

Caches II

CSE 351 Autumn 2017

Instructor:

Justin Hsia

Teaching Assistants:

Lucas Wotton

Michael Zhang

Parker DeWilde

Ryan Wong

Sam Gehman

Sam Wolfson

Savanna Yee

Vinny Palaniappan



Administrivia

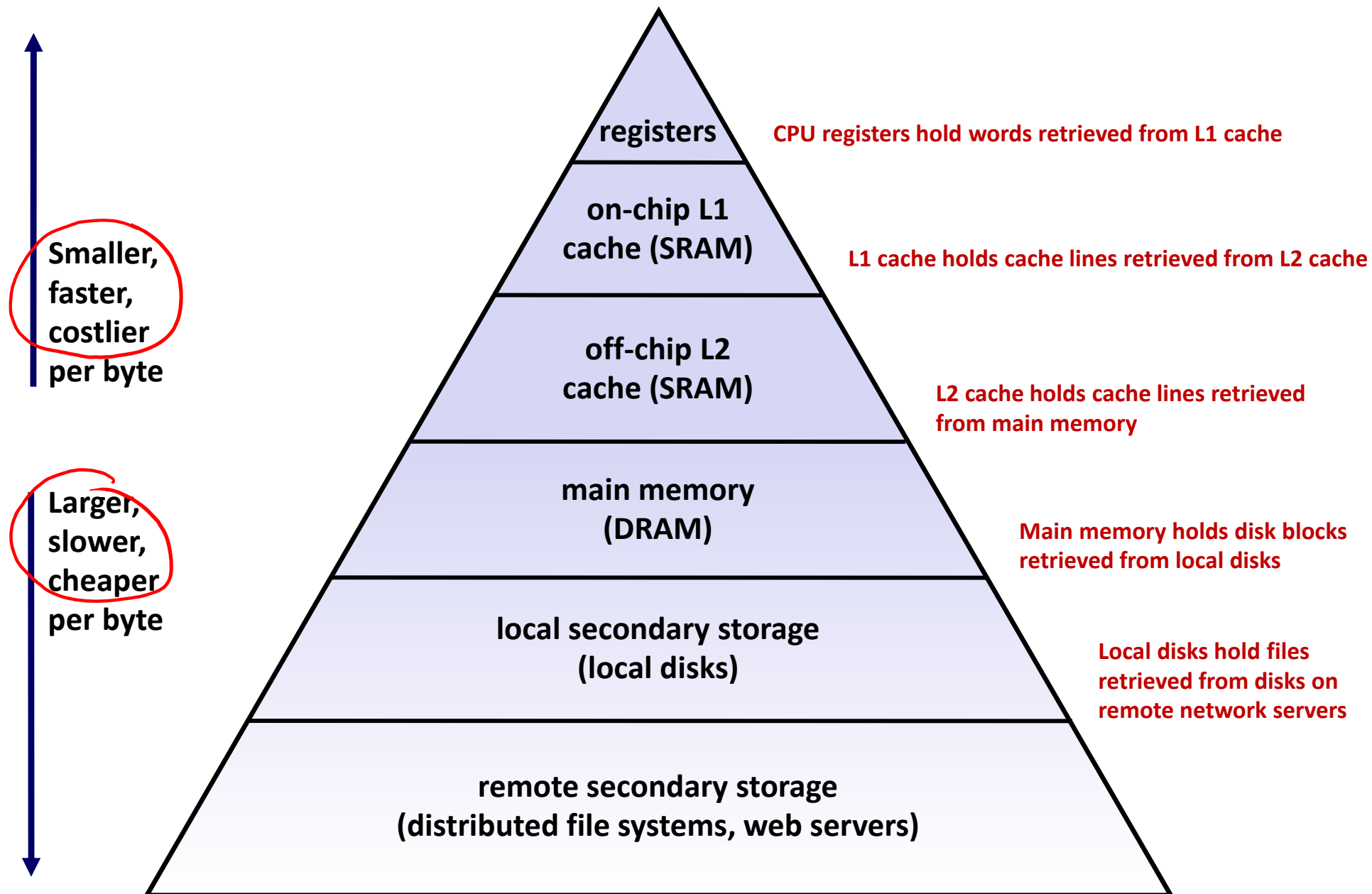
- ❖ Homework 4 released tomorrow (Structs, Caches)
- ❖ Midterm Regrade Requests due Wednesday (11/8)
- ❖ Lab 3 due *Friday* (11/10)

- ❖ **Mid-Quarter Survey Feedback**
 - Pace is “moderate” to “a bit too fast”
 - You talk too fast in lecture (or rush at the end) and I wish there were more peer instruction questions
 - Canvas quiz answer keys are annoying, but instant homework feedback is great

