

Jin-Soo Kim  
([jinsoo.kim@snu.ac.kr](mailto:jinsoo.kim@snu.ac.kr))

Systems Software &  
Architecture Lab.

Seoul National University

Fall 2020

# Introduction to Computer Architecture

Chap. 1.1 – 1.5, 1.7 – 1.8



# Computer Systems



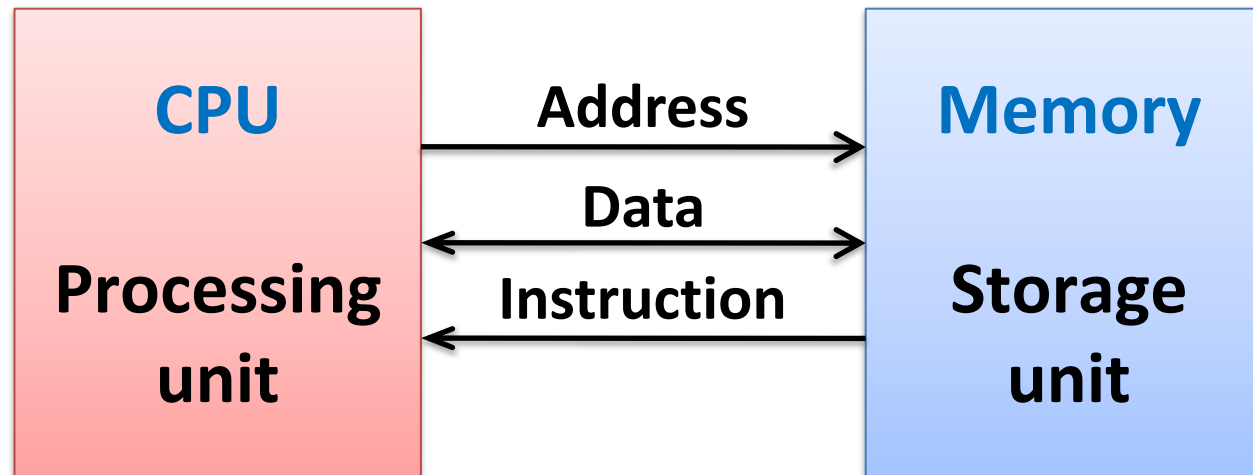
# Classes of Computers

- **Personal computers**
  - General purpose, variety of software
  - Subject to cost/performance tradeoff
- **Server computers**
  - Network based
  - High performance, capacity, reliability
  - Range from small servers to large data centers
- **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

A computer is a  machine.

# von Neumann Architecture

- By John von Neumann, 1945

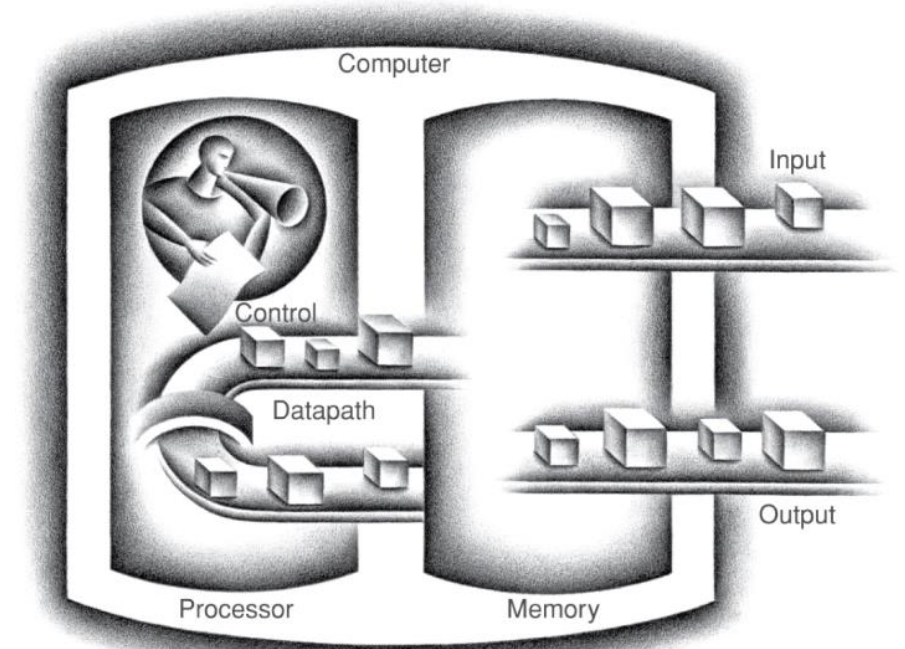
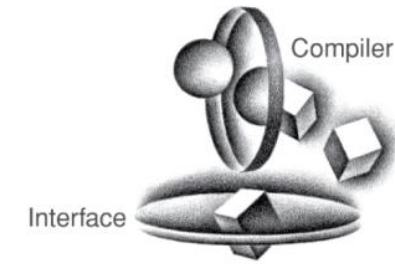


