

# Computer Abstractions and Technology

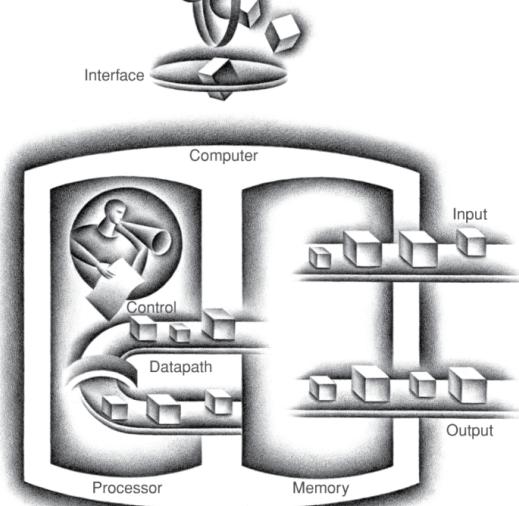
CS 154: Computer Architecture

Lecture #2

Winter 2020

Ziad Matni, Ph.D.

Dept. of Computer Science, UCSB



# A Word About Registration for CS154

---

## **FOR THOSE OF YOU NOT YET REGISTERED:**

- This class is **FULL**
- **If you want to add this class AND you are on the waitlist, see me after lecture**

# Lecture Outline

---

- Tech Details
  - Trends
  - Historical context
  - The manufacturing process of Ics
- Important Performance Measures
  - CPU time
  - CPI
  - Other factors (power, multiprocessors)
  - Pitfalls

# Parts of the CPU

- The **Datapath**, which includes the **Arithmetic Logic Unit (ALU)** and other items that perform operations on data
- **Cache Memory**, which is small & fast RAM memory for immediate access to data. Resides inside the CPU.  
(other types of memory are outside the CPU, like DRAM, etc...)
- The **Control Unit (CU)**  
which sequences how Datapath + Memory interact

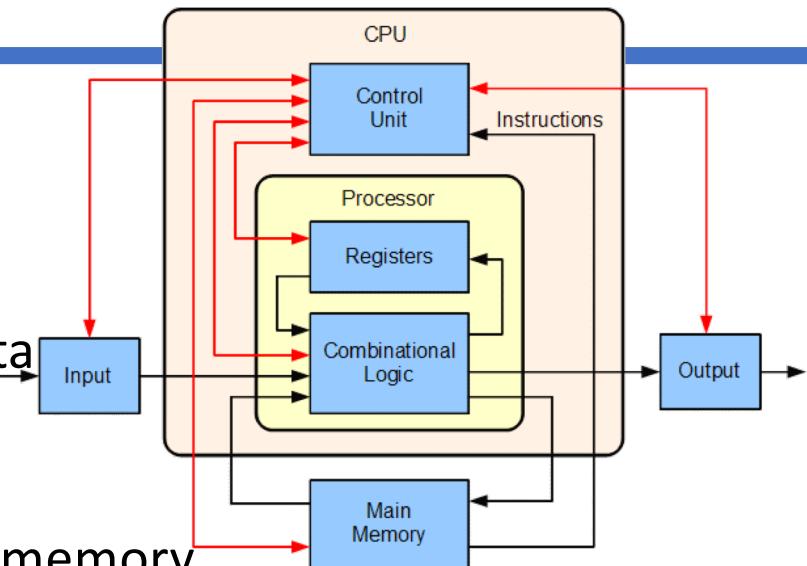
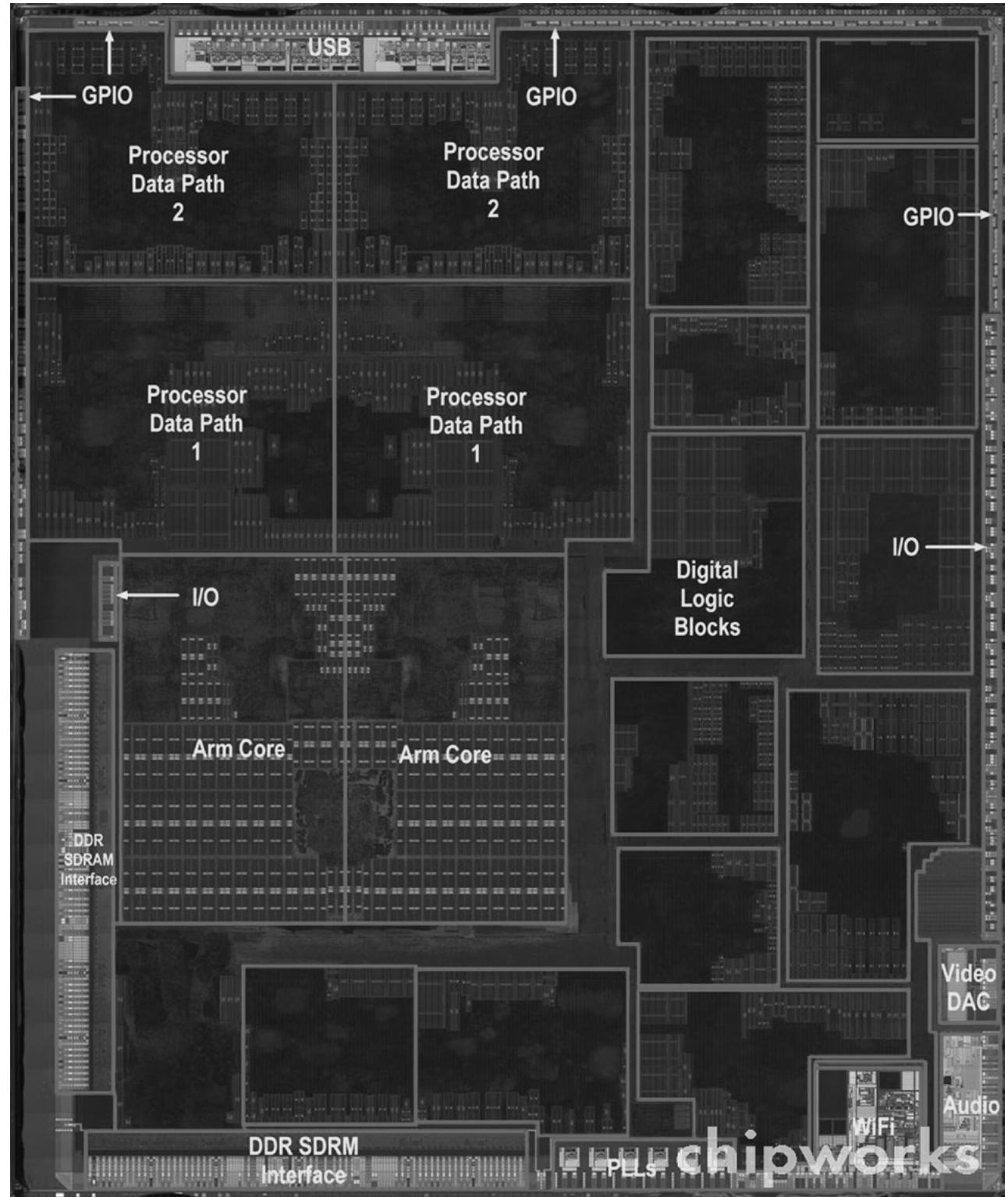


Image from wikimedia.org

# Inside the Apple A5 Processor

Manufactured in 2011 – 2013  
32 nm technology  
37.8 mm<sup>2</sup> die size



# The CPU's Fetch-Execute Cycle

---

- **Fetch** the next instruction
- **Decode** the instruction
- **Get data** if needed
- **Execute** the instruction
  - Maybe access mem again and/or write back to reg.

