

# **ECE 485/585**

## **Computer Organization and Design**

### **Lecture 1: Introduction**

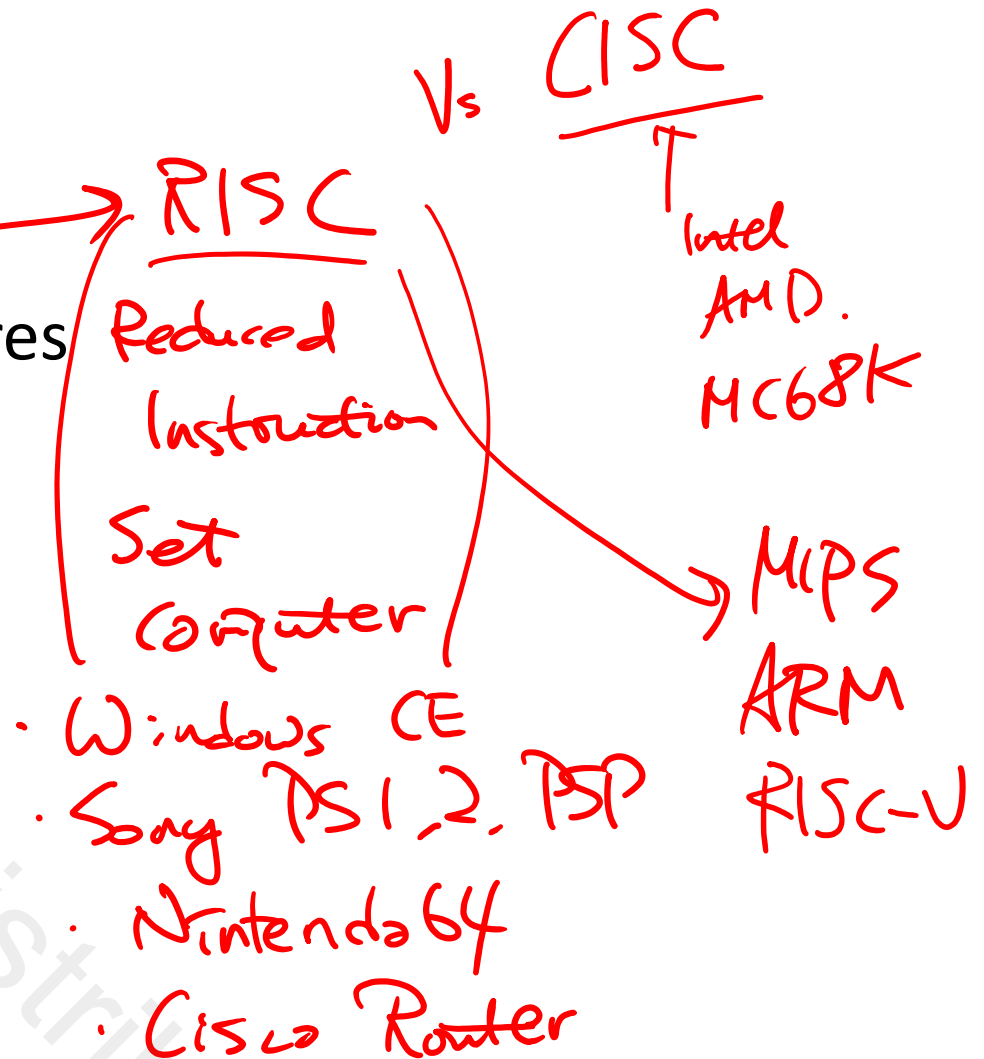
#### **Fall 2022**

**Won-Jae Yi, Ph.D.**

**Department of Electrical and Computer Engineering**  
**Illinois Institute of Technology**

# Topics

- Introduction to Computers
- Computer Performance Measures
- MIPS Instruction Set
- ALU Design
- Datapath Design
- Control Unit design
- Exceptions
- Pipelined Datapath Design
- Hazards
- Cache
- Memory
- Introduction to Parallel Processors



# Introduction

- This course is all about how computers operate
- What do we mean by a computer??
  - Different types: desktop, server, embedded system (e.g., smartphone)
  - Different usage: automobiles, graphics, finance, shopping...
  - Different manufacturers: Intel, AMD, Apple, Microsoft, HP, Samsung...
  - Different underlying technologies and different costs
- Best way to learn
  - Focus on specific instance and learn how it works
  - While learning general principles and historical perspectives

# What will you learn?

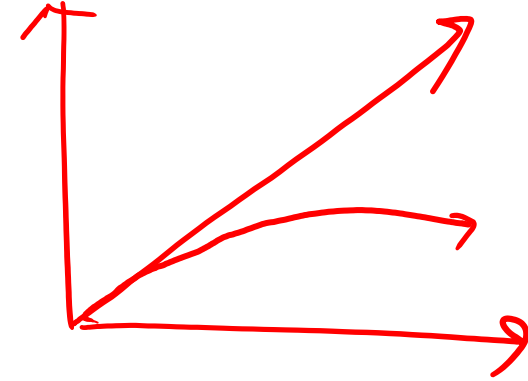
- How programs written in a high-level language are executed in hardware
- Interface between software and hardware
- What determines the performance of a program and how can a programmer improve it?
- How can a hardware engineer improve the performance?