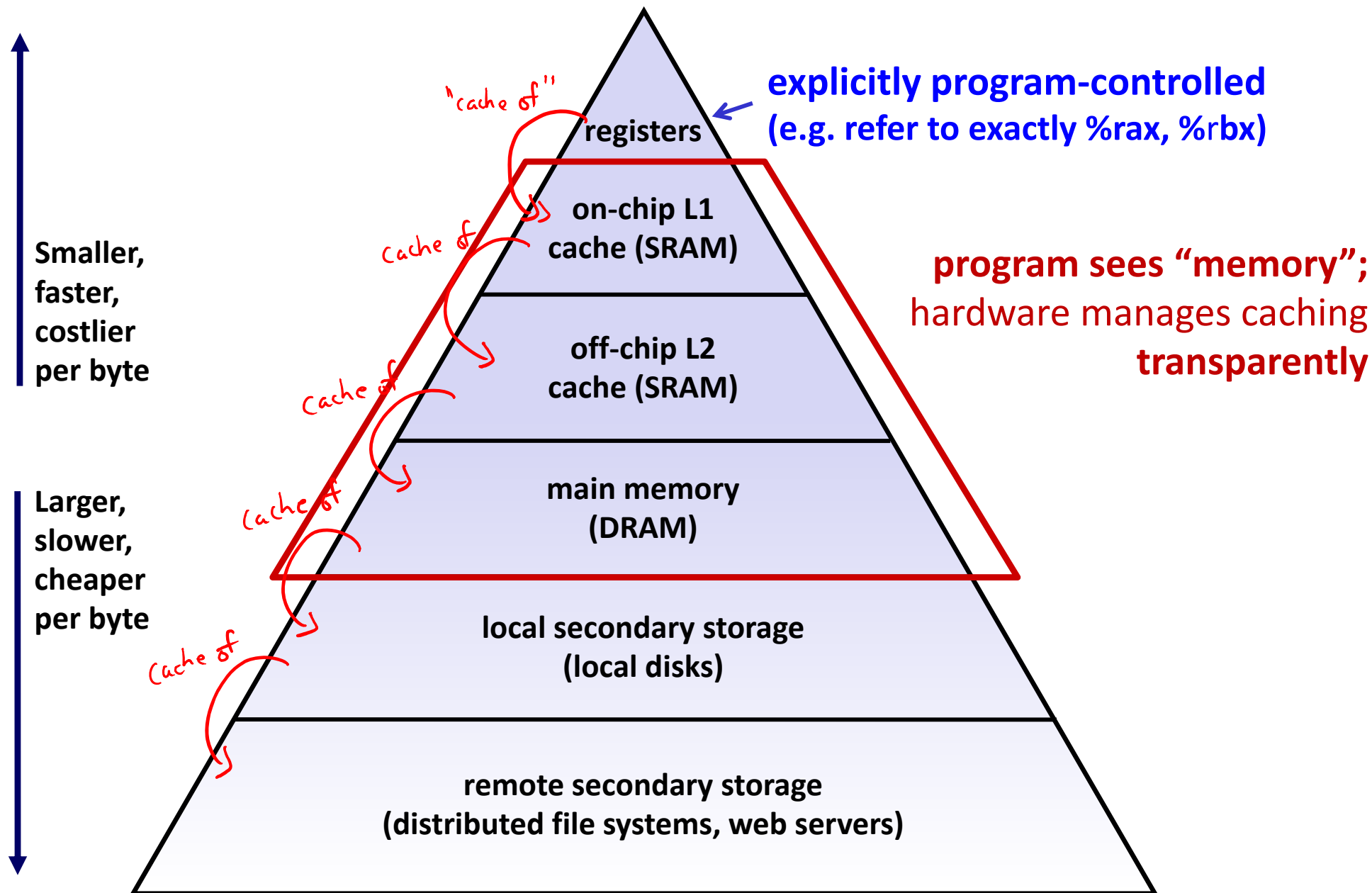
Memory Hierarchies

- ❖ 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

An Example Memory Hierarchy



An Example Memory Hierarchy

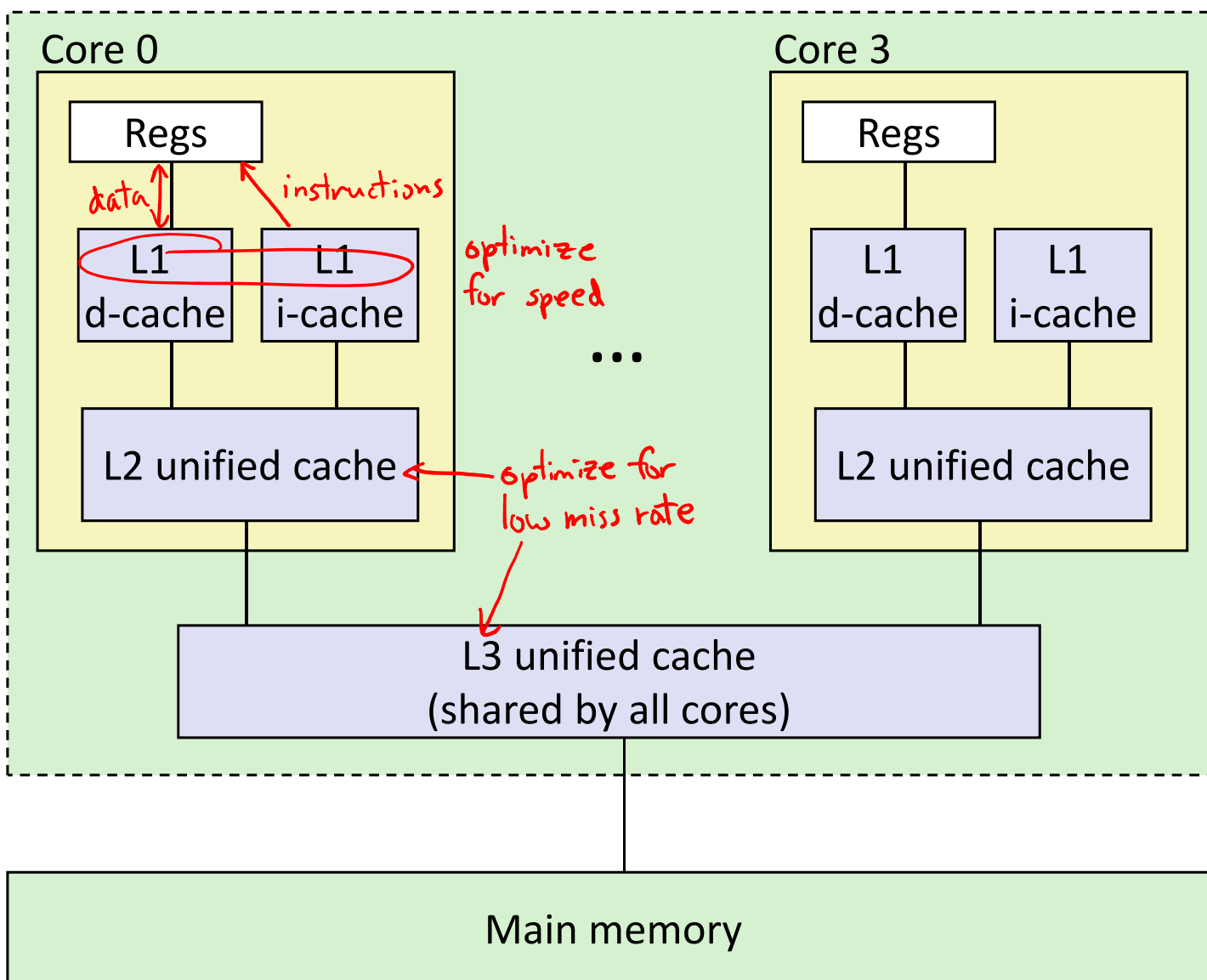


Memory Hierarchies

- ❖ Fundamental idea of 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 $k+1$
- ❖ Why do memory hierarchies work?
 - Because of locality, programs tend to access the data at level k more often than they access the data at level $k+1$
 - Thus, the storage at level $k+1$ can be slower, and thus larger and cheaper per bit
- ❖ *Big Idea*: The memory hierarchy creates a large pool of storage that costs as much as the cheap storage near the bottom, but that serves data to programs at the rate of the fast storage near the top

