

# CSE261: Computer Architecture

## 15. Memory Hierarchy (1): Introduction

Seongil Wi

# HW2

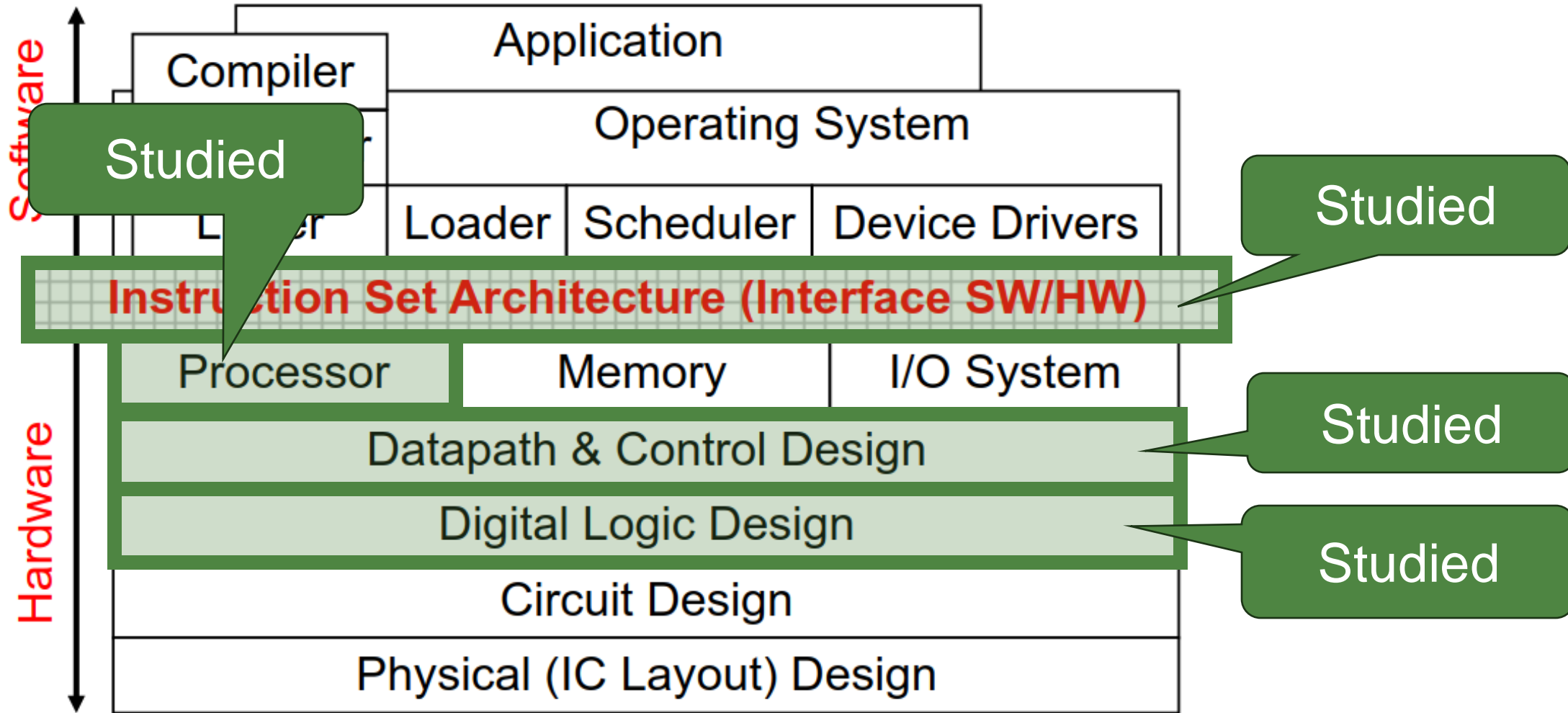
---

2

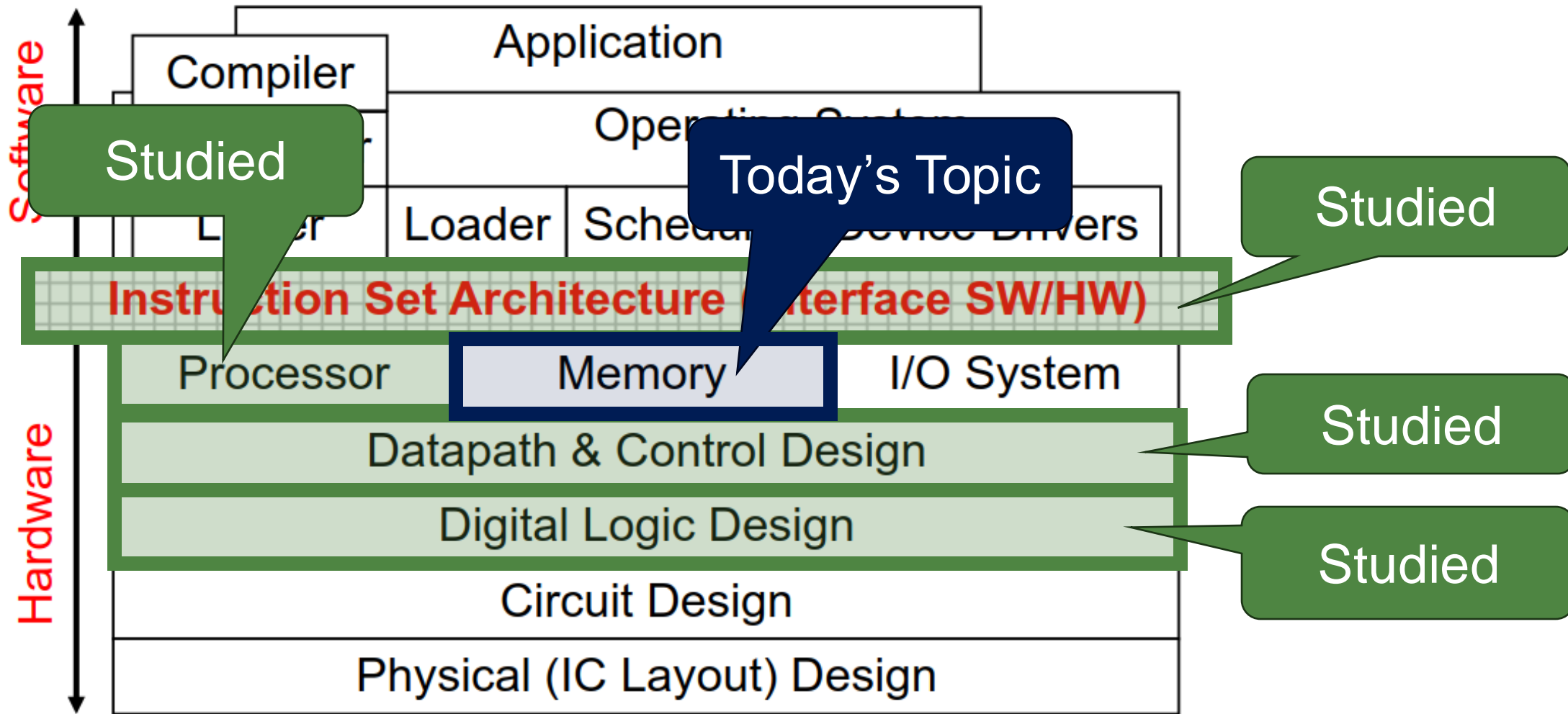


Due date: **11/26, 11:59PM**

# Where Are We?



# Today's Topic – Memory

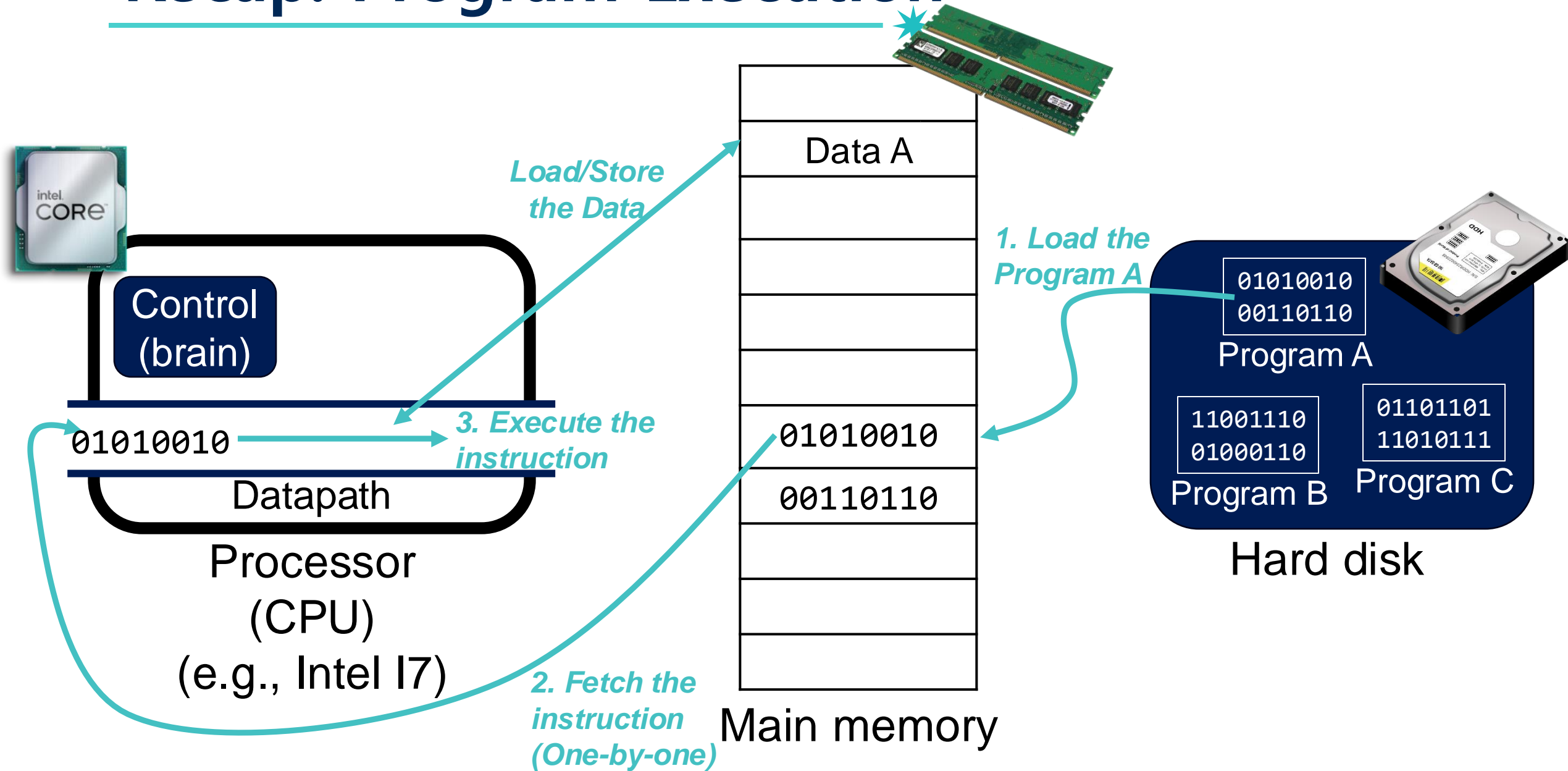




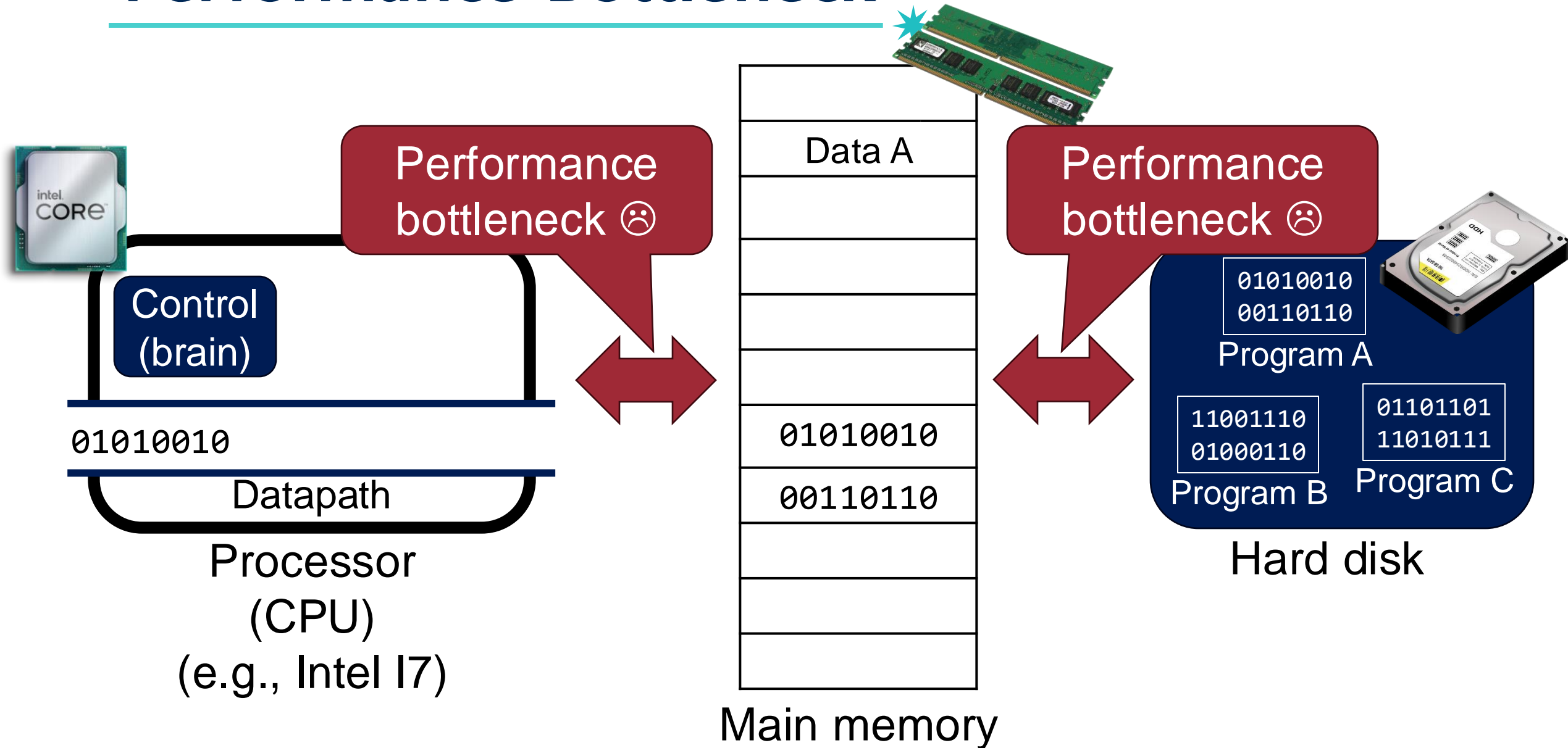
## Memory

**Accessing memory is one of the biggest performance bottlenecks 😞**

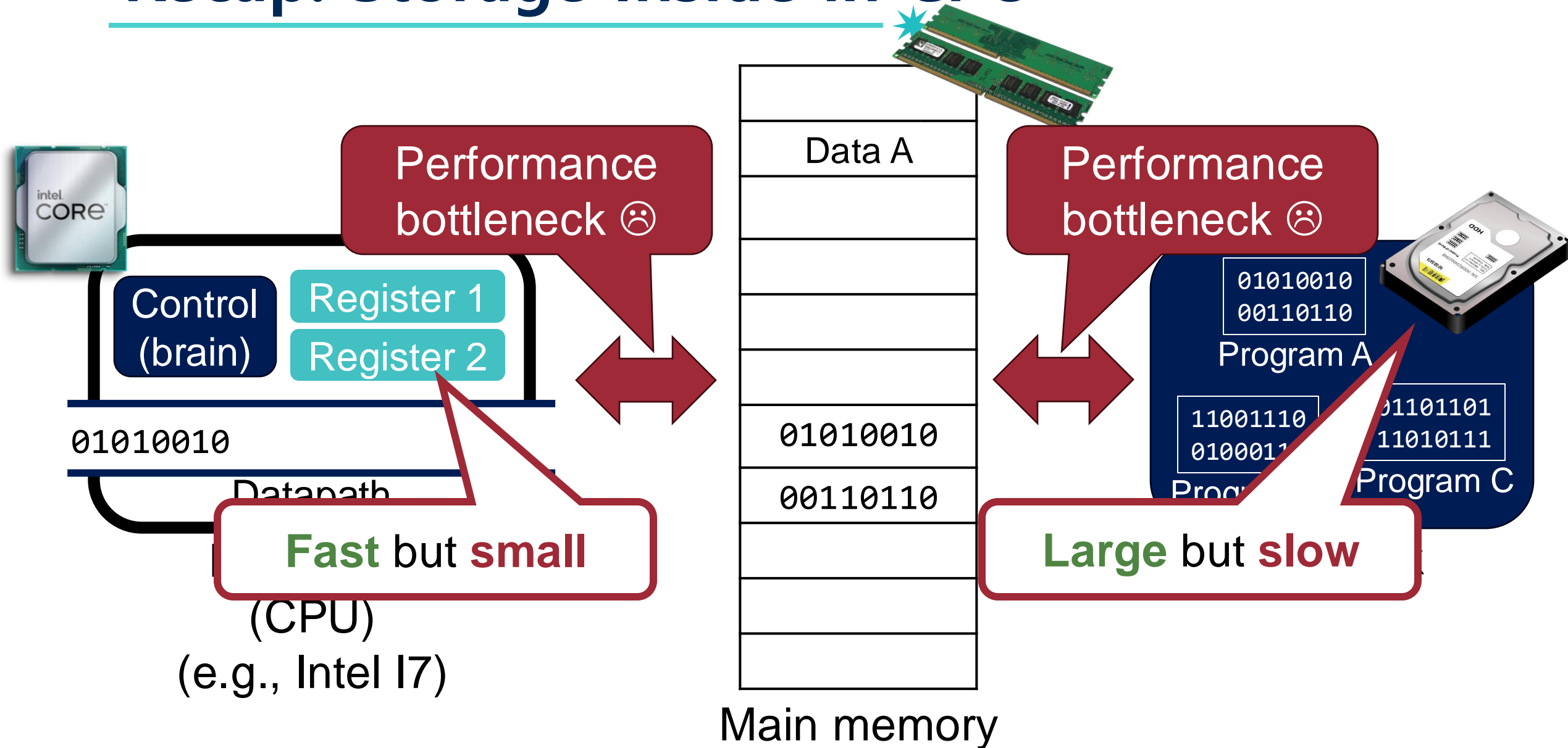
