

# Caches II

CSE 351 Autumn 2024

**Instructor:**

Ruth Anderson

**Teaching Assistants:**

Alexandra Michael

Connie Chen

Chloe Fong

Chendur Jayavelu

Joshua Tan

Nikolas McNamee

Nahush Shrivatsa

Naama Amiel

Neela Kausik

Renee Ruan

Rubee Zhao

Samantha Dreussi

Sean Siddens

Waleed Yagoub



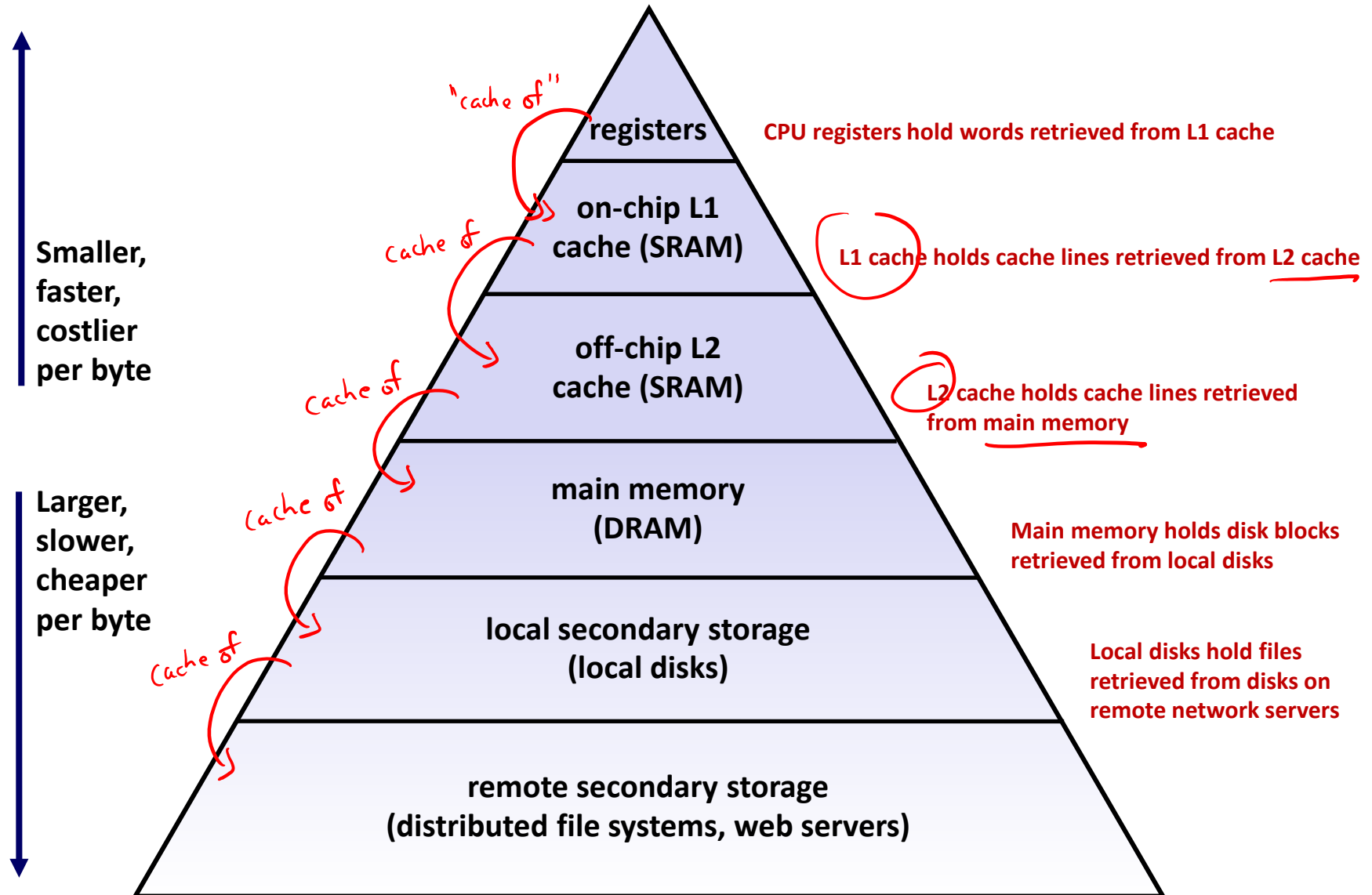
# Relevant Course Information

- ❖ HW15 due TONIGHT, Monday (11/04) @ 11:59 pm
- ❖ HW16 due Wednesday (11/06) @ 11:59 pm
- ❖ Lab 3 due Mon 11/11
  - Encouraged to aim for Fri 11/08, actual deadline Mon 11/11
  - You have everything you need to do the lab as of 10/28
  - Last part of HW15 is useful for Lab 3
- ❖ Mid-quarter Survey due Saturday (11/09)

