

Caches II

CSE 351 Summer 2020

Instructor:

Porter Jones

Teaching Assistants:

Amy Xu

Callum Walker

Sam Wolfson

Tim Mandzyuk



Administrivia

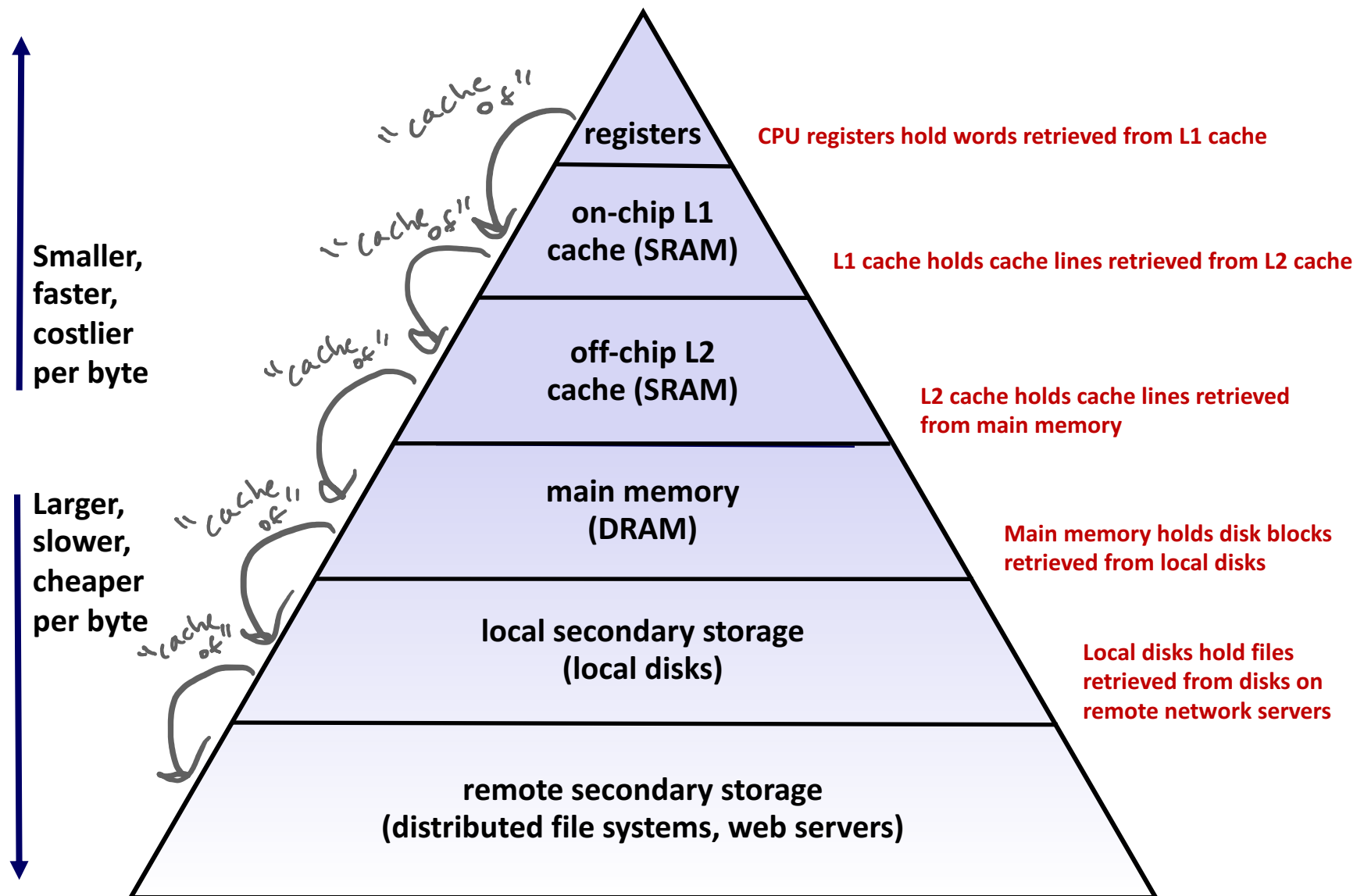
- ❖ Questions doc: <https://tinyurl.com/CSE351-7-29>
- ❖ hw15 due Friday (7/31) – 10:30am
- ❖ No homework due Monday!
- ❖ Lab 3 due Friday (7/31) – 11:59pm
 - You get to write some buffer overflow exploits!
- ❖ Unit Summary 2 Due next Wednesday (8/5) – 11:59pm

Memory Hierarchies

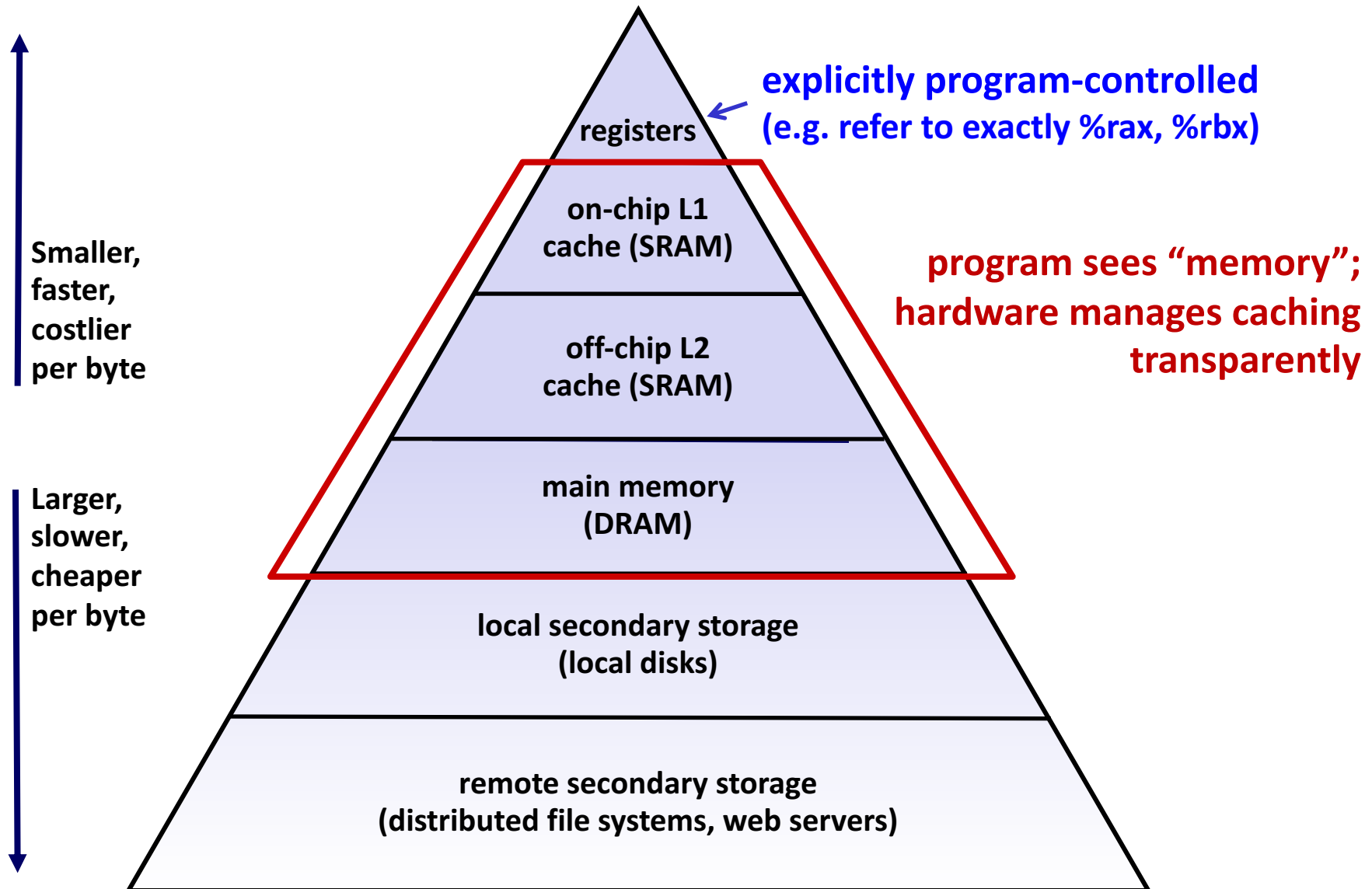
- ❖ Some fundamental and enduring properties of hardware and software systems:
 - Faster storage technologies almost always cost more per byte and have lower capacity
 - The gaps between memory technology speeds are widening
 - True for: registers \leftrightarrow cache, cache \leftrightarrow DRAM, DRAM \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 k , the faster, smaller device at level k serves as a cache for the larger, slower device at level $k+1$

