# Caches II

## CSE 351 Summer 2020

**Instructor:**

Porter Jones

**Teaching Assistants:**

Amy Xu

Callum Walker

Sam Wolfson

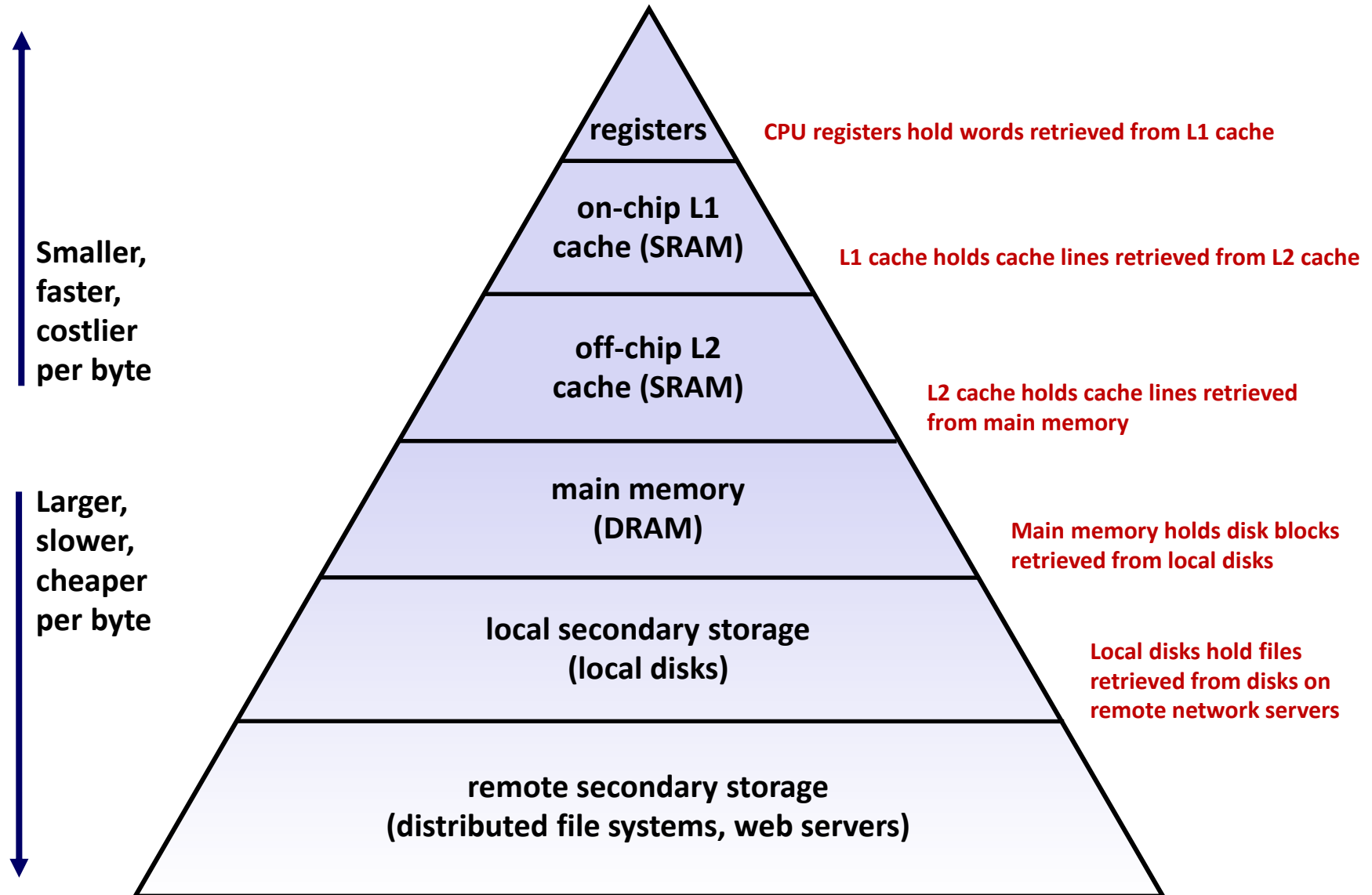Tim Mandzyuk