# Why learn this topic?

- You want to call yourself a “computer engineer”
- You want to build hardware or software people use (need performance)
- You need to make a purchasing decision or offer “expert” advice

# 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



# 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

Edge Computing.

Vs.

Cloud Computing



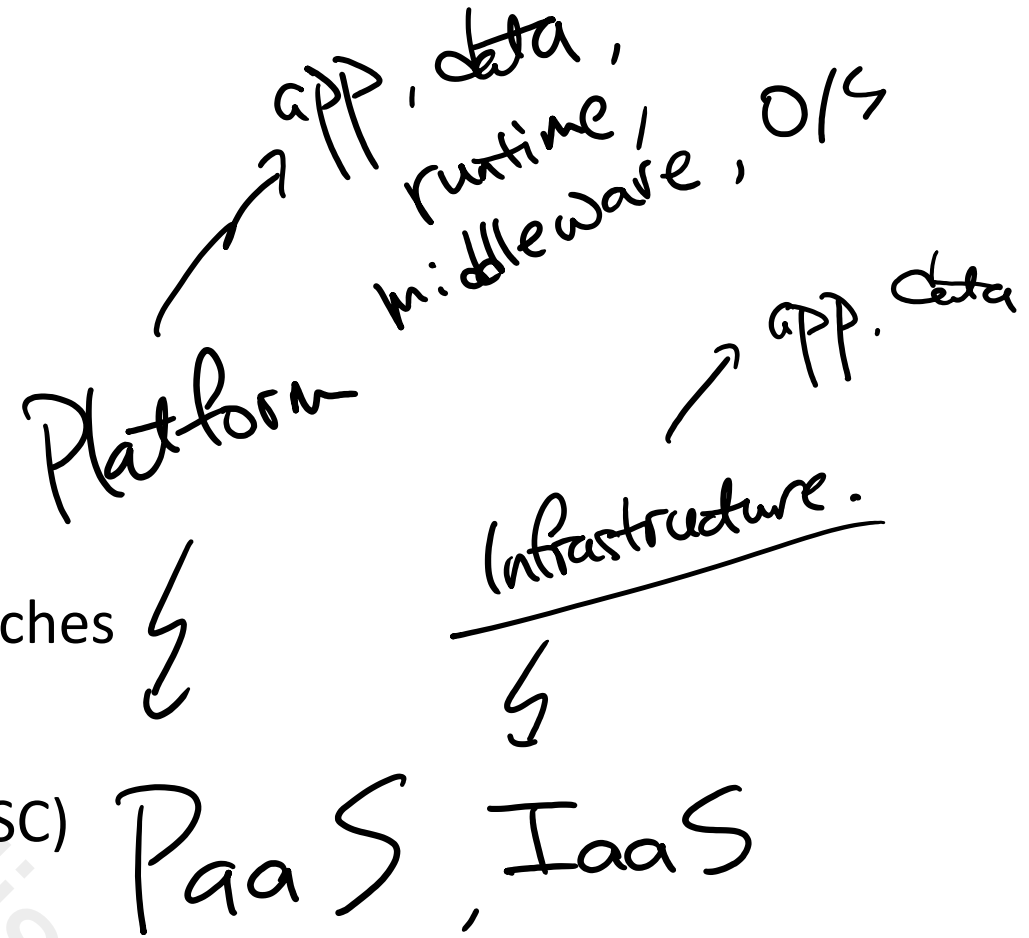
# Classes of Computers

- Supercomputers
  - High-end scientific and engineering calculations
  - Highest capability but represent a small fraction of the overall computer market
- Embedded systems/computers/devices
  - Hidden as components of systems
  - Stringent power/performance/cost constraints



# PostPC Era

- Personal Mobile Device (PMD)
  - Battery operated
  - Connects to the Internet
  - Hundreds of dollars
  - Smartphones, tablets, smart watches
- Cloud computing
  - Warehouse Scale Computers (WSC)
  - Software as a Service (SaaS)
  - Portion of software run on a PMD and a portion run in the Cloud
  - Amazon, Google, Microsoft