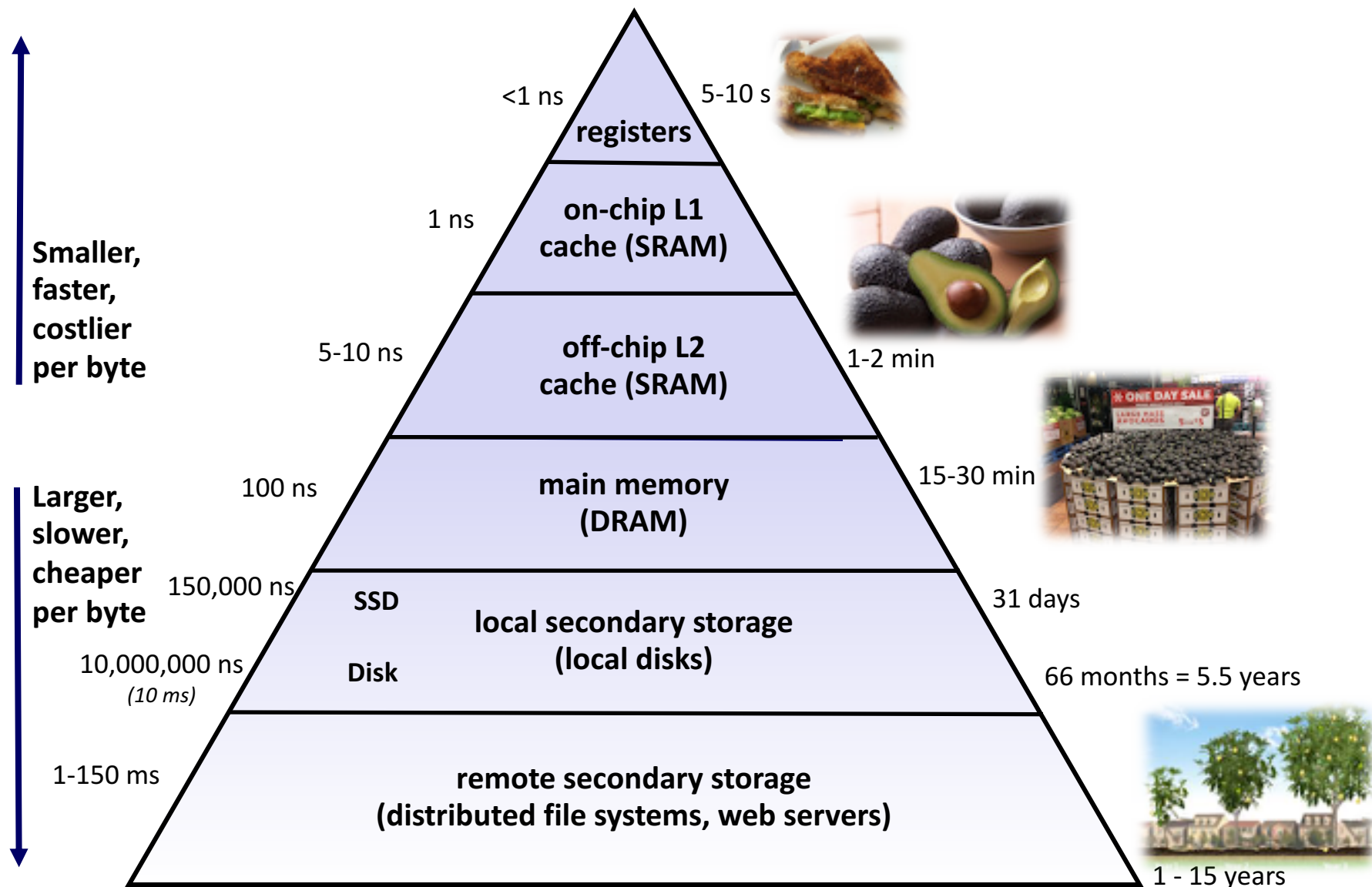
An Example Memory Hierarchy



An Example Memory Hierarchy



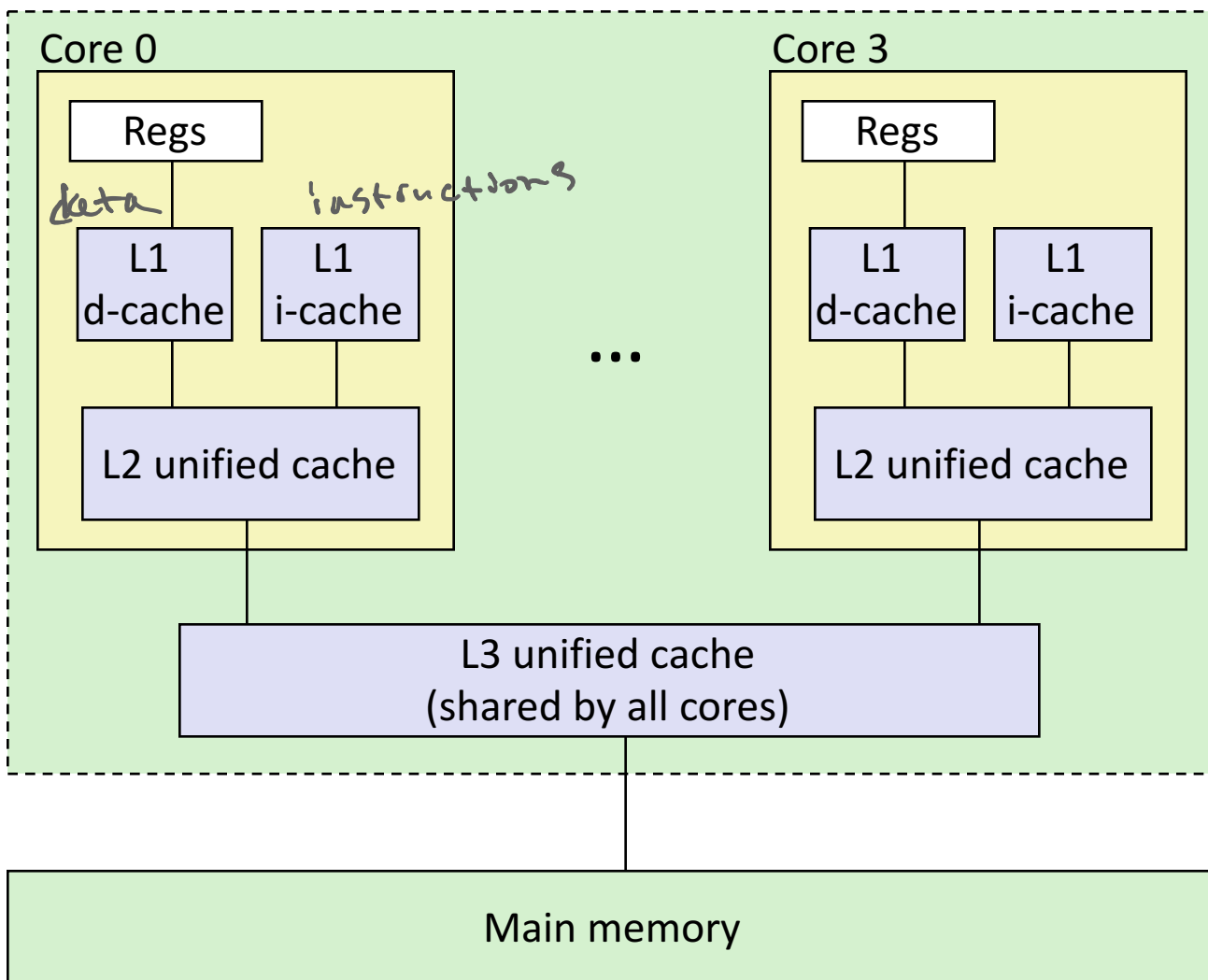
An Example Memory Hierarchy



Intel Core i7 Cache Hierarchy

Core: what executes instructions

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!