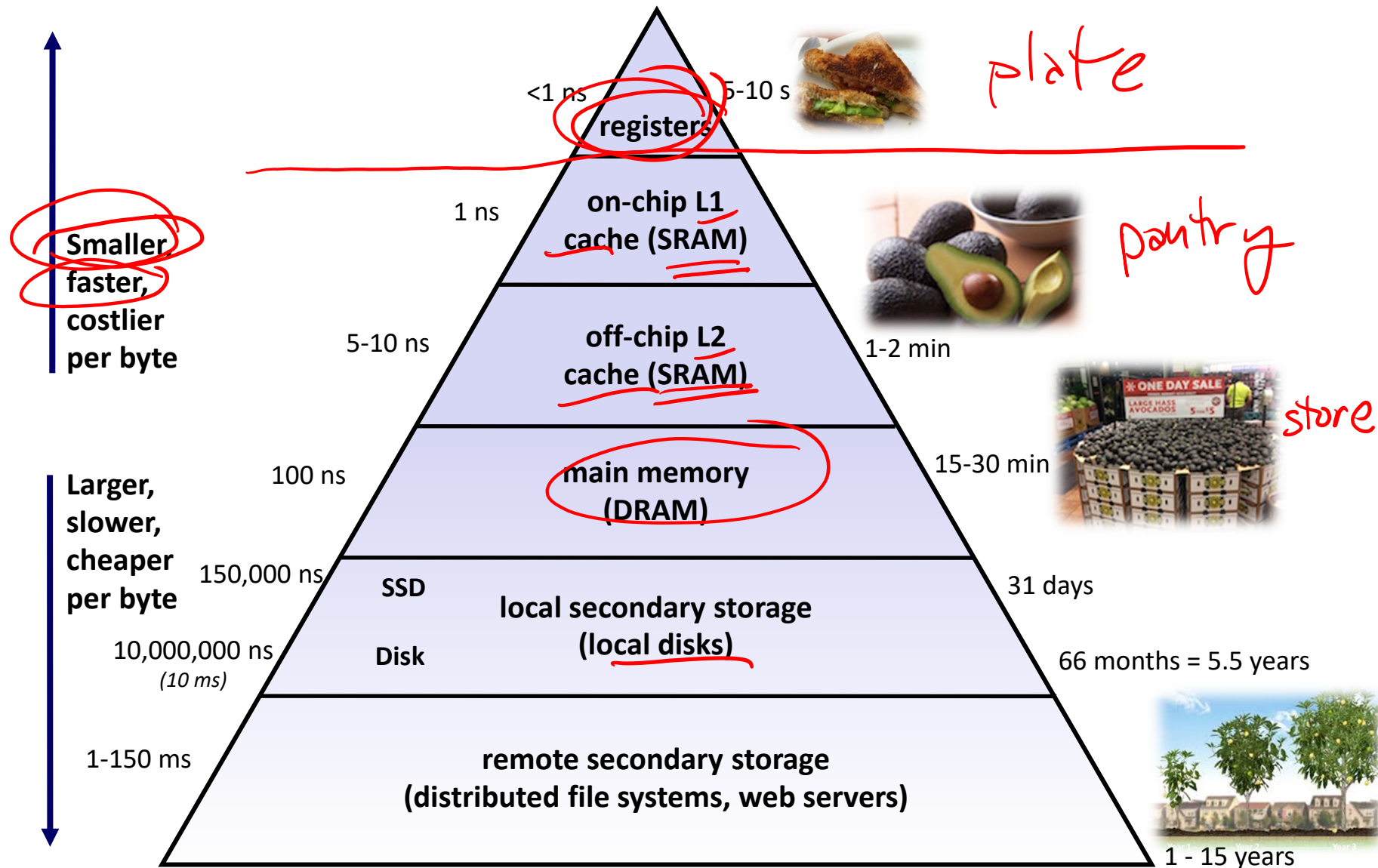
# Memory Hierarchies (Review)

- ❖ Some fundamental and enduring properties of hardware and software systems:
  - Faster storage technologies almost always cost more per byte and have lower capacity
  - The gaps between memory technology speeds are widening
    - True for: registers  $\leftrightarrow$  cache, cache  $\leftrightarrow$  DRAM, DRAM  $\leftrightarrow$  disk, etc.
  - Well-written programs tend to exhibit good locality
- ❖ These properties complement each other beautifully
  - They suggest an approach for organizing memory and storage systems known as a memory hierarchy
    - For each level  $x$ , the faster, smaller device at level  $x$  serves as a cache for the larger, slower device at level  $x+1$

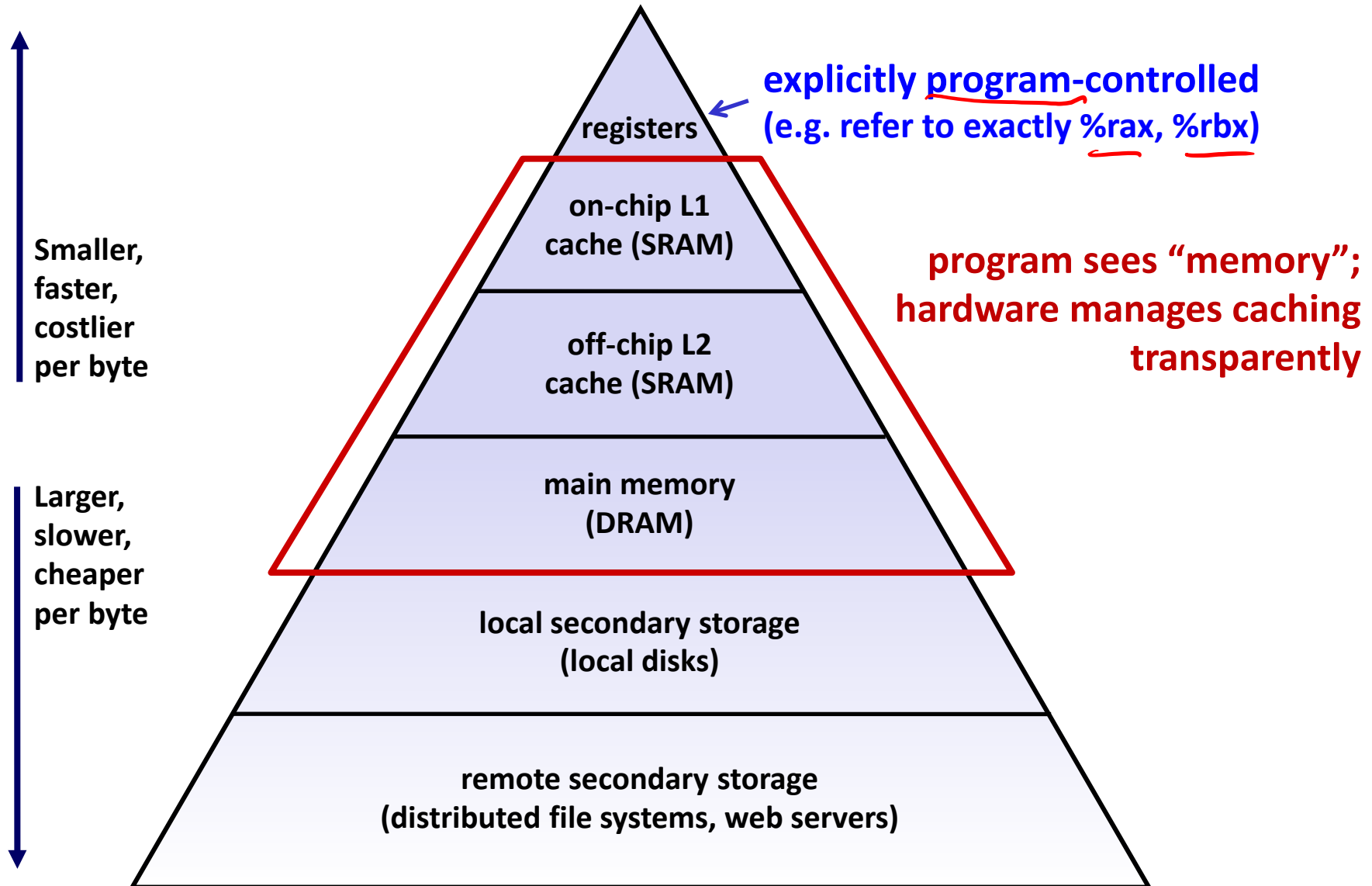
# An Example Memory Hierarchy (1)