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

Cache Organization (1)

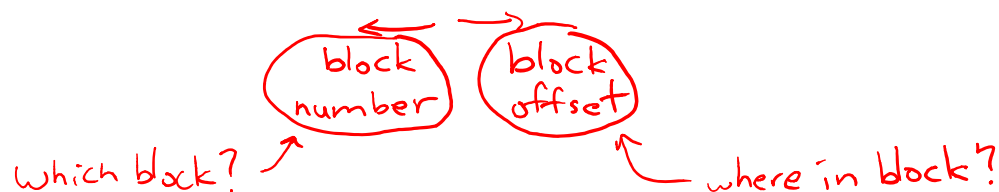
Note: The textbook uses "B" for block size

- ❖ **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 1: within Block



0: 0b 0...00 | 000 000
 63: 0b 0...00 | 111 111
 64: 0b 0...01 | 000 000
 127: 0b 0...01 | 111 111



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! *each bit in address has value 2^i*

$$2^i \bmod 2^n = \begin{cases} 0, & n \leq i \text{ (discard upper bits)} \\ 2^i, & n > i \text{ (keep lower bits)} \end{cases}$$

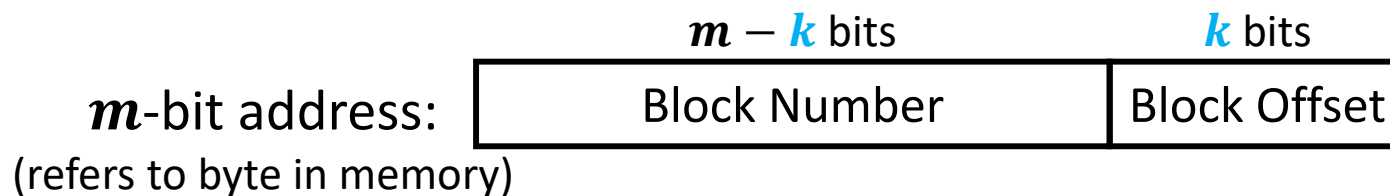
❖ **Offset field**

- 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?