Data movement  
Arithmetic & logical ops  
Control transfer

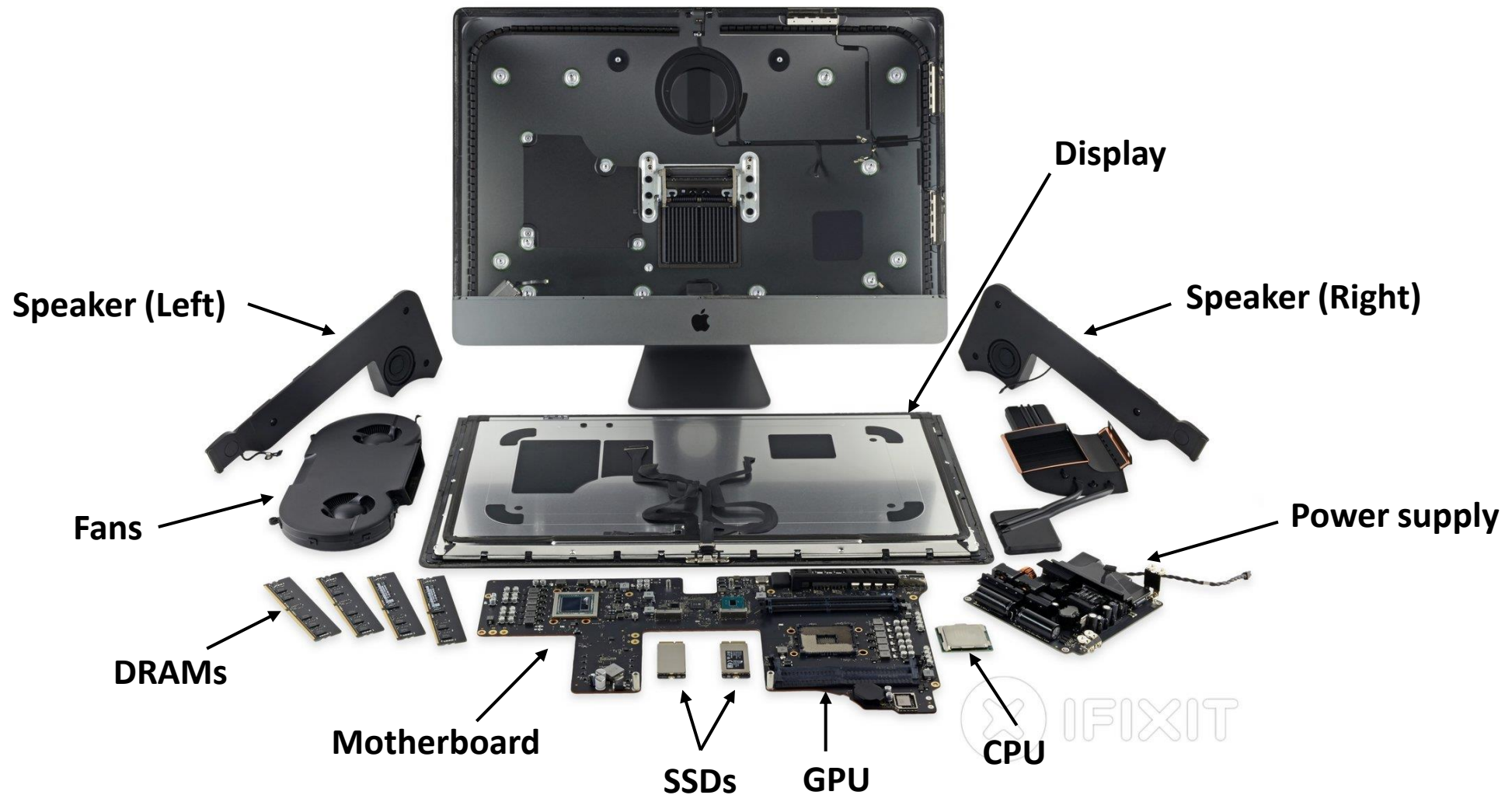
Byte addressable array  
Code + data (user program, OS)  
Stack to support procedures

# Components of a Computer

- CPU: Control + Datapath
- Memory
- I/Os
  - GPUs
  - User-interface devices:  
Display, keyboard, mouse, sound, ...
  - Storage devices:  
HDD, SSD, CD/DVD, ...
  - Network adapters:  
Ethernet, 3G/4G/5G, WiFi, Bluetooth, ...
- Same components for all kinds of computer



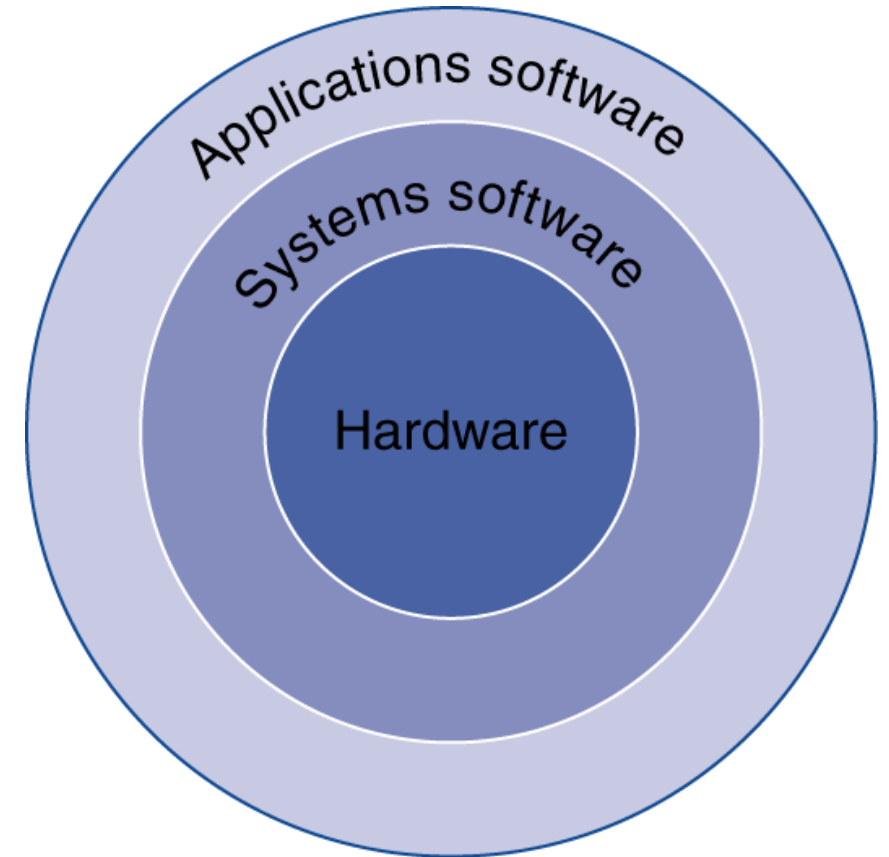
# Opening the Box (iMac Pro)