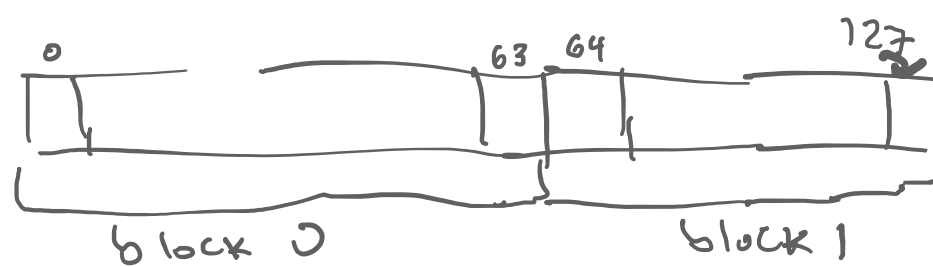
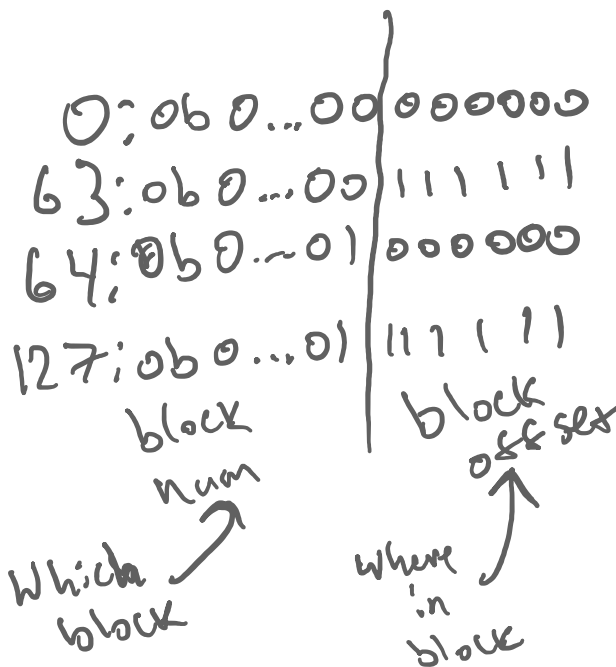
- ❖ Cache basics
- ❖ Principle of locality
- ❖ Memory hierarchies
- ❖ **Cache organization**
 - **Direct-mapped (*sets*; index + tag)**
 - **Associativity (*ways*)**
 - Replacement policy
 - Handling writes
- ❖ Program optimizations that consider caches

Cache Organization (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!

Lab 1a within same block



start of block

$$\text{addr} = \text{block size} * \text{block num} + \text{block offset}$$

Cache Organization (1)

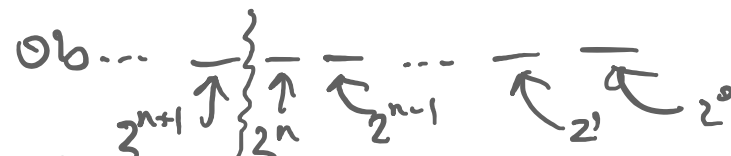
Note: The textbook uses "b" for offset bits

- ❖ **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!

$x \% 2^n = \text{value of lowest } n \text{ bits}$

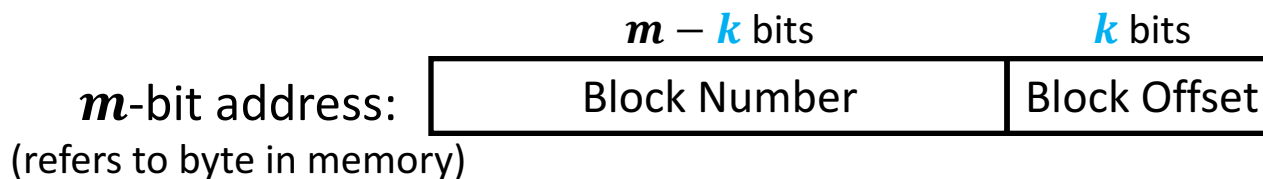
- ❖ **Offset field**

64
 2^6



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

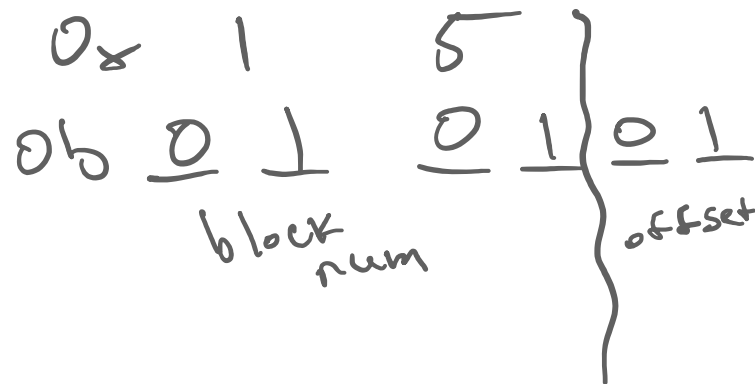
How many bits do I need to specify every byte in a block?



Polling Question [Cache II-a]

- ❖ If we have 6-bit addresses and block size $K = 4$ B, which block and byte does 0x15 refer to?
 - Vote at: <http://pollev.com/pbjones>

	Block Num	Block Offset
A.	1	1
B.	1	5
C.	5	1
D.	5	5
E.	We're lost...	



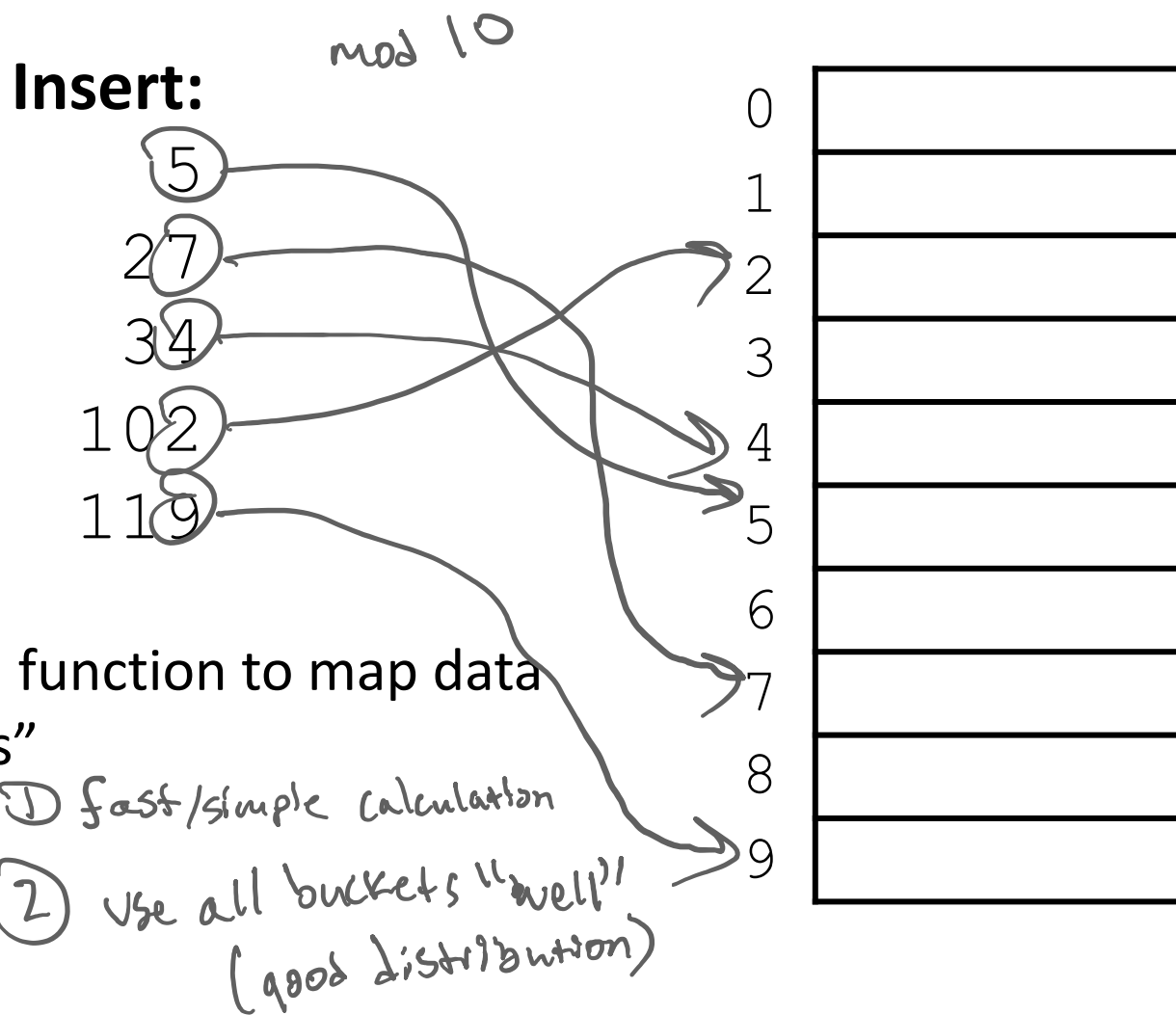
$$\text{offset width}(K) = \log_2(K) \\ = \log_2(4) = 2 \text{ bits}$$