*This is what happens inside a computer interacting with a program at the “lowest” level*

# Pipelining (Parallelism) in CPUs

---

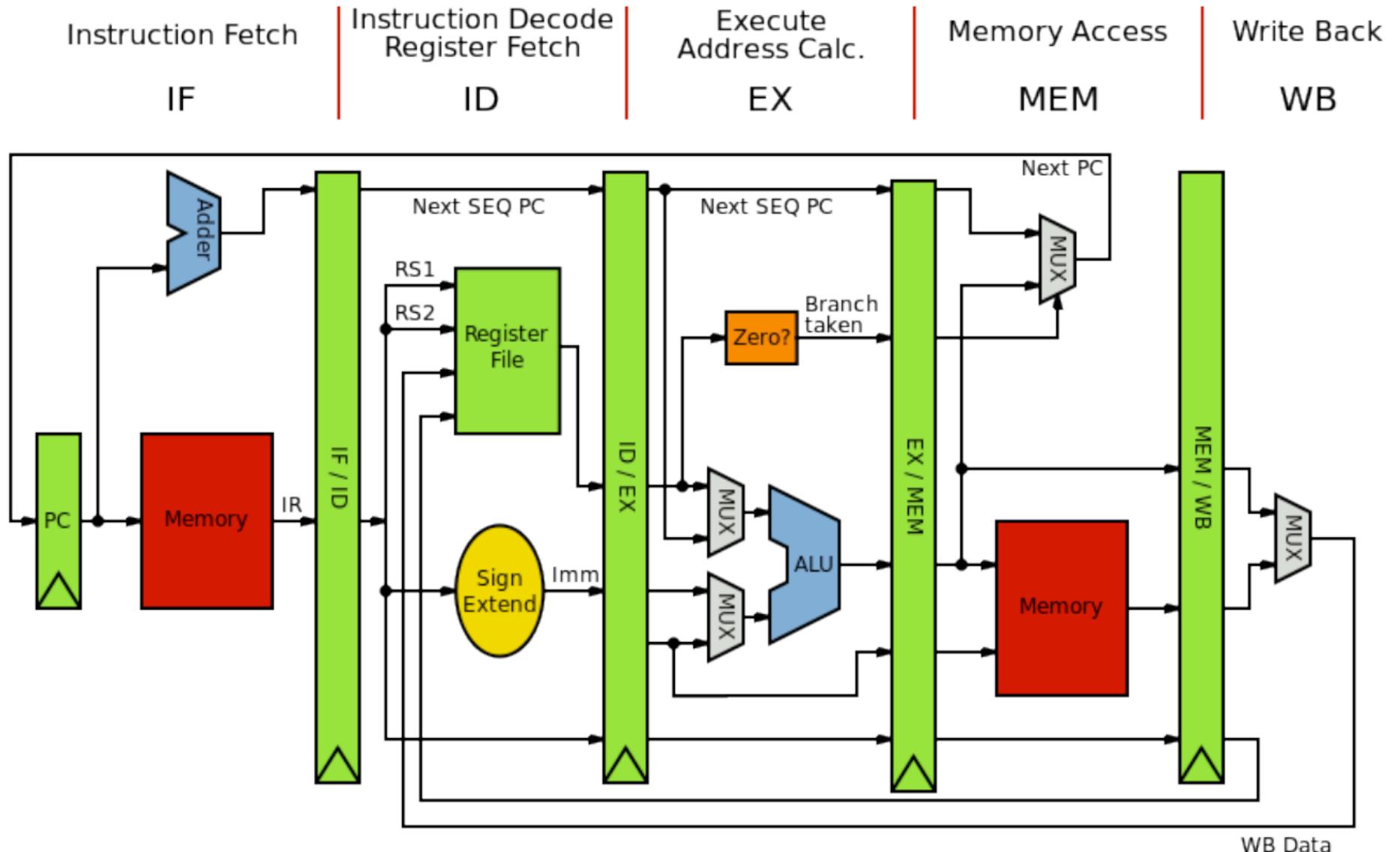
- Pipelining is a fundamental design in CPUs
- Allows multiple instructions to go on at once
  - a.k.a instruction-level parallelism

**Basic five-stage pipeline**

Instr. No.	Clock cycle	1	2	3	4	5	6	7
1		IF	ID	EX	MEM	WB		
2			IF	ID	EX	MEM	WB	
3				IF	ID	EX	MEM	WB
4					IF	ID	EX	MEM
5						IF	ID	EX

(IF = Instruction Fetch, ID = Instruction Decode, EX = Execute,  
MEM = Memory access, WB = Register write back).

# Digital Design of a CPU (Showing Pipelining)



# Computer Languages and the F-E Cycle

---

- Instructions get executed in the CPU in machine language (i.e. all in “1”s and “0”s)
- Even *small* instructions, like  
“add 2 to 3 then multiply by 4”,  
need *multiple* cycles of the CPU to get fully executed
- But **THAT’S OK!** Because, typically,  
**CPUs can run *many millions* of instructions per second**

# Computer Languages and the F-E Cycle

---

- But **THAT'S OK!** Because, typically,  
CPUs can run *many millions* of instructions per second
- In *low-level languages* (like assembly or machine lang.) you need to spell those parts of the cycles one at a time
- In *high-level languages* (like C, Python, Java, etc...) you don't
  - 1 HLL statement, like " $x = c*(a + b)$ " is enough to get the job done
  - This would translate into multiple statements in LLLs
  - **What translates HLL to LLL?**
  - **What translates LLL to ML?**

# Machine vs. Assembly Language

- **Machine language (ML)** is the actual 1s and 0s

Example:

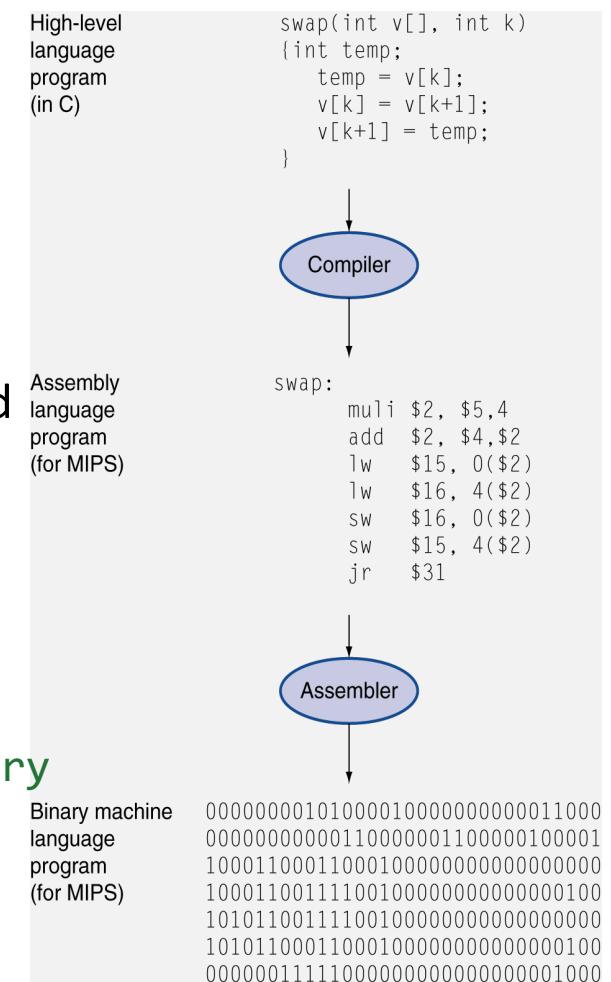
```
1011110111011100000101010101000
```

- **Assembly language** is one step above ML

- Instructions are given **mnemonic codes** but still displayed one step at a time
- Advantage? Better human readability