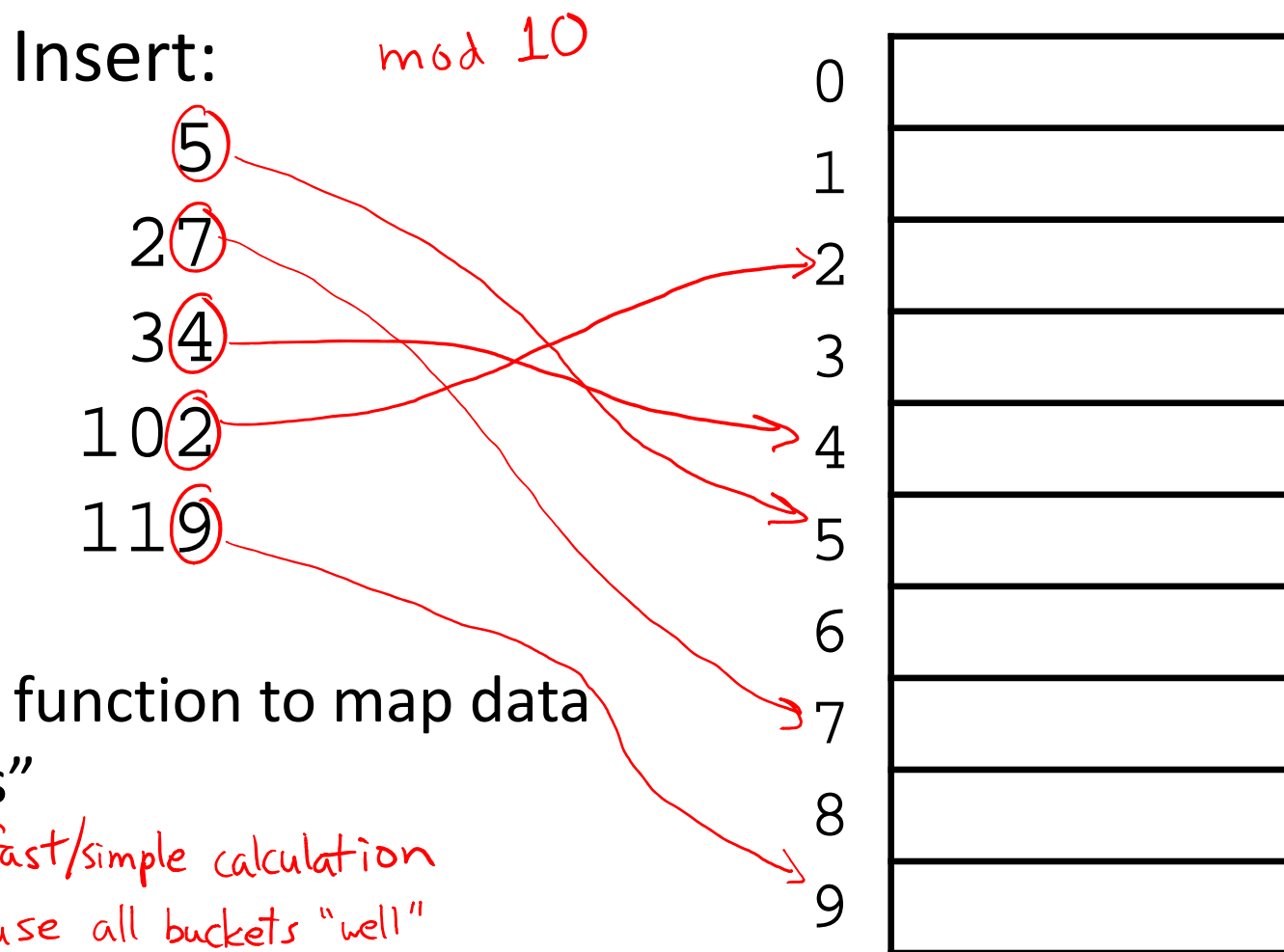


Cache Organization (2)

- ❖ **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

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

Review: Hash Tables for Fast Lookup

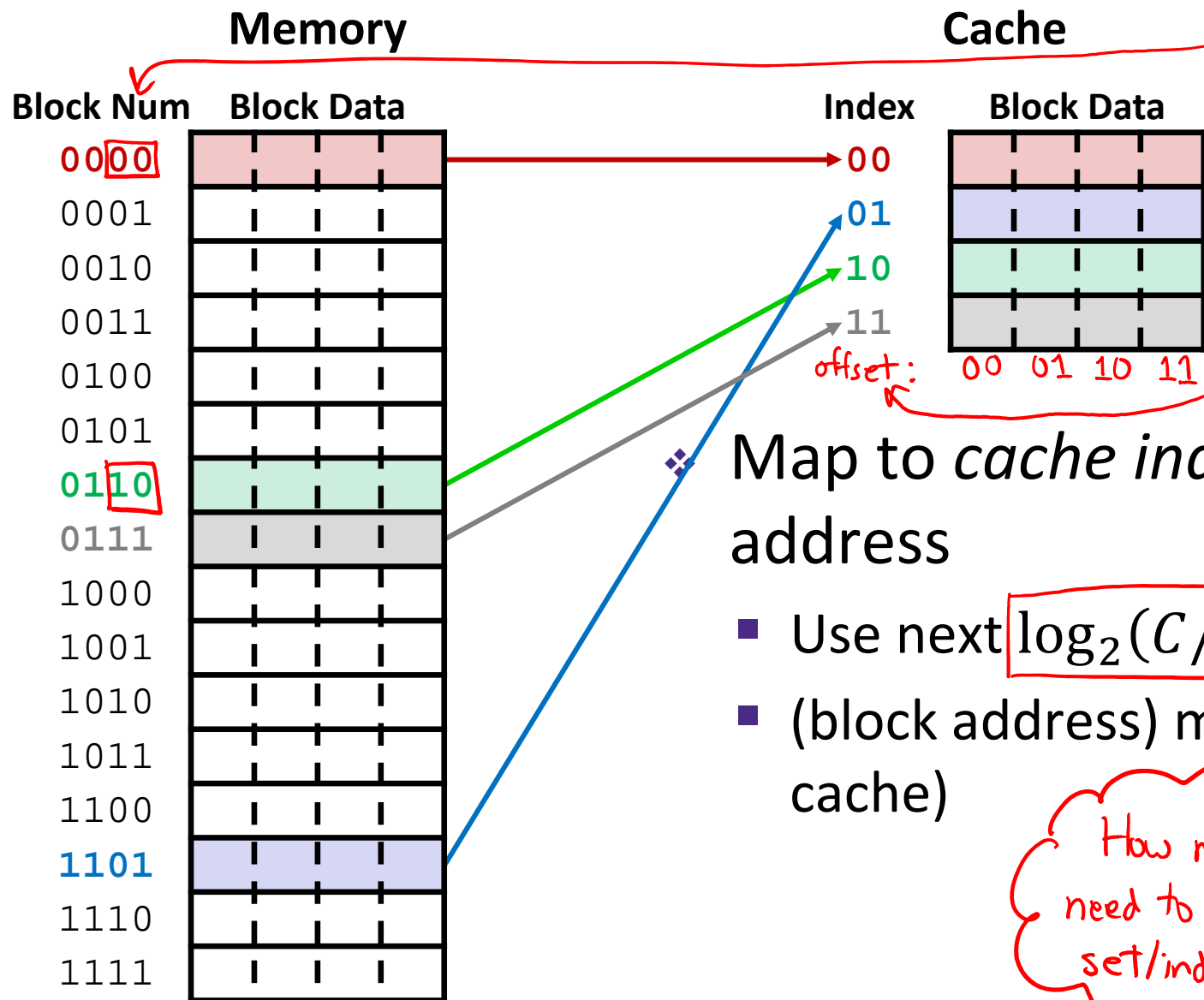


Apply hash function to map data to "buckets"

Goals: ① fast/simple calculation
 ② use all buckets "well"