# An Example Memory Hierarchy (2)

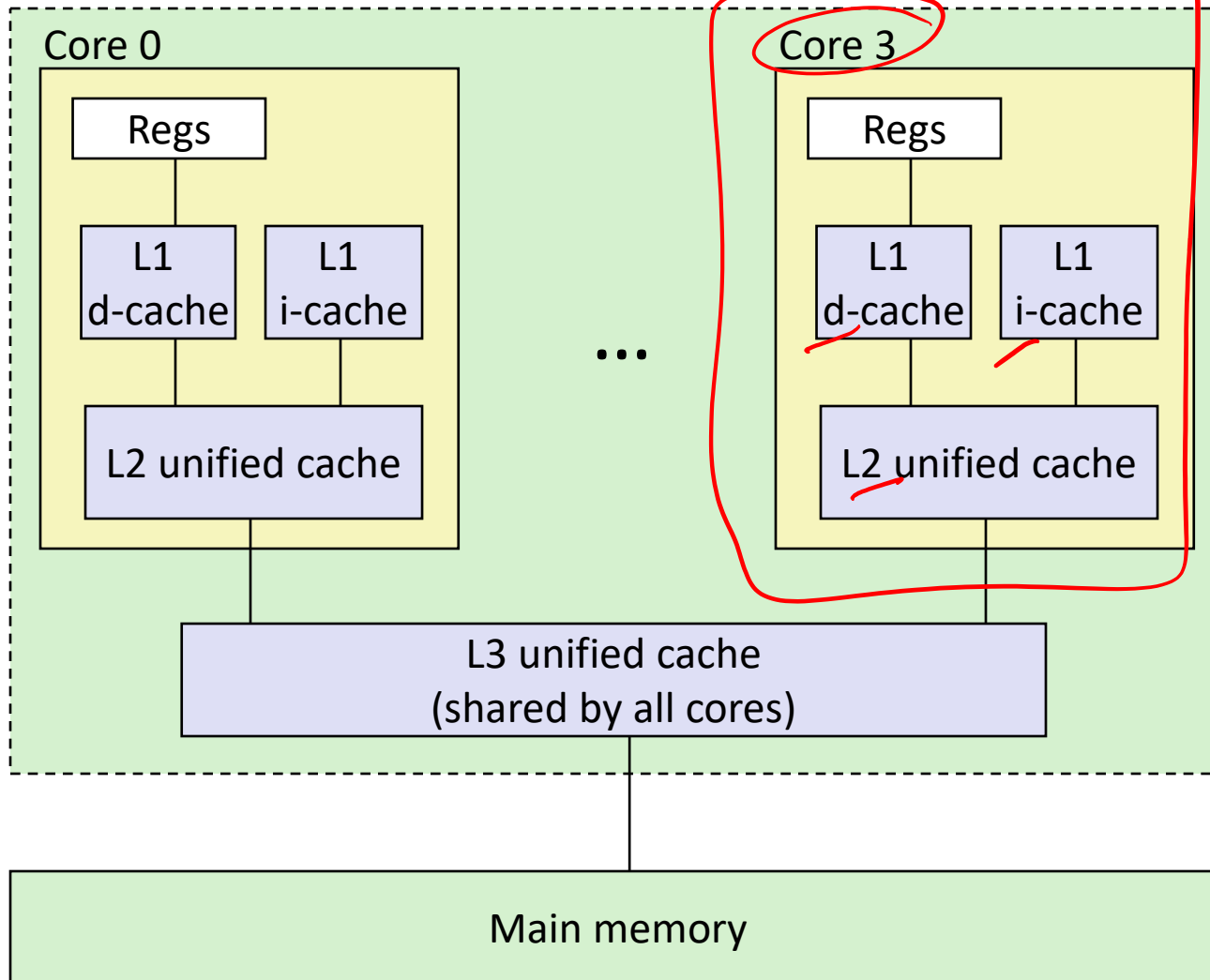


# Registers vs. Memory



# Intel Core i7 Cache Hierarchy

## Processor package



### Block size:

64 bytes for all caches

### L1 i-cache and d-cache:

32 KiB, 8-way,  
Access: 4 cycles

### L2 unified cache:

256 KiB, 8-way,  
Access: 11 cycles

### L3 unified cache:

8 MiB, 16-way,  
Access: 30-40 cycles

# Making memory accesses fast!

- ❖ Cache basics
- ❖ Principle of locality
- ❖ Memory hierarchies
- ❖ **Cache organization**
  - Direct-mapped (sets; index + tag)
  - Associativity (ways)
  - Replacement policy
  - Handling writes
- ❖ Program optimizations that consider caches



# Reading Review

## ❖ Terminology:

- Memory hierarchy
- Cache parameters: block size ( $K$ ), cache size ( $C$ )
- Addresses: block offset field ( $k$  bits wide)
- Cache organization: direct-mapped cache, index field

# Review Questions

- ❖ We have a direct-mapped cache with the following parameters:

- Block size of 8 bytes  $K = 2^3 \text{ B}$
- Cache size of 4 KiB  $C = 2^{12} \text{ B}$

- ❖ How many blocks can the cache hold?  $C/K = \frac{2^{12}}{2^3} = 2^9 = 512 \text{ blocks}$
- ❖ How many bits wide is the block offset field?  $\log_2 K = 3 \text{ bits}$
- ❖ Which of the following addresses would fall under block number 3?