Example:

```
lw    $t0, 4($sp)    # fetch N from someplace in memory  
add  $t0, $t0, $t0    # add N to itself  
                      # and store the result in N
```



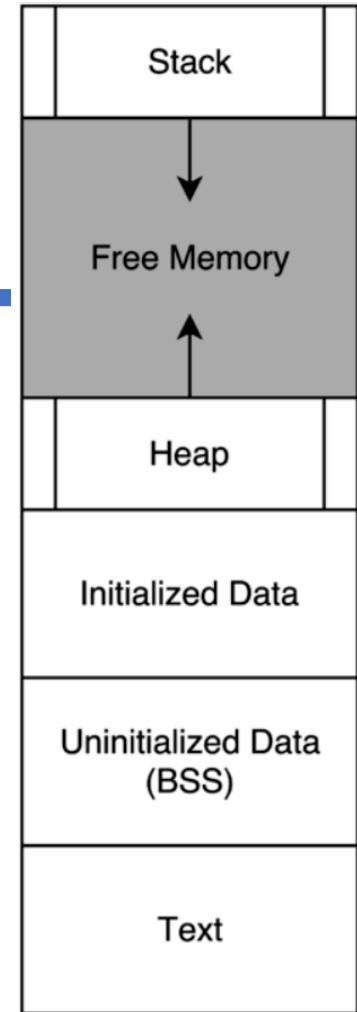
# Computer Memory

Usually organized in two parts:

- **Address:** *Where* can I find my data?
- **Data (payload):** *What* is my data?

Recall:

- A bit (b) is \_\_\_\_\_
- A byte (B) is \_\_\_\_\_
- MIPS CPUs operate ***instructions*** that are \_\_\_\_\_ bits long
- MIPS CPUs organize ***memory*** in units called \_\_\_\_\_ that are \_\_\_\_\_ bits long
- MIPS memory is addressable in \_\_\_\_\_  ***endian***



# Reminder of some MIPS instructions

---

- **add** vs **addi** vs **addu** vs **addui**
- **mult** and **mflo**
- **sll**
- **srl** vs **sra**
- **la** vs **li** vs **lw** vs **sw**



# Eight Great Ideas in Computer Architecture

---

- 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

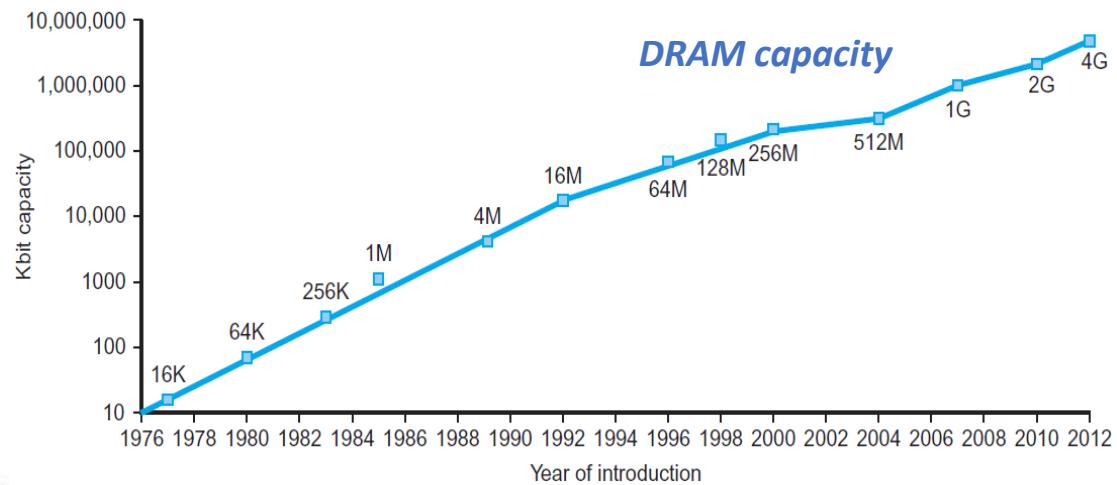
# Electronic Circuitry Tech Trends

- Electronics technology continues to evolve
  - Increased memory capacity (at same price/size)
  - Increased CPU performance
  - Reduced costs overall

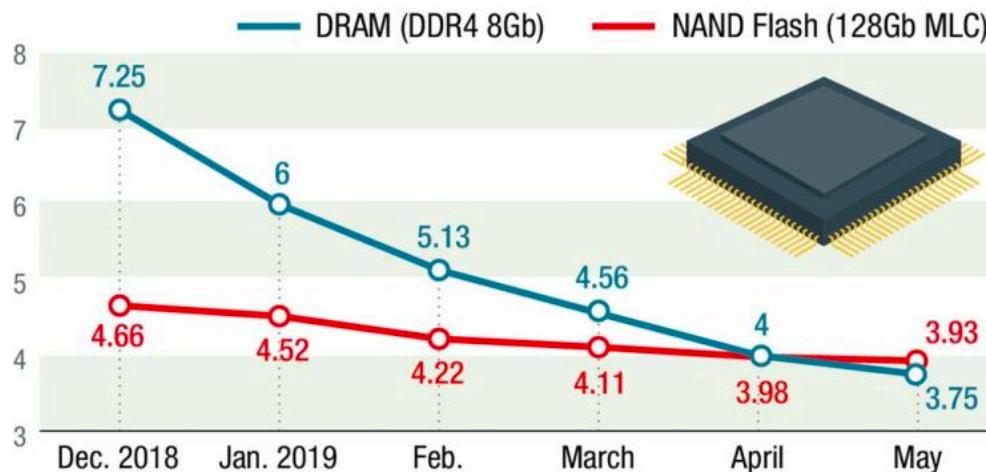
Year	Technology	Relative Performance
1951	Vacuum tube	1
1965	Transistor	35
1975	Integrated circuit (IC)	900
1995	Very large scale IC (VLSI)	2.4 million
2013	Application Specific IC or ASIC (ultra-large scale)	250 million