Place Data in Cache by Hashing Address

addresses are 6 bits: $0b\ XX\ XX\ /\ XX$
 block num / offset



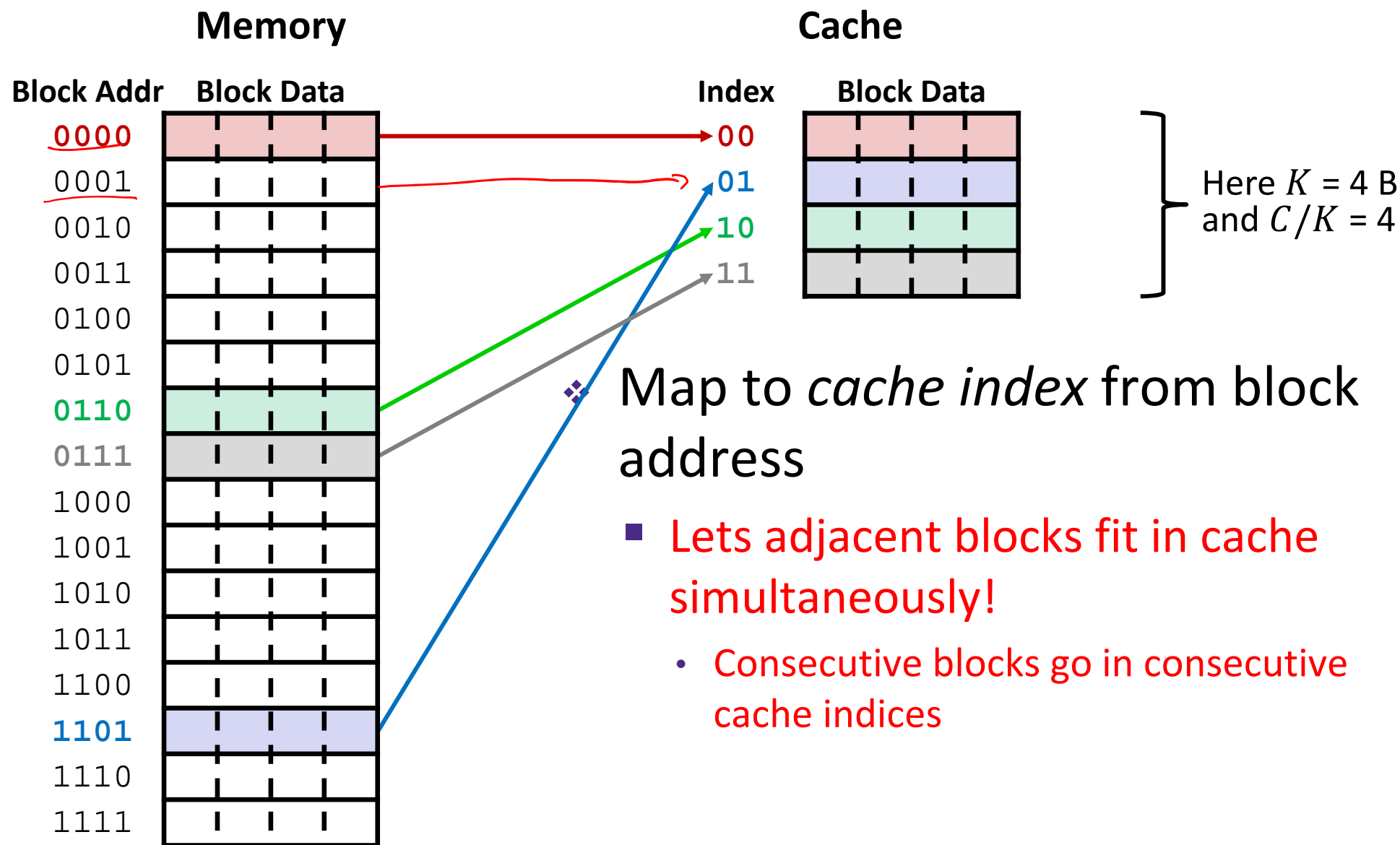
Here $K = 4\text{ B}$
 and $C/K = 4$
 blocks

Map to *cache index* from block address

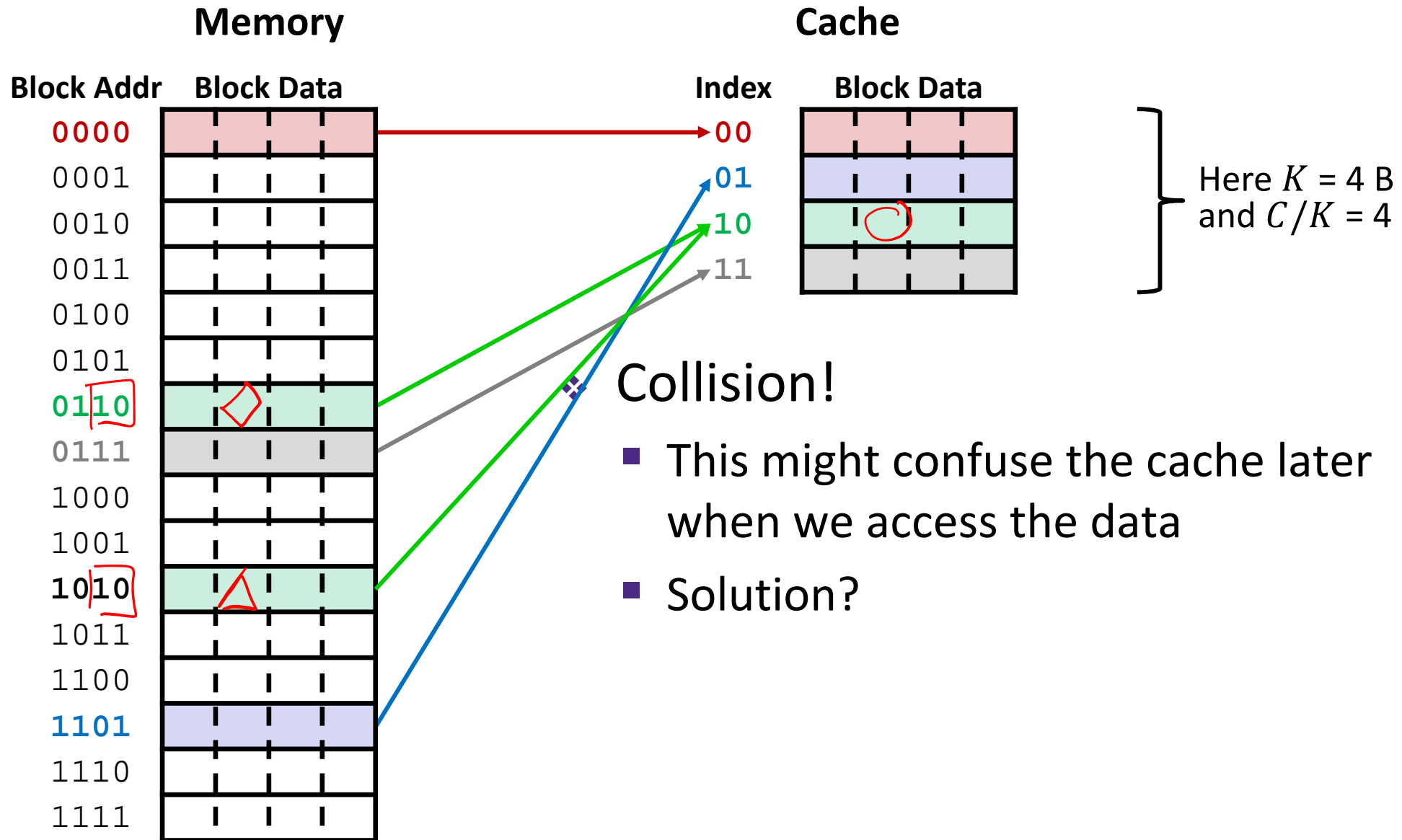
- Use next $\log_2(C/K) = s$ bits
- (block address) mod (# blocks in cache)

