CMPE 200
Computer Architecture & Design
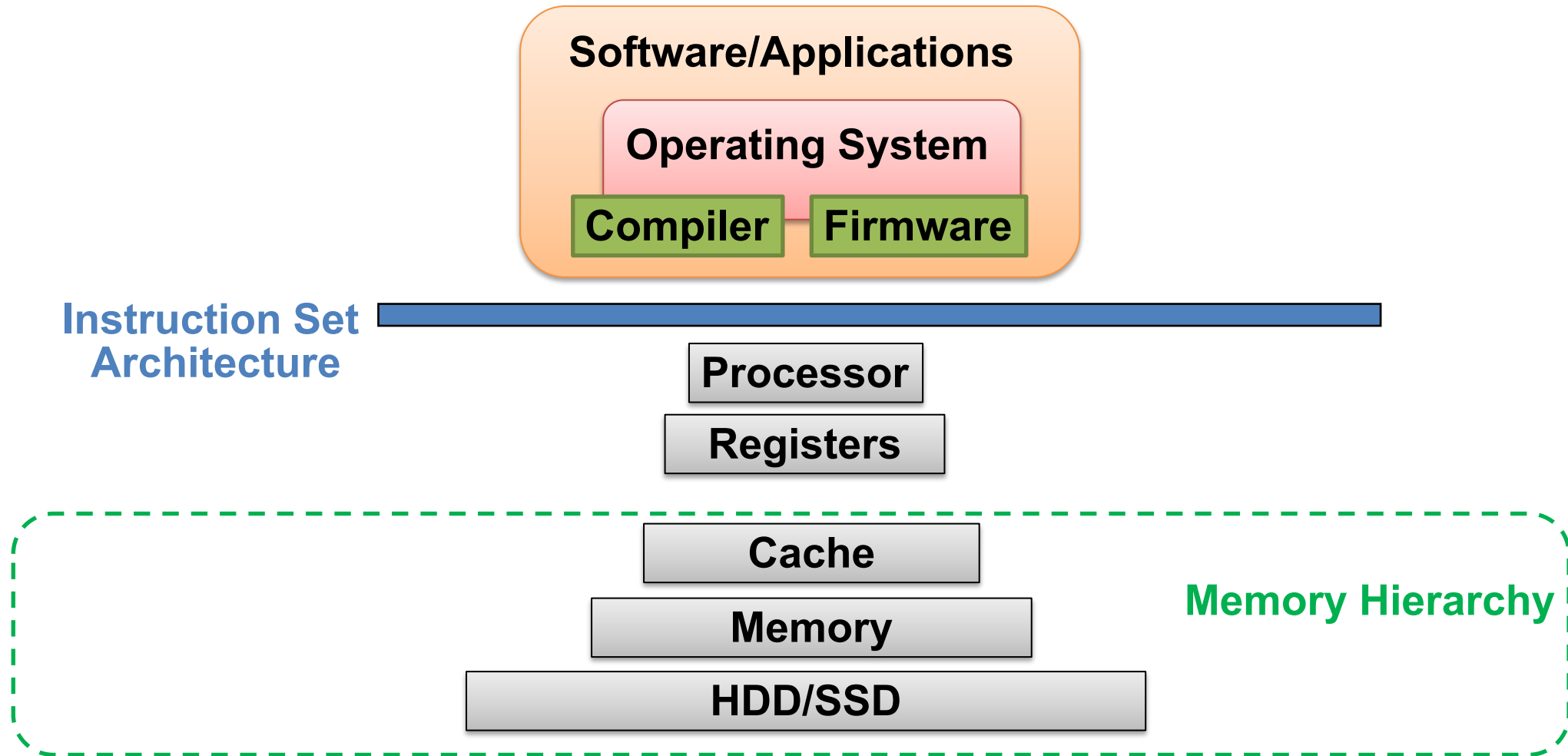
# Lecture 4.
# Memory Hierarchy (1)