# DRAM capacity goes up and the prices come down...

- DRAM = Dynamic RAM
- Very common tech used for computer memory

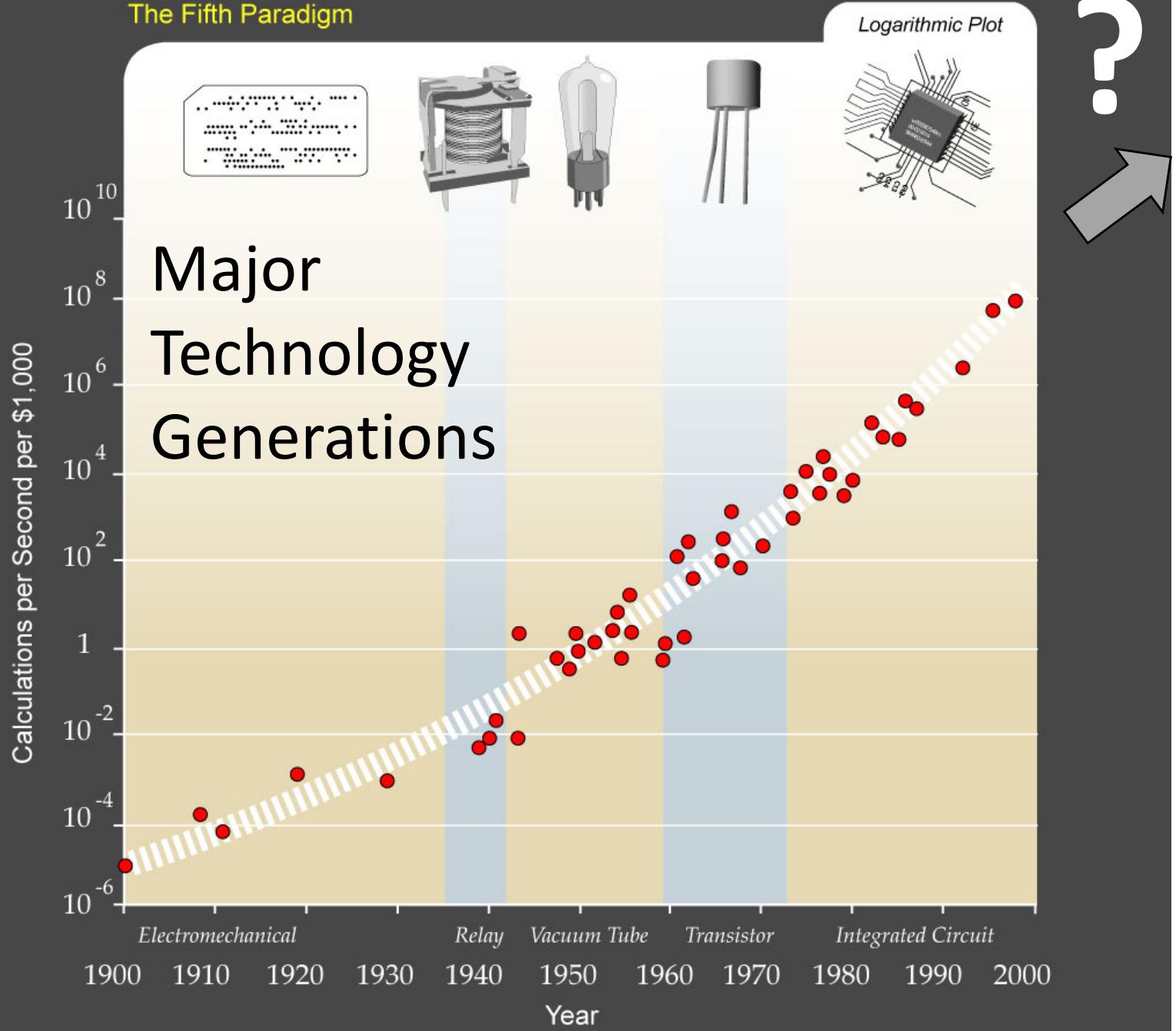


## Decreasing memory chip prices (Unit: dollar)



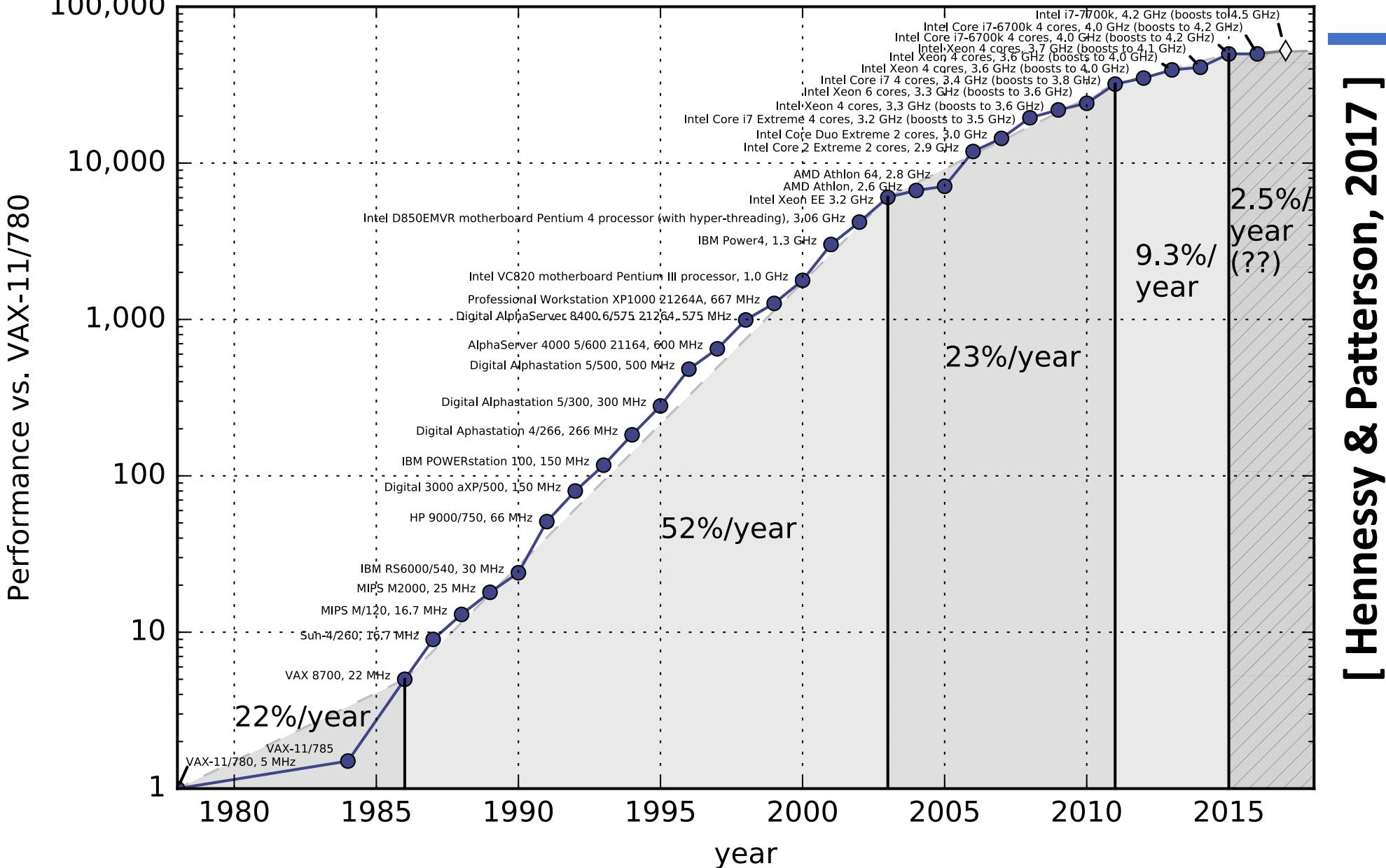
Source: DRAMEXchange  
Matni, CS154, Wi20

**Moore's Law**  
The Fifth Paradigm



# Single-Thread Processor Performance

## 40 years of Processor Performance



# Computer Architecture: A Little History

---

Throughout the course we'll use a historical narrative to help understand why certain ideas arose

Why worry about old ideas?

- Helps to illustrate the design process, and explains why certain decisions were taken
- Because future technologies might be as constrained as older ones
- Those who ignore history are doomed to repeat it
  - Every mistake made in mainframe design was also made in minicomputers, then microcomputers, where next?