# Below Your Program

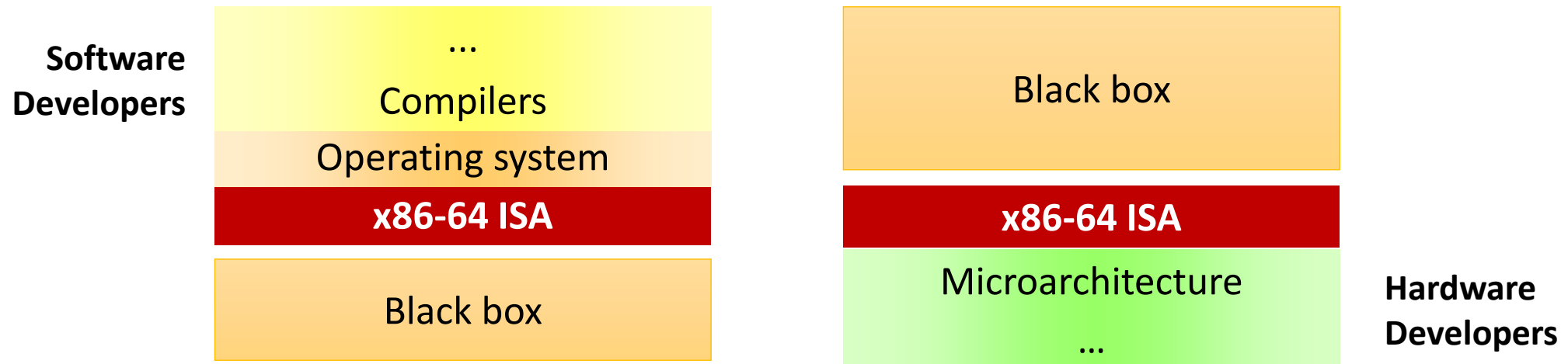
- **Application software**
  - Written in high-level language
- **System software**
  - Compiler: translates HLL code to machine code
  - Operating system: service code
    - Handling input/output
    - Managing memory and storage
    - Scheduling tasks & sharing resources
- **Hardware**
  - Processor,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





# Instruction Set Architecture (ISA)

- The hardware/software interface
  - Hardware abstraction visible to software (OS, compilers, ...)
  - Instructions and their encodings, registers, data types, addressing modes, etc.
  - Written documents about how the CPU behaves
  - e.g., All 64-bit Intel CPUs follow the same x86-64 (or Intel 64) ISA

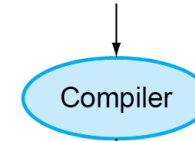


# Levels of Program Code

- High-level language
  - Level of abstraction closer to problem domain
  - Provides for productivity and portability
- Assembly language
  - Textual representation of instructions
  - For humans
- Machine language
  - Hardware representation
  - Binary digits (bits)
  - Encoded instructions and data

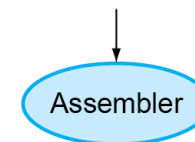
High-level  
language  
program  
(in C)

```
swap(int v[], int k)
{int temp;
  temp = v[k];
  v[k] = v[k+1];
  v[k+1] = temp;
}
```



Assembly  
language  
program  
(for RISC-V)

```
swap:
  slli x6, x11, 3
  add x6, x10, x6
  ld x5, 0(x6)
  ld x7, 8(x6)
  sd x7, 0(x6)
  sd x5, 8(x6)
  jalr x0, 0(x1)
```



Binary machine  
language  
program  
(for RISC-V)

```
0000000001101011001001100010011
00000000011001010000001100110011
00000000000000110011001010000011
00000000100000110011001110000011
00000000011100110011000000100011
00000000010100110011010000100011
0000000000000000100000001100111
```

# What Happened:

1997

2017

104 cabinets  
(76 computes,  
8 switches,  
20 disks)