# Recap: Program Execution



# Performance Bottleneck



# Recap: Storage Inside in CPU





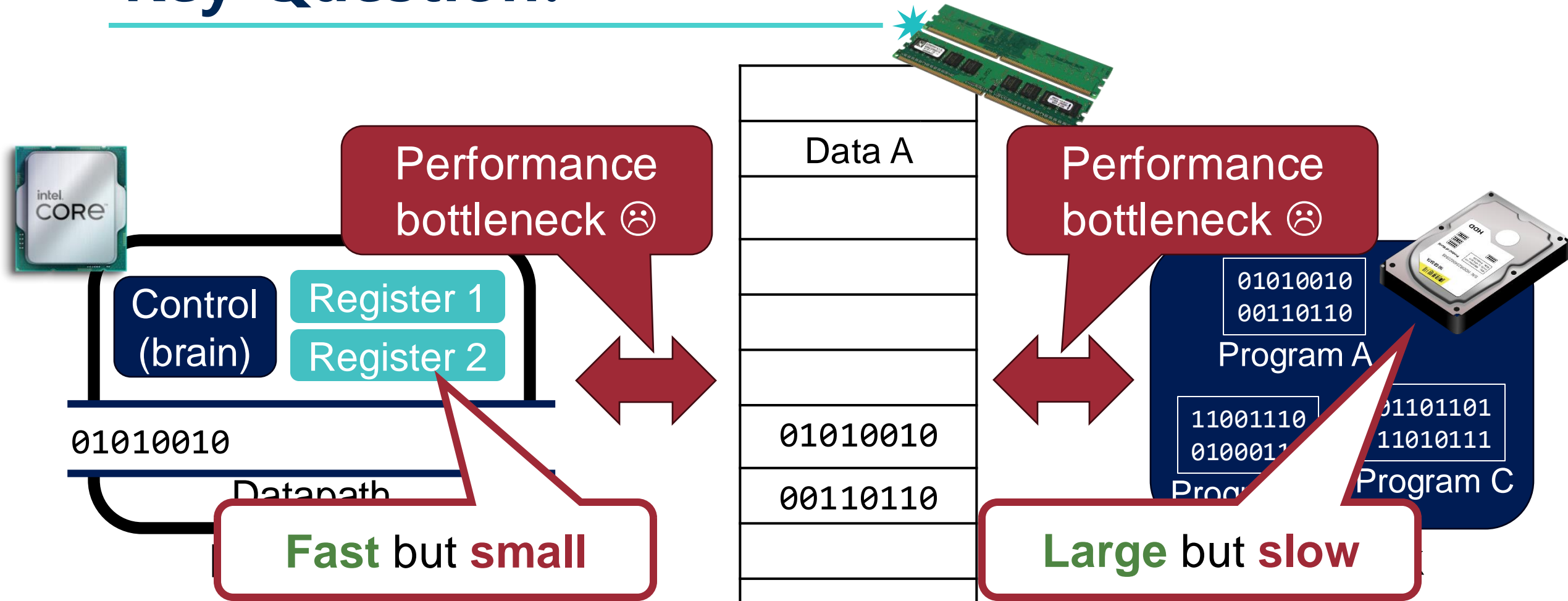
# Recall that 'Smaller is Faster!'



The faster memories are *more expensive per bit* than the slower memories and thus are smaller!

Memory	Typical access time	\$ per GB (in 2020)
Static RAM (SRAM)	0.5ns – 2.5ns	\$2,000 – \$5,000
Dynamic RAM (DRAM)	50ns – 70ns	\$20 – \$75
Magnetic Disk	5ms – 20ms	\$0.20 – \$2

# Key Question!



There is a *conflict* having large and fast memory ☹️  