# Digital Computers

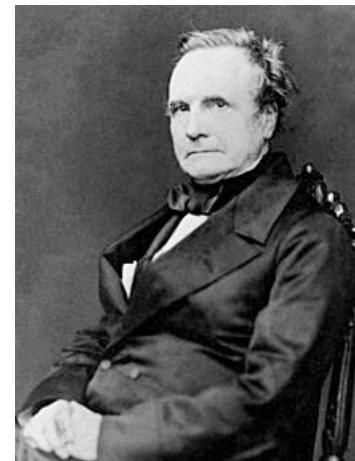
---

- An improvement over Analog Computers...
- Represent problem variables as numbers encoded using discrete steps
  - Discrete steps provide noise immunity
- Enables accurate and deterministic calculations
  - Same inputs give same outputs exactly

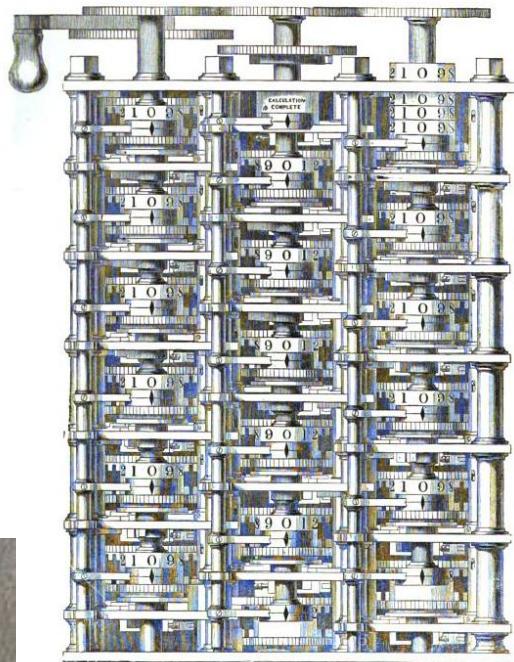
# Computing Devices for General Purposes

- **Charles Babbage (UK)**

- *Analytical Engine* could calculate polynomial functions and differentials
- Inspired by older generation of calculating machines made by Blaise Pascal (1623-1662, France)
- Calculated results, but also *stored intermediate findings* (i.e. precursor to computer memory)
- “**Father of Computer Engineering**”



C. Babbage (1791 – 1871)



Part of Babbage's Analytical Engine

- **Ada Byron Lovelace (UK)**

- Worked with Babbage and foresaw computers doing much more than calculating numbers
- Loops and Conditional Branching
- “**Mother of Computer Programming**”



A. Byron Lovelace (1815 – 1852)

Images from Wikimedia.org

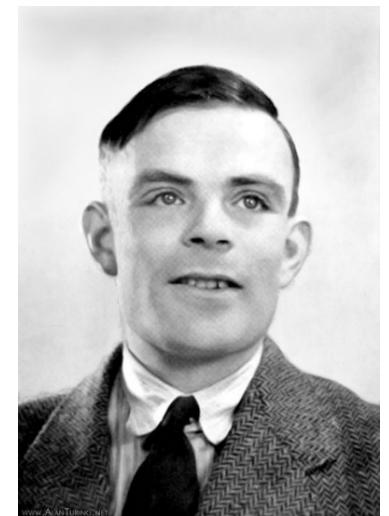
# The Modern Digital Computer

---

- Calculating machines kept being produced in the early 20<sup>th</sup> century (IBM was established in the US in 1911)
- Instructions were very simple, which made hardware implementation easier, but this hindered the creation of complex programs.

## Alan Turing (UK)

- Theorized the possibility of computing machines capable of performing *any* conceivable mathematical computation as long as this was representable as an *algorithm*
  - Called “*Turing Machines*” (1936) – ideas live on today...
  - Lead the effort to create a machine to successfully decipher the German “Enigma Code” during World War II



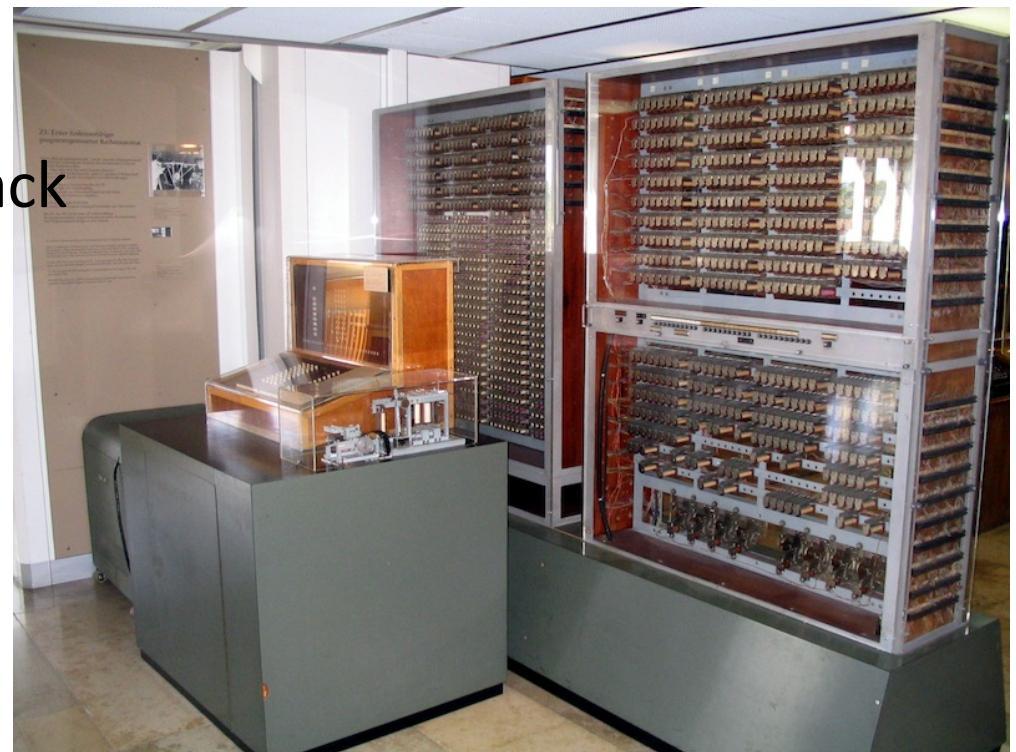
A. Turing (1912 – 1954)

# Zuse Z3 (1941)

---

- Built by Konrad Zuse in wartime Germany using 2000 relays
- Could do *floating-point* arithmetic with hardware
- 22-bit word length ; clock frequency of about 4–5 Hz!!
- 64 words of memory!!!
- Two-stage pipeline
  - 1) fetch & execute, 2) writeback
- No conditional branch
- Programmed via paper tape