# **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!
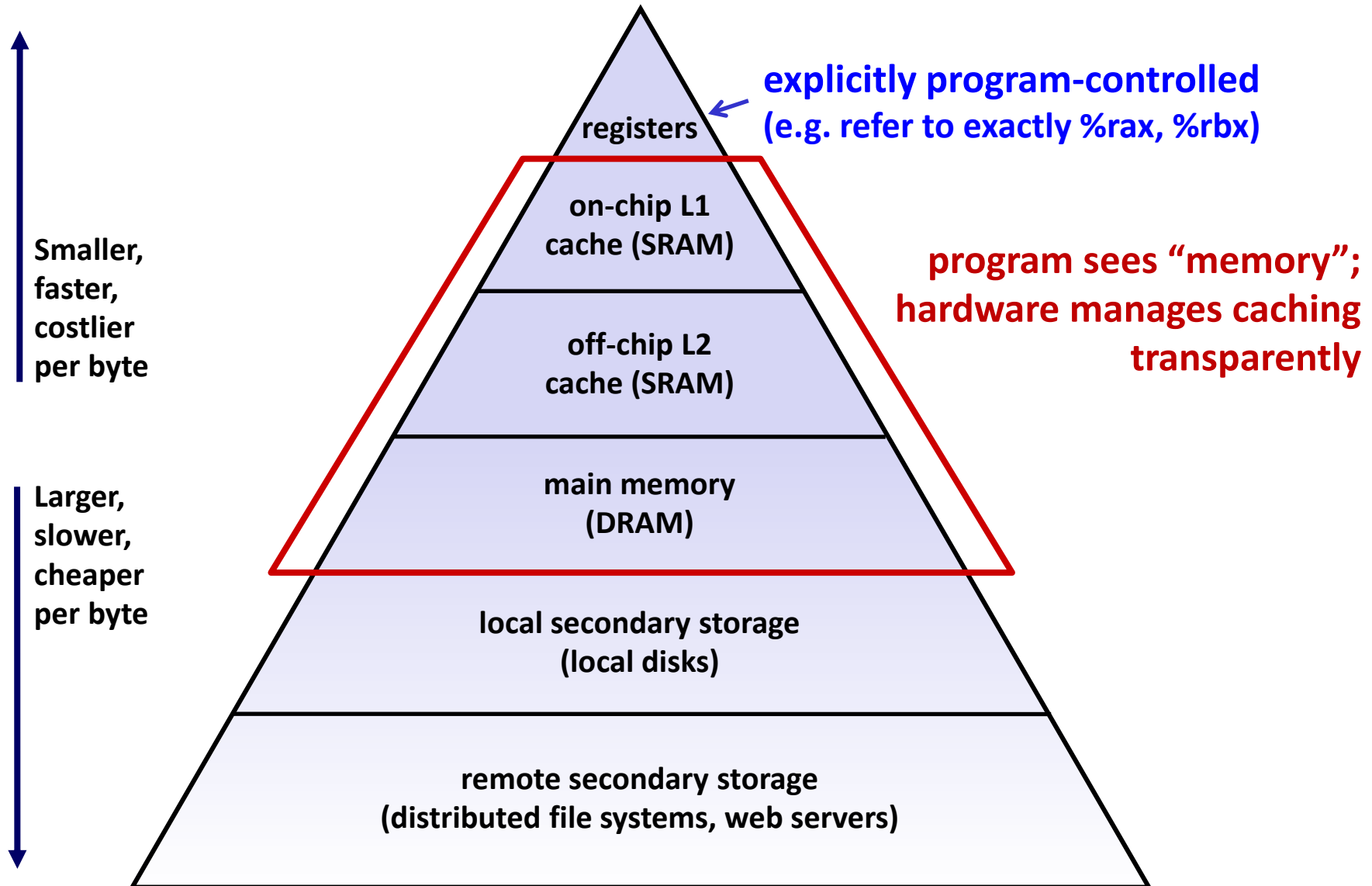
❖ Unit Summary 2 Due next Wednesday (8/5) – 11:59pm