9298 cores

150m<sup>2</sup>



ASCI Red at Sandia

1.3 TF/s, 850 KW

Cavium ThunderX2

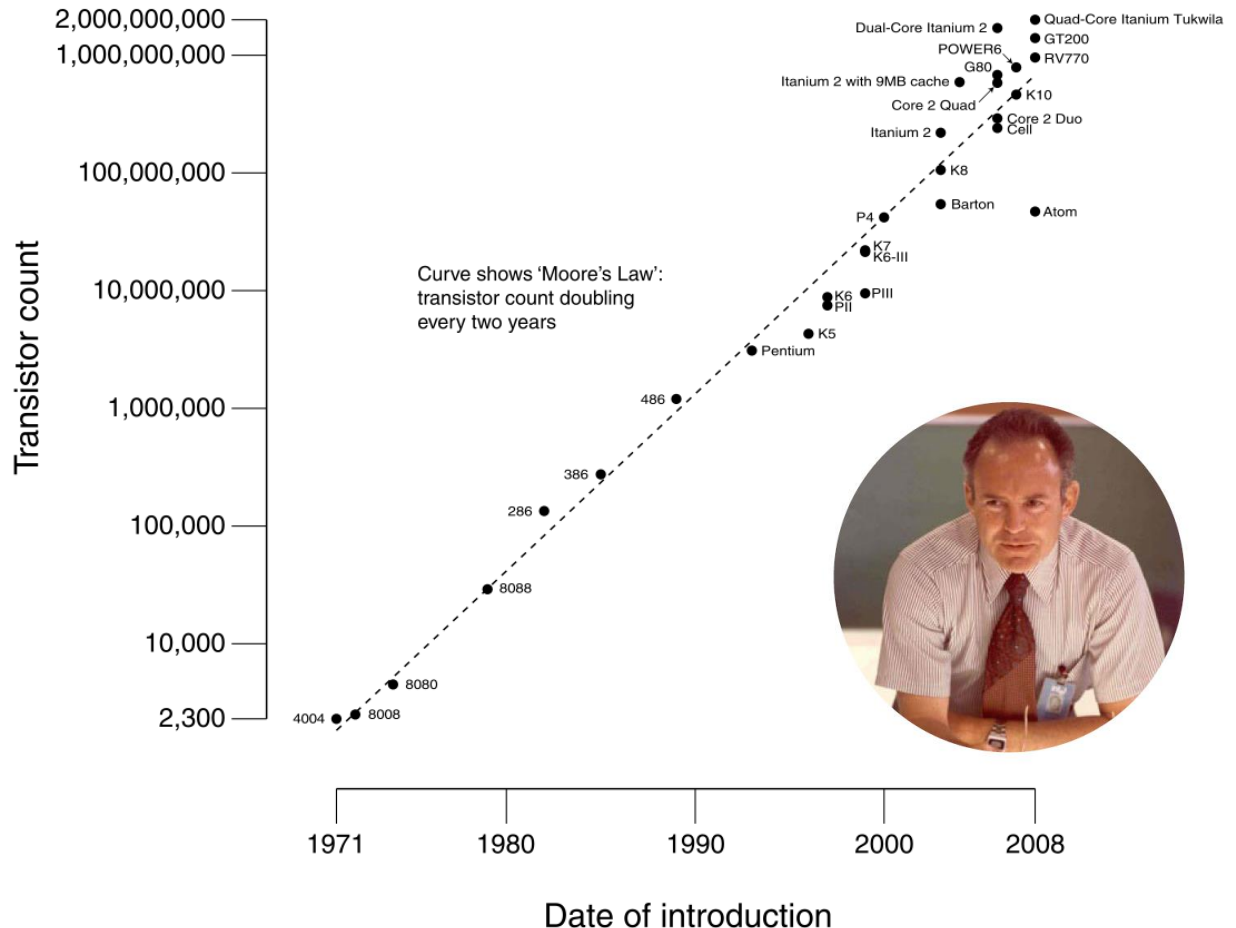
~ 1.1 TF/s, ~ 0.2 KW

3.5 orders of  
magnitude

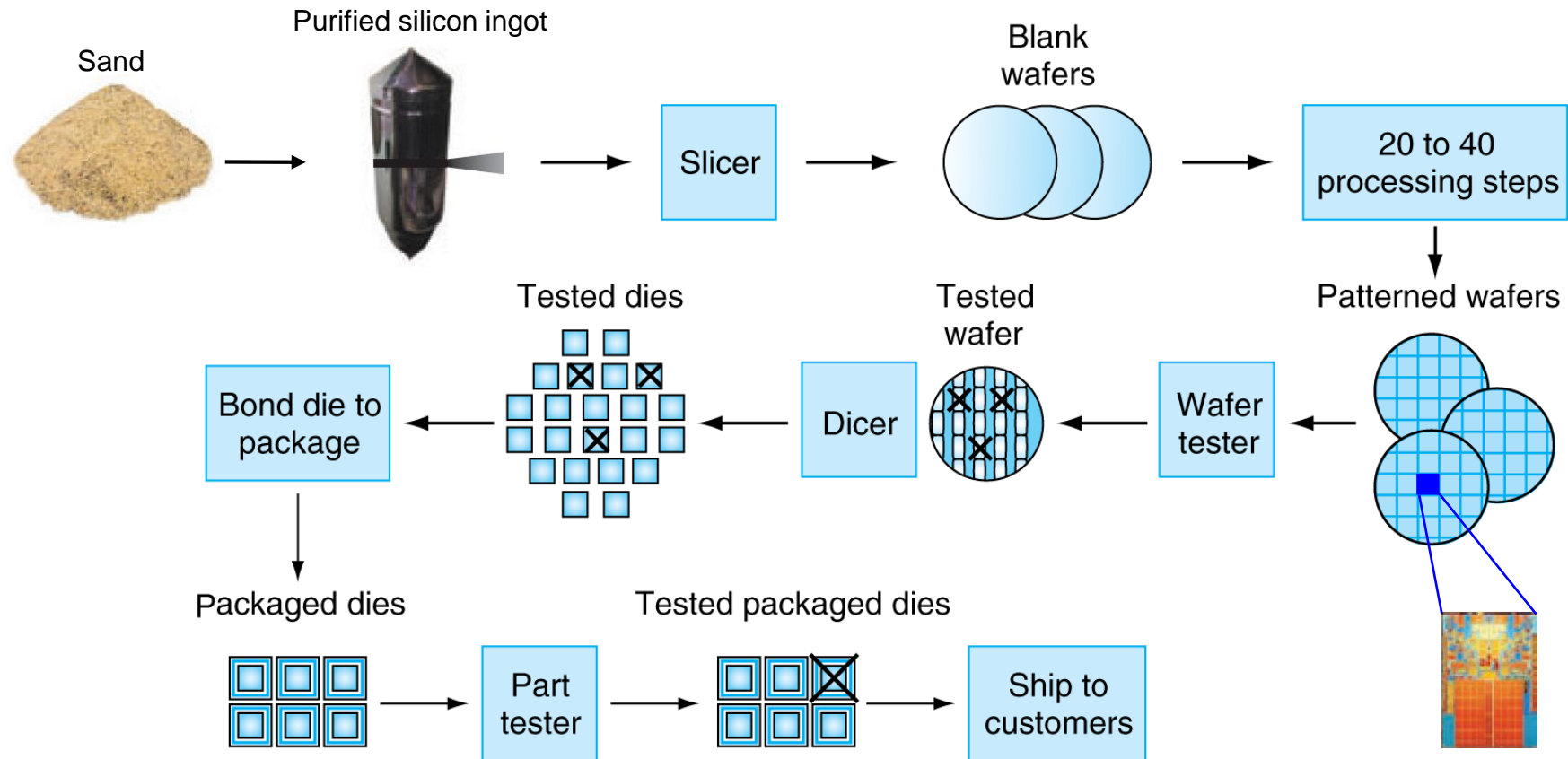
# The Moore's Law

- *"The number of transistors incorporated in a chip will approximately double every 24 months"*  
(Gordon Moore, Intel Co-founder, 1965)
- **Makes novel applications feasible**
  - WWW, search engines,
  - Smartphones, VR/AR,
  - AI, Self-driving cars,
  - Human genome project, ...
- **Computers are pervasive**

CPU Transistor Counts 1971-2008 & Moore's Law

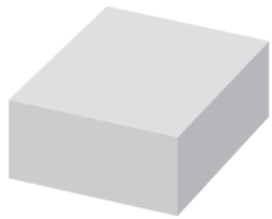


# From Sand to Circuits

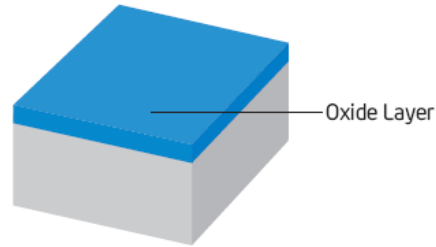


**Yield: proportion of working dies per wafer**

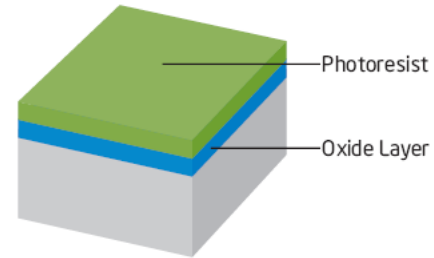
# Processing ICs



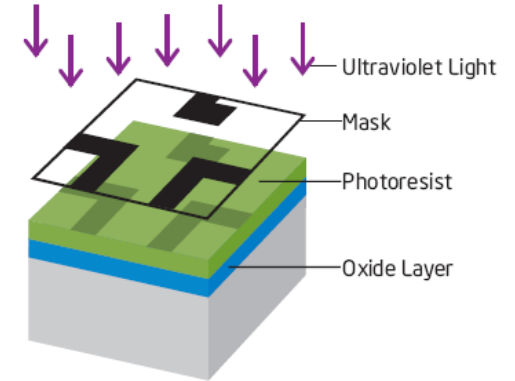
1. Start with a partially processed die on a silicon wafer.