How can we achieve large and fast memory?

# Solution: Memory Hierarchy

A structure that uses multiple levels of memories

# Exploiting Hierarchy Example: Book Storage



# Exploiting Hierarchy Example: Book Storage<sup>12</sup>



Bookshelves: Level 2  
Slower but bigger

Desk: Level 1  
Fastest but smallest

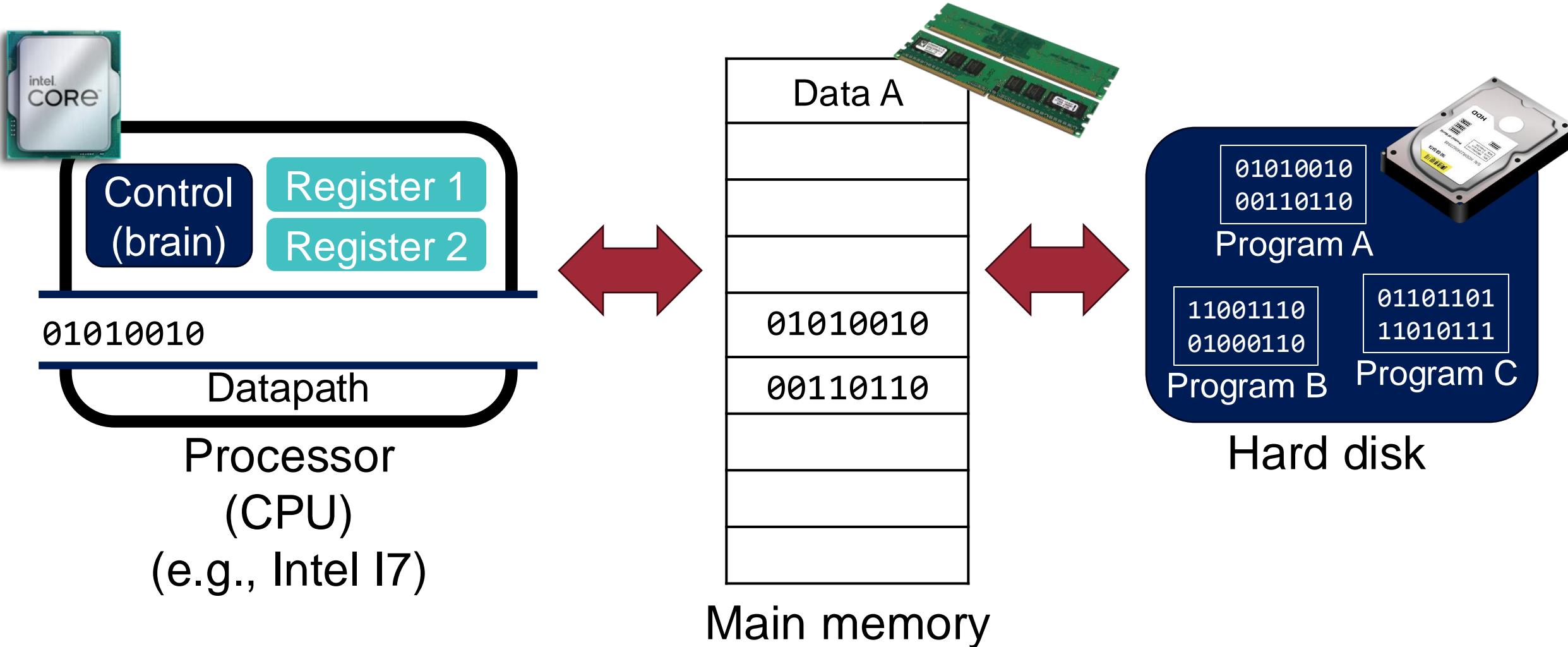
How about Level 3?





# Solution: Memory Hierarchy

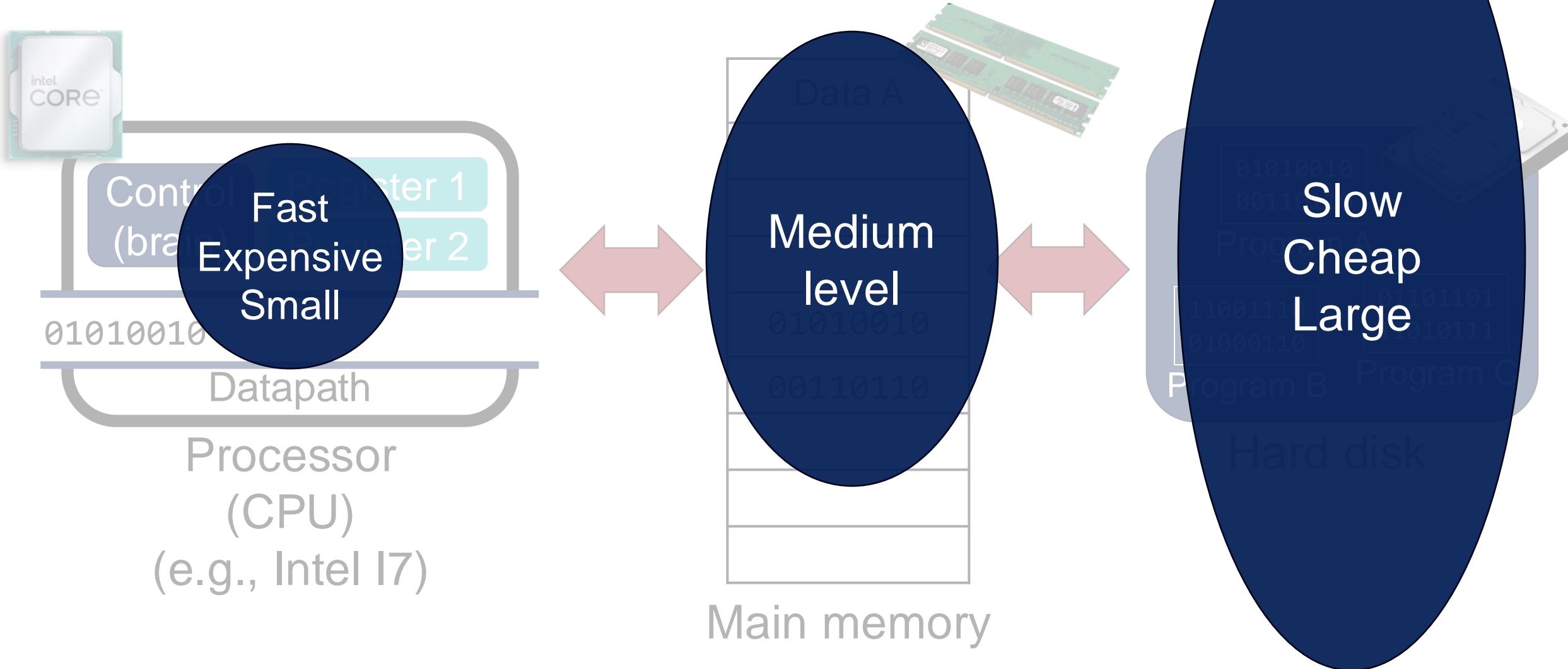
A structure that uses multiple levels of memories