# 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 ↔ cache, cache ↔ DRAM, DRAM ↔ 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 <u>memory hierarchy</u>
  - For each level k, the faster, smaller device at level k serves as a cache for the larger, slower device at level k+1
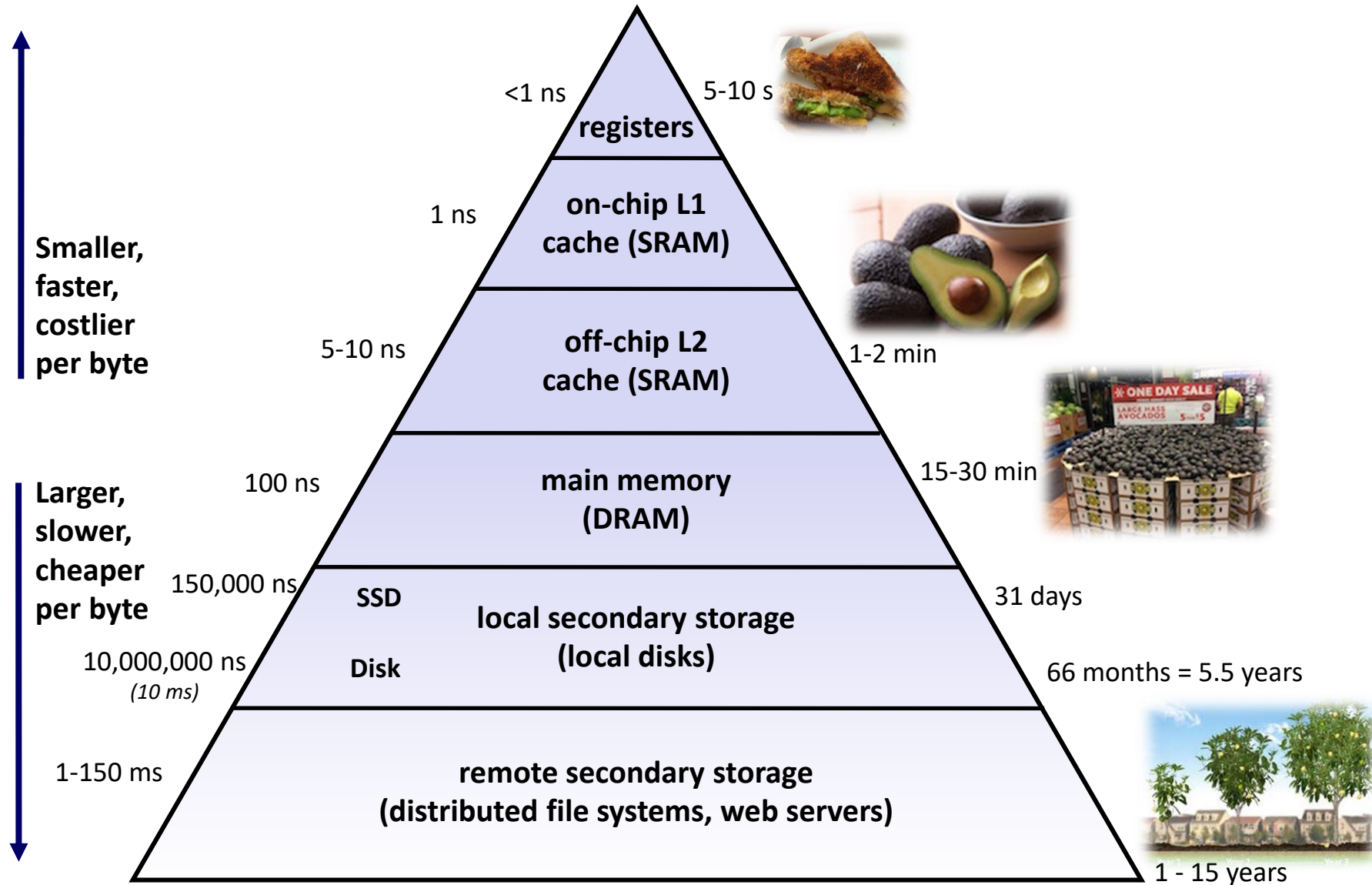
3

# An Example Memory Hierarchy

**Smaller,
faster,
costlier
per byte**

**registers**    <span style="color:darkred">**CPU registers hold words retrieved from L1 cache**</span>

**on-chip L1
cache (SRAM)**    <span style="color:darkred">**L1 cache holds cache lines retrieved from L2 cache**</span>

**off-chip L2
cache (SRAM)**    <span style="color:darkred">**L2 cache holds cache lines retrieved
from main memory**</span>

**Larger,
slower,
cheaper
per byte**

**main memory
(DRAM)**    <span style="color:darkred">**Main memory holds disk blocks
retrieved from local disks**</span>

**local secondary storage
(local disks)**    <span style="color:darkred">**Local disks hold files
retrieved from disks on
remote network servers**</span>

**remote secondary storage
(distributed file systems, web servers)**

# An Example Memory Hierarchy

explicitly program-controlled
(e.g. refer to exactly %rax, %rbx)

**Smaller, faster, costlier per byte**

**Larger, slower, cheaper per byte**

registers

on-chip L1 cache (SRAM)

off-chip L2 cache (SRAM)

main memory (DRAM)

program sees "memory"; hardware manages caching transparently

local secondary storage (local disks)

remote secondary storage (distributed file systems, web servers)

5

# An Example Memory Hierarchy

**Smaller,
faster,
costlier
per byte**

**Larger,
slower,
cheaper
per byte**

<1 ns — registers — 5-10 s

1 ns — on-chip L1 cache (SRAM)