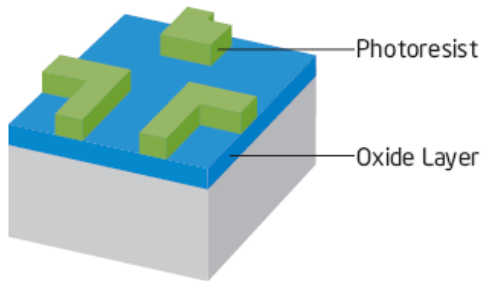
2. Deposit oxide layer.



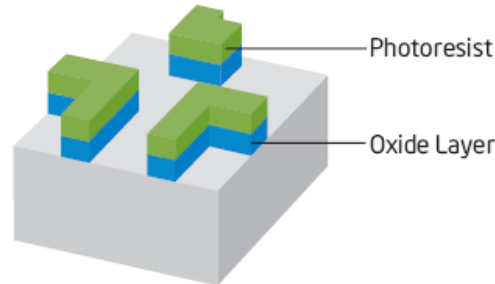
3. Coat with photoresist.



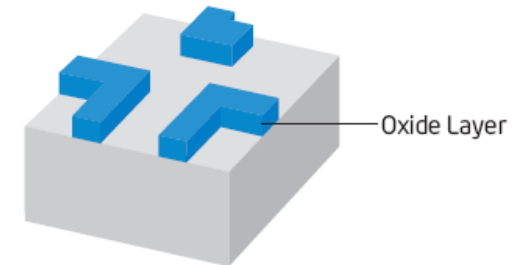
4. Position mask and flash ultraviolet light.



5. Rinse with solvent.



6. Etch with acid.



7. Remove remaining photoresist.



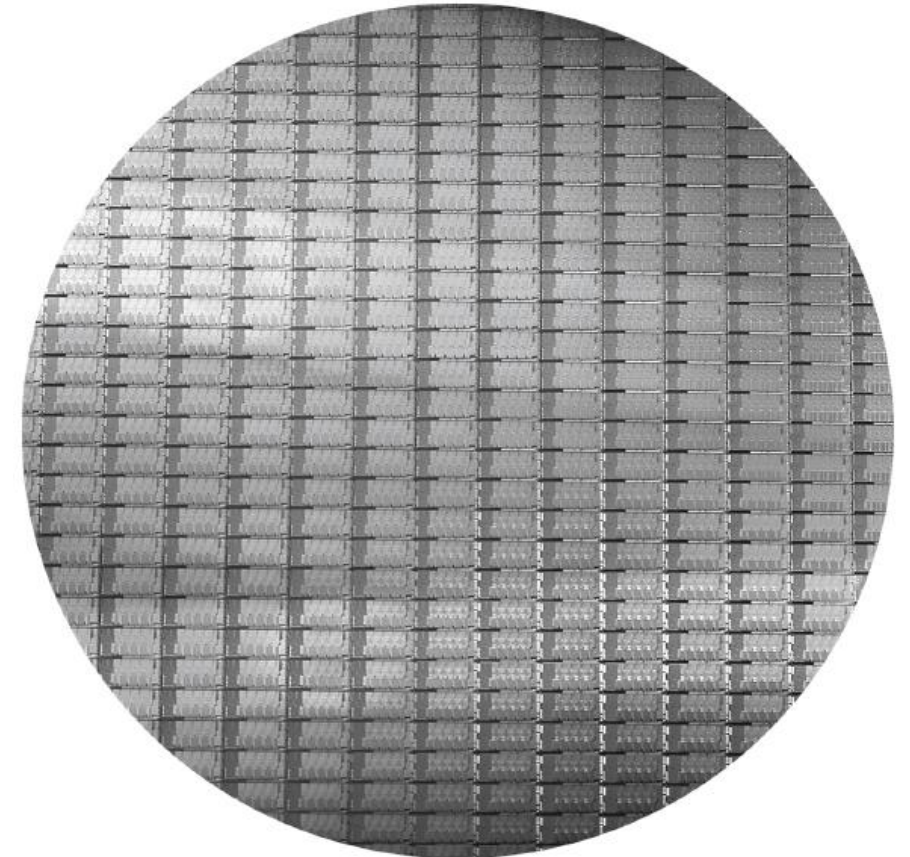
# Intel Core i7 Wafer

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

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

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

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



- Wafer cost and area are fixed
- Defect rate determined by manufacturing process
- Die area determined by architecture and circuit design