How many bits do I need to specify every set/index in my cache?

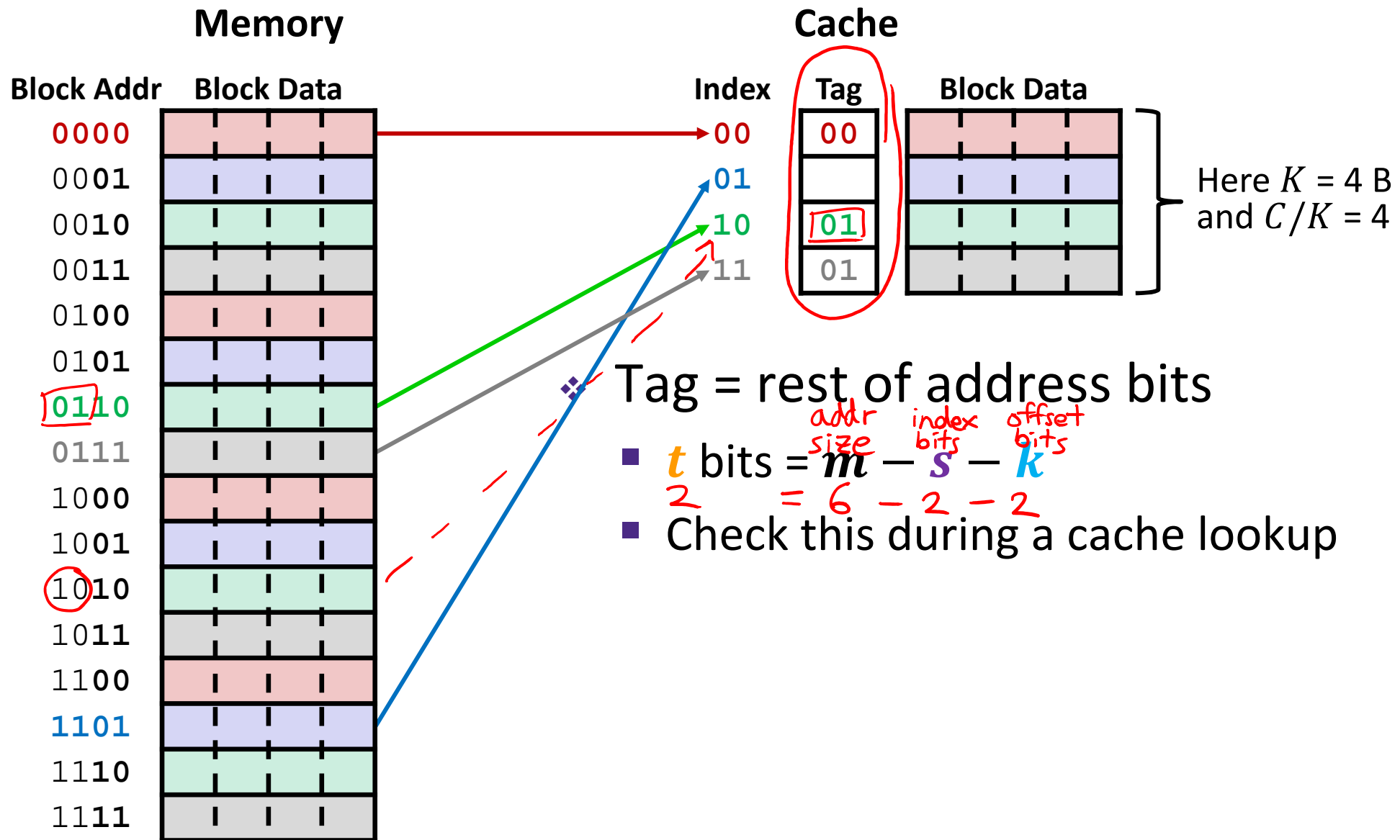
Place Data in Cache by Hashing Address



Place Data in Cache by Hashing Address



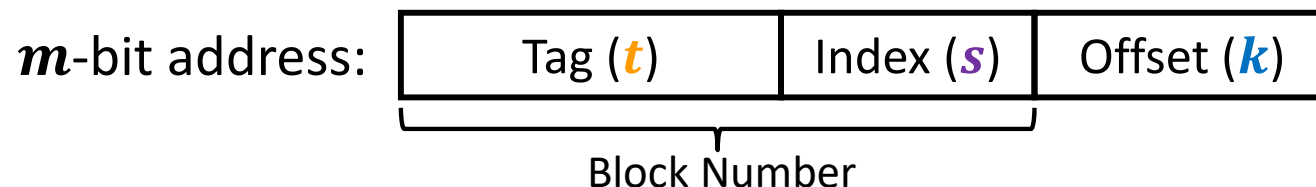
Tags Differentiate Blocks in Same Index



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 his or her phone number

- ❖ TIO address breakdown:



- 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 #1

Vote at <http://PollEv.com/justinh>

❖ 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

① pulls block containing 14 into \$
 ② 14 & 15 are in the same block
 ③

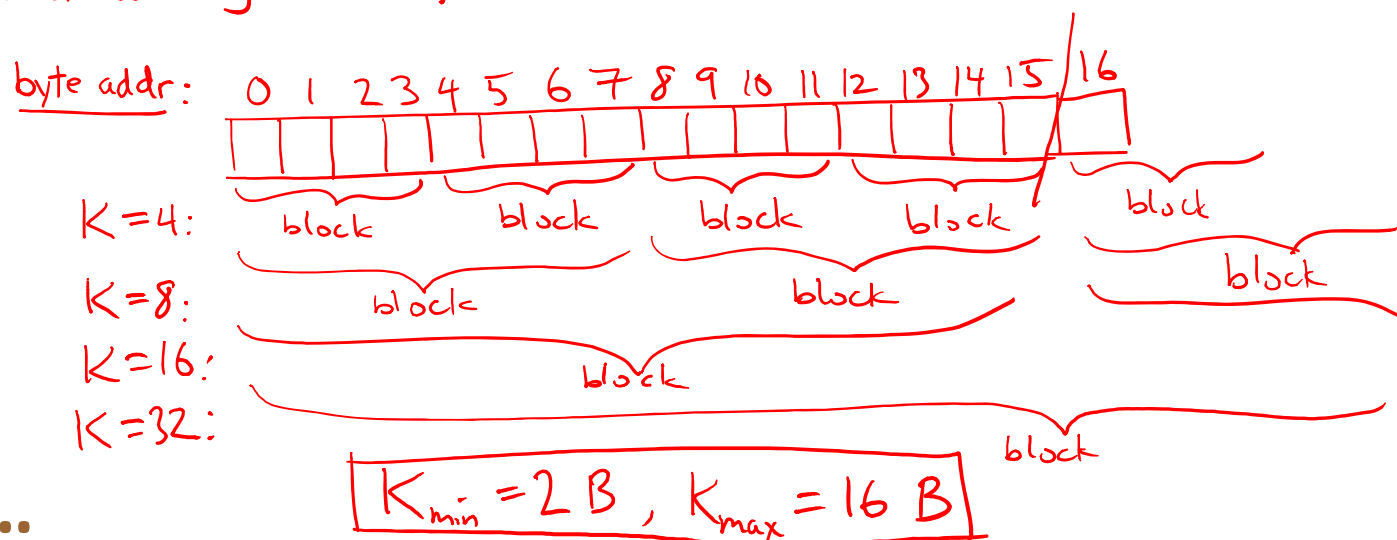
A. 4 bytes

B. 8 bytes

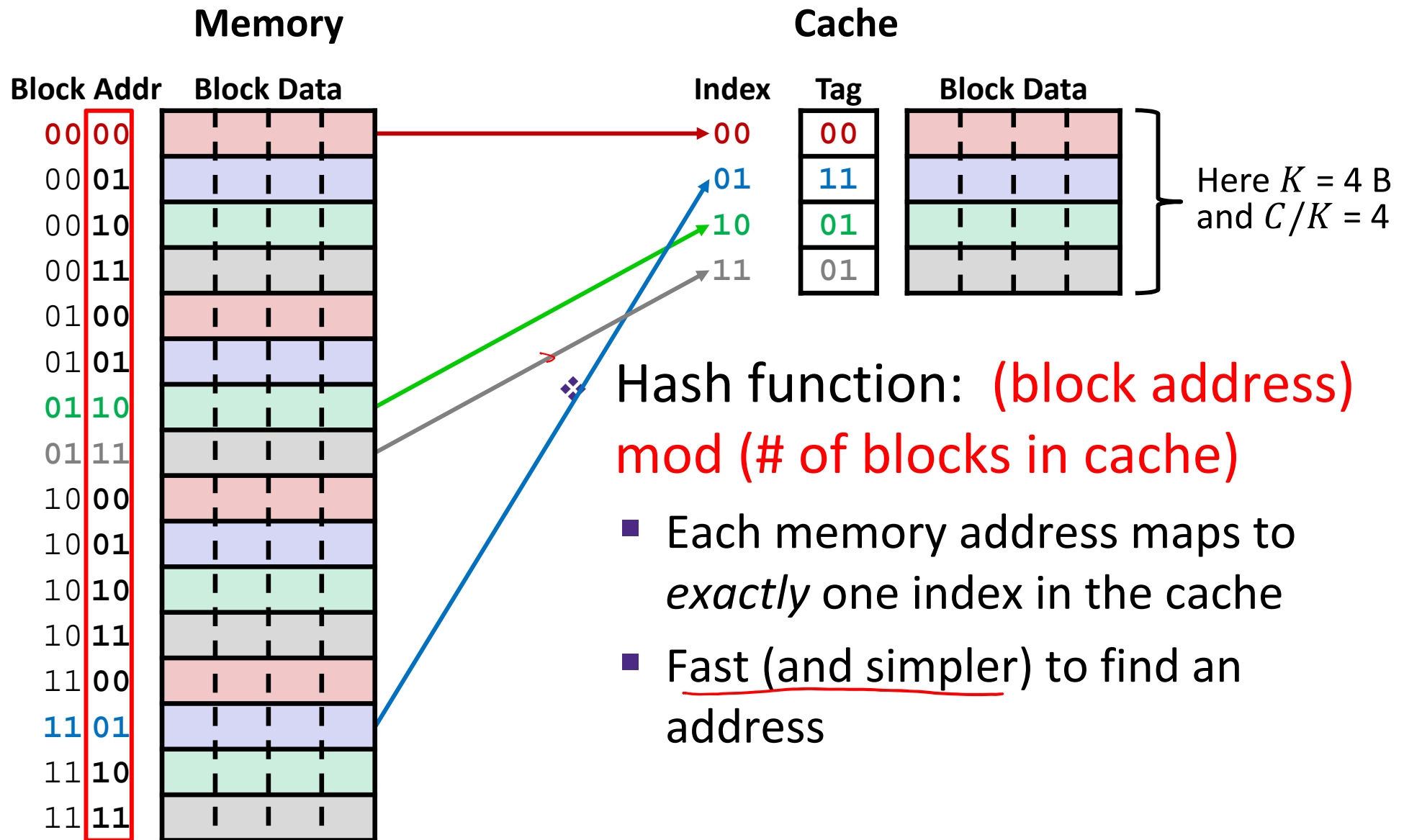
C. 16 bytes

D. 32 bytes

E. We're lost...