5-10 ns — off-chip L2 cache (SRAM) — 1-2 min

100 ns — main memory (DRAM) — 15-30 min

150,000 ns — SSD — local secondary storage (local disks) — 31 days

10,000,000 ns *(10 ms)* — Disk — 66 months = 5.5 years

1-150 ms — remote secondary storage (distributed file systems, web servers) — 1 - 15 years

# 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
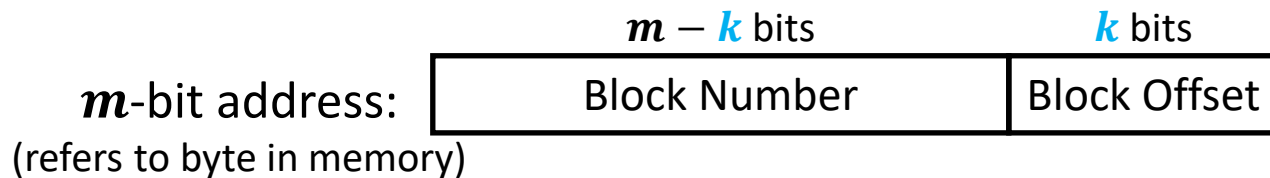
# **Cache Organization (1)**

**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!

# 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!

❖ 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)

$m$-bit address: (refers to byte in memory)

| $m - k$ bits | $k$ bits |
|---|---|
| Block Number | Block Offset |

# 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… | |

# 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 ($C/K$)
- <u>Example</u>:  $C$ = 32 KiB = 512 blocks if using 64-B 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