Cache Organization (2)

in bytes

- ❖ **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: $C = 32 \text{ KiB} = 512 \text{ blocks}$ if using 64-B blocks
$$2^5 * 2^{10} / 2^6 = 2^9 \text{ blocks}$$
- ❖ 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!

Review: Hash Tables for Fast Lookup

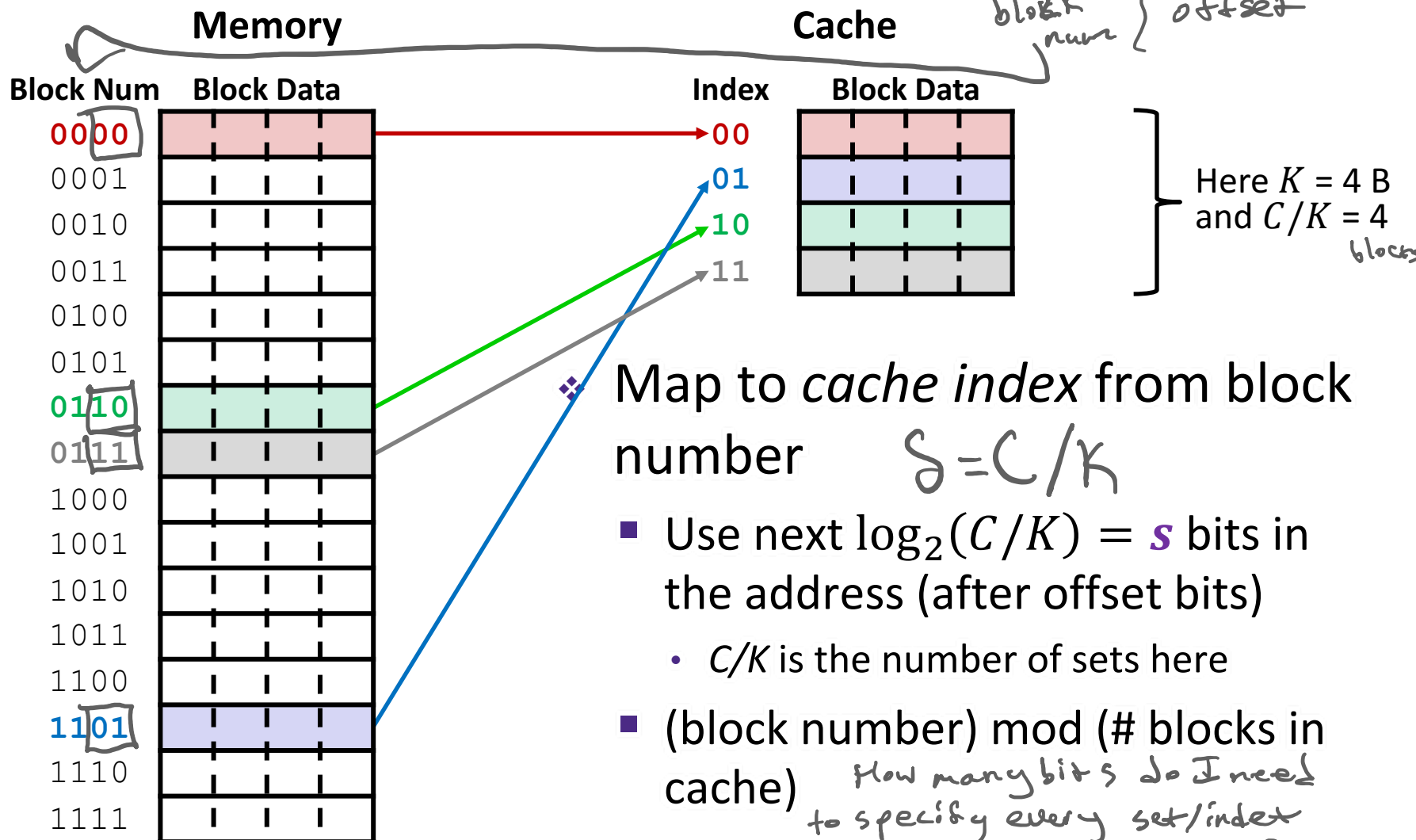


Apply hash function to map data to "buckets"

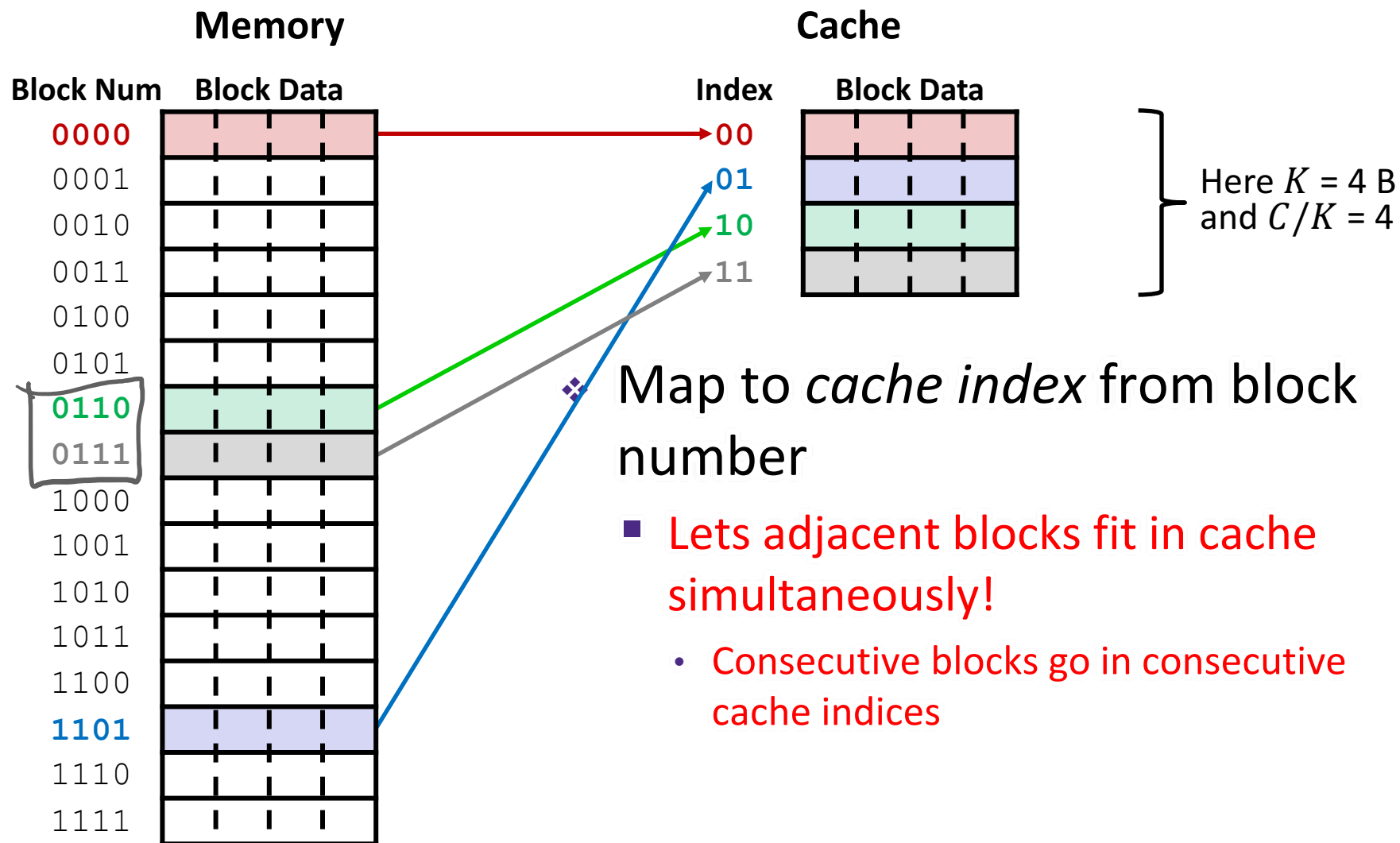
- Goals:
- ① fast/simple calculation
 - ② use all buckets "well" (good distribution)