# Solution: Memory Hierarchy

15

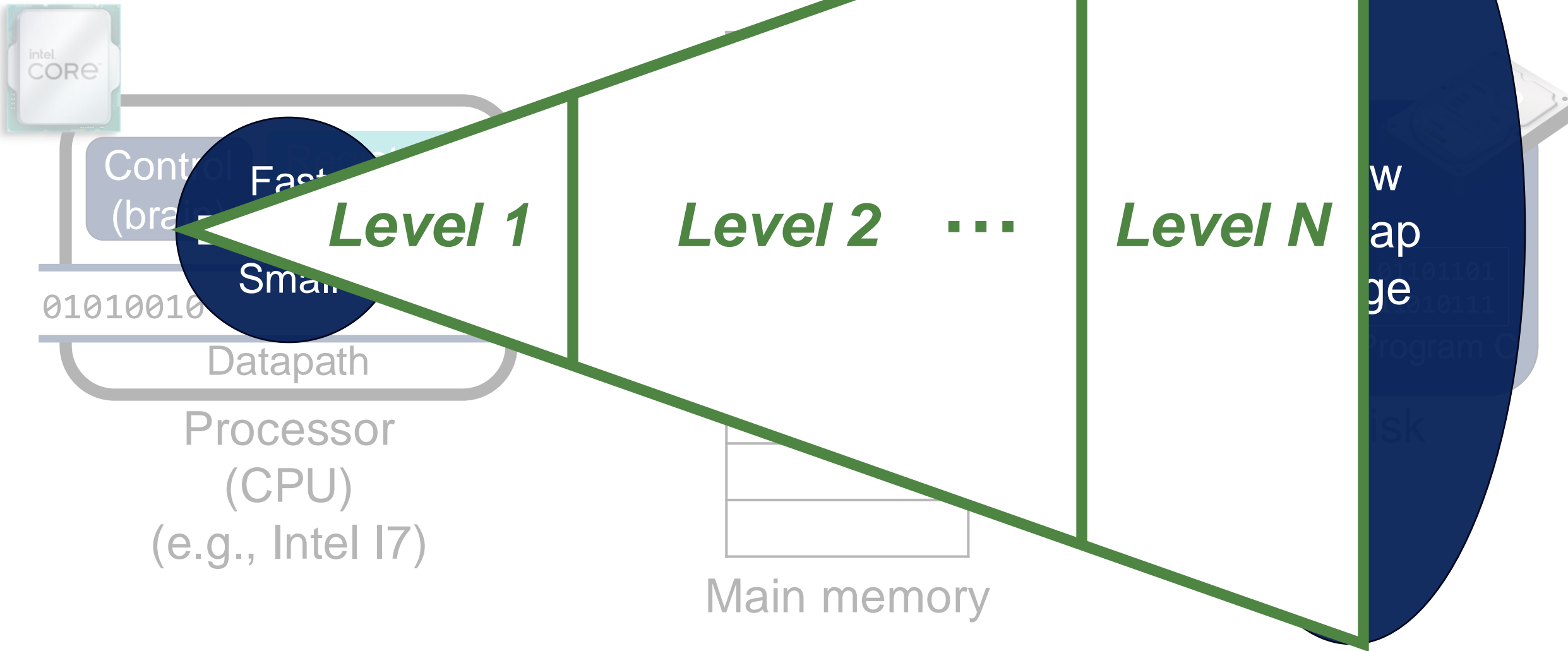
A structure that uses multiple levels of memories



# Solution: Memory Hierarchy

16

A structure that uses multiple levels of mem

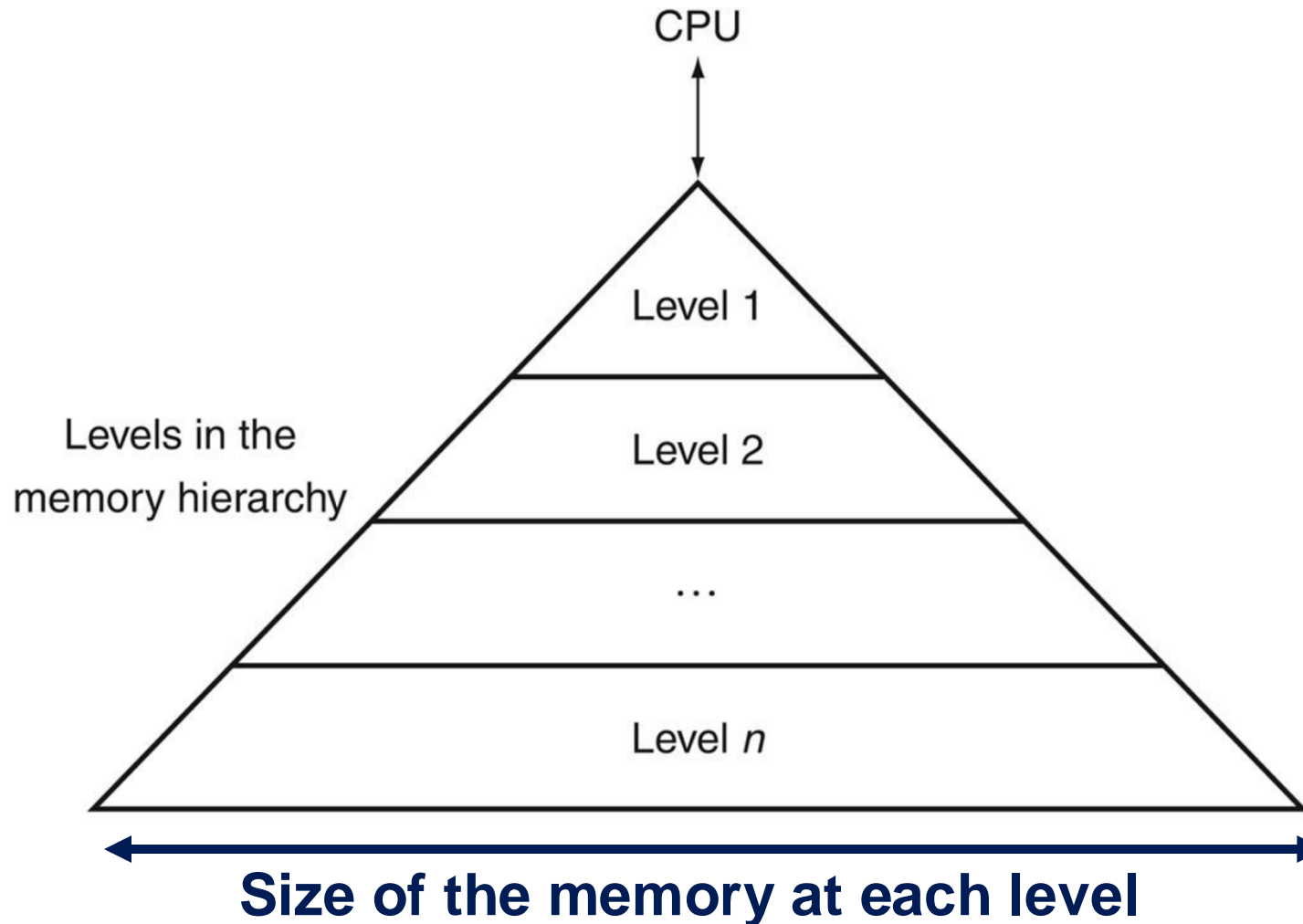




# Solution: Memory Hierarchy

A structure that uses multiple levels of memories

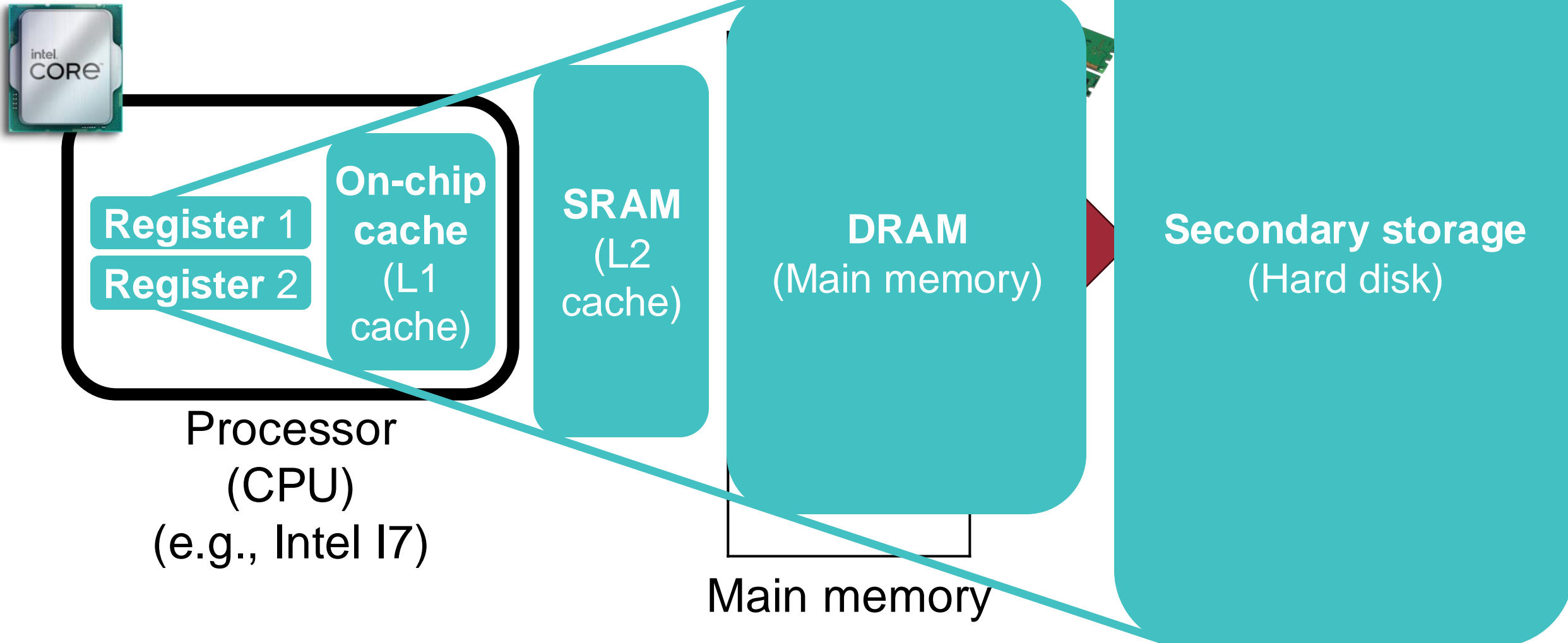
**Faster,  
expensive,  
smaller**



**Increasing distance  
from the CPU in  
access time**

# Memory Hierarchy Details

18



# Why does the Memory Hierarchy Work?

---

19

Because of the principle of *locality*!

