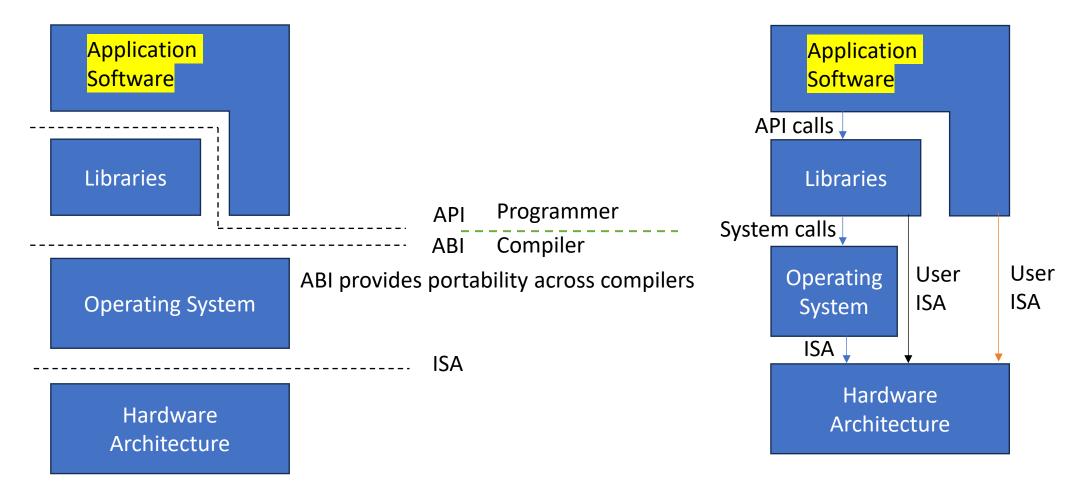


Key:

RF = Register File
Mem = Memory
HLL = High Level Language
OS = Operating System

Operating System Interface Software / Hardware





API = Application Programming Interface – Through which software interfaces with another software at source level

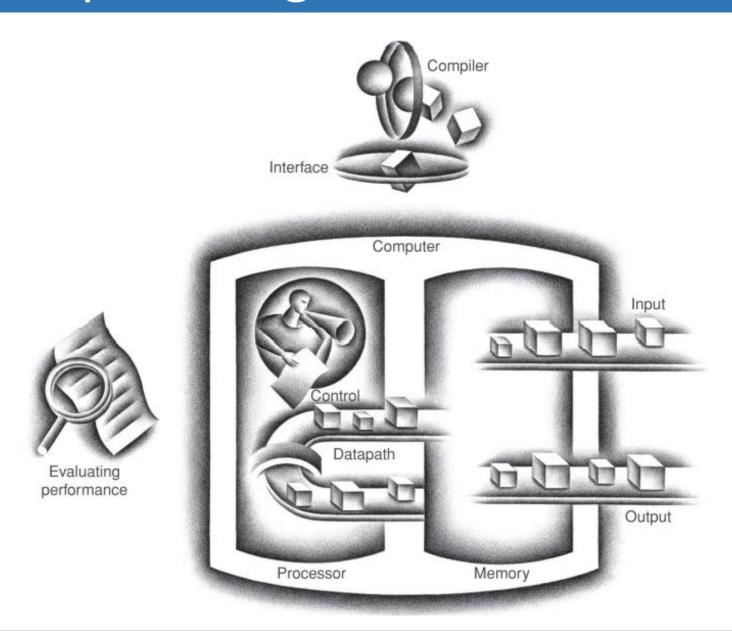
ABI = Application Binary Interface – low level binary interface between two pieces of software on an architecture