Haonan Wang

SJSU

# Computer Architecture Overview

**Software/Applications**

**Operating System**

**Compiler**  **Firmware**

**Instruction Set Architecture**

**Processor**

**Registers**

**Cache**

**Memory**

**HDD/SSD**

**Memory Hierarchy**

SJSU   SAN JOSÉ STATE UNIVERSITY

# Memory Hierarchy

**Example:** a C program that reads two integer values from "file.txt" file and prints the sum of them.

```c
#include <stdio.h>
#include <string.h>

int numbers[2];

void myfunction(void)
{
    FILE *fp;
    int size = 2;
    int sum = 0;

    /* Open file for reading */
    fp = fopen("mynumbers.txt", "r");

    /* Read and display data */
    fread(numbers, sizeof(int), size, fp);
    fclose(fp);
    sum = numbers[0] + numbers[1];
    printf("Sum = %d\n", sum);
}

int main (void) {
    myfunction();
    return(0);
}
```

**Processor (CPU)**

Understands and executes each line of the code.

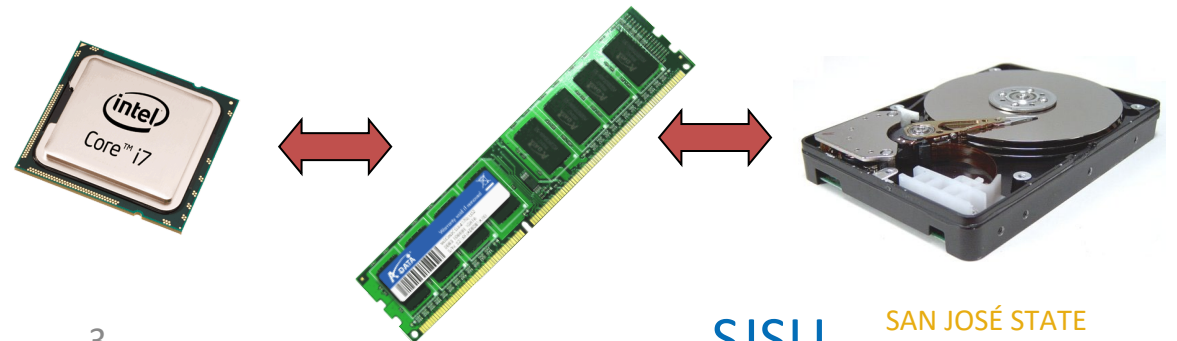Uses fast on-chip memories

**Memory (DRAM)**

Provides operands to CPU
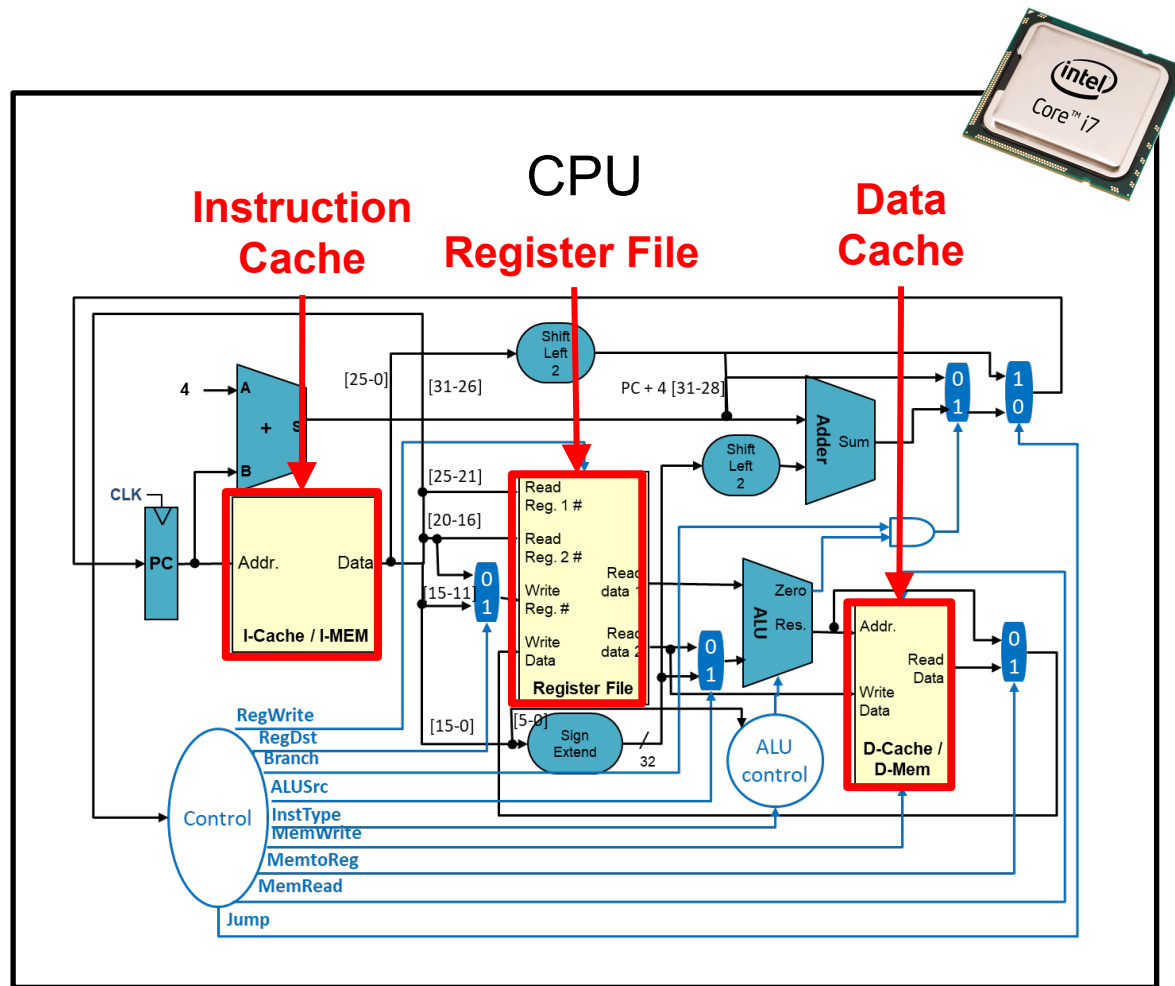
**(*fp, size, sum, numbers[2])**