*Replica of the Zuse Z3 in the  
Deutsches Museum, Munich*



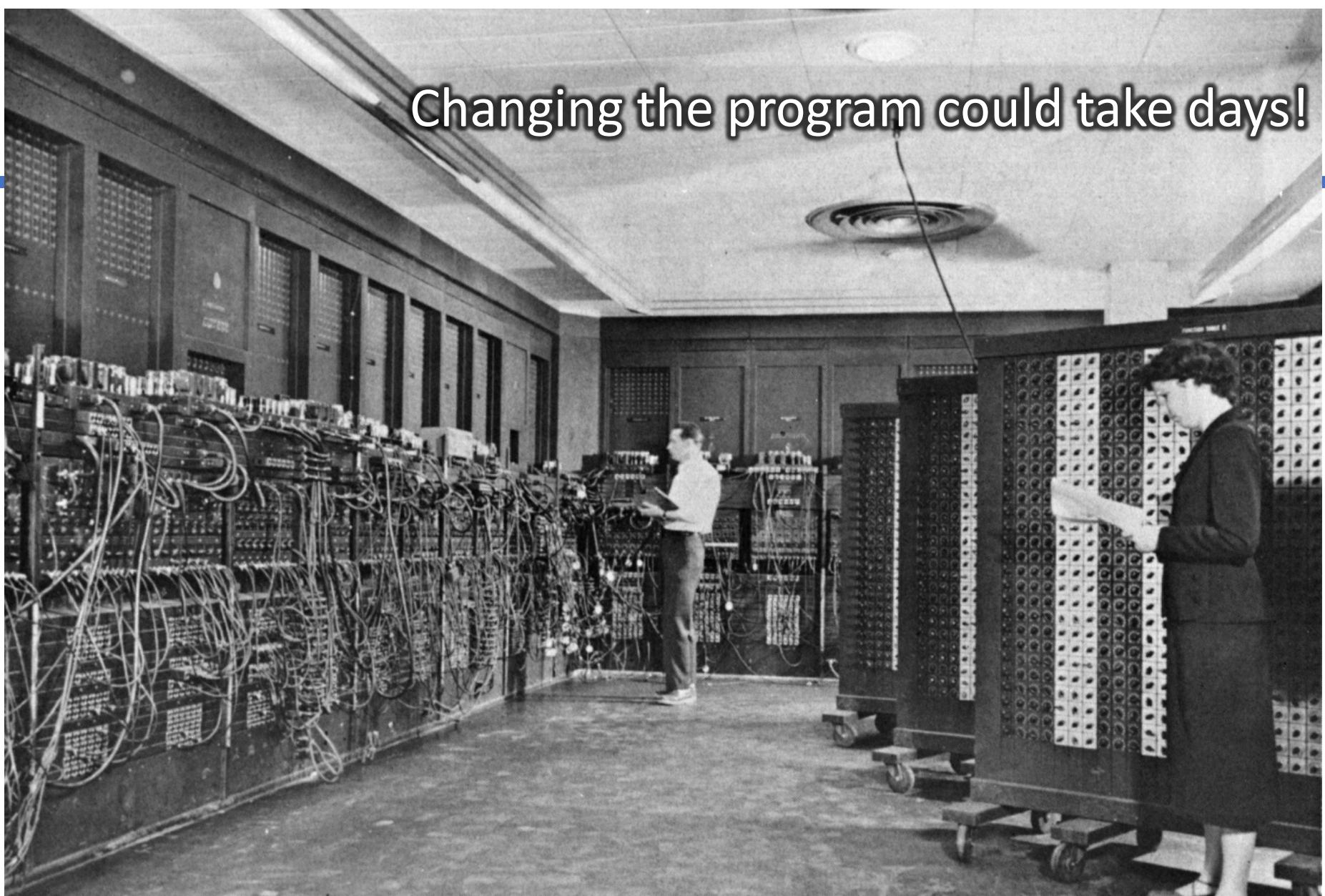
[Venusianer, Creative Commons BY-SA 3.0 ]

# ENIAC (1946)

---

- First electronic general-purpose computer
- Constructed during WWII to calculate firing tables for US Army
  - Trajectories (for bombs) computed in 30 seconds instead of 40 hours
  - Was very fast for its time – started to replace human “computers”
- Used vacuum tubes (transistors hadn’t been invented yet)
- Weighed **30 tons**, occupied **1800 sq ft**
- It used **160 kW** of power (about 3000 light bulbs worth)
- It cost **\$6.3 million** in today’s money to build.
- Programmed by plugboard and switches, time consuming!
- As a result of large number of tubes, it was often broken  
(5 days was longest time between failures!)

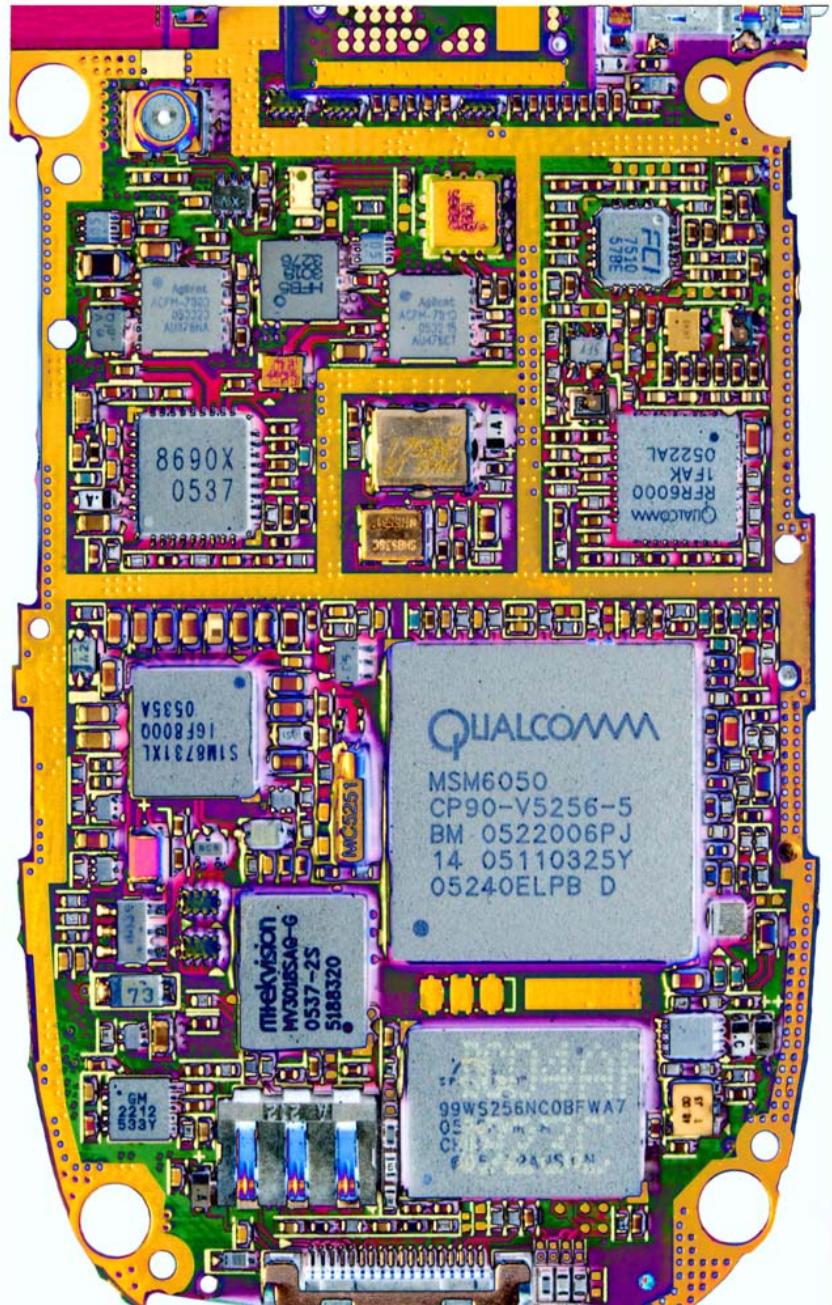
Changing the program could take days!