Place Data in Cache by Hashing Address

addresses are 6-bits; 0b $\overset{s}{\times \times \times \times}$ } $\times \times$
 block num } offset



Place Data in Cache by Hashing Address



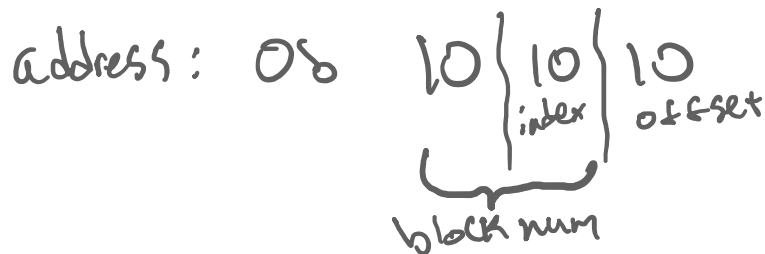
Practice Question

$\text{num offset bits} = k = \log_2(4) = 2$

❖ 6-bit addresses, block size $K = 4$ B, and our cache holds $S = 4$ blocks. $= C / K$ ($C = 16$ B) $\text{num index bits} = s = \log_2(4) = 2$

❖ A request for address **0x2A** results in a cache miss. Which set index does this block get loaded into and which 3 other addresses are loaded along with it?

■ No voting for this question

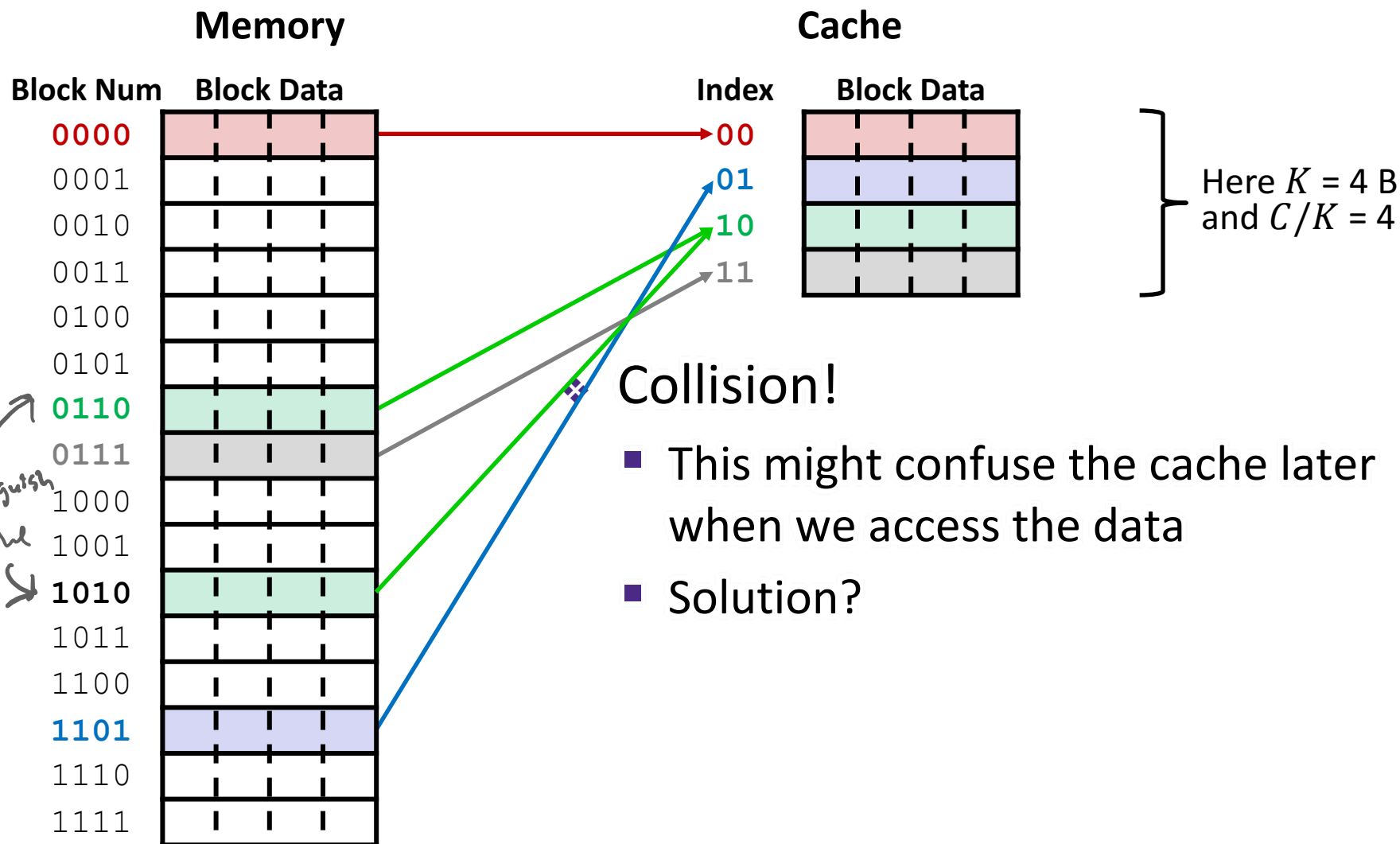


Index = 0b10 = 2

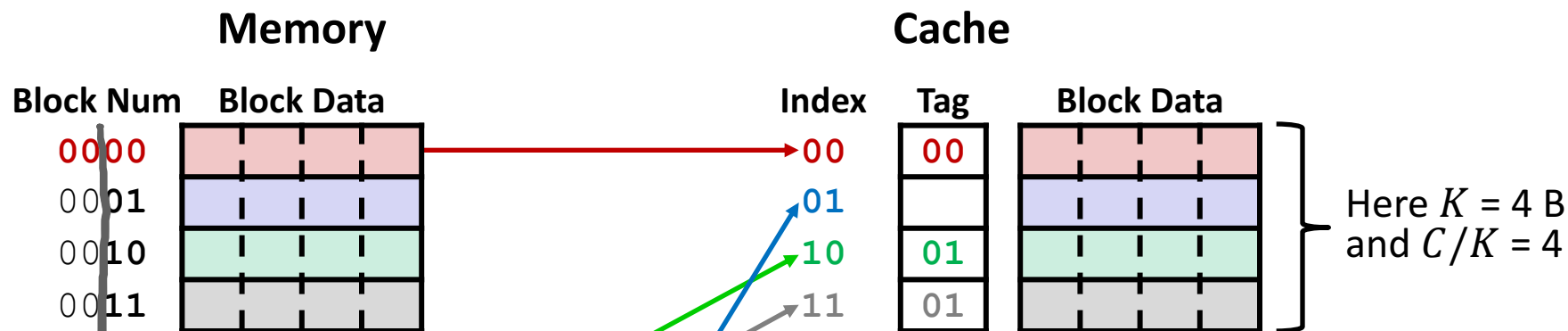
addresses w/ block num 1010

- 0b 101000 = 28
 - 0b 101001 = 29
 - 0b 101010 = 2A
 - 0b 101011 = 2B
- These are loaded into cache