**Storage (HDD)**

Provides file inputs and program code
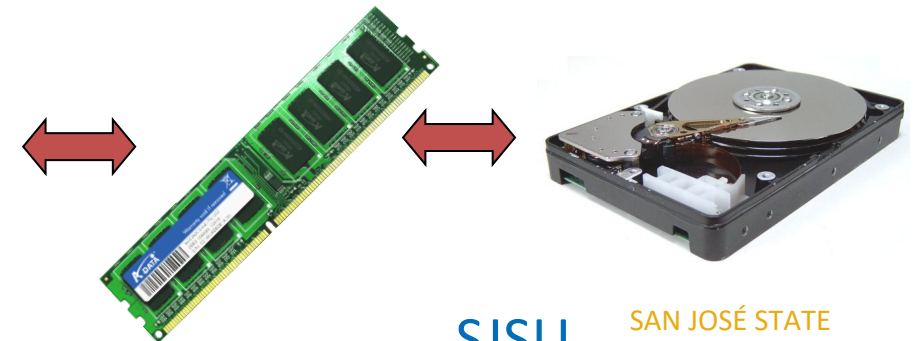
**(file.txt)**

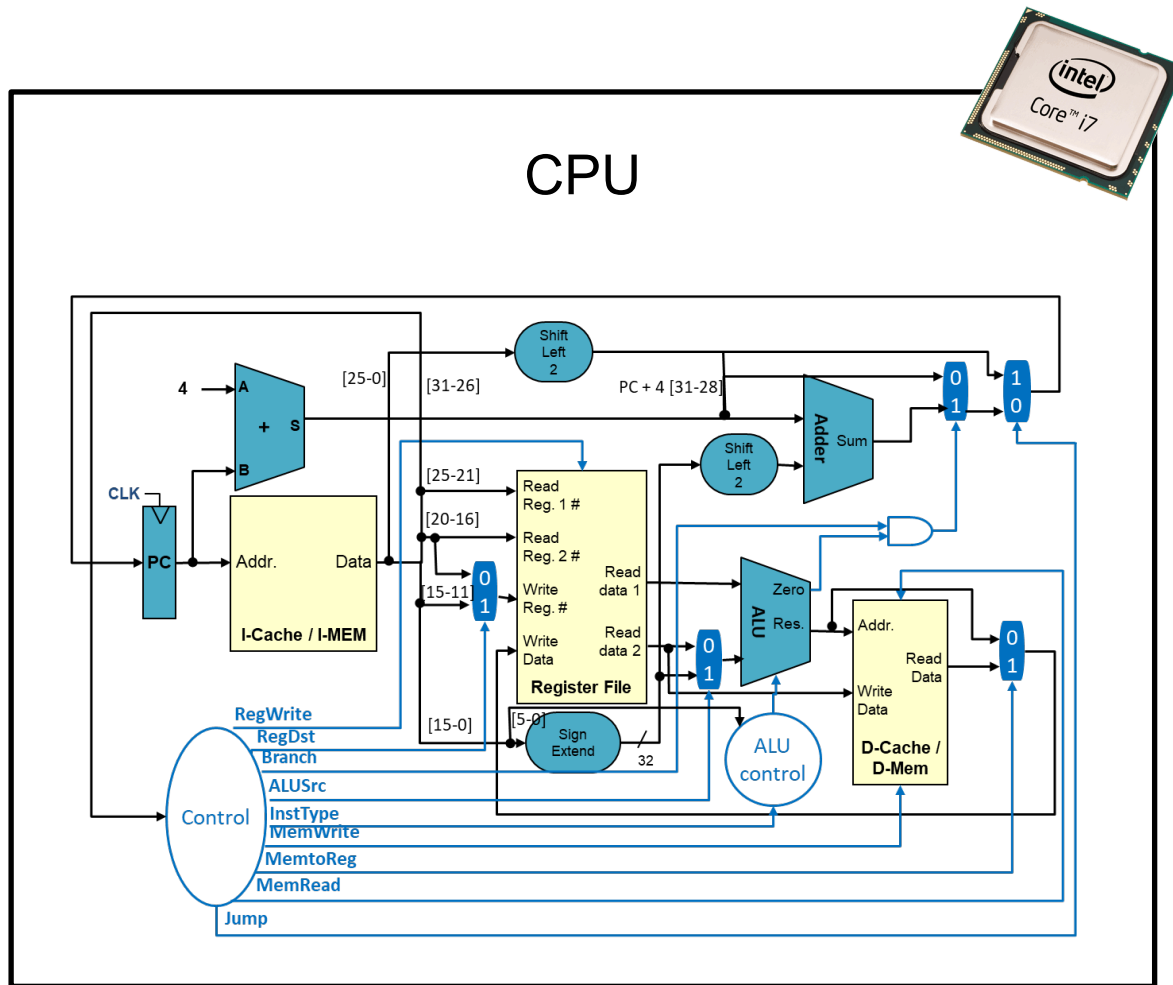SJSU   SAN JOSÉ STATE UNIVERSITY

# Memory Hierarchy



- **On-chip Memories (Memories inside of CPU)**
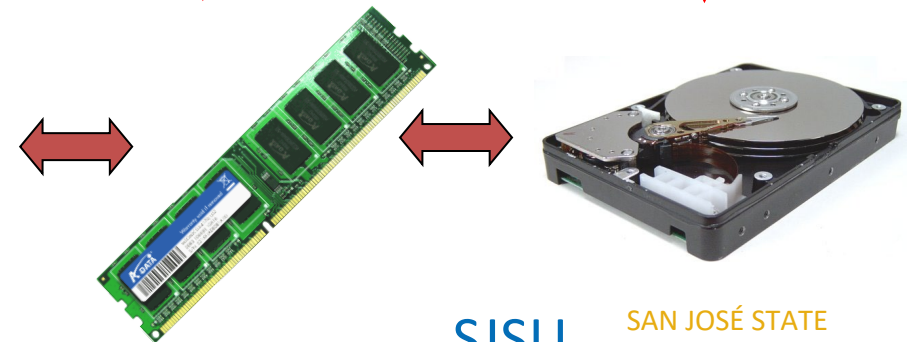  - Register file, Caches
  - Small but Fast