[Public Domain, US Army Photo]

Comparing today's cell phones  
(with dual CPUs), with ENIAC,  
we see they

cost 17,000X less  
are 40,000,000X smaller  
use 400,000X less power  
are 120,000X lighter  
AND...  
**are 1,300X more powerful.**



# EDVAC (1951)

---

- ENIAC team started discussing *stored-program concept* to speed up programming and simplify machine design
- Based on ideas by John von Neumann & Herman Goldstine
- Still the basis for our general CPU architecture today

# Commercial computers: BINAC (1949) and UNIVAC (1951) at EMC

---

- Eckert and Mauchly left academia and formed the Eckert-Mauchly Computer Corporation (EMC)
- World's first commercial computer was BINAC which didn't work...
- Second commercial computer was UNIVAC
  - Famously used to predict presidential election in 1952
  - Eventually 46 units sold at >\$1M each

# IBM 650 (1953)

- The first mass-produced computer
- Low-end system aimed at businesses rather than scientific enterprises
- Almost 2,000 produced



[Cushing Memorial Library and Archives, Texas A&M,  
Creative Commons Attribution 2.0 Generic ]

# Improvements in C.A.

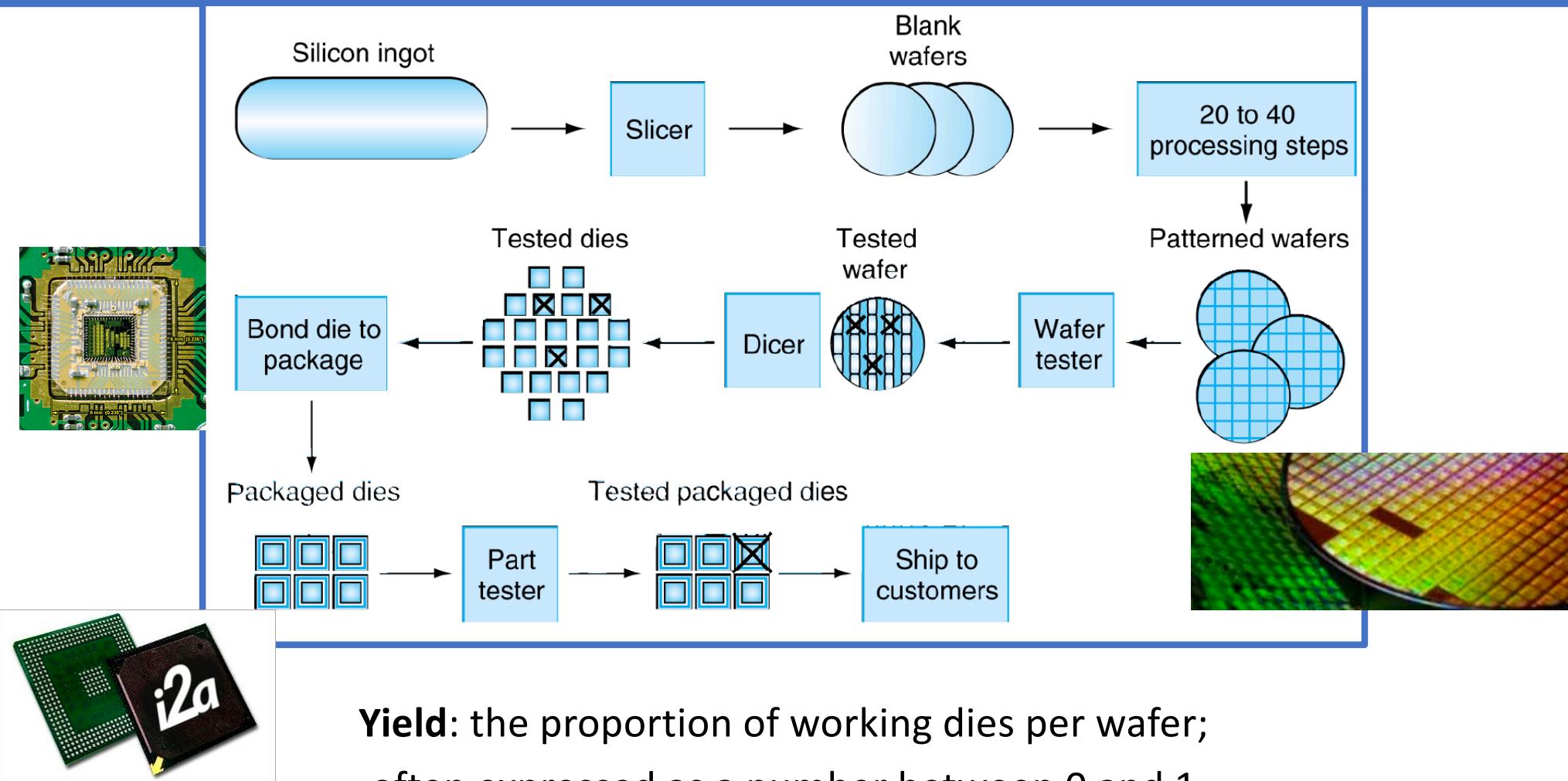
---

- IBM 650's instruction set architecture (ISA)
  - 44 instructions in base instruction set, expandable to 97 instructions
- Hiding instruction set completely from programmer using the concept of ***high-level languages*** like Fortran (1956), ALGOL (1958) and COBOL (1959)
  - Allowed the use of stack architecture, nested loops, recursive calls, interrupt handling, etc...

*Adm. Grace Hopper (1906 – 1992),  
inventor of several High-level language concepts*



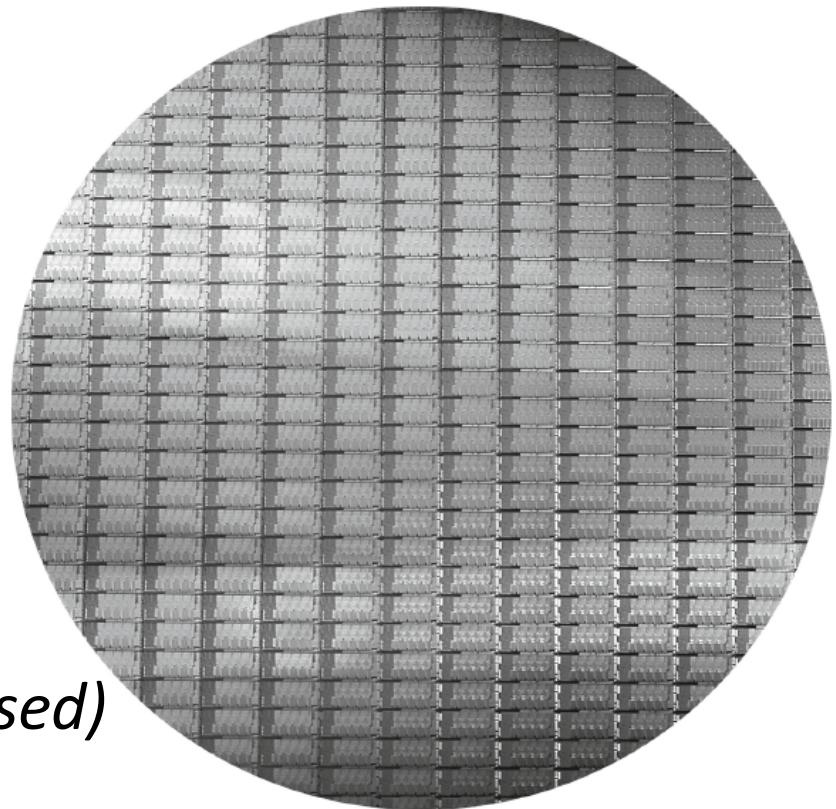
# Manufacturing ICs



# Example: Intel Core i7 Wafer

---

- 300mm (diameter) wafer
- 280 chips
- Each chip is 20.7 mm x 10.5 mm
- 32nm CMOS technology  
*(the size of the smallest piece of logic  
and the type of Silicon semiconductor used)*