Exploring Computer Organization

Typical Computer Organization Includes





From Patterson Hennessy book

LUMS

A Computer has these fundamental Main Components

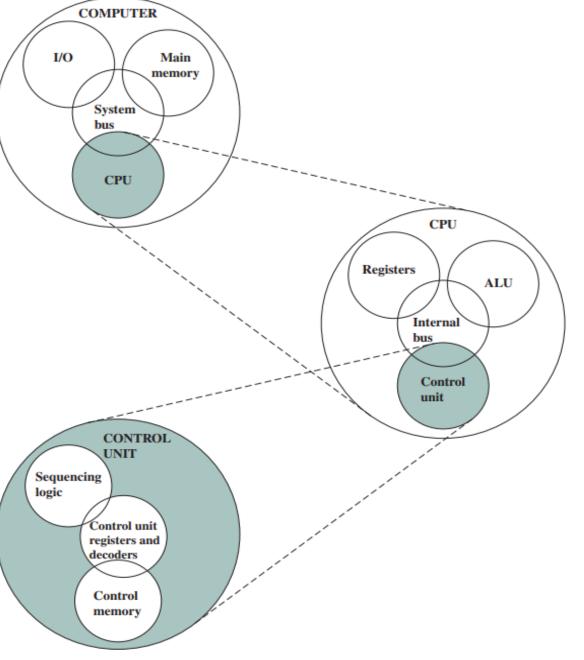


Figure 1.1 The Computer: Top-Level Structure

A Multi-core Computer View

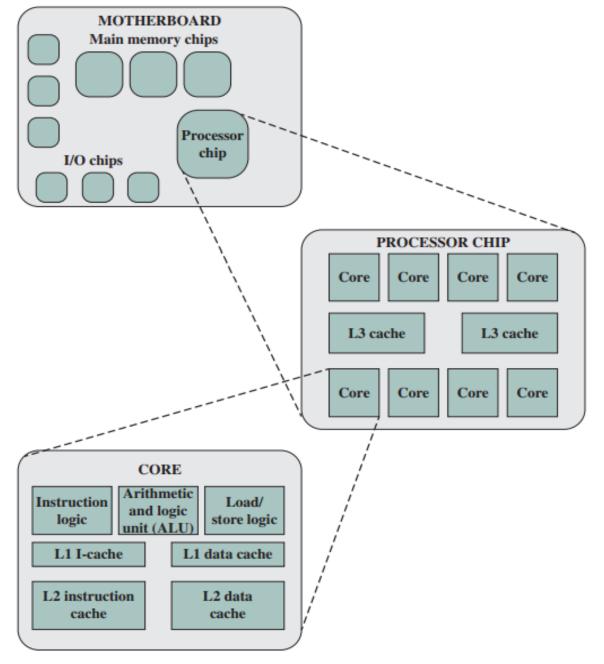
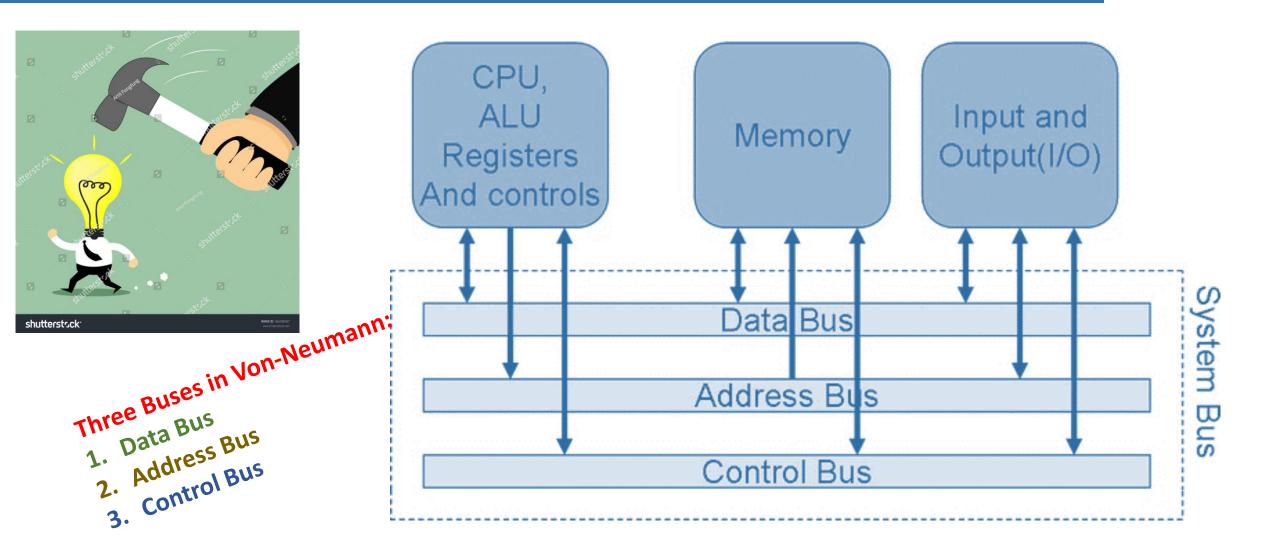