# Memory Hierarchy



CPU

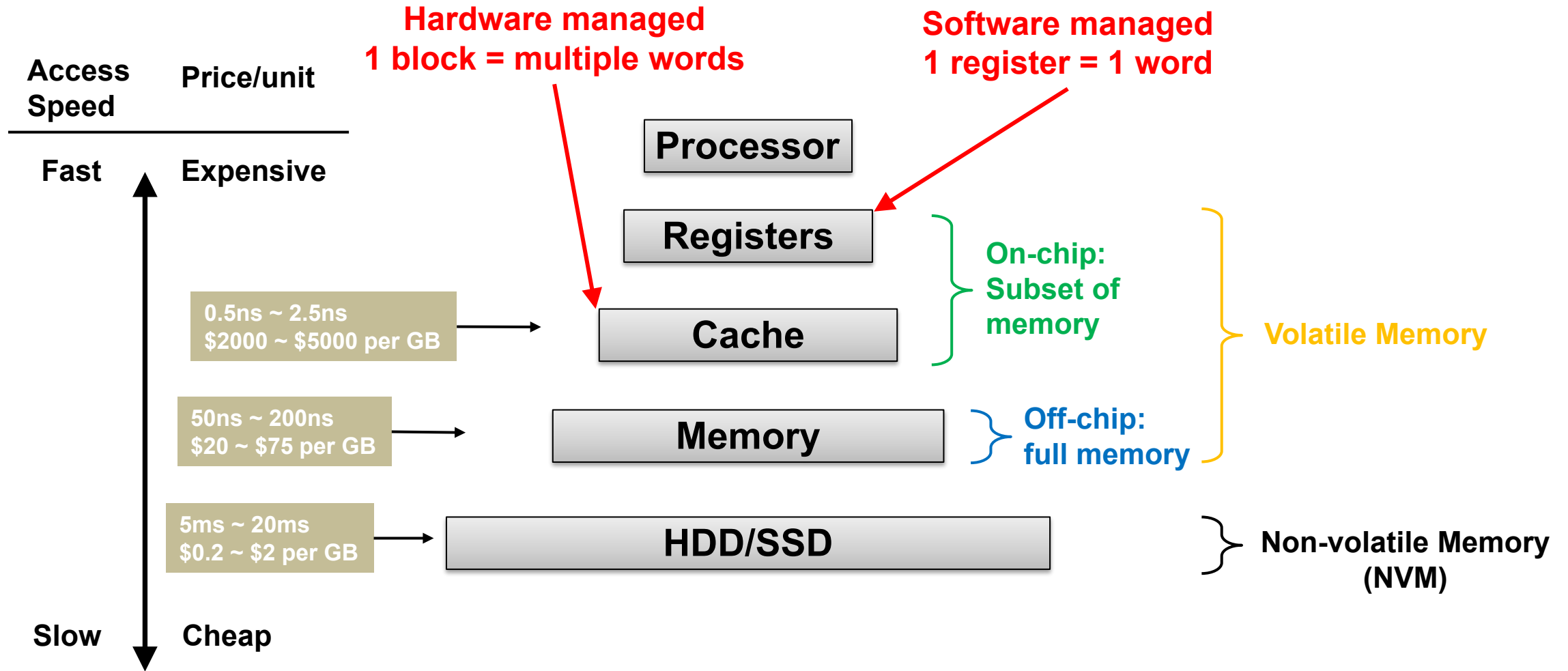- **Off-chip Memories (Memories outside CPU)**
  - System Memory, Storage
  - Large but Slow

**System Memory**

**Storage**

SJSU    SAN JOSÉ STATE UNIVERSITY

# Memory Hierarchy

**Hardware managed**
**1 block = multiple words**

**Software managed**
**1 register = 1 word**

**Access Speed**  **Price/unit**

**Fast**  **Expensive**

**Processor**

**Registers**

**On-chip: Subset of memory**

0.5ns ~ 2.5ns
$2000 ~ $5000 per GB

**Cache**

**Volatile Memory**

50ns ~ 200ns
$20 ~ $75 per GB

**Memory**

**Off-chip: full memory**

5ms ~ 20ms
$0.2 ~ $2 per GB

**HDD/SSD**

**Non-volatile Memory (NVM)**

**Slow**  **Cheap**

SJSU  SAN JOSÉ STATE UNIVERSITY

# Memories in Your PC

- **Windows**
  - This PC → Properties
  - cmd window → wmic
  - 3<sup>rd</sup> party tool like CPU-Z



**System Memory**



**Storage**



**Caches**

- **Linux**
  - lscpu
  - cat /proc/cpuinfo
  - etc.