# Technology Trends

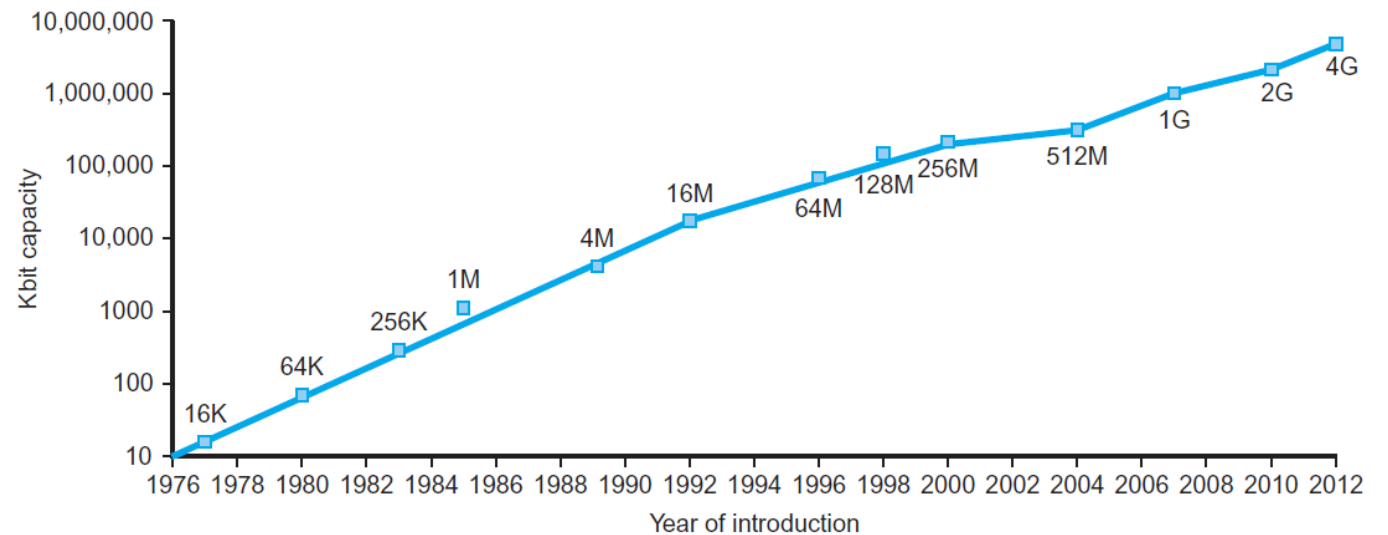
## ■ CPU

- Logic capacity:  $\sim 30\%$  / year
- Clock rate:  $\sim 20\%$  / year

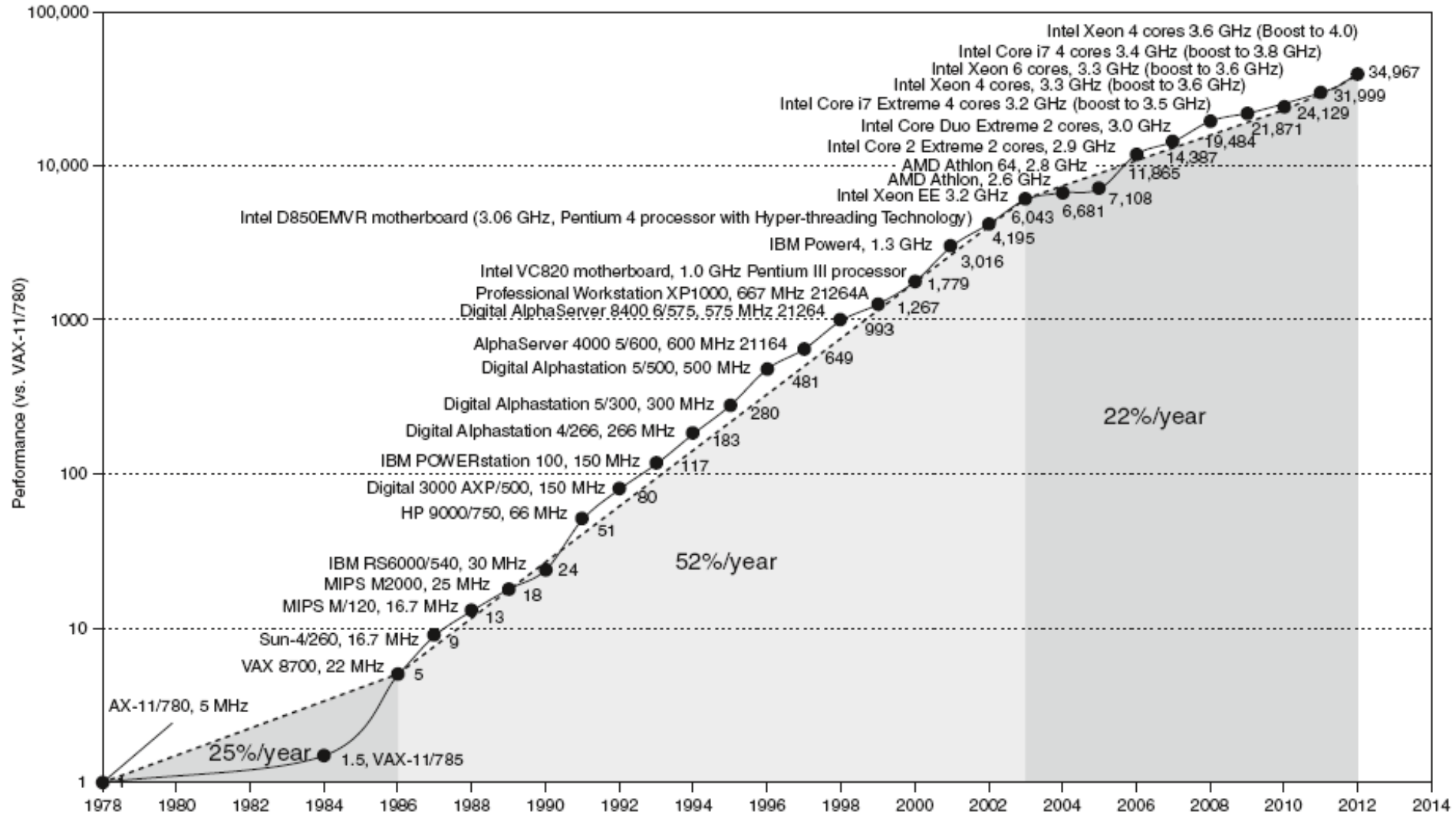
## ■ Memory

- DRAM capacity:  $\sim 60\%$  / year  
(4x every 3 years)
- DRAM speed:  $\sim 10\%$  / year