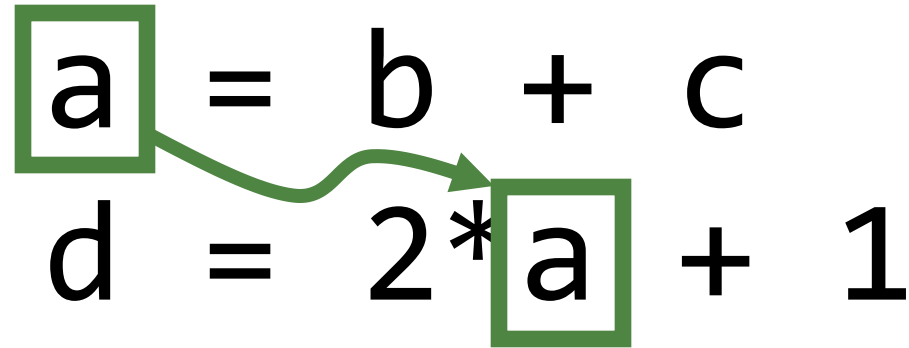
# Important Principle: Locality (지역성)



The tendency to access the same set of memory locations *repetitively* over a short period of time

1. **Temporal locality** (locality in time)
2. **Spatial locality** (locality in space)

# Temporal Locality Example


$$\begin{array}{l} \boxed{a} = b + c \\ d = 2 * \boxed{a} + 1 \end{array}$$

**Temporal locality:** items accessed recently are likely to be accessed again soon!

# Important Principle: Locality (지역성)



The tendency to access the same set of memory locations repetitively over a short period of time

## 1. Temporal locality (locality in time)

- If an item is referenced, the same item will tend to be referenced again soon

## 2. Spatial locality (locality in space)

# Spatial Locality Example



```
for (i=0; i<10; i++)  
    sum = sum + a[i]
```

**Spatial locality:** Items near those accessed recently are likely to be accessed soon

# Important Principle: Locality (지역성)



The tendency to access the same set of memory locations repetitively over a short period of time

## 1. Temporal locality (locality in time)

- If an item is referenced, the same item will tend to be referenced again soon

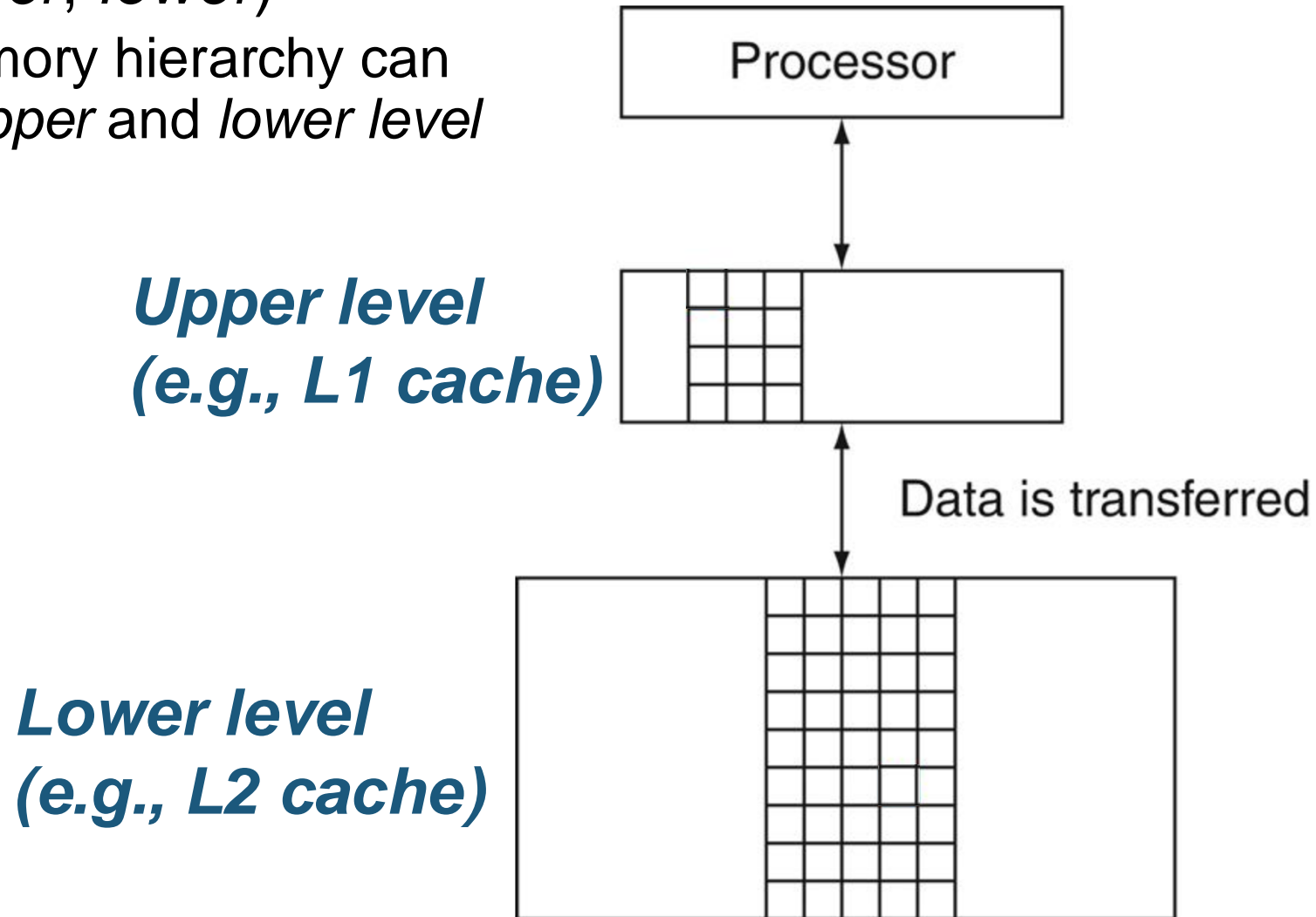
## 2. Spatial locality (locality in space)

- If an item is referenced, nearby items will tend to be referenced soon



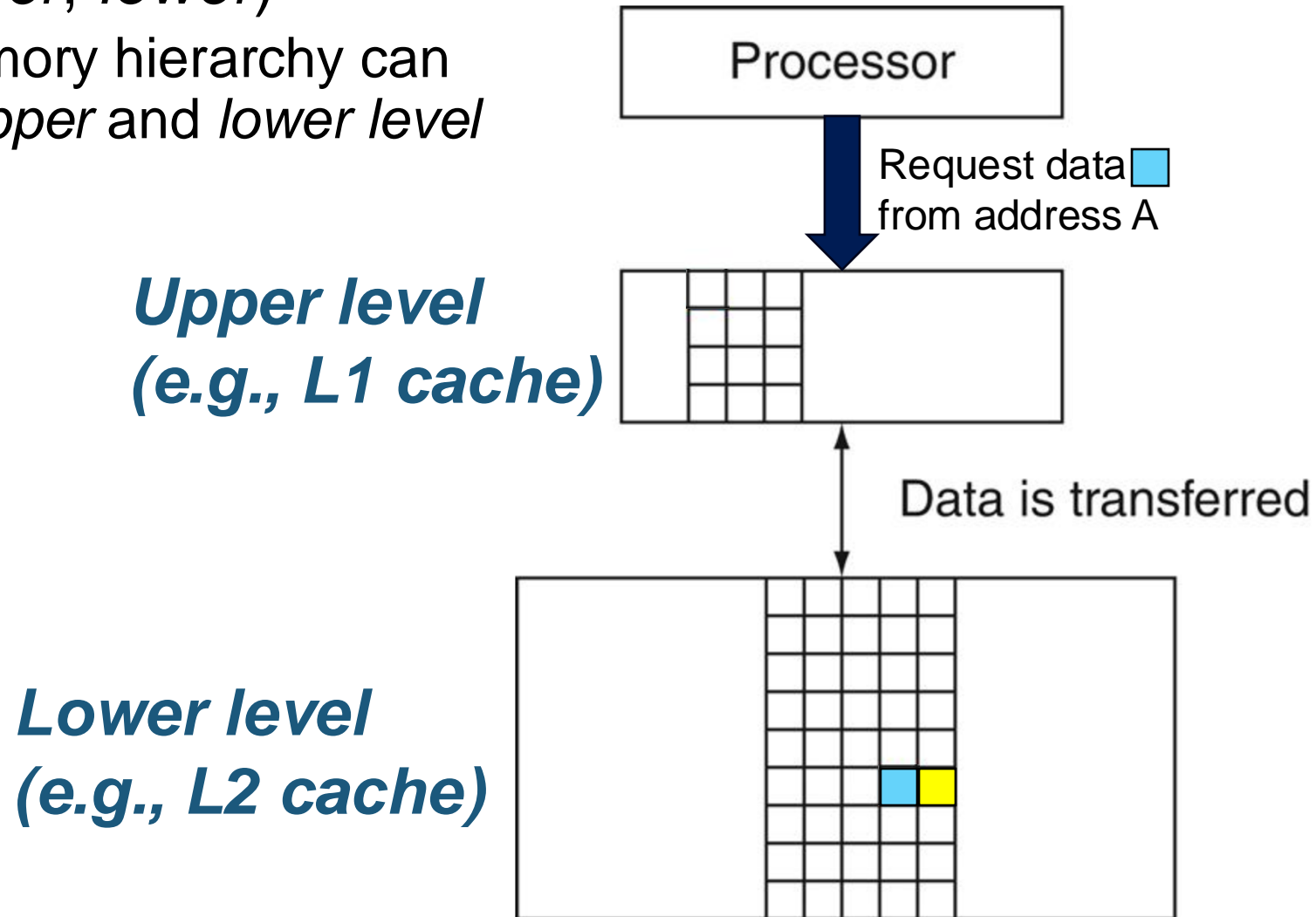
# Locality Example: 1st Access

- Assumption: two levels (*upper*, *lower*)
  - Each pair of levels in the memory hierarchy can be thought of as having an *upper* and *lower level*