SJSU SAN JOSÉ STATE UNIVERSITY

# Discussion

**How would you design the memory system?**

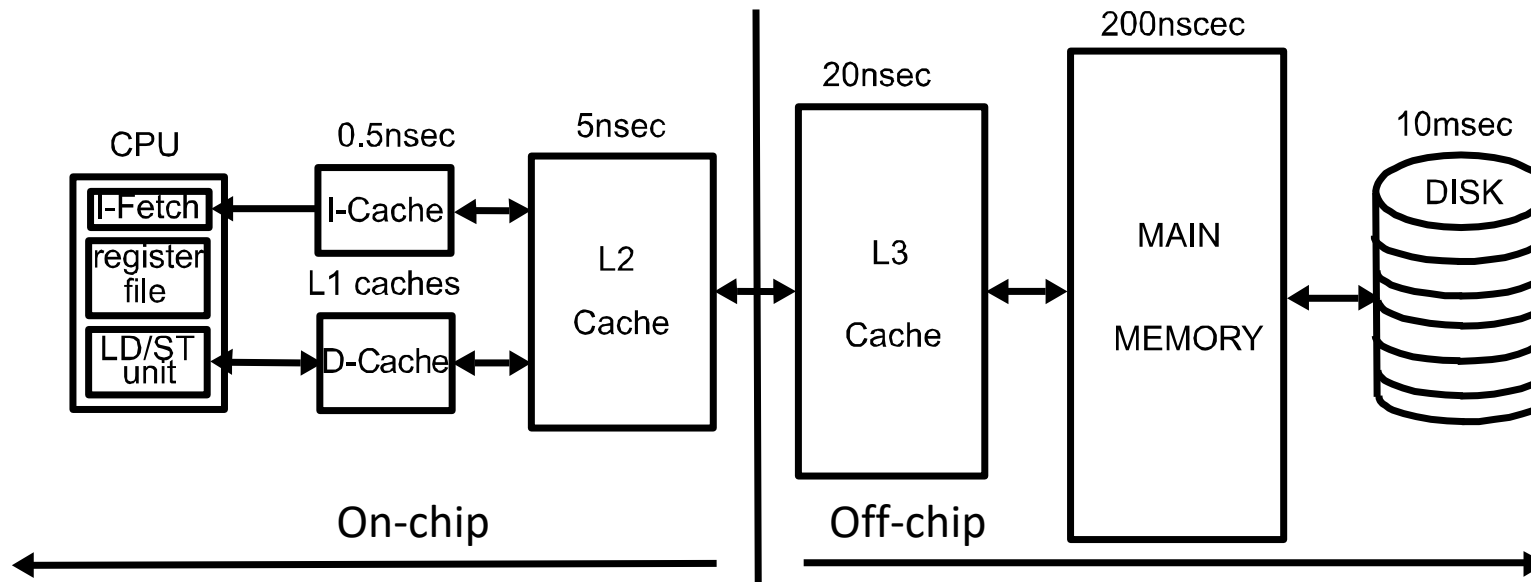**Single piece of memory