$\lfloor 3/8 \rfloor = 0$   
**A. 0x3**  
 0b...000011  
 block 0

$\lfloor 31/8 \rfloor = 3$   
**B. 0x1F**  
 0b...0111111  
 block 3

$\lfloor 48/8 \rfloor = 6$   
**C. 0x30**  
 0b...0110000  
 block 6

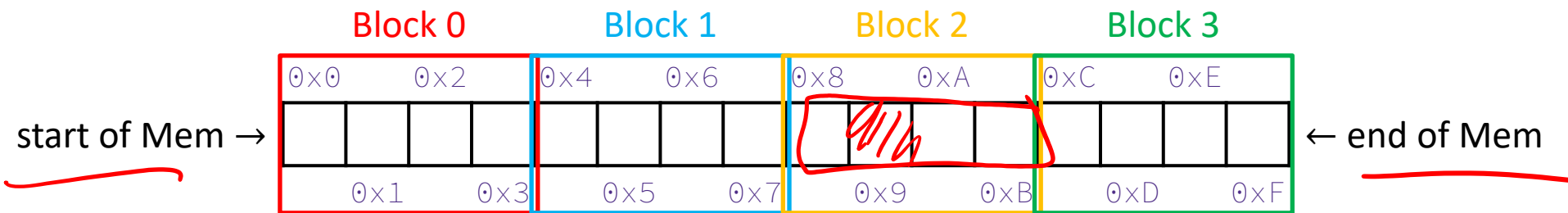
$\lfloor 56/8 \rfloor = 7$   
**D. 0x38**  
 0b...0111000  
 block 7

# Block Size (1)

**Note:** The textbook uses “B” for block size

- ❖ **Block Size ( $K$ ):** unit of transfer between \$ and Mem
  - Given in bytes and always a power of 2 (e.g., 64 B)
  - Blocks consist of adjacent bytes (differ in address by 1)
    - Spatial locality!

- Small example ( $K = 4$  B):



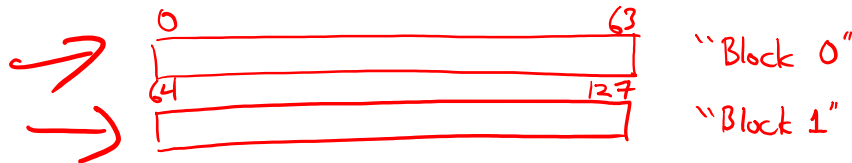
Address Space = 16 bytes

# Block Size (2)

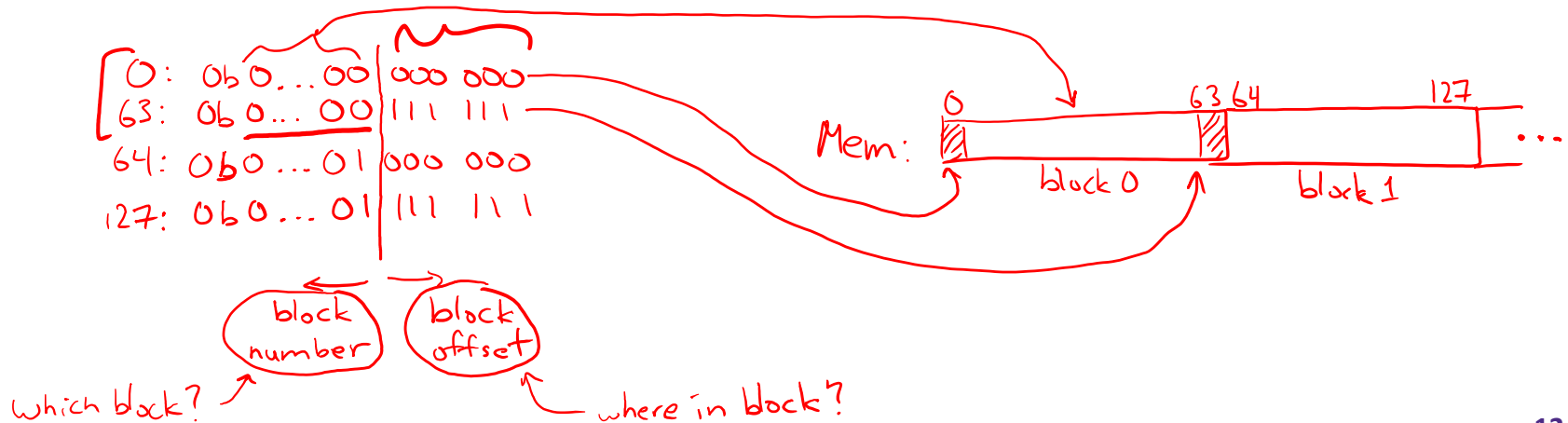
**Note:** The textbook uses "B" for block size

- ❖ **Block Size ( $K$ ):** unit of transfer between \$ and Mem
  - Given in bytes and always a power of 2 (e.g., 64 B)
  - Blocks consist of adjacent bytes (differ in address by 1)
    - Spatial locality!

Lab 1a: within Same Block



64 Byte Block  
 $2^6$  6 bits for offset



# Block Size (3)

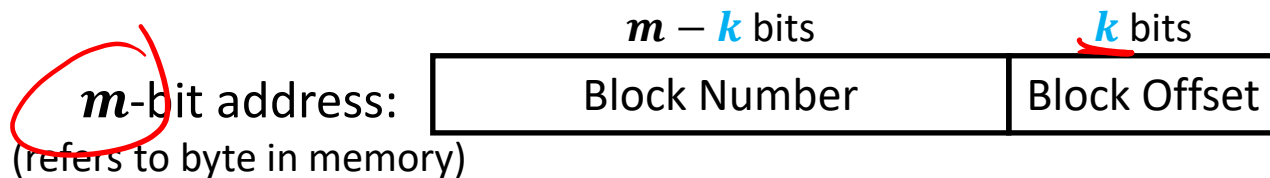
**Note:** The textbook uses “b” for offset bits

- ❖ **Block Size ( $K$ ):** unit of transfer between \$ and Mem
  - Given in bytes and always a power of 2 (e.g., 64 B)
  - Blocks consist of adjacent bytes (differ in address by 1)
    - Spatial locality!

## ❖ Offset field

- Low-order  $\log_2(K) = k$  bits of address tell you which byte within a block
  - $(\text{address}) \bmod 2^n = n$  lowest bits of address
- $(\text{address}) \bmod (\# \text{ of bytes in a block})$

How many bits do I need to specify every byte in a block?



# Block Size (4)

**Note:** The textbook uses “b” for offset bits

❖ **Block Size ( $K$ ):** unit of transfer between \$ and Mem


- Given in bytes and always a power of 2 (e.g., 64 B)
- Blocks consist of adjacent bytes (differ in address by 1)
  - Spatial locality!