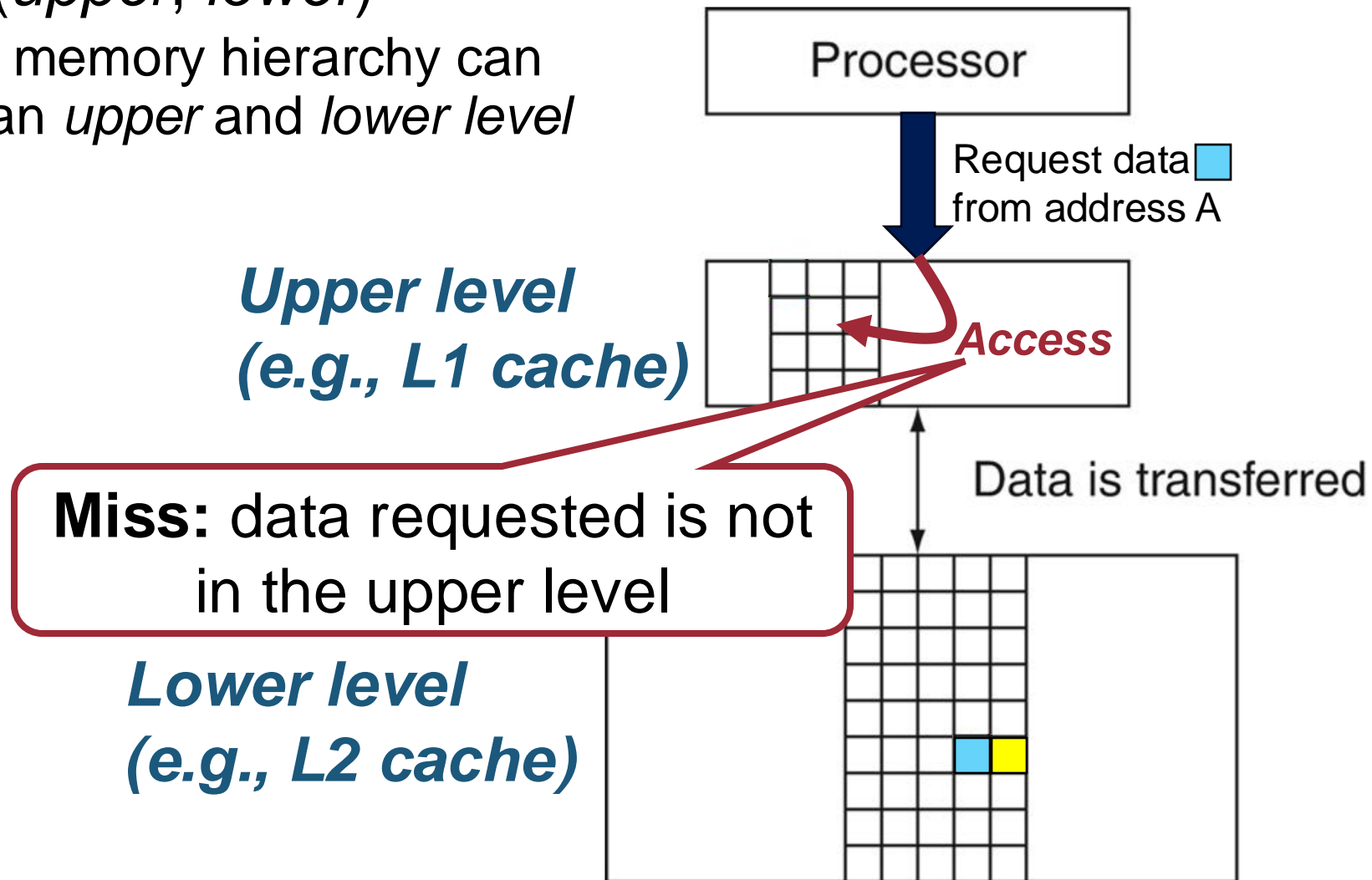
# Locality Example: 1st Access

- Assumption: two levels (*upper*, *lower*)
  - Each pair of levels in the memory hierarchy can be thought of as having an *upper* and *lower level*



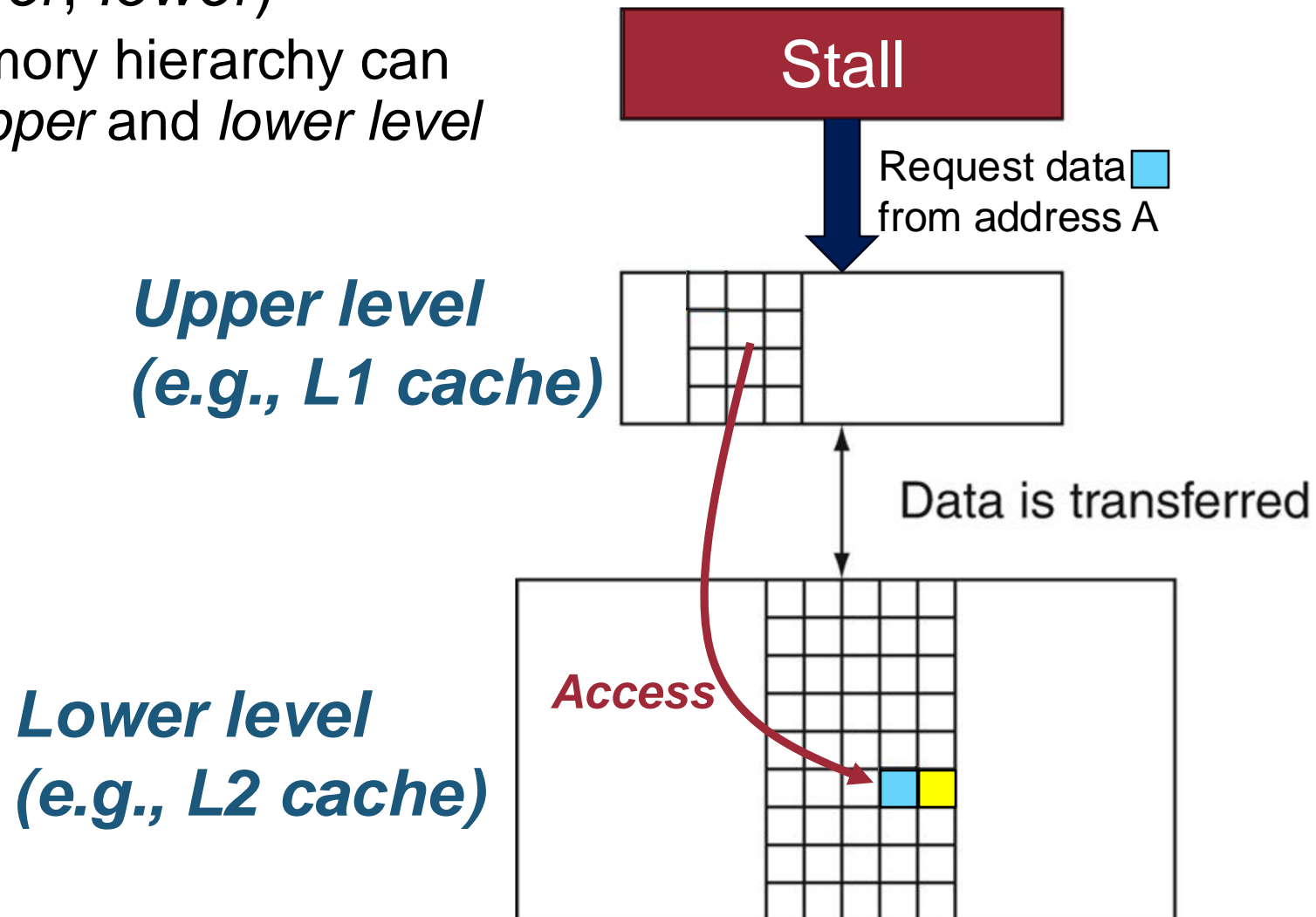
# Locality Example: 1st Access

- Assumption: two levels (*upper*, *lower*)
  - Each pair of levels in the memory hierarchy can be thought of as having an *upper* and *lower level*



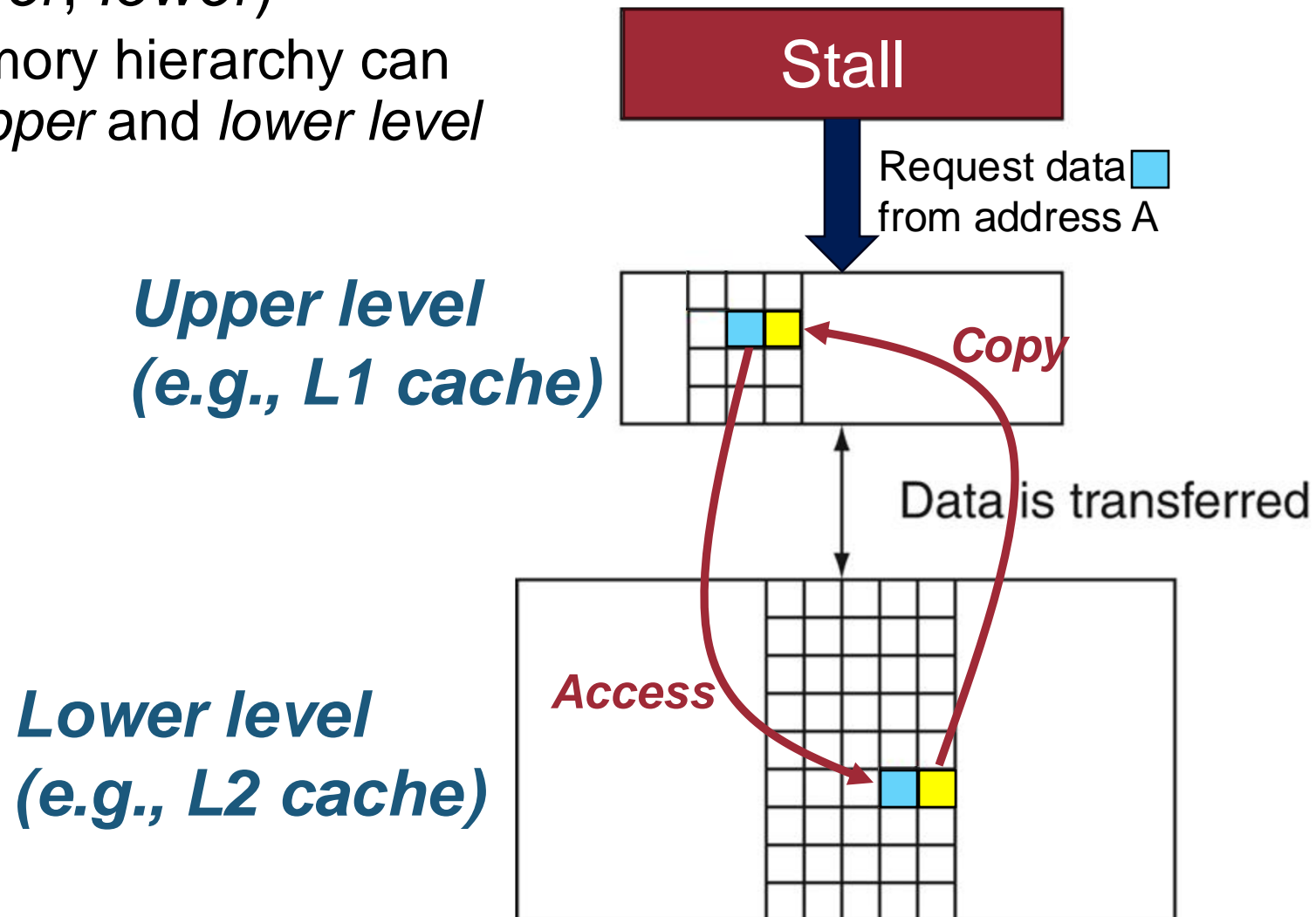
# Locality Example: 1st Access

- Assumption: two levels (*upper*, *lower*)
  - Each pair of levels in the memory hierarchy can be thought of as having an *upper* and *lower level*