Place Data in Cache by Hashing Address



Tags Differentiate Blocks in Same Index



Tag = rest of address bits

t bits = m - s - k

addr size index bits offset bits

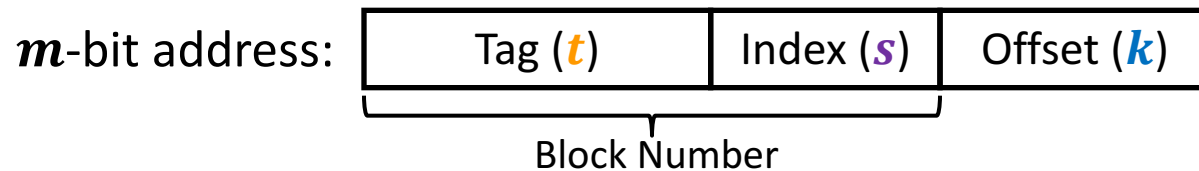
- Check this during a cache lookup

tag bits
index bits

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

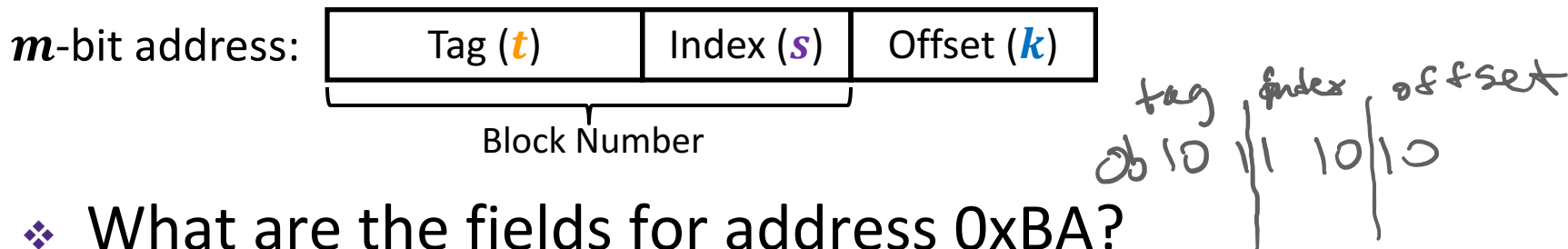
- ❖ TIO address breakdown:



- **Index** field tells you where to look in cache
- **Tag** field lets you check that data is the block you want
- **Offset** field selects specified start byte within block
- **Note:** *t* and *s* sizes will change based on hash function

Checking for a Requested Address Example

- ❖ Using 8-bit addresses.
- ❖ Cache Params: block size (K) = 4 B, cache size (C) = 32 B (which means number of sets is $C/K = 8$ sets).
 - Offset bits (k) = $\log_2(K) = \log_2(4) = 2$ bits
 - Index bits (s) = $\log_2(\text{num sets}) = \log_2(8) = 3$ bits
 - Tag bits (t) = Rest of the bits in the address = $8 - 2 - 3 = 3$ bits



- ❖ What are the fields for address 0xBA?
 - Tag bits (unique id for block): 0b101 = 5
 - Index bits (cache set block maps to): 0b110 = 6
 - Offset bits (byte offset within block): 0b10 = 2

Cache Puzzle [Cache II-b] [Vote at http://pollev.com/pbjones](http://pollev.com/pbjones)

❖ 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)

① pulls block w/14 into cache

② 15 is in the same block as 14

③ 16 is not in block w/14 and 15

hit: block is already in cache!
miss: block is not in cache, pulls block from memory and puts it in cache