❖ Example:

- If we have 6-bit addresses and block size  $K = 4$  B, which block and byte does  $0 \times 15$  refer to?

address:  $0b \underbrace{0101}_{\text{block num (value 5)}} \underbrace{01}_{\text{offset (value 1)}}$

offset width =  $\log_2(K) = \log_2(4) = 2 \text{ bits}$

block number 5:   
offset: 00 01 10 11

# Cache Size

❖ **Cache Size ( $C$ )**: amount of *data* the \$ can store

- Cache can only hold so much data (subset of next level)

- Given in bytes ( $C$ ) or number of blocks ( $C/K$ )

- Example:  $C = 32 \text{ KiB} = 512$  blocks if using 64-B blocks

$$2^5 \times 2^{10} = 2^{15} \text{ B} \times \frac{1 \text{ block}}{2^6 \text{ B}} = 2^9 \text{ blocks}$$

❖ Where should data go in the cache?

- We need a mapping from memory addresses to specific locations in the cache to make checking the cache for an address **fast**

❖ What is a data structure that provides fast lookup?

- Hash table!

# Hash Tables for Fast Lookup

mod by table size

Insert:

5

27

34

102

119

Block  
#s

32

Apply hash function to map data  
to "buckets"

0	
1	
2	102
3	
4	
5	5
6	
7	27
8	
9	119

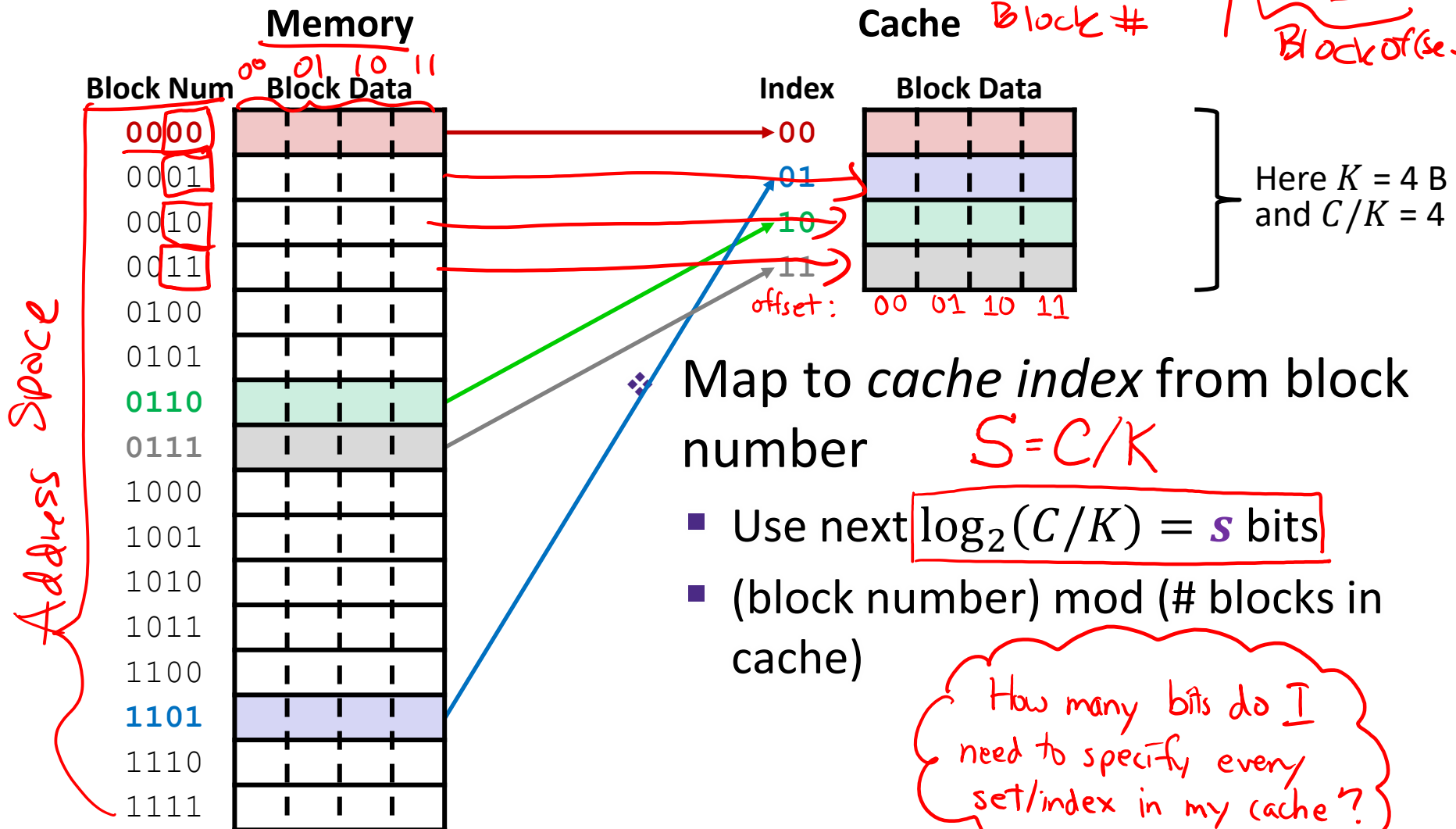
Table  
Size 10



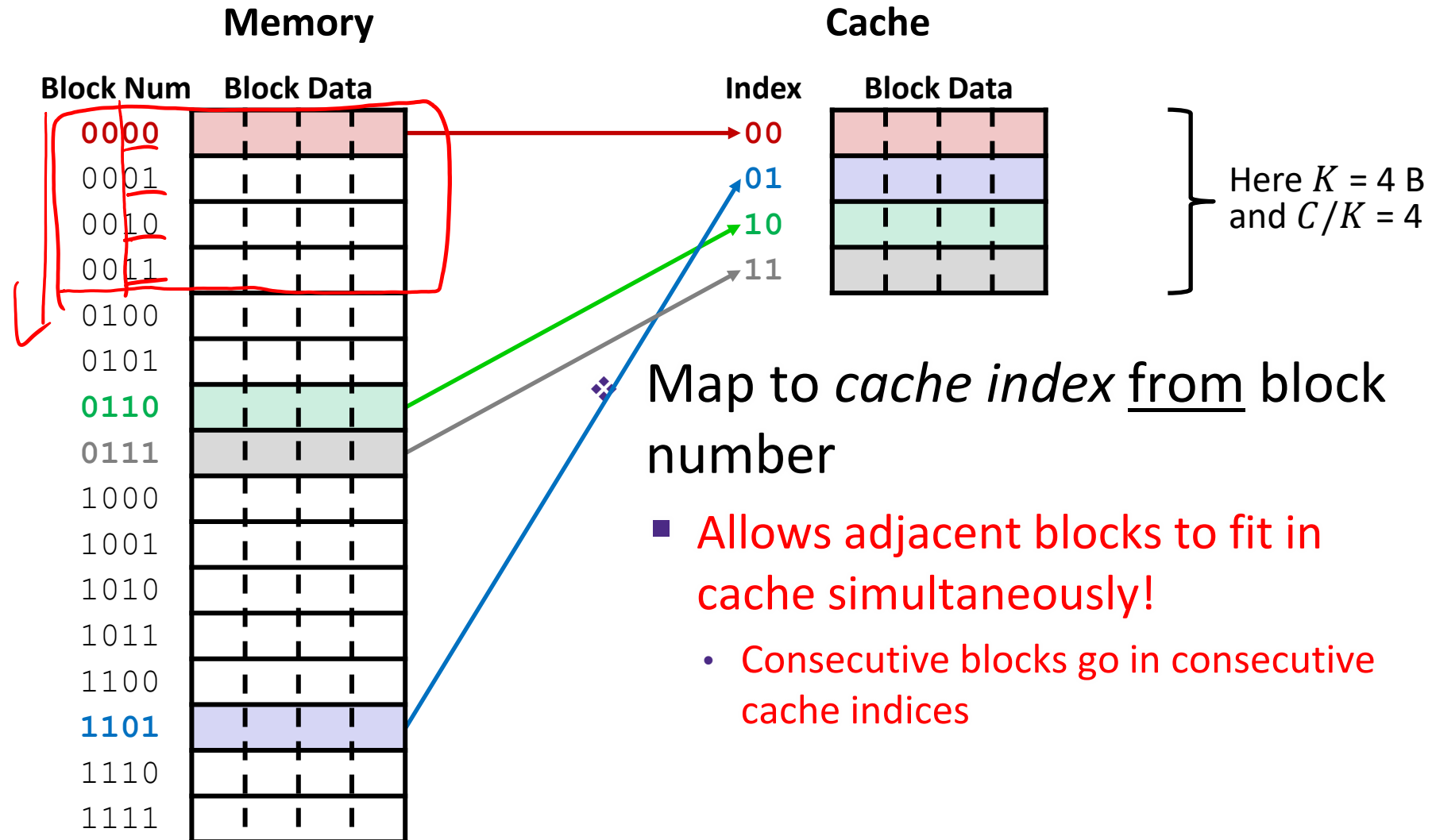
# Place Data in Cache by Hashing Address

Cache Size  $C = 16$  Bytes  $K = \text{Block Size} = 4$  Bytes  
 Addresses are 6 bits

Cache Block # — | — Block of set



# Place Data in Cache by Hashing Address



# Polling Question

- ❖ <sup>m</sup> 6-bit addresses, block size  $K = 4$  B, and our cache holds  $S = 4$  blocks. =  $C/K$  ,  $s = \log_2(4) = 2$  bits
- ❖ A request for address **0x2A** results in a cache miss. Which index does this block get loaded into and which 3 other addresses are loaded along with it?
- Vote on Ed Lessons

address: 0x 0b 10 / 10 / 10

index      offset

(value 2) (value 2)

cache:

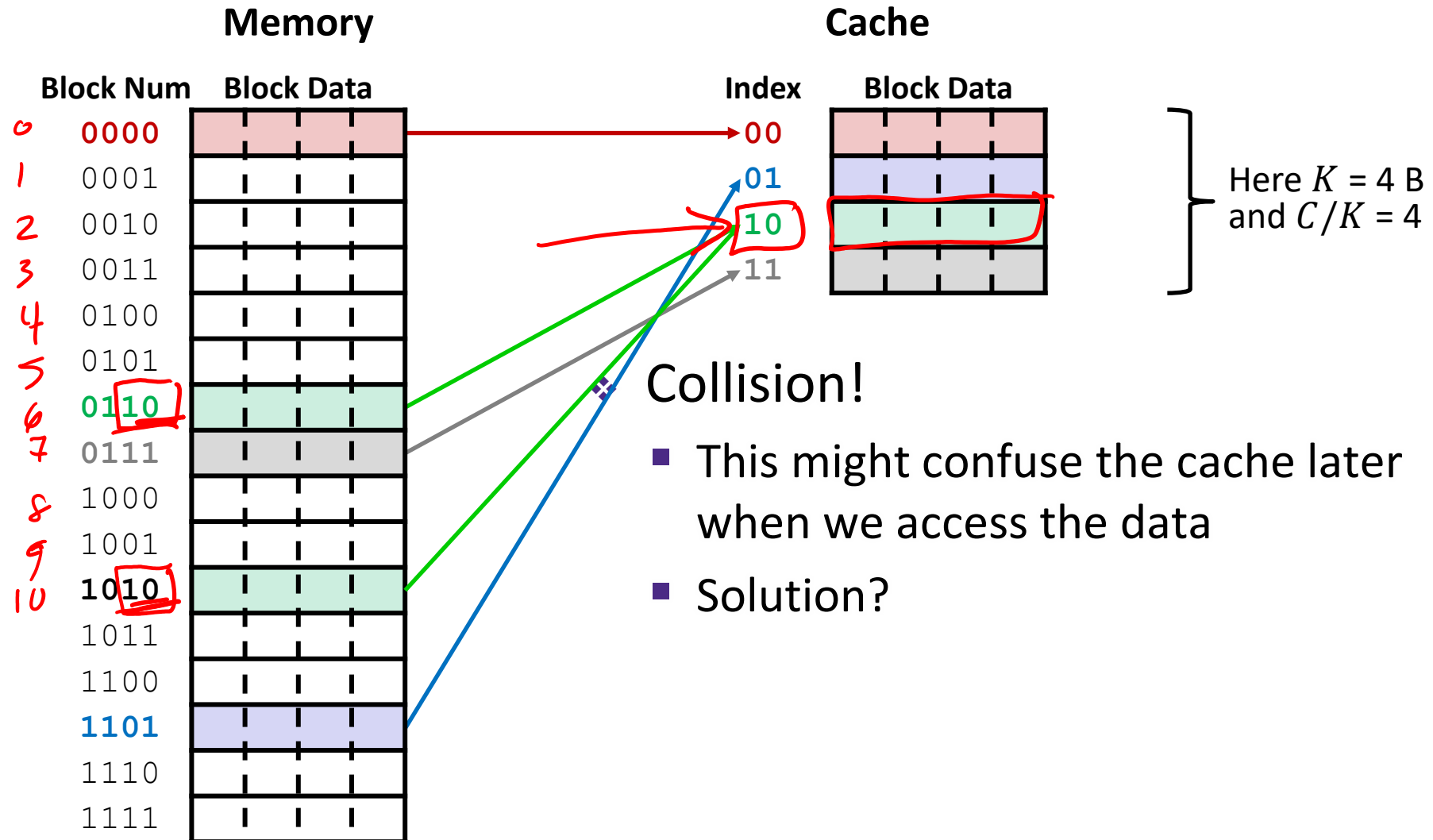
Index	00	01	10	11
00				
01				
→ 10	38	29	2A	2B
11				