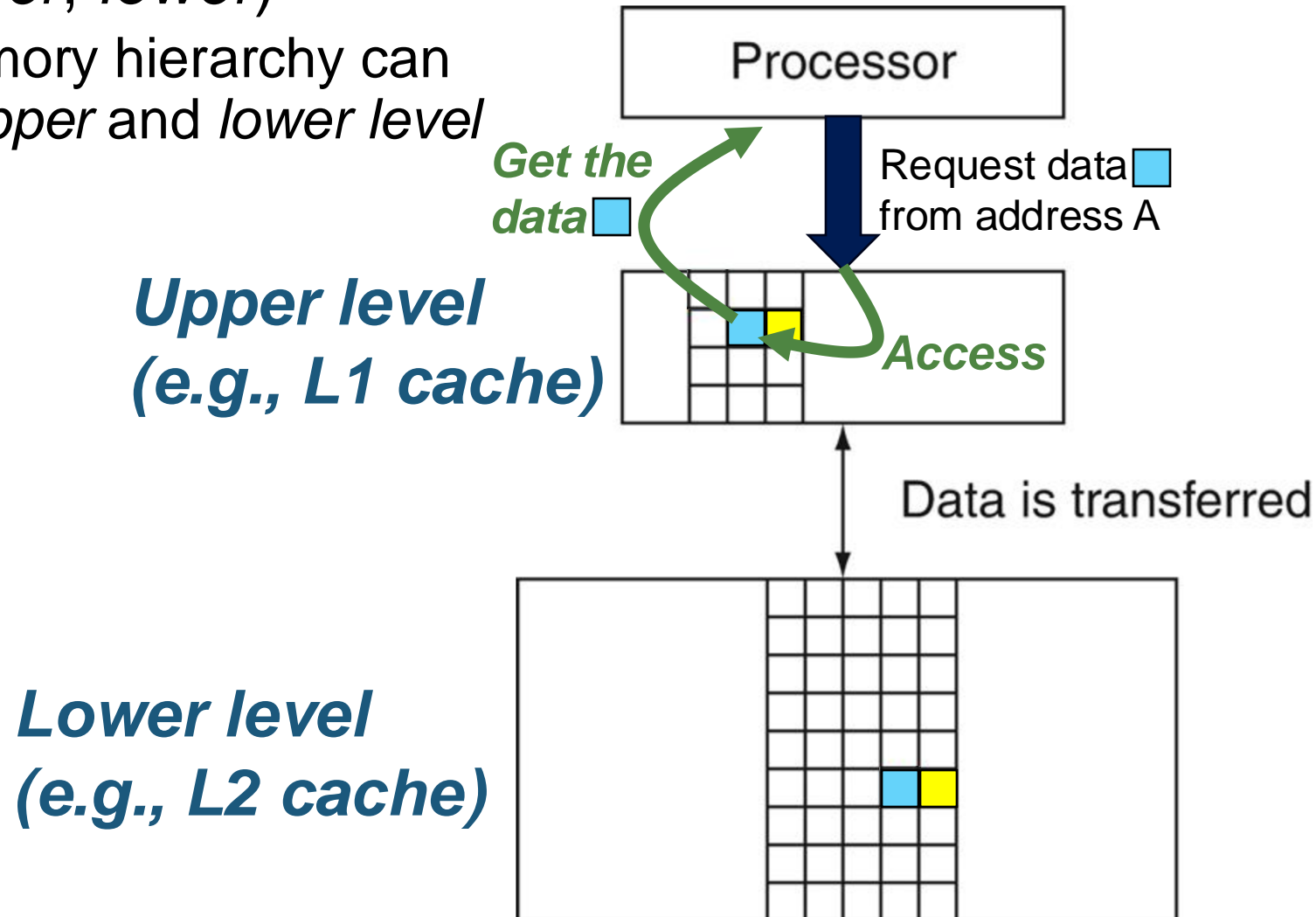
# Locality Example: 1st Access

- Assumption: two levels (*upper*, *lower*)
  - Each pair of levels in the memory hierarchy can be thought of as having an *upper* and *lower level*



# Locality Example: 1st Access

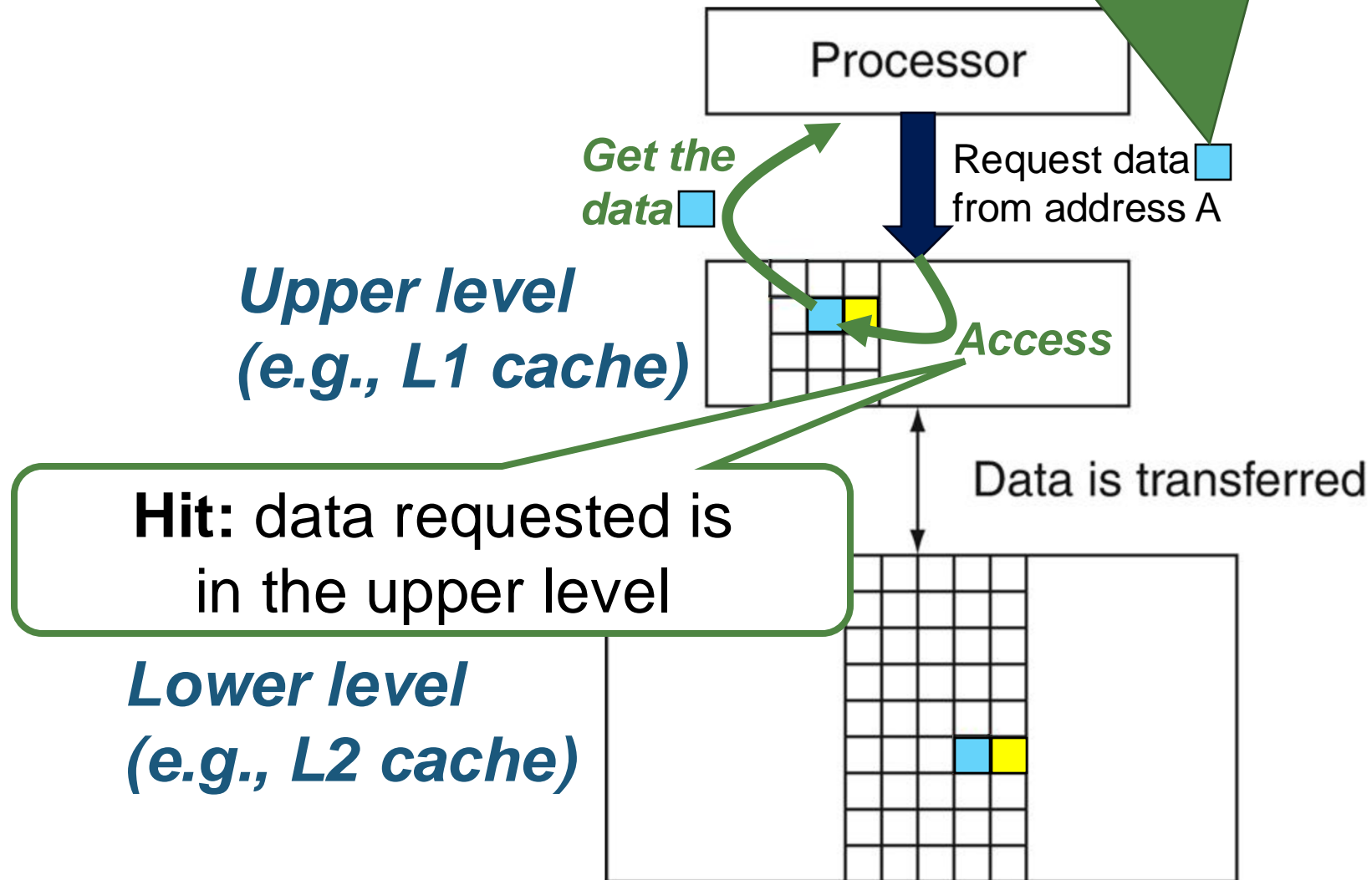
- Assumption: two levels (*upper*, *lower*)
  - Each pair of levels in the memory hierarchy can be thought of as having an *upper* and *lower* level