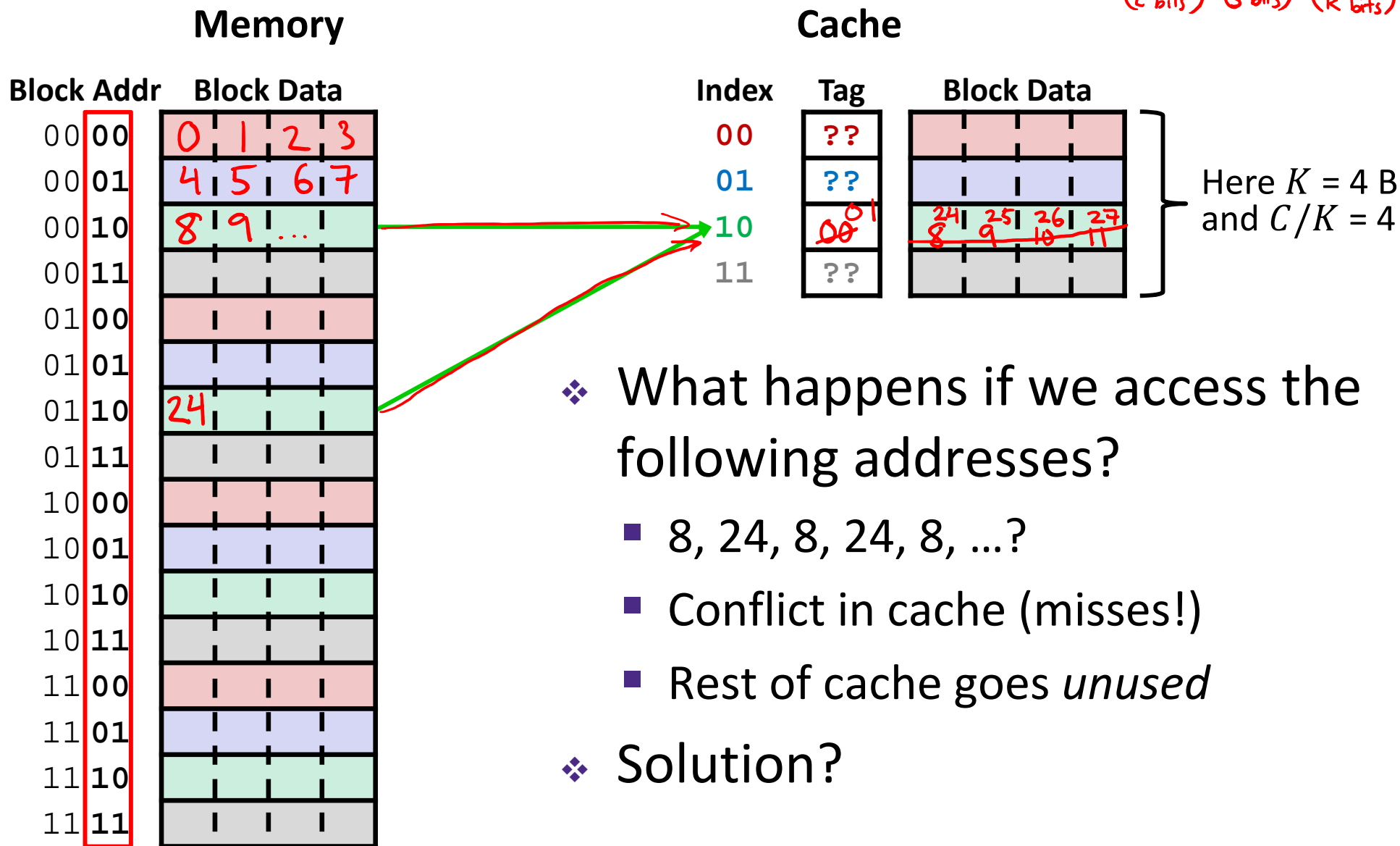


Direct-Mapped Cache



Direct-Mapped Cache Problem

8: 0b 000|10|00
 24: 0b 01|10|00
 Tag (t bits) | Index (s bits) | Offset (k bits)



- ❖ What happens if we access the following addresses?
 - 8, 24, 8, 24, 8, ...?
 - Conflict in cache (misses!)
 - Rest of cache goes *unused*
- ❖ Solution?

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**