↓      ↓      ↓      ↓

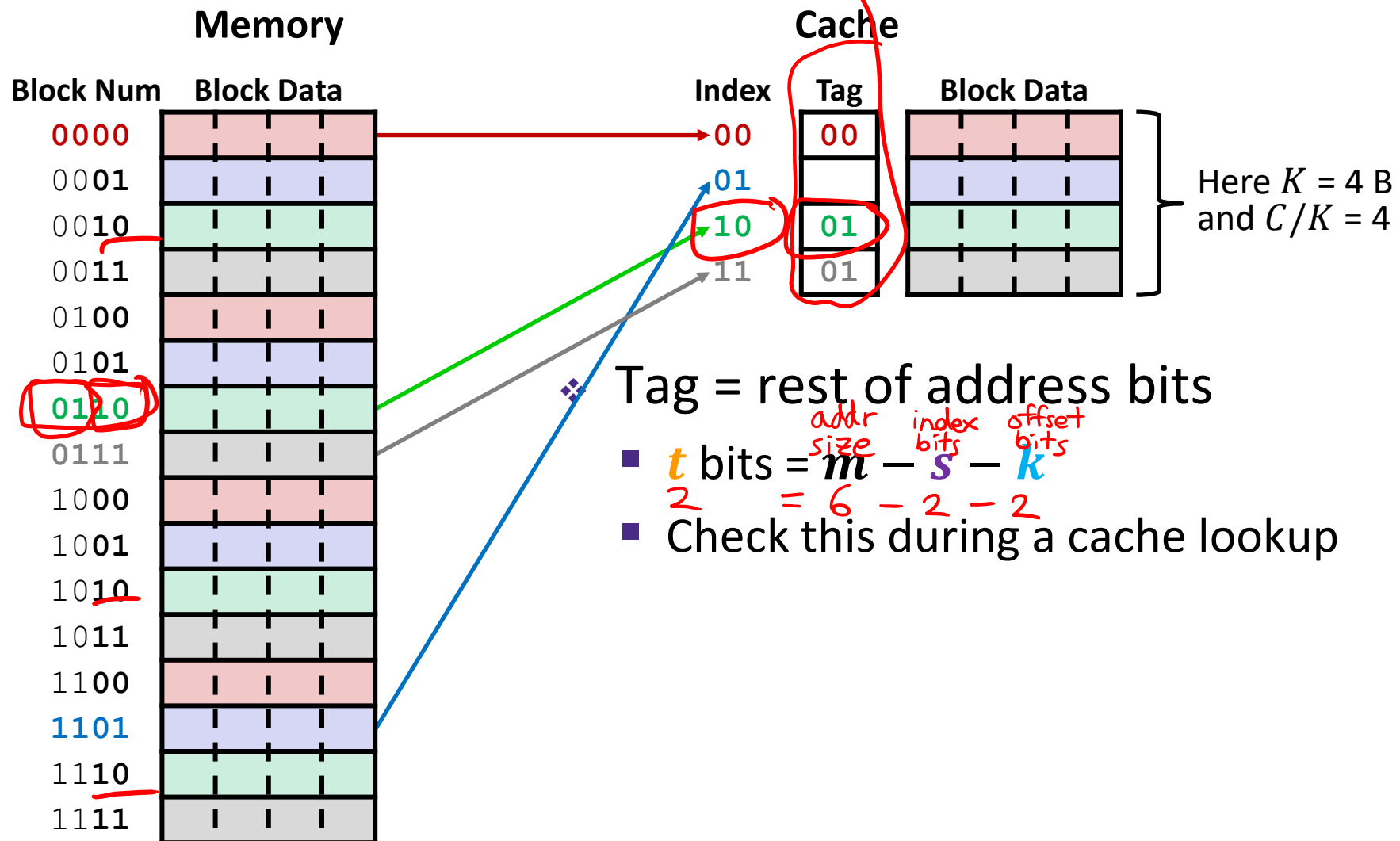
Index 2 (0b10)

along with: 0b1010 00 = 0x28  
 01 = 0x29  
 11 = 0x2B

# Place Data in Cache by Hashing Address



# Tags Differentiate Blocks in Same Index

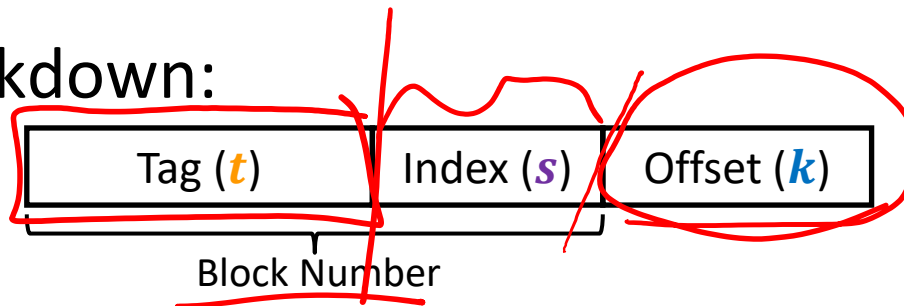


# Checking for a Requested Address

- ❖ CPU sends address request for chunk of data
  - Address and requested data are not the same thing!
    - Analogy: your friend  $\neq$  their phone number

- ❖ TIO address breakdown:

*m*-bit address:



- ① ■ **Index** field tells you where to look in cache
- ② ■ **Tag** field lets you check that data is the block you want
- ③ ■ **Offset** field selects specified start byte within block

- **Note:** *t* and *s* sizes will change based on hash function

# Cache Puzzle

- ❖ Based on the following behavior, which of the following block sizes is NOT possible for our cache?
  - Cache starts *empty*, also known as a **cold cache**
  - Access (addr: hit/miss) stream:
    - (14: miss), (15: hit), (16: miss)

hit: block with data already in \$  
miss: data not in \$, pulls block containing data from Mem