When a miss occurs, performance degrades significantly 😞

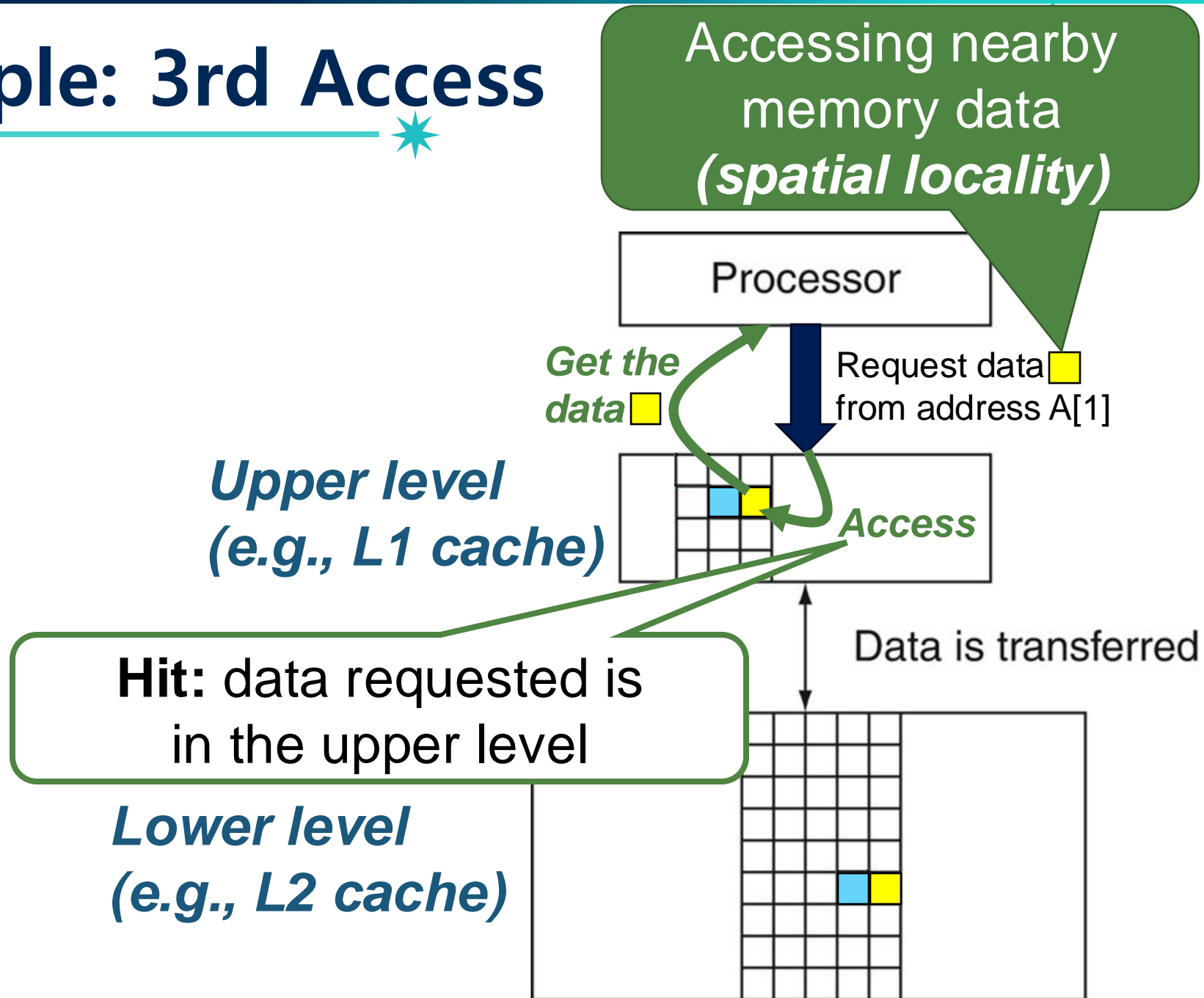
# Locality Example: 2nd Access

Accessing the same memory data again  
(*temporal locality*)





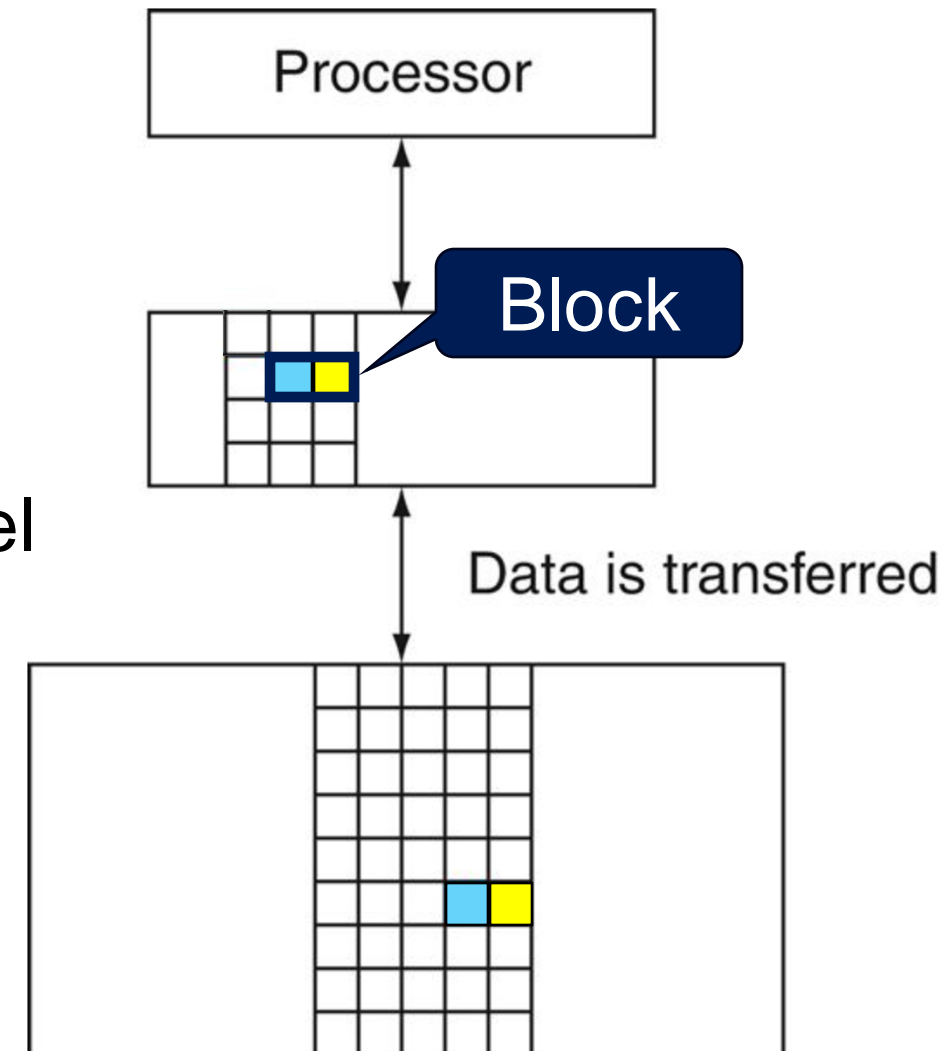
# Locality Example: 3rd Access



When a miss occurs, performance degrades significantly 😞  
However, once a miss happens, many hits can be expected  
due to **locality**, leading to performance improvement!

# Terms

- Block (a.k.a., line): unit of copying
  - Several words in cache memory
- **Hit**: data requested is in the upper level
  - Hit ratio:  $\text{hits}/\text{accesses}$
- **Miss**: data requested is not in the upper level
  - Block copied from lower level
  - Miss penalty: time taken to resolve miss
  - Miss ratio:  $\text{misses}/\text{accesses}$   
 $= 1 - \text{hit ratio}$



# Memory Hierarchy

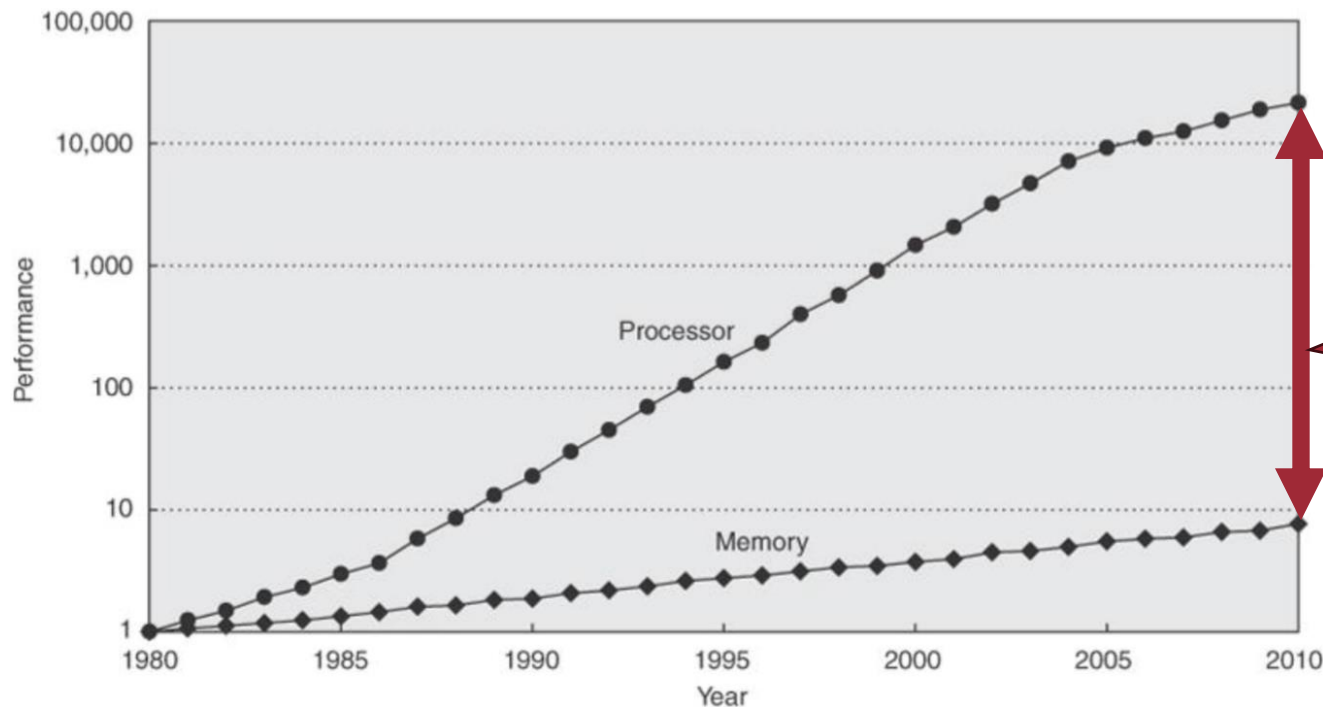
---



- Taking advantage of **locality**
  - Store everything on disk
  - Copy recently accessed (and nearby) items from disk to smaller DRAM memory
  - Copy more recently accessed (and nearby) items from DRAM to smaller SRAM memory

# Why Memory Hierarchy Crucial?

- Processor speed improvement: 60%/ year (2x/1.5 year)
- Memory latency improvement: 9% / year (2x/10 years)



Processor-memory performance gap: (grows 50% / year)

- Still it takes **around 100+ CPU cycles** for DRAM access
  - Hierarchy (계층구조) is the key!

**Question?**