that does everything**
**Vs.**
**Multiple levels of memories**
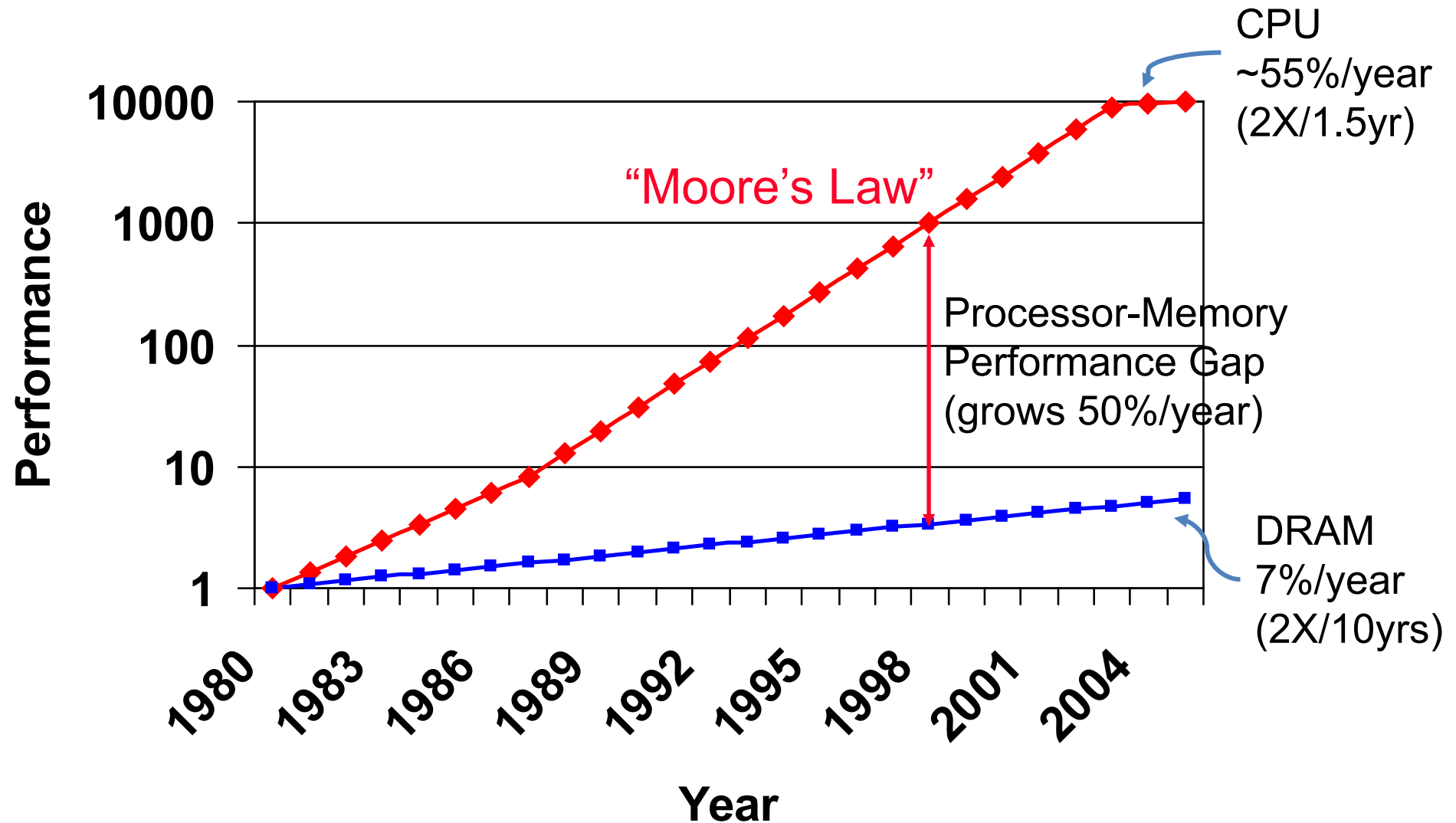
**Why should we use this hierarchy design?**