# 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

) ch 2, 3

) ch 4, 5, 6

# 8 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

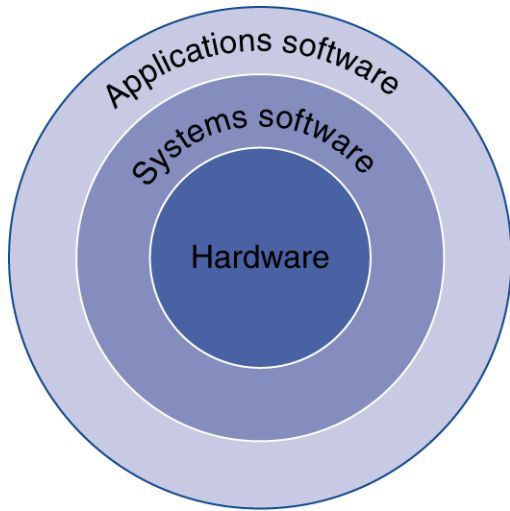


L1 L2 L3  
Cache

RAM

HDD/SSD

# 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

- 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

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  
program  
(for MIPS)

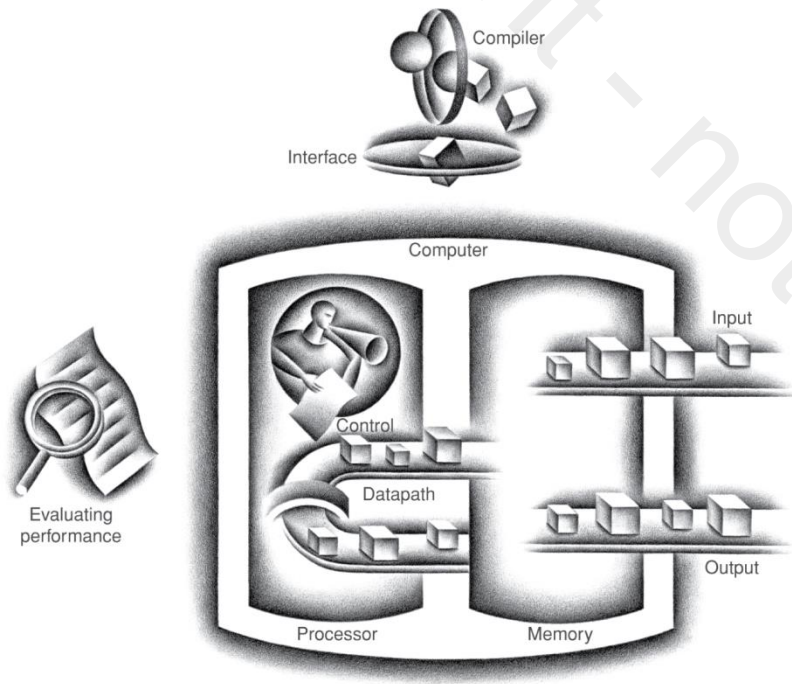
```
swap:
  muli $2, $5, 4
  add  $2, $4, $2
  lw   $15, 0($2)
  lw   $16, 4($2)
  sw   $16, 0($2)
  sw   $15, 4($2)
  jr   $31
```

Assembler

Binary machine  
language  
program  
(for MIPS)

```
00000000101000010000000000011000
00000000000110000001100000100001
10001100011000100000000000000000
100011001111001000000000000000100
10101100111100100000000000000000
10101100011000100000000000000100
0000001111100000000000000001000
```