**Insert:**

5

27
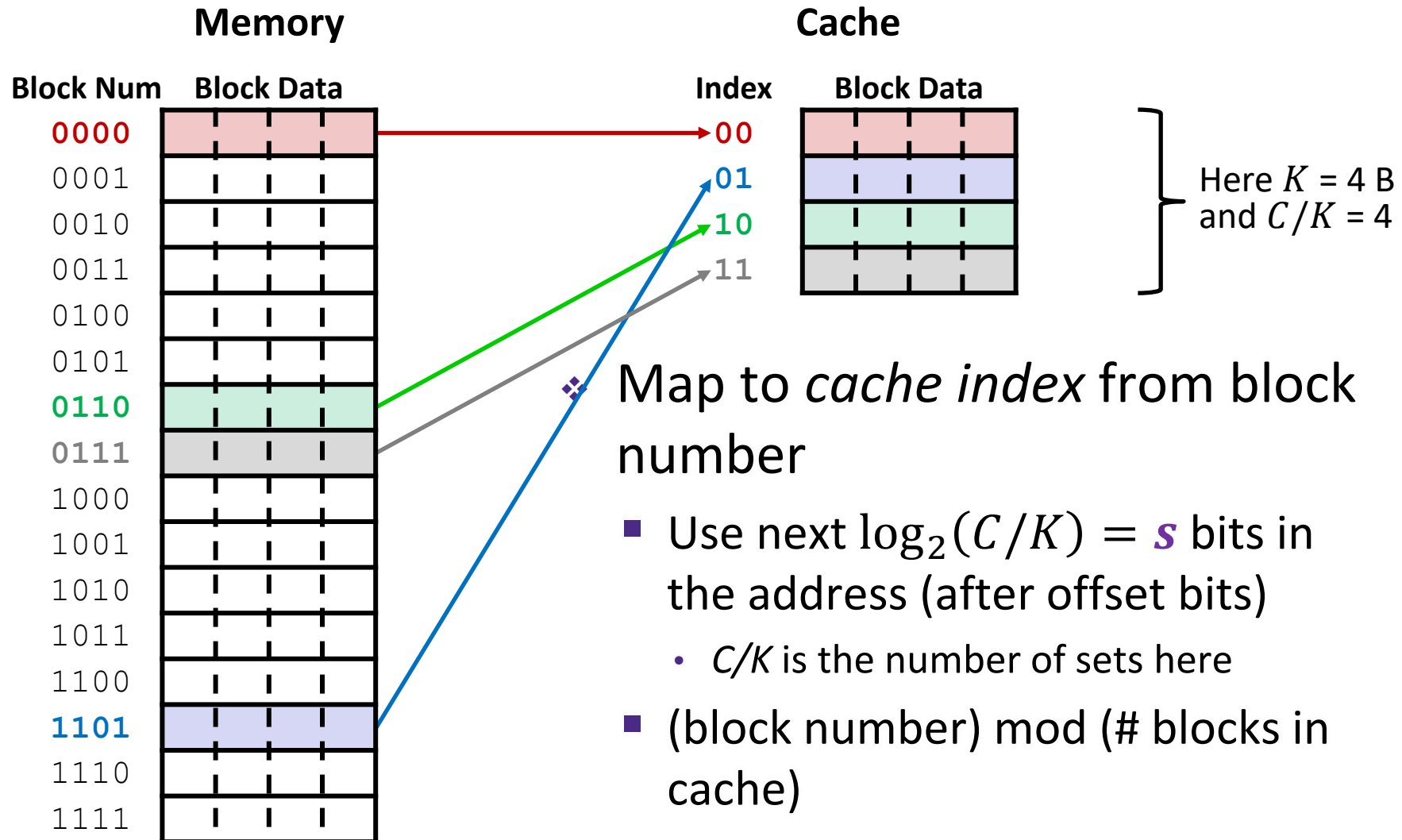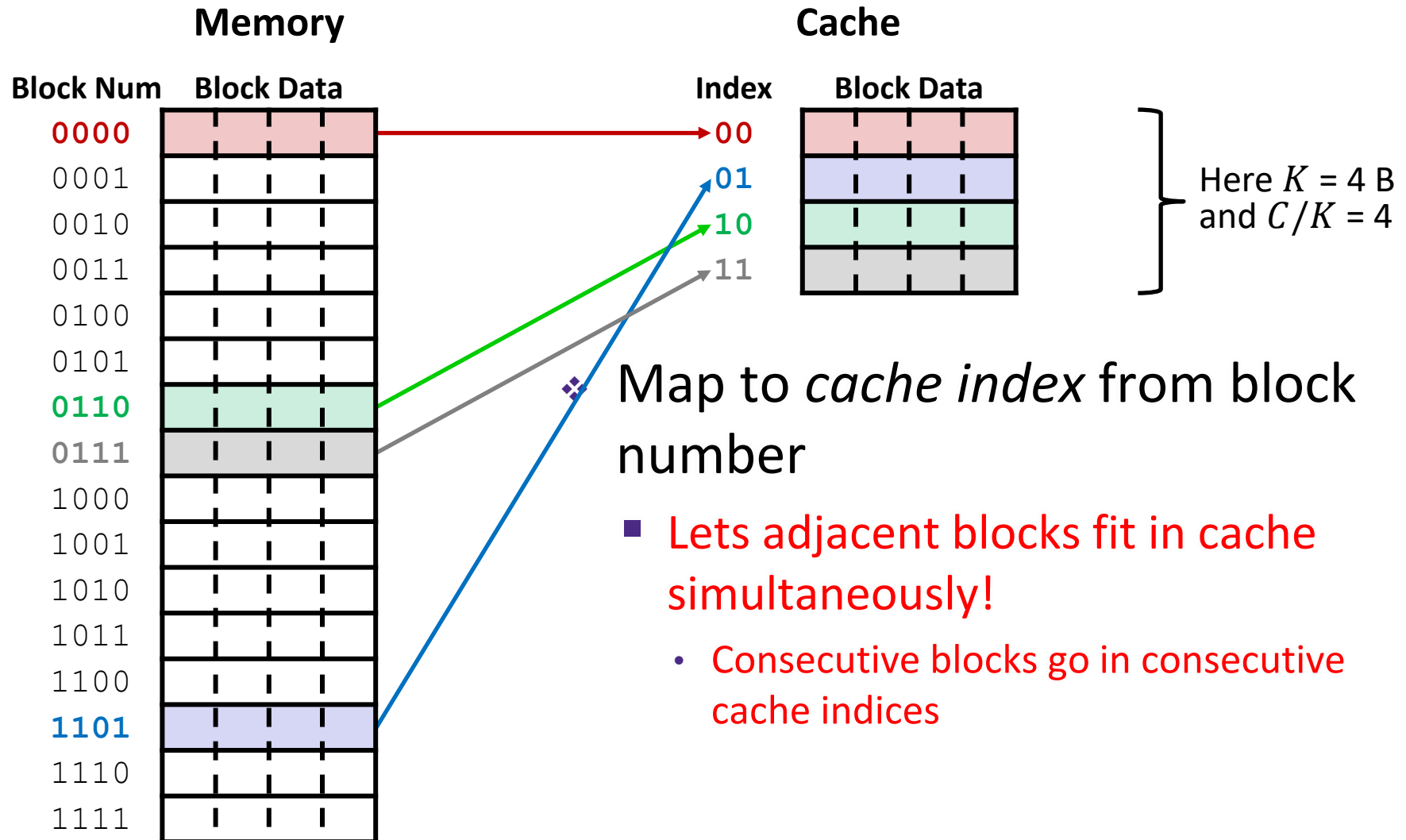
34