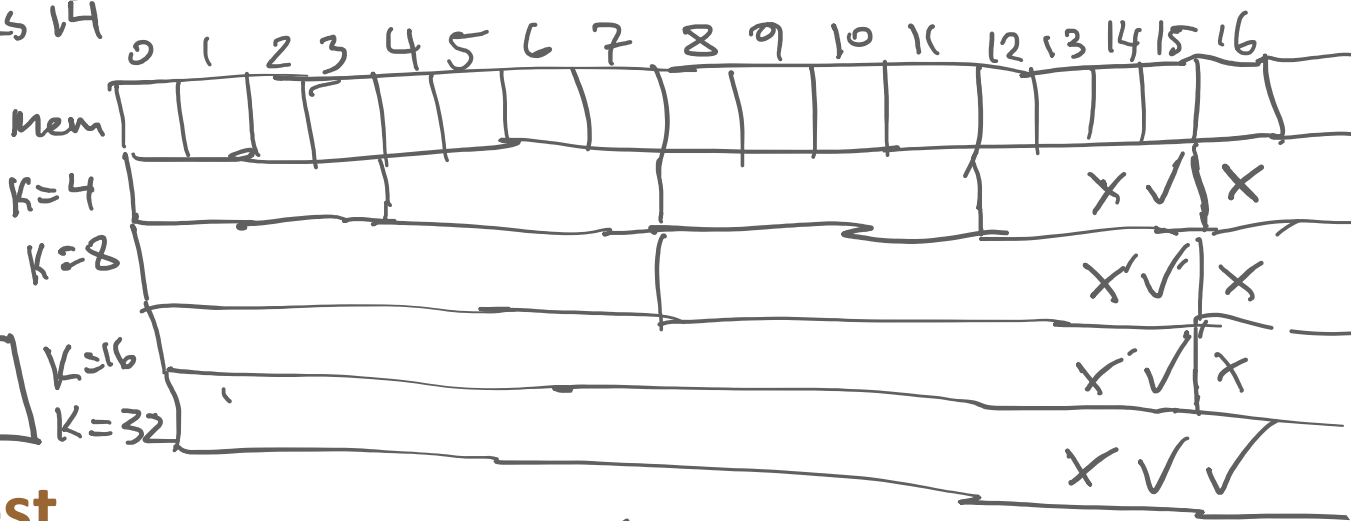
A. 4 bytes

B. 8 bytes

C. 16 bytes

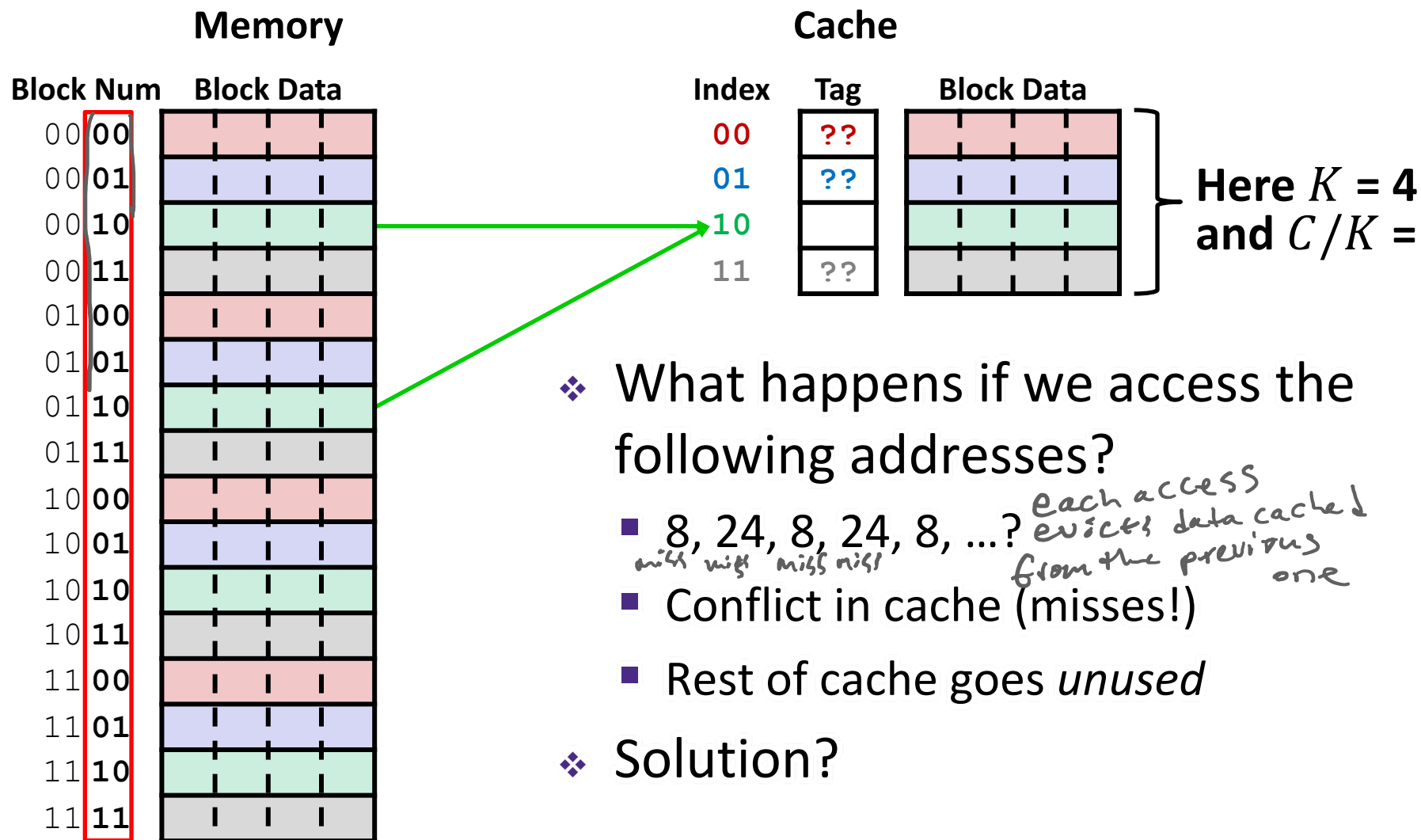
D. 32 bytes

E. We're lost...



✓ = hit x = miss

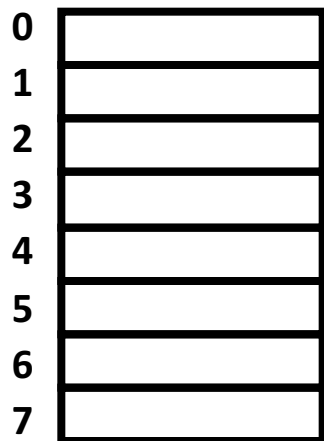
Direct-Mapped Cache Problem



Associativity

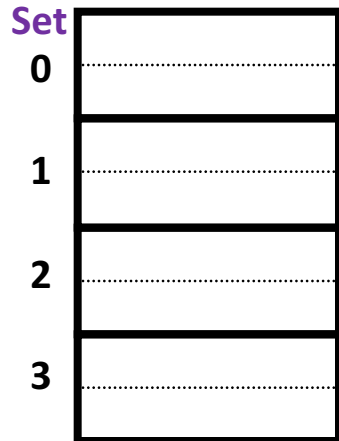
- ❖ What if we could store data in any place in the cache?
 - More complicated hardware = more power consumed, slower
- ❖ So we *combine* the two ideas:
 - Each address maps to exactly one **set**
 - Each set can store block in more than one **way**

1-way:
8 sets,
1 block each

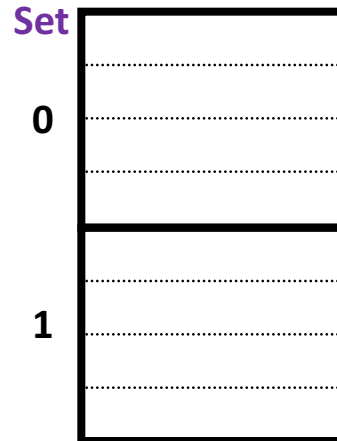


direct-mapped

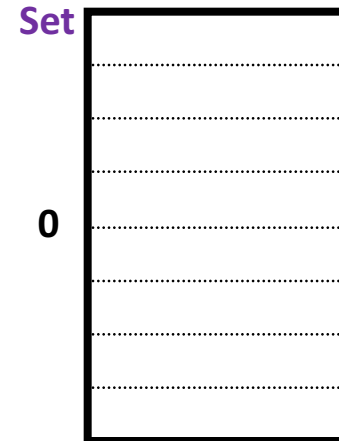
2-way:
4 sets,
2 blocks each



4-way:
2 sets,
4 blocks each



8-way:
1 set,
8 blocks



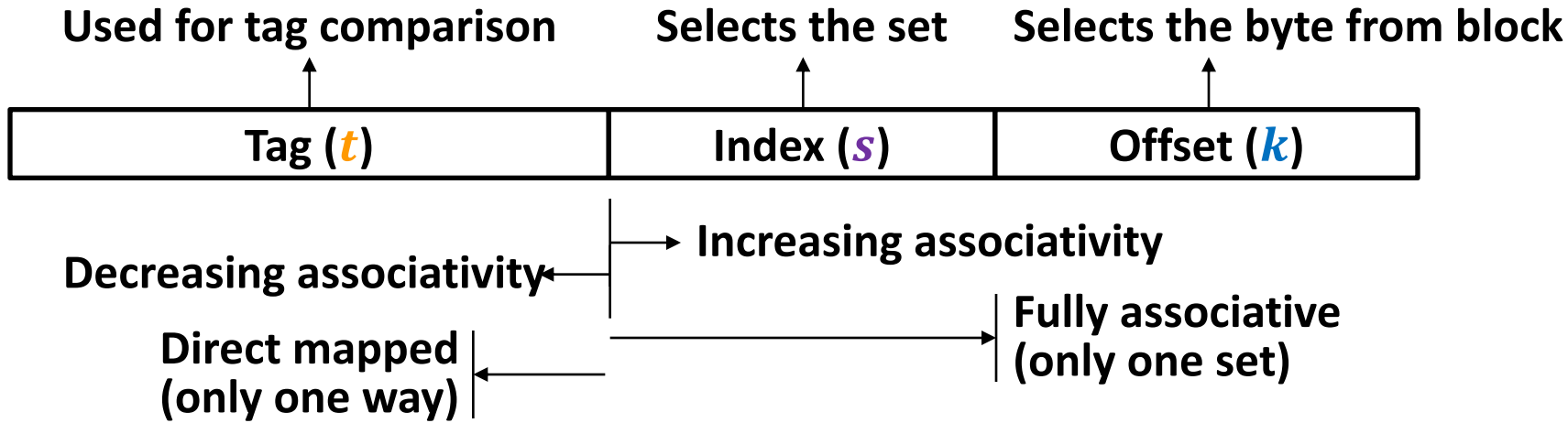
fully associative

Cache Organization (3)

Note: The textbook uses "b" for offset bits

- ❖ **Associativity (E):** # of ways for each set
 - Such a cache is called an "*E-way set associative cache*"
 - We now index into cache *sets*, of which there are $S = C/K/E$

\uparrow *num blocks* \uparrow *blocks per set*
 - Use lowest $\log_2(C/K/E) = s$ bits of block address
 - Direct-mapped: $E = 1$, so $s = \log_2(C/K)$ as we saw previously
 - Fully associative: $E = C/K$, so $s = 0$ bits



Example Placement

block size:	16 B	K
capacity:	8 blocks	C/K
address:	16 bits	M

❖ Where would data from address 0×1833 be placed?