# Why?

- **Two Types of Locality:**
  - **Temporal Locality** (Locality in Time): If an address is referenced, it tends to be referenced again (e.g., loops, variable reuse)

  - **Spatial Locality** (Locality in Space): If an address is referenced, neighboring addresses tend to be referenced (e.g., array, stack, etc.)

# The "Memory Wall"

SJSU  SAN JOSÉ STATE UNIVERSITY

# The Memory Hierarchy Goal

- **How do we create a memory system that gives the illusion of being large, cheap and fast (most of the time)?**
  - With hierarchy
    - try the fast parts first -- most of the time, this works well
    - if not, move the data so it works well the next time
  - With parallelism
    - use multiple identical parts operating simultaneously
    - for large quantities of data, this works well

- **Example – keep a subset of the data in fast memory**
- **Example – 1-byte-wide memory ➔ 4 × 1-byte-wide memory ➔ 4-byte-wide memory**
  - load word takes 4 memory cycles vs. 1 memory cycle

SJSU SAN JOSÉ STATE UNIVERSITY

# Let us Conclude

**What is cache?**

**Fast memory**

**What are the two types of data localities?**

**Temporal & Spatial**

**What are the two memory design directions inspired by the principle of locality?**

**Hierarchy & parallelism**

SJSU   SAN JOSÉ STATE UNIVERSITY

SAN JOSÉ STATE UNIVERSITY *powering* SILICON VALLEY