102

119

Apply hash function to map data to "buckets"

| | |
|---|---|
| 0 | |
| 1 | |
| 2 | |
| 3 | |
| 4 | |
| 5 | |
| 6 | |
| 7 | |
| 8 | |
| 9 | |

# Place Data in Cache by Hashing Address

**Memory**

**Cache**

| Block Num | Block Data |
|-----------|------------|
| **0000** | |
| 0001 | |
| 0010 | |
| 0011 | |
| 0100 | |
| 0101 | |
| **0110** | |
| 0111 | |
| 1000 | |
| 1001 | |
| 1010 | |
| 1011 | |
| 1100 | |
| **1101** | |
| 1110 | |
| 1111 | |

| Index | Block Data |
|-------|------------|
| **00** | |
| **01** | |
| **10** | |
| **11** | |

Here $K$ = 4 B and $C/K$ = 4

❖ Map to *cache index* from block number

- Use next $\log_2(C/K) = \boldsymbol{s}$ bits in the address (after offset bits)
  - *C/K* is the number of sets here
- (block number) mod (# blocks in cache)

# Place Data in Cache by Hashing Address

**Memory**

**Block Num**   **Block Data**

| Block Num | Block Data |
|---|---|
| 0000 | |
| 0001 | |
| 0010 | |
| 0011 | |
| 0100 | |
| 0101 | |
| 0110 | |
| 0111 | |
| 1000 | |
| 1001 | |
| 1010 | |
| 1011 | |
| 1100 | |
| 1101 | |
| 1110 | |
| 1111 | |

**Cache**

**Index**   **Block Data**

| Index | Block Data |
|---|---|
| 00 | |
| 01 | |
| 10 | |
| 11 | |

Here $K$ = 4 B
and $C/K$ = 4

❖ Map to *cache index* from block number

- Lets adjacent blocks fit in cache simultaneously!
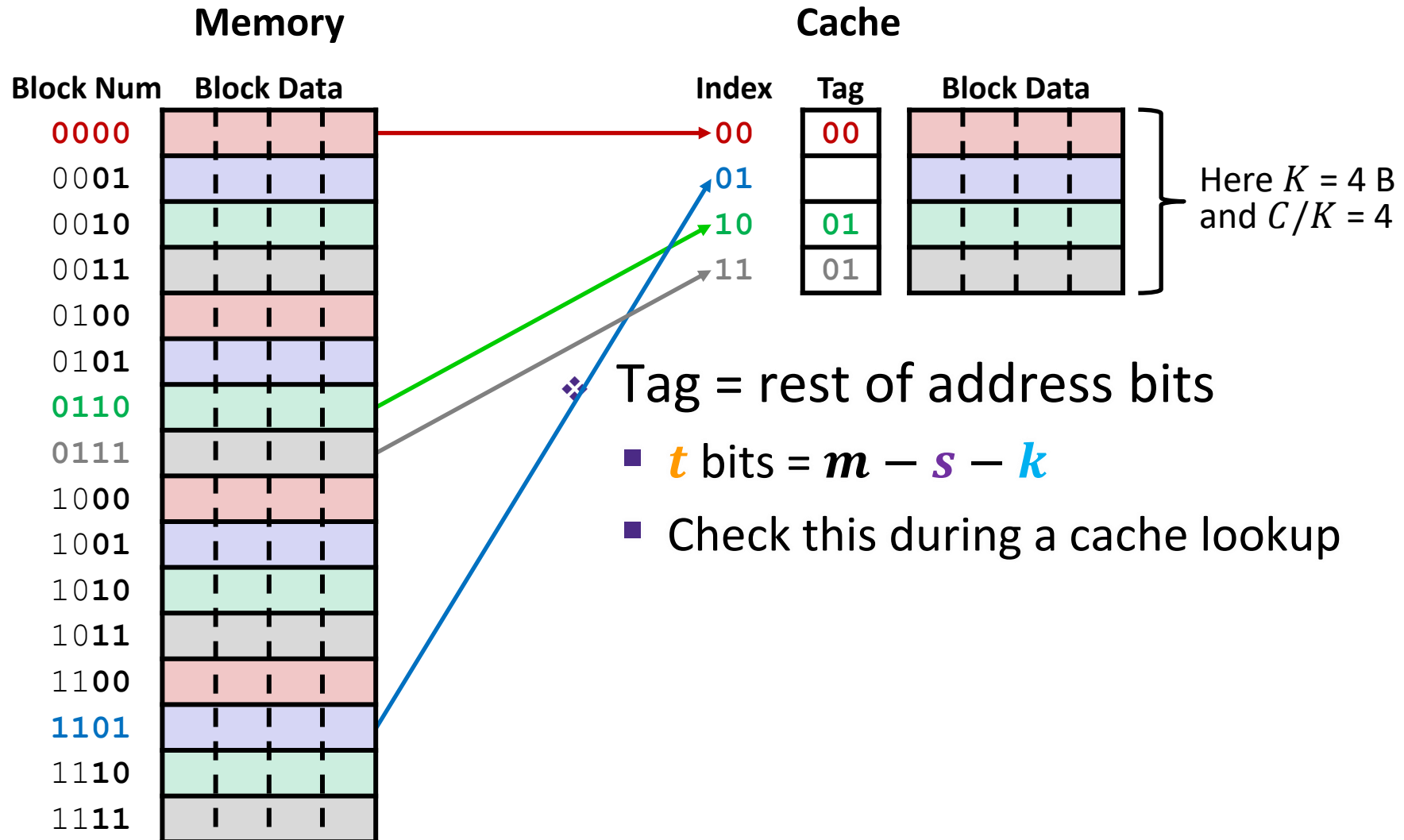  - Consecutive blocks go in consecutive cache indices