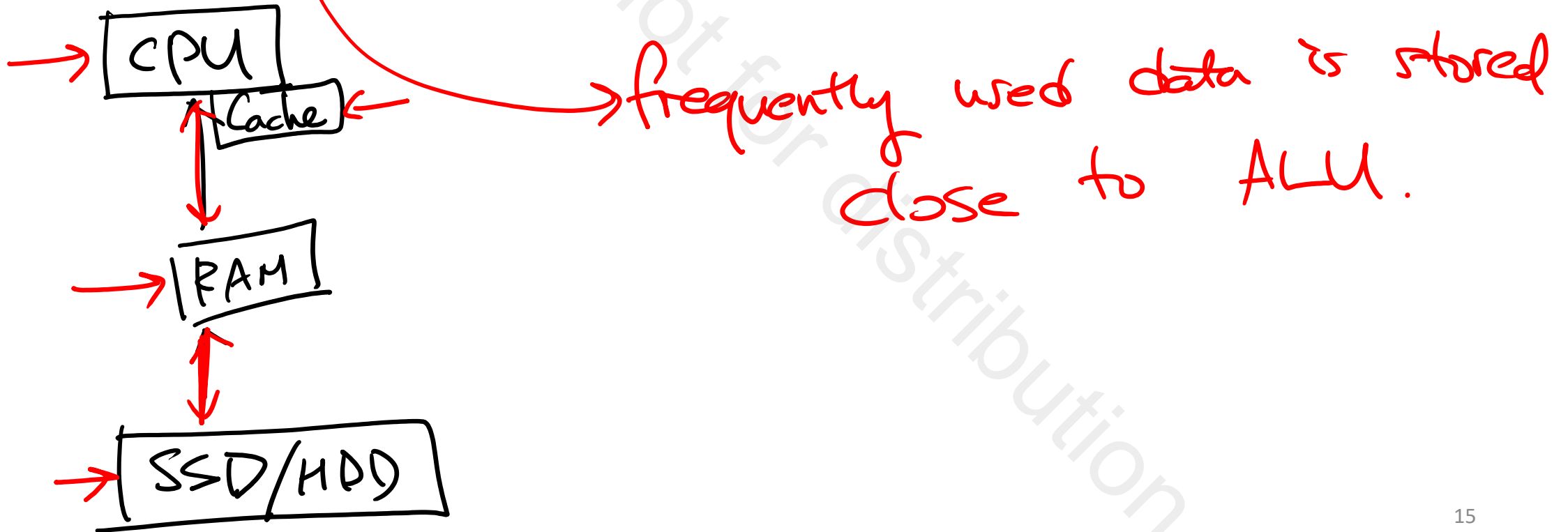
# Components of a Computer



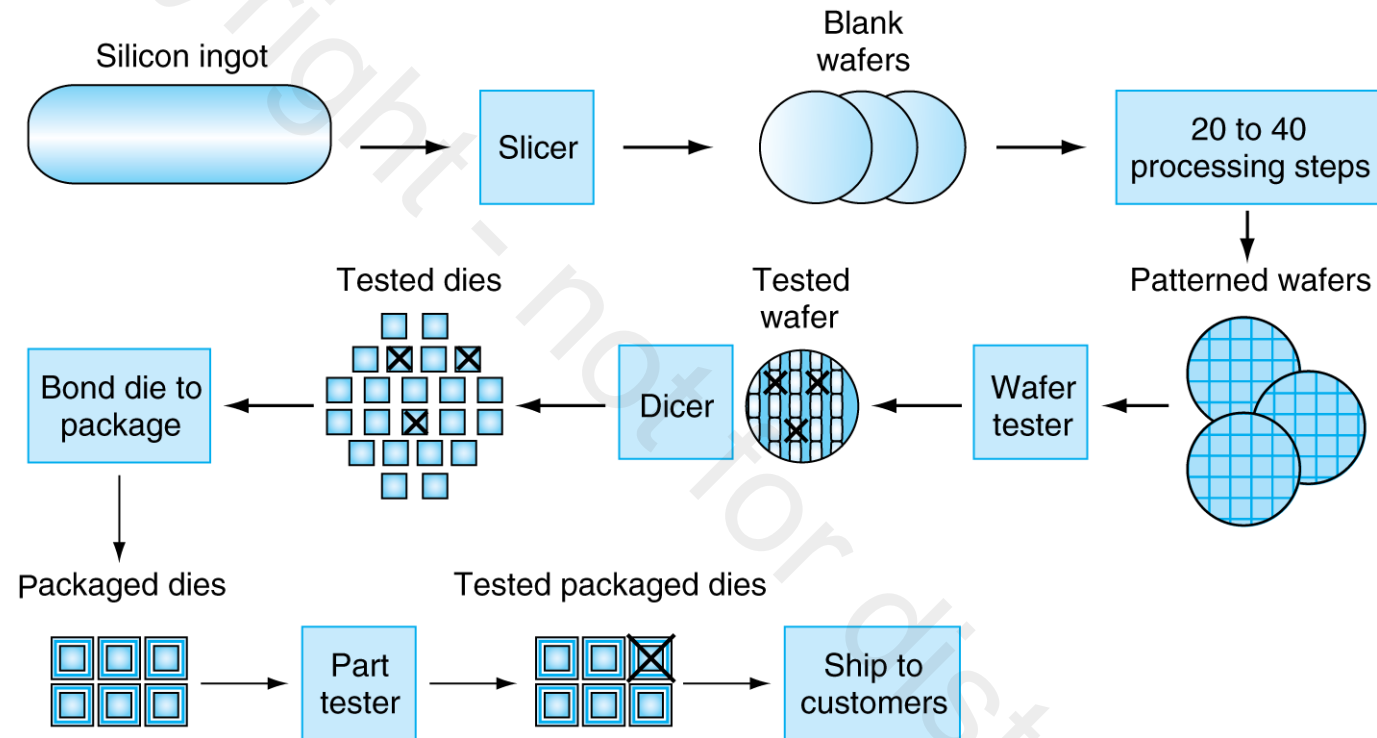
- Same components for all kinds of computer
  - Desktop, server, embedded
- Input/output includes
  - User-interface devices
    - Display, keyboard, mouse
  - Storage devices
    - Hard disk, SSD, Magnetic tapes
  - Network adapters
    - For communicating with other computers (machines or human)

# Inside the Processor (CPU)

- Datapath: performs operations on data
- Control: sequences datapath, memory, ...
- Cache memory
  - Small fast SRAM memory for immediate access to data