Figure 1.2 Simplified View of Major Elements of a Multicore Computer

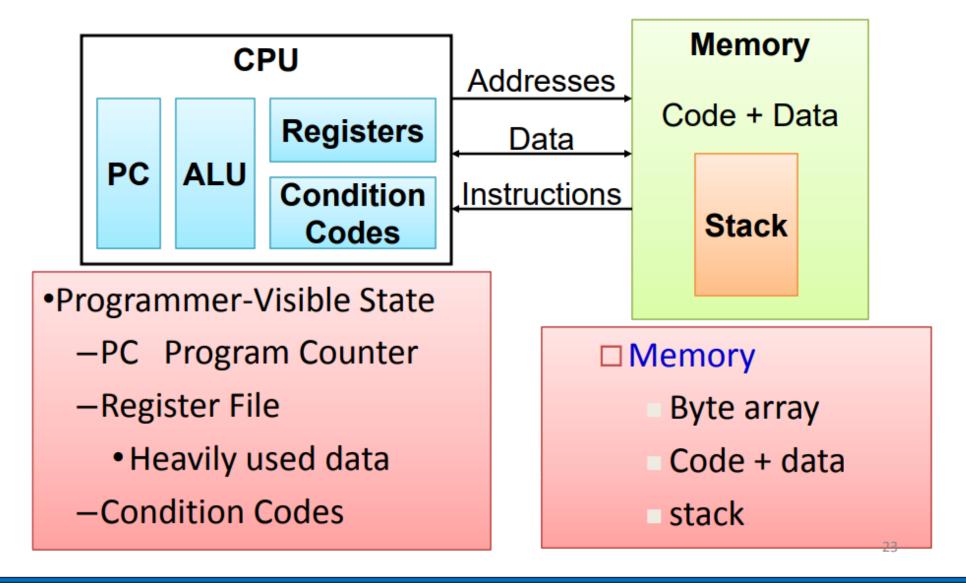
The Von-Neumann Architecture





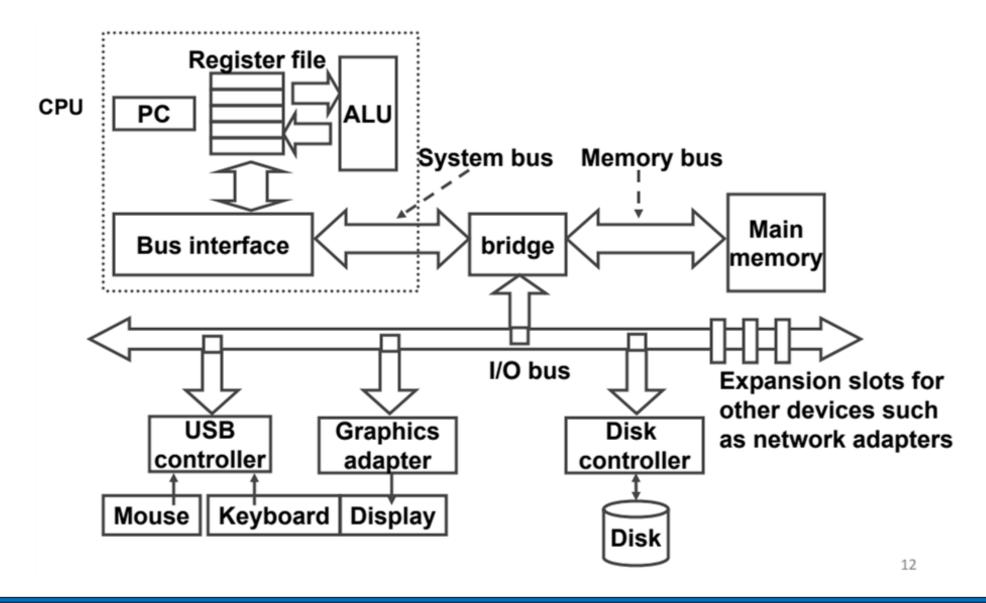
Computer Abstraction





Hardware Abstraction Example









Simple Computer Performance Measures



- Time to execute a program in seconds
- Clock Speed of Computer MHz or GHz
- In terms of the width of Data Bus (16 bit, 32 bit, 64 bit)
- In terms of the size of different memory (8GB / 256GB)
 - Cache size
 - RAM
 - Hard disk
 - SSD, etc.
- In terms of multiple cores and multiple I/O available
- Software
 - Algorithm
 - Language / Compiler
 - Optimization Parallel, Vector, etc.

What is MIPS?



- MIPS is a leading example of RISC Architecture
- Used in many household products such as Nintendo, Sony, Mobile phones, Routers, etc.
- **Simple and Small Instruction Set**

MIPS = Microprocessor without Interlock Pipeline Stages ()



Vs

MIPS = Million Instructions Per Second

What is MFLOPS

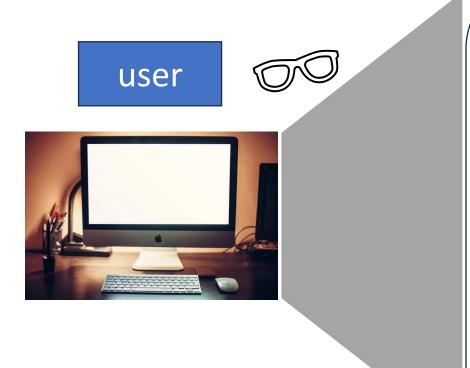


FLOPS = Floating Point Linear Operations Per Second

 $= \frac{Number\ of\ executed\ floating\ point\ operations\ in\ a\ program}{Execution\ time\ \times 10^6}$

Overall Computer Performance





Software Performance

Programming Languages
Operating System
Compilers
Algorithms
Parallelism / Threads

optimization

Software



<u>Hardware Performance</u>

Processor Architecture
Vector Processing
Memory Hierarchy
Multicores
Network support

Hardware optimization

Make Common Case FAST

Use Domain Specific Concepts for Special Tasks

Conclusion

- Complexity Wall
- Power Wall
- Power Calculations
- Abstraction
- Von-Neumann Computer Design
- Hardware / Software Interface
- Introducing Quantitative Performance