# Practice Question

❖ 6-bit addresses, block size $K$ = 4 B, and our cache holds $S$ = 4 blocks.

❖ 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
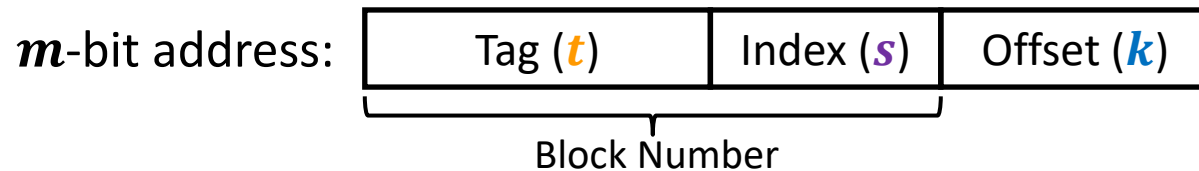
# Place Data in Cache by Hashing Address

**Memory**                                    **Cache**

| Block Num | Block Data |
|-----------|------------|

**Index**   **Block Data**

Here $K$ = 4 B
and $C/K$ = 4

0000
0001
0010
0011
0100
0101
0110
0111
1000
1001
1010
1011
1100
1101
1110
1111

00
01
10
11

## Collision!

- This might confuse the cache later when we access the data
- Solution?

# Tags Differentiate Blocks in Same Index

**Memory**

**Cache**

| Block Num | Block Data |
|---|---|
| 0000 | |
| 0001 | |
| 0010 | |
| 0011 | |
| 0100 | |
| 0101 | |
| 0110 | |
| 0111 | |
| 1000 | |
| 1001 | |
| 1010 | |
| 1011 | |
| 1100 | |
| 1101 | |
| 1110 | |
| 1111 | |

| Index | Tag | Block Data |
|---|---|---|
| 00 | 00 | |
| 01 | | |
| 10 | 01 | |
| 11 | 01 | |

Here $K$ = 4 B and $C/K$ = 4

❖ Tag = rest of address bits

- $t$ bits $= m - s - k$

- Check this during a cache lookup

# **Checking for a Requested Address**

❖ CPU sends address request for chunk of data
  ▪ Address and requested data are not the same thing!
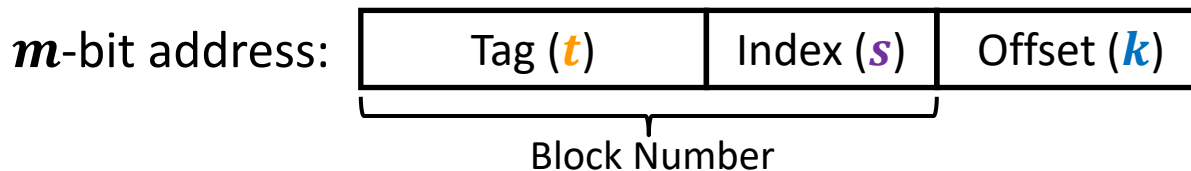    • Analogy: your friend ≠ their phone number


❖ TIO address breakdown:

$m$-bit address:

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

Block Number

  ▪ **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) = $\log_2(K)$ =

  ▪ Index bits (s) = $\log_2(num\ sets)$ =

  ▪ Tag bits (t) = Rest of the bits in the address =

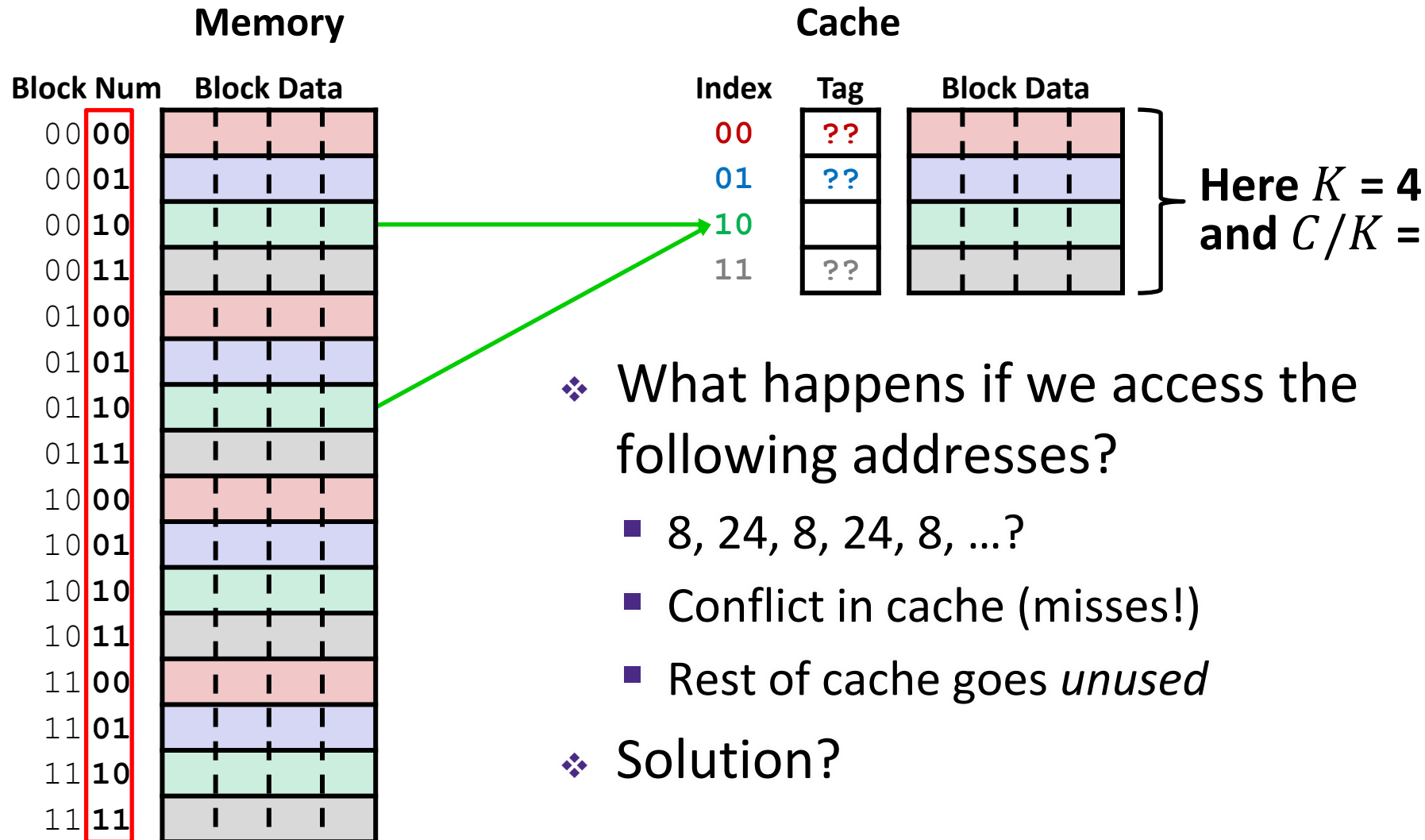$m$-bit address:

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

Block Number

❖ What are the fields for address 0xBA?

  ▪ Tag bits (unique id for block):

  ▪ Index bits (cache set block maps to):

  ▪ Offset bits (byte offset within block):

20

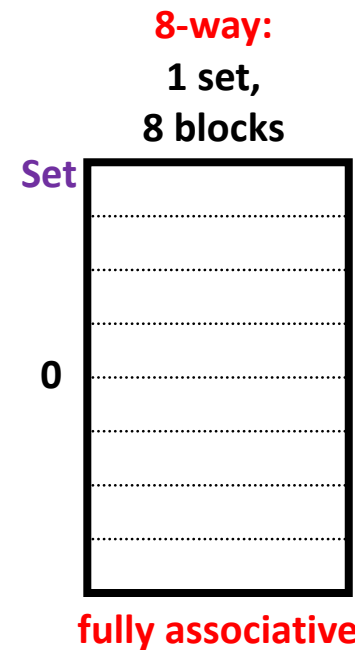