# Manufacturing ICs

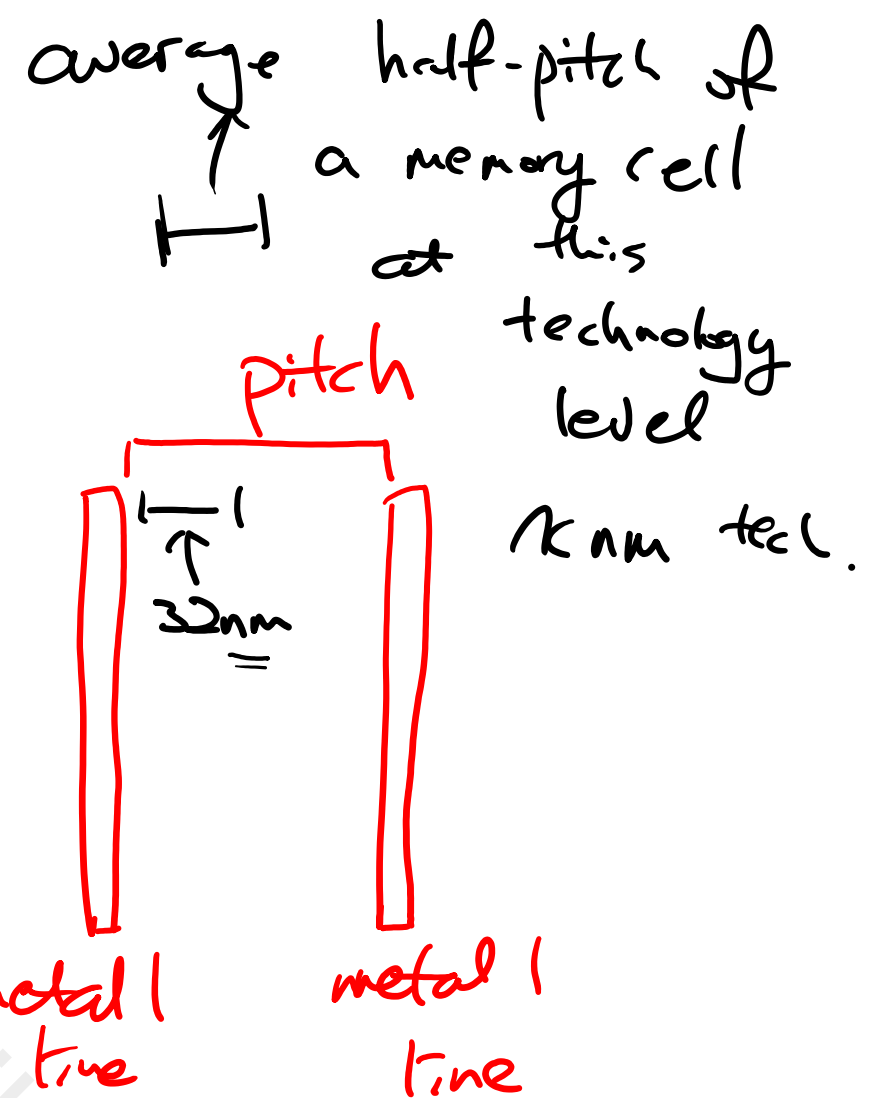
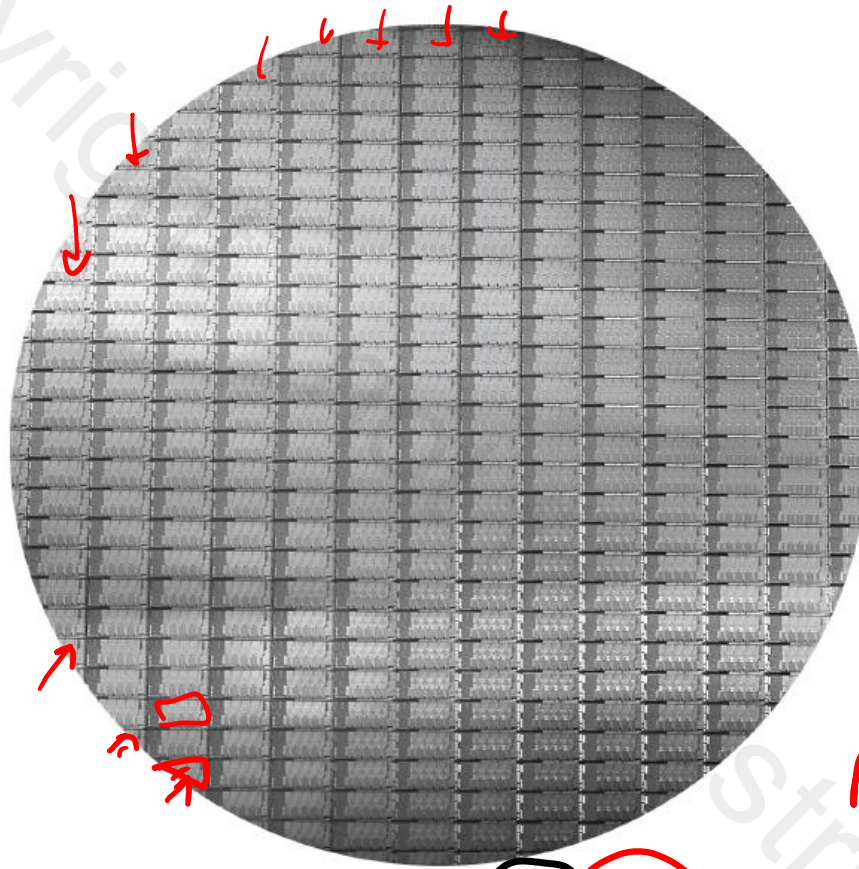