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



# Uniprocessor Performance



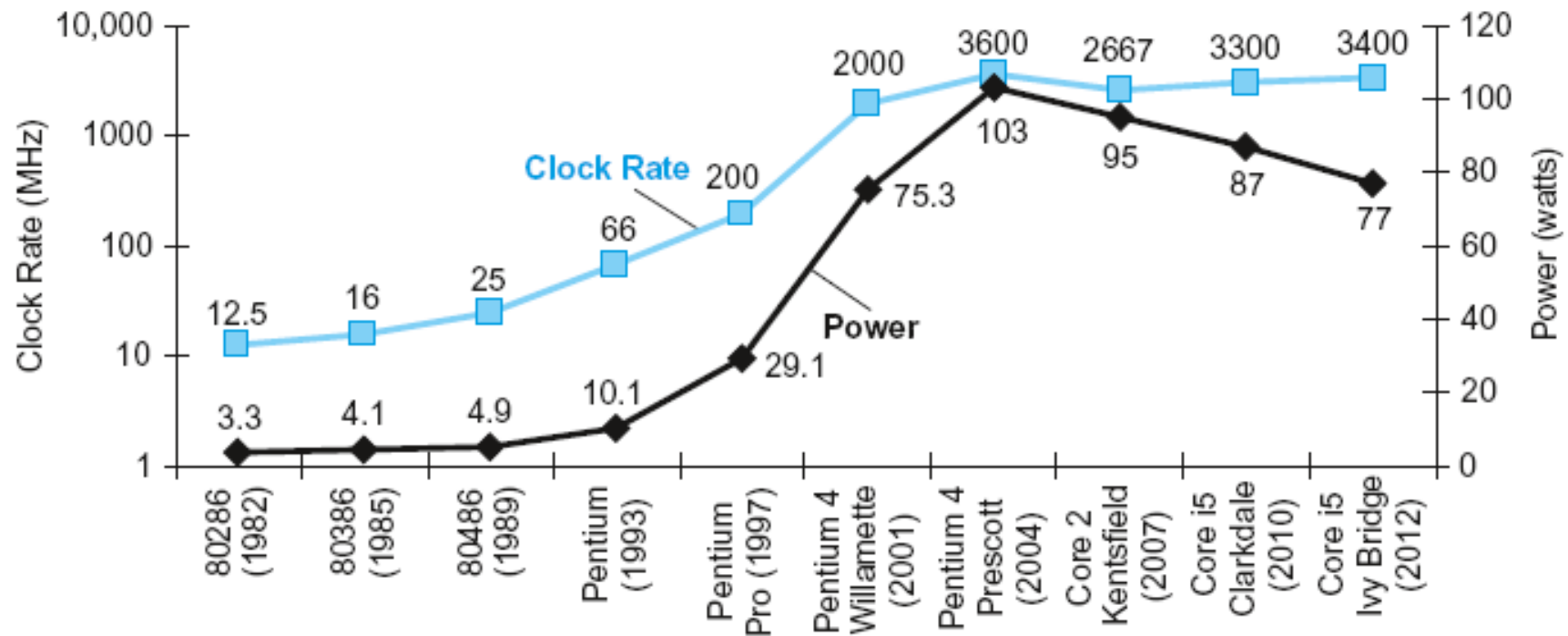
# Architecture: Exploiting Parallelism

- **Instruction level parallelism (ILP)**
  - Pipelining
  - Superscalar
  - Out-of-order execution
  - Branch prediction
  - VLIW (Very Long Instruction Word)
- **Data level parallelism (DLP)**
  - SIMD / Vector instructions
- **Task level parallelism (TLP)**
  - Simultaneous multithreading (Hyperthreading)
  - Multicore

# Power Trends

- In CMOS IC technology,

$$Power = Capacitive Load \times Voltage^2 \times Frequency$$



# Reducing Power

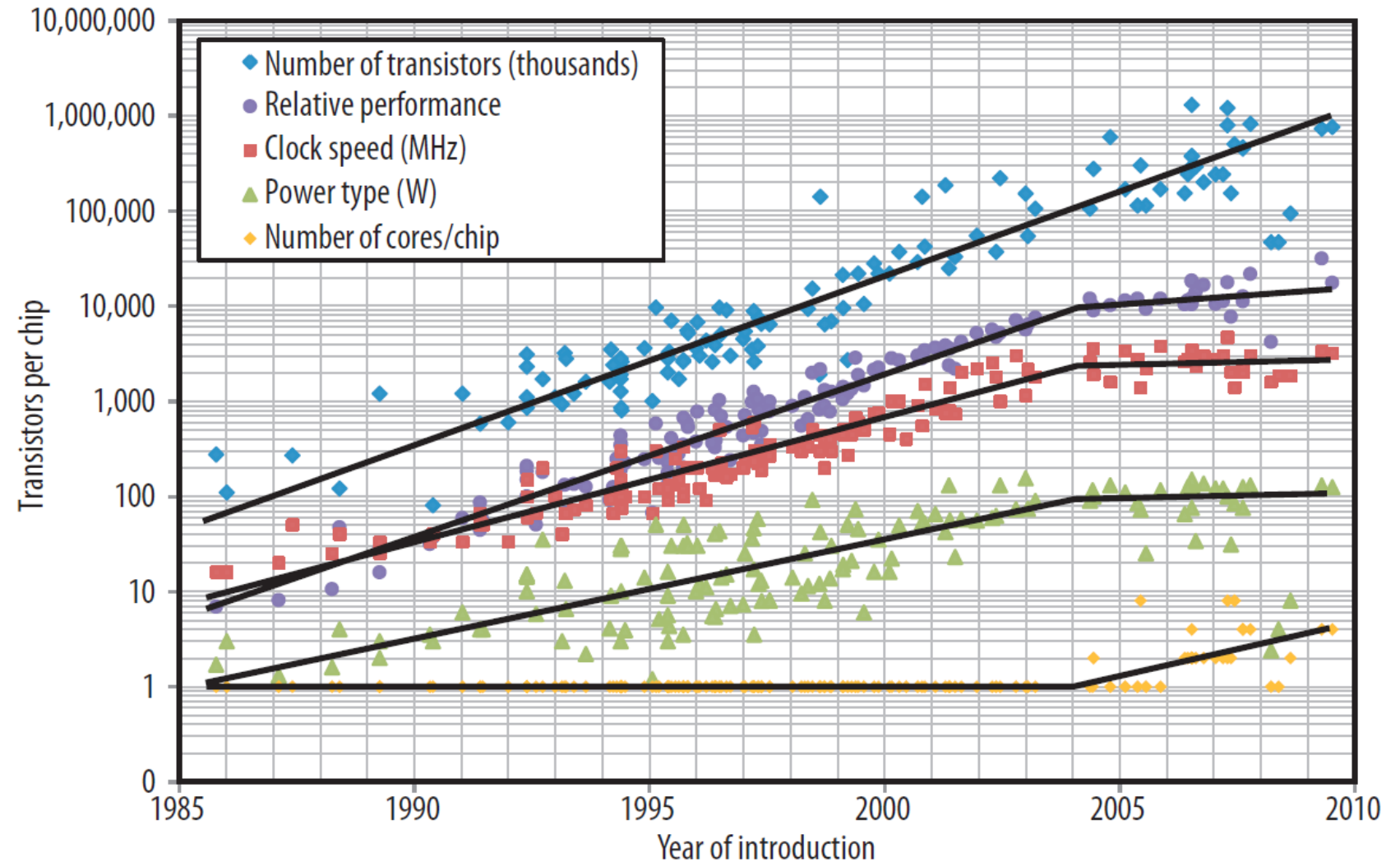
- Suppose a new CPU has
  - 85% of capacitive load of old CPU
  - 15% voltage and 15% frequency reduction