# Costs of Manufacturing ICs

---

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

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

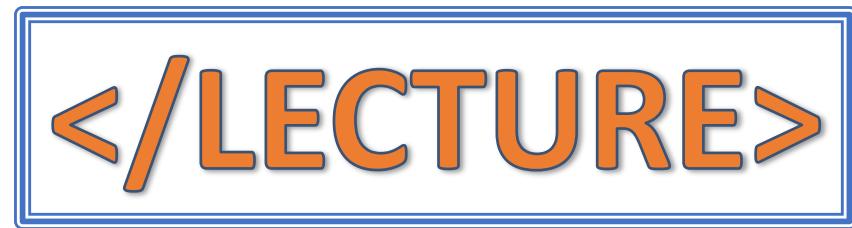
$$\text{Yield} = \frac{1}{(1 + (\text{Defects per area} \times \text{Die area}/2))^2}$$

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

# YOUR TO-DOs for the Week

---

- Do your reading for next class (see syllabus)
- Work on Assignment #1 for lab (*lab01*)
  - Meet up in the lab this Friday
  - Do the lab assignment
  - You have to submit it as a PDF using *Gradescope*
  - Due on **Wednesday, 1/15, by 11:59:59 PM**



# 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?

# Performance Measures

---

- Execution Time: Total response time, including **EVERYTHING**
  - CPU time (processing), I/O use, OS overhead, any idle time
  - This determines **system performance**
- CPU time:
  - Time spent just processing a given job (discounts I/O time, OS time, etc...)
  - CPU time = *user* CPU time + *system* CPU time
- Define Performance =  $1/\text{Execution Time}$
- Relative performance
  - The performance of system A vs performance of system B, ie.  $P_A / P_B$

# CPU Clocking

---

- Most digital hardware today operates to 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**: clock *rate* or *cycles per second*
  - e.g.  $4.0 \text{ GHz} = 4000 \text{ MHz} = 4.0 \times 10^9 \text{ Hz}$
- Hertz (Hz) is “cycles per second”, so  
**clock freq. = 1 / clock period**

# Useful Prefixes (Multipliers) to Know

Prefix	Symbol	Multiplier	In words...	Scientific Notation
Kilo	k	1,000	thousand	$10^3$
Mega	M	1,000,000	million	$10^6$
Giga	G	1,000,000,000	billion	$10^9$
Tera	T	1,000,000,000,000	trillion	$10^{12}$
Peta	P	1,000,000,000,000,000	quadrillion	$10^{15}$

Prefix	Symbol	Multiplier	In words...	Scientific Notation
milli	m	0.001	thousandth	$10^{-3}$
micro	m	0.000001	millionth	$10^{-6}$
nano	n	0.000000001	billionth	$10^{-9}$
pico	p	0.000000000001	trillionth	$10^{-12}$

# CPU Time

CPU Time = CPU Clock Cycles × Clock Cycle Time

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

- Performance improved (i.e. make CPU Time **less**) by
  - Reducing number of clock cycles
  - Increasing clock rate
  - Hardware designer must often trade off clock rate against cycle count
- Example: it took the CPU 1000 cycles to run the program. The clock cycle time (i.e. period) is 10 ns, so the CPU time is:  
 $1000 \times 10 \text{ ns} = 10000 \text{ ns} = 10 \mu\text{s}$ , or  $10 \times 10^{-6} \text{ s}$

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

# 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 (**CPI**)
  - Determined by CPU hardware
  - If different instructions have different CPI, then *Average CPI* is affected by instruction mix
- Example: *next slide*

# CPI Example

---

- Computer A: Cycle Time = 250 ps, CPI = 2.0
- Computer B: Cycle Time = 500 ps, CPI = 1.2
- Same Instruction Set Architecture (ISA)
- **Which is faster?**
  - CPU Time = Instruction Count x CPI x Cycle Time
  - $\text{CPU\_Time\_A} = \text{NI} \times 2.0 \times 250 \times 10^{-12} \text{ s} = \text{NI} \times 500 \times 10^{-12} \text{ s}$
  - $\text{CPU\_Time\_B} = \text{NI} \times 1.2 \times 500 \times 10^{-12} \text{ s} = \text{NI} \times 600 \times 10^{-12} \text{ s}$
  - **So, CPU A is faster than CPU B**
- **By how much is it faster?**
  - Relative Performance =  $\text{NI} \times 600 \times 10^{-12} \text{ s} / \text{NI} \times 500 \times 10^{-12} \text{ s} = \underline{\underline{1.2}}$
  - **So, CPU A is 1.2 times faster than B (or you could say it's 20% faster)**

# CPI Example using Weighted Classes

---

- An instruction class = instruction type
  - e.g. arithmetic type vs. branching type vs. jump type, etc...
- A CPU compiles 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, so Clock Cycles =  $2 \times 1 + 1 \times 2 + 2 \times 3 = 10$
- So, **Avg. CPI = 10/5 = 2.0**
- **Sequence 2:** IC = 6, so Clock Cycles =  $4 \times 1 + 1 \times 2 + 1 \times 3 = 9$
- So, **Avg. CPI = 9/6 = 1.5**

# Other Factors to CPU Performance: Power Consumption

---

- Market trends DEMAND that power consumption of CPUs keep decreasing
- **Power = Capacitive Load x Voltage<sup>2</sup> x Clock Frequency**
- So:
  - Decreasing Voltage helps power, but can make individual logic go slower!
  - Increasing clock frequency helps performance, but increases power!
- It's a dilemma that has contributed to Moore's Law "plateau"

# Other Factors to CPU Performance: Multiple Processors

---

- 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

# Pitfalls: Amdahl's Law

---

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

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

- Your benchmark time is 100, 80 of which comes from a part of the CPU that you want to improve by a factor of  $n$ , so:

$$T_{\text{improved}} = (80 / n) + 20$$

- If you wanted to improve your overall  $T$  by a factor of 2 (i.e. drop total from 100 to 50), then you'd need to make  $n = 2.7$

because  $50 = (80 / 2.7) + 20$     ...ok...

- Keep that up! Let's go for  $n = 5$ , so drop total from 100 to 20:

$$20 = (80 / 5) + 20$$
    ...uh... can't do that... 😞

# Pitfalls: Idle Power

---

- Simply put, CPUs will still draw *disproportionate* power when idling.
- Example, even when operating at 10% load, the i7 will draw 47% of the power
- Becomes a problem when dealing with large scale implementations, like data centers (Google, Facebook, Amazon, etc...)
- Design challenge: design processors to draw power more proportional to load (requires Physics-level approach, tho...)

# Pitfall: MIPS as a Performance Metric

---

- Note: We're NOT talking about **MIPS the processor type!!!!**
- **MIPS** (millions of instructions per second) is a popular performance metric, *HOWEVER...*
- Doesn't account for
  - Differences in ISAs between computers  
(some ISAs may be more efficient than others)
  - Differences in complexity between instructions (weighted CPIs)

# YOUR TO-DOs for the Week

---

- Do your reading for next class (see syllabus)
- Work on Assignment #1 for lab (*lab01*)
  - Meet up in the lab this Friday
  - Do the lab assignment
  - You have to submit it as a PDF using *Gradescope*
  - Due on **Wednesday, 1/15, by 11:59:59 PM**