■ Binary: 0b 0001 1000 0011 0011



$$t = m - s - k \quad s = \log_2(C/K/E) \quad k = \log_2(K)$$

m -bit address:

Tag (t)	Index (s)	Offset (k)
-----------------------------	-------------------------------	--------------------------------

$\log_2(8)$
 $S = ?$ 3 bits
 Direct-mapped

$\log_2(4)$
 $S = ?$ 2 bits
 2-way set associative

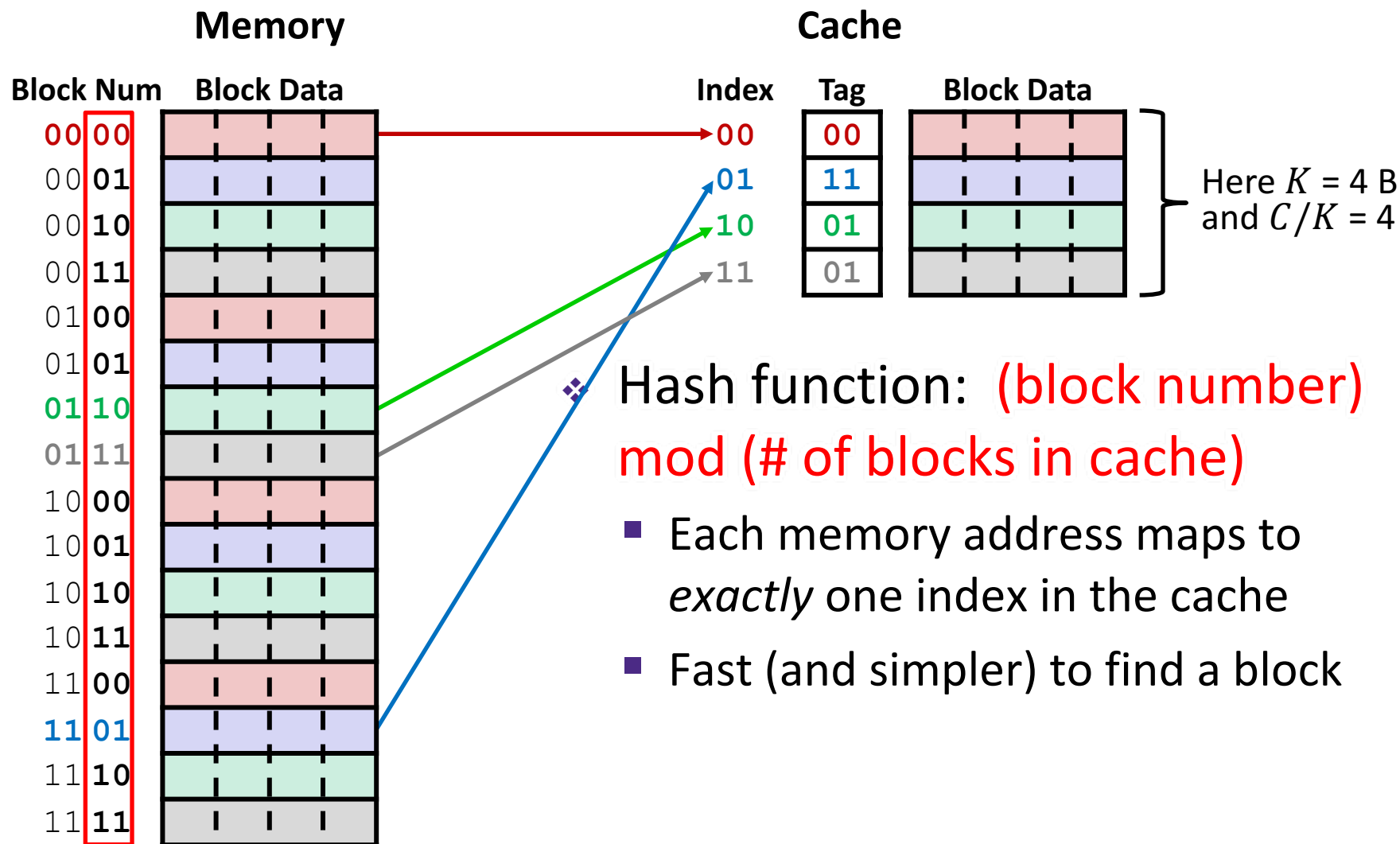
$\log_2(2)$
 $S = ?$ 1 bit
 4-way set associative

	Set	Tag	Data
000	0		
001	1		
010	2		
011	3		✓
100	4		
101	5		
110	6		
111	7		

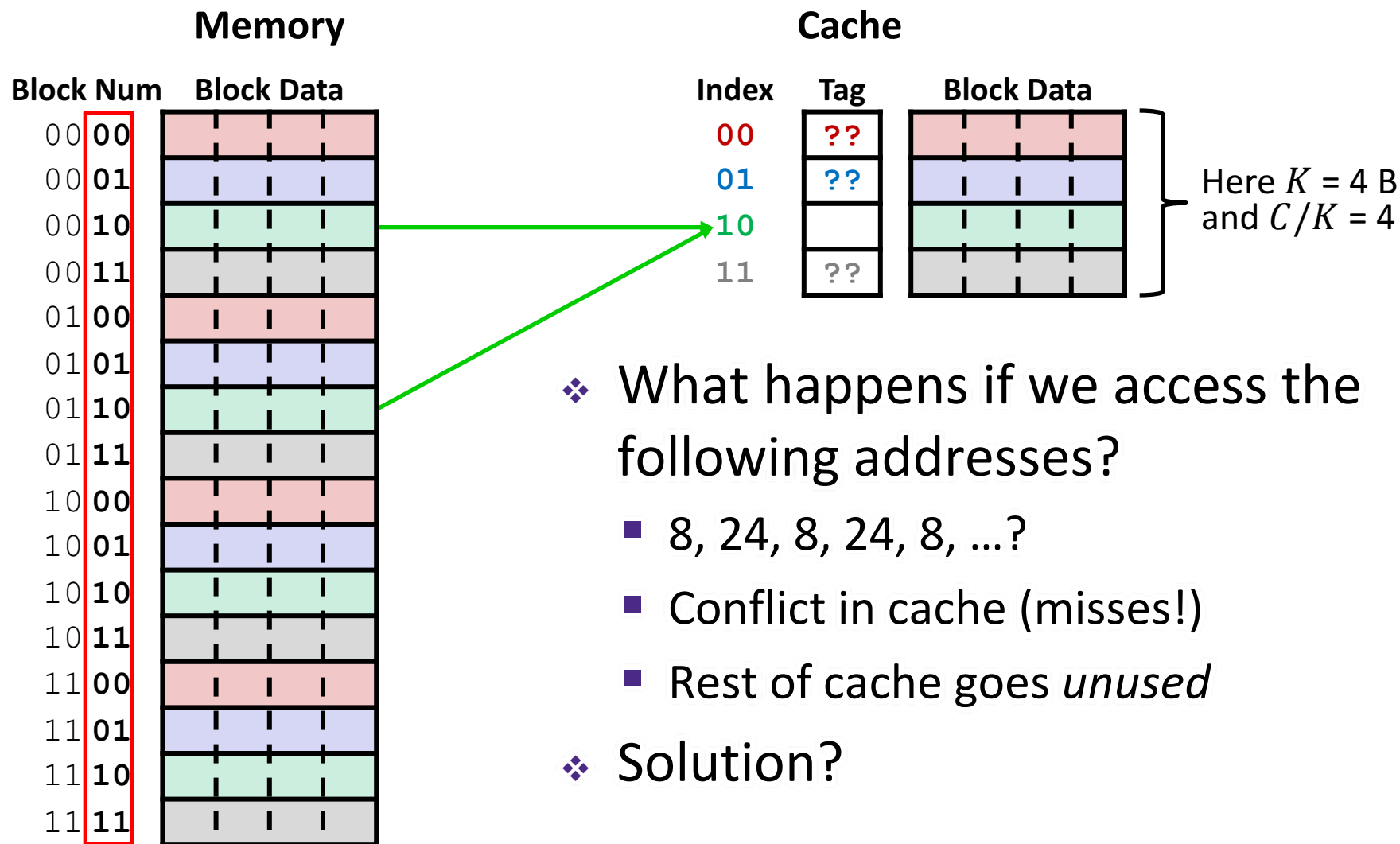
	Set	Tag	Data
00	0		
01	1		
10	2		
11	3		✓
			✓

	Set	Tag	Data
0			
1			✓
			✓
			✓
			✓

Direct-Mapped Cache



Direct-Mapped Cache Problem



Notes Diagrams