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

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

# The Shift to Multicores

- **ILP wall**
  - Control dependency
  - Data dependency
- **Memory wall**
  - Memory latency improved by 10% / year
  - Cache shows diminishing returns
- **Power wall**

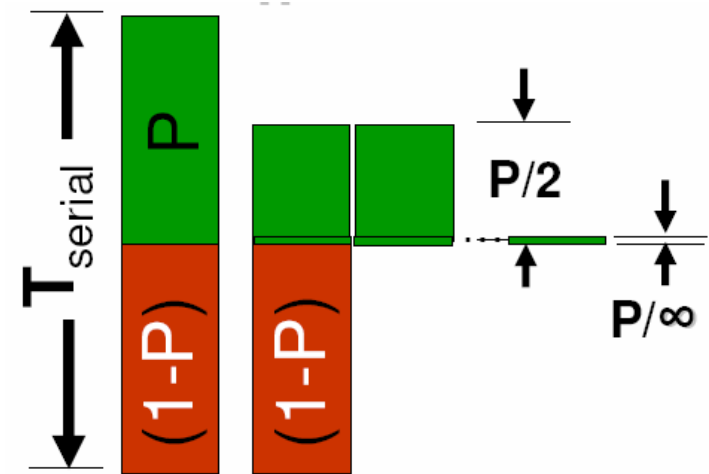


# Amdahl's Law

- The theoretical upper limit of speed up is limited by the serial portion of the code

$$Speedup = \frac{1}{(1 - P) + \frac{P}{n}}$$

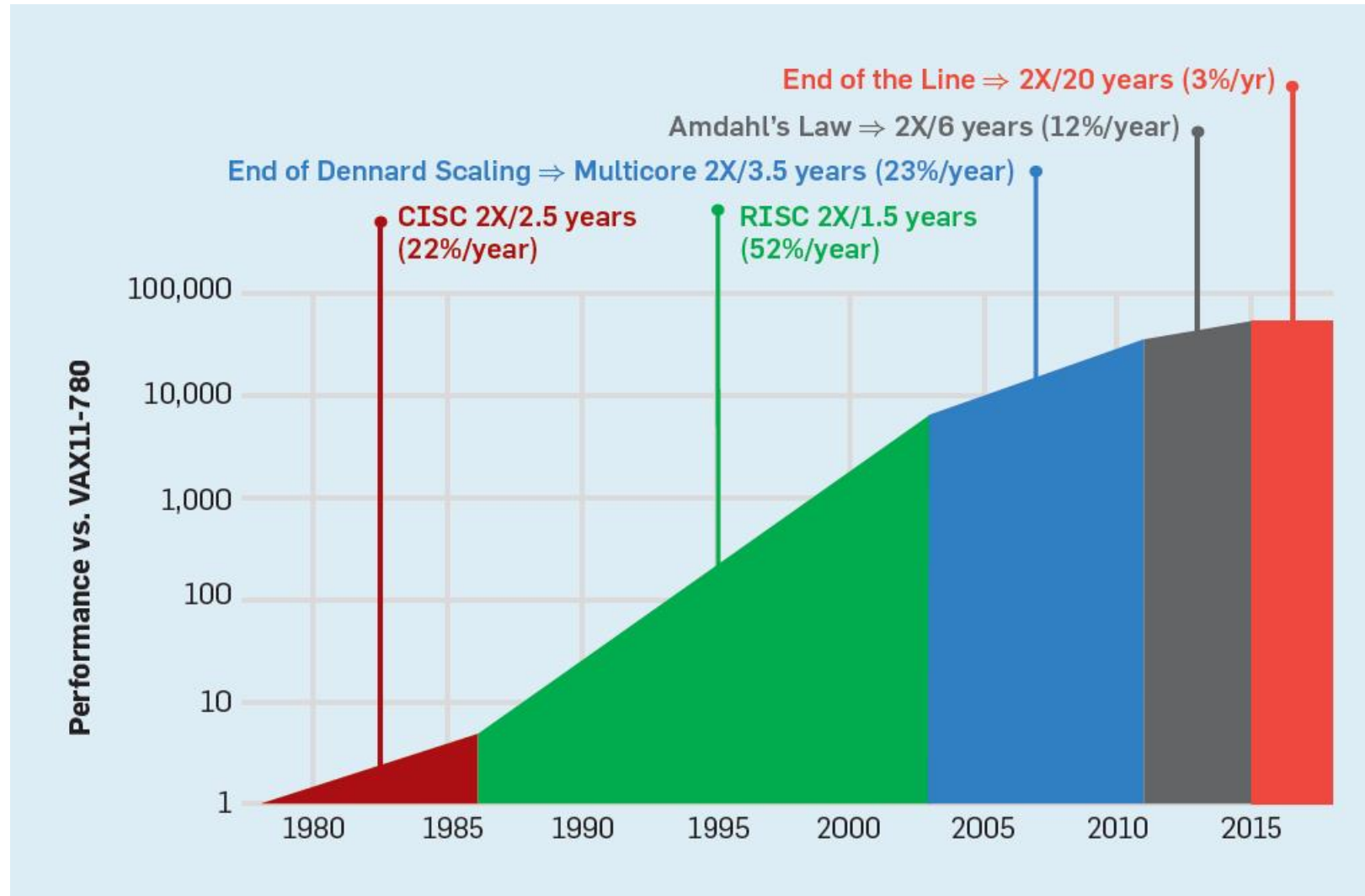
- $S = (1 - P)$ : the time spent executing the serial portion
- $n$ : the number of processor cores



- Corollary: make the common case fast



# Summary: A New Golden Age?



**End of  
Moore's law**

**A New Golden Age  
for Computer  
Architecture  
(CACM, Feb. 2019)**