- Yield: proportion of working dies per wafer

<https://youtu.be/Q5paWn7bFg4>



# Intel Core i7 Wafer



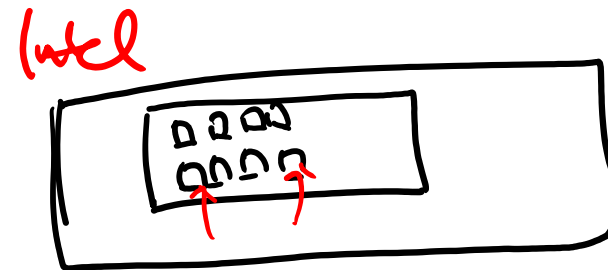
- 300mm wafer, 280 chips, 32nm technology
- Each chip is 20.7 x 10.5 mm

Intel : one chip.

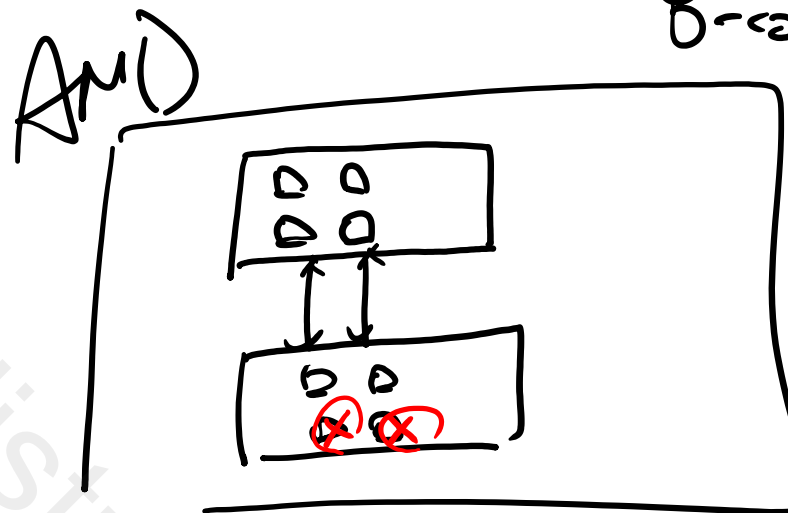
AMD :

4  
6  
8

8-core



8-cores.