- ❖ [Not in Ed Lessons]
  - ③ 16 is in a different block
  - ② 14 & 15 are in the same block
  - ① pulls block containing 14 into \$

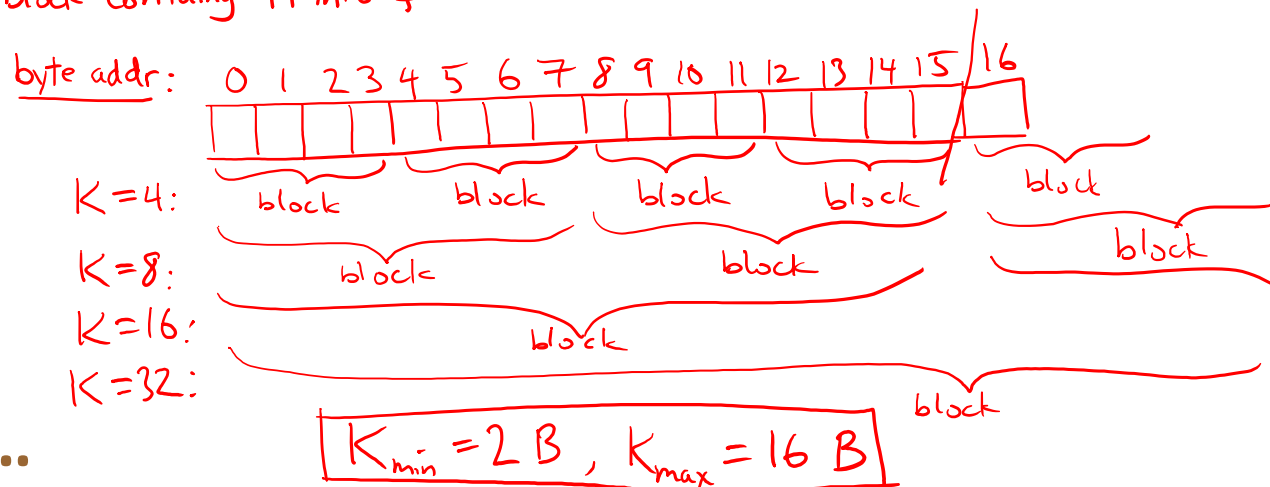
A. 4 bytes

B. 8 bytes

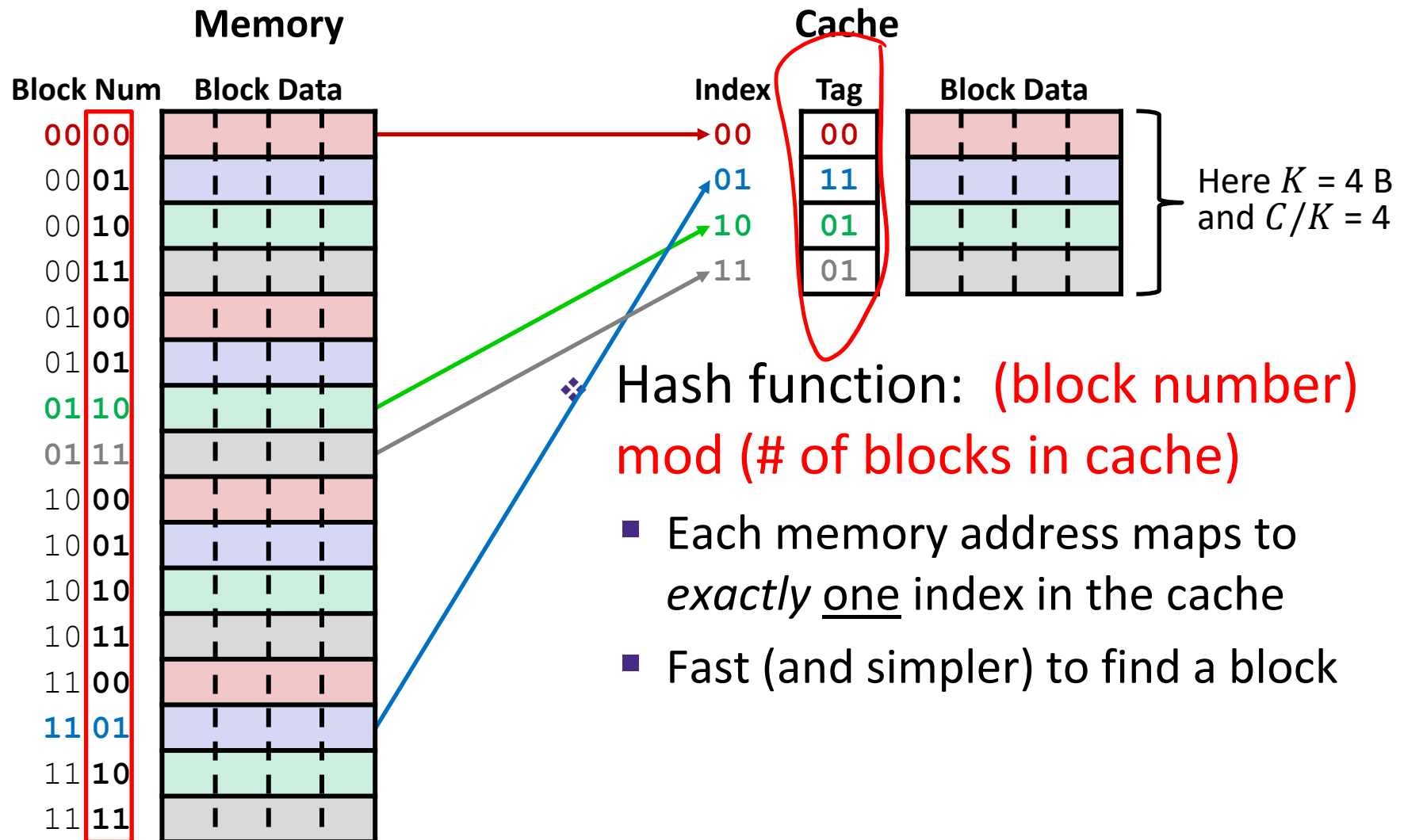
C. 16 bytes

D. 32 bytes

E. We're lost...



# Summary: Direct-Mapped Cache





# Direct-Mapped Cache Problem