6 core

# Integrated Circuit Cost

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

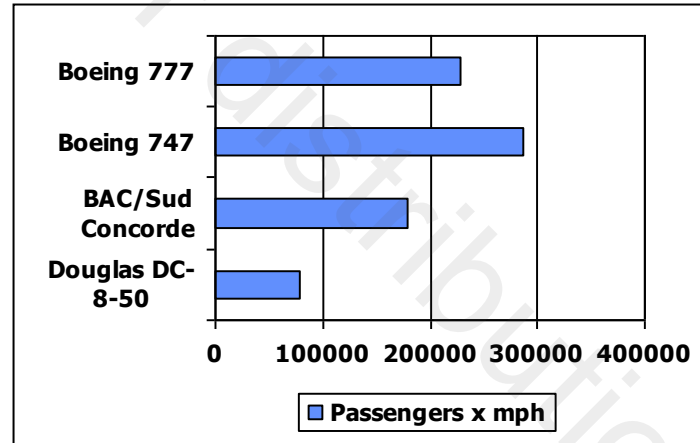
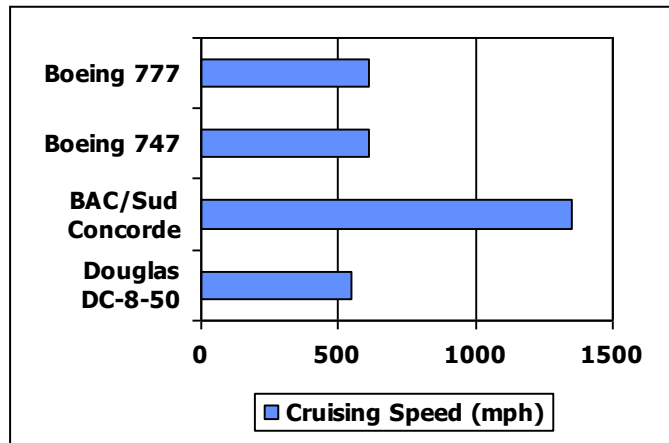
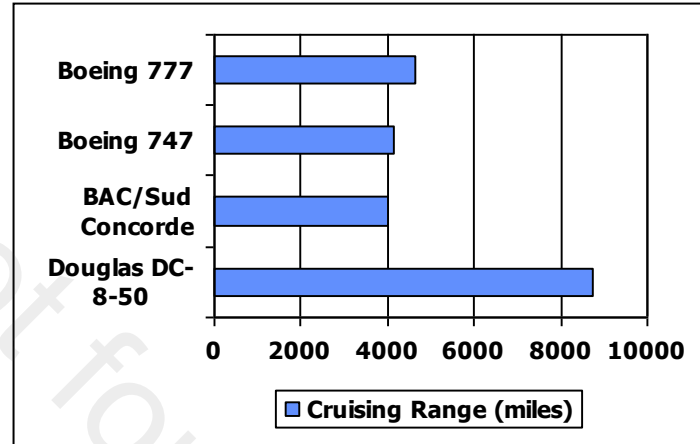
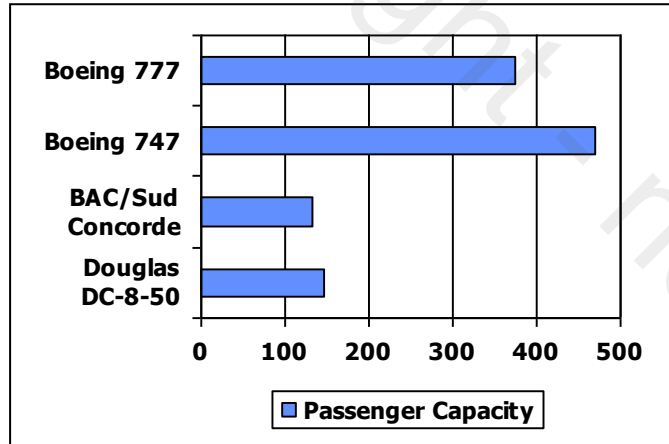
$$\text{Dies per wafer} \approx \text{Wafer area} / \text{Die area}$$

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

- 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

# Defining Performance

- Which airplane has the best performance?



# Performance Metrics

- Response time
- Throughput
- Relative performance
- Execution time
- CPU time
- • Instruction count IC
- • Clocks per Instruction (CPI)

# Response Time and Throughput

- Processor times
- I/O times
- Memory times
- Propagation delays
- Response time <sup>"minimized"</sup> → time difference between submission of a request until the response begins to be received.
  - How long it takes to do a task
- Throughput <sup>"maximized"</sup> → number of processes completed per unit time.
  - 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...

# Relative Performance

- Define Performance =  $1/\text{Execution Time}$
- "X is  $n$  time faster than Y"

→ time spent by the job/task actively using processor resources

→ 
$$\frac{\text{Performance}_X^A}{\text{Performance}_Y^B} = \frac{\text{Execution time}_Y^B}{\text{Execution time}_X^A} = n$$

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

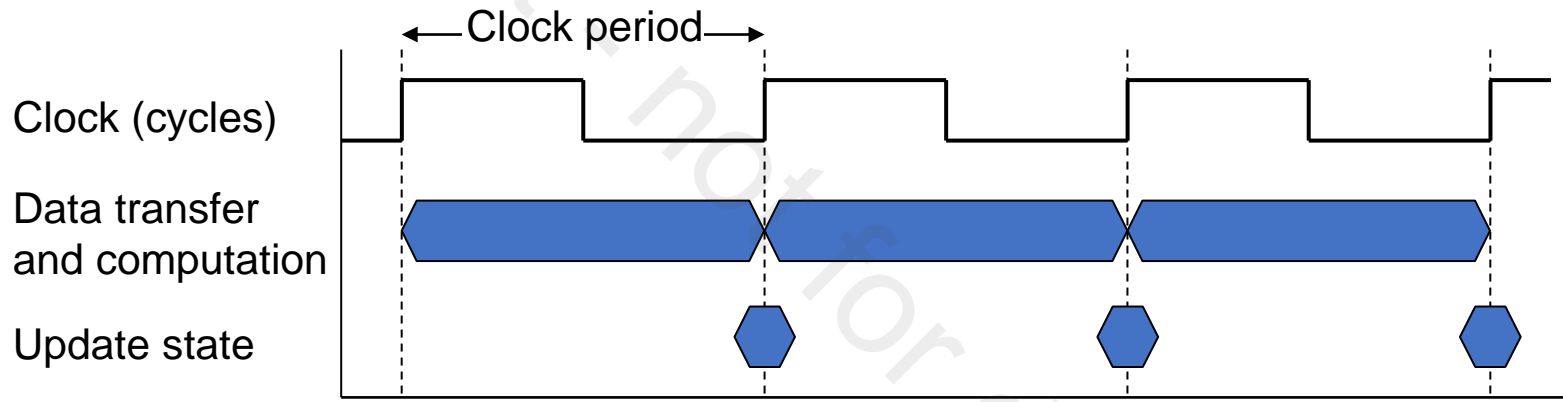
# Measuring Execution Time

- Elapsed time = Total Execution Time.
  - Total response time, including all aspects
    - **Processing, I/O, OS overhead, idle time**
  - Determines system performance
- CPU time ⇒ ONLY TIME SPENT AT CPU
  - Time spent processing a given job
    - **Does not count I/O time, other jobs' shares**
  - Comprises user CPU time and system CPU time
  - Different programs are affected differently by CPU and system performance
    - Program
    - OS



# CPU Clocking

- Operation of digital hardware governed by a constant-rate clock



- Clock period: duration of a clock cycle

- e.g., 250ps = 0.25ns =  $250 \times 10^{-12}\text{s}$

- Clock frequency (rate): cycles per second

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

$$T = \frac{1}{f}$$

$$f = \frac{1}{T}$$

# CPU Time (CPU Exec. Time)

amount of time  
for one  
clock  
period

$$\text{CPU Time} = \text{CPU Clock Cycles} \times \text{Clock Cycle Time}$$

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

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

# 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?

$$\text{CPU time} = \frac{\text{clock cycle}}{\text{clock rate}}$$

$$\text{clock cycle}_A = \text{CPU time}_A \times \text{clock rate}_A$$

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

$$\text{Clock Cycles}_A = \text{CPU Time}_A \times \text{Clock Rate}_A$$

$$= 10s \times 2\text{GHz} = 20 \times 10^9$$

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

# Instruction Count and CPI

$\text{clock cycle time} \times \text{CPI}$

$$\text{Clock Cycles} = \text{Instruction Count} \times \text{Cycles per Instruction}$$
$$\Rightarrow \text{CPU Time} = \text{Instruction Count} \times \text{CPI} \times \text{Clock Cycle Time}$$
$$= \frac{\text{Instruction Count} \times \text{CPI}}{\text{Clock Rate}}$$

$$\text{CPU time} = IC \times CPI \times T_c$$
$$= \frac{IC \times CPI}{f_c}$$

- Instruction Count for a program
  - Determined by program, ISA (Instruction Set Architecture) and compiler
- Average cycles per instruction  $CPI_{avg}$ 
  - Determined by CPU hardware
  - If different, instructions have different CPI
    - Average CPI affected by instruction mix

# CPI Example

$$\Rightarrow \text{CPU}_{\text{time}} = \text{IC} \times \text{CPI} \times T_c$$

↑                      ↑

- Computer A: Cycle Time = 250ps, CPI = 2.0
- Computer B: Cycle Time = 500ps, CPI = 1.2
- Same ISA    IC is same
- Which is faster, and by how much?

$$\text{CPU Time}_A = \text{Instruction Count} \times \text{CPI}_A \times \text{Cycle Time}_A$$

$$= 1 \times 2.0 \times 250\text{ps} = 1 \times 500\text{ps}$$

A is faster...

$$\text{CPU Time}_B = \text{Instruction Count} \times \text{CPI}_B \times \text{Cycle Time}_B$$

$$= 1 \times 1.2 \times 500\text{ps} = 1 \times 600\text{ps}$$

$$\frac{\text{CPU Time}_B}{\text{CPU Time}_A} = \frac{1 \times 600\text{ps}}{1 \times 500\text{ps}} = 1.2$$

...by this much

# CPI in More Detail

- <sup>CPI</sup> 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

# CPI Example

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

Class	A	B	C
CPI for class	<u>1</u>	<u>2</u>	<u>3</u>
<del>IC</del> IC in sequence 1	<u>2</u>	<u>1</u>	<u>2</u>
IC in sequence 2	<u>4</u>	<u>1</u>	<u>1</u>

## ■ Sequence 1: IC = 5

### ■ Clock Cycles

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

A      B      C

$$= 10$$

$$\text{Avg. CPI} = 10/5 = 2.0$$

## ■ Sequence 2: IC = 6

### ■ Clock Cycles

$$= 4 \times 1 + 1 \times 2 + 1 \times 3$$

A      B      C

$$= 9$$

$$\text{Avg. CPI} = 9/6 = 1.5$$

# Performance Summary

$$\text{CPU Time} = \frac{\overset{\checkmark}{\text{IC}} \times \text{Instruction ns}}{\text{Program}} \times \frac{\text{Clock cycles}}{\overset{\text{CPI}}{\text{Instruction n}}} \times \frac{\text{Seconds}}{\overset{T_c}{\text{Clock cycle}}}$$

- Performance depends on
  - Algorithm: affects IC, possibly CPI
  - Programming language: affects IC, CPI
  - Compiler: affects IC, CPI
  - Instruction set architecture: affects IC, CPI,  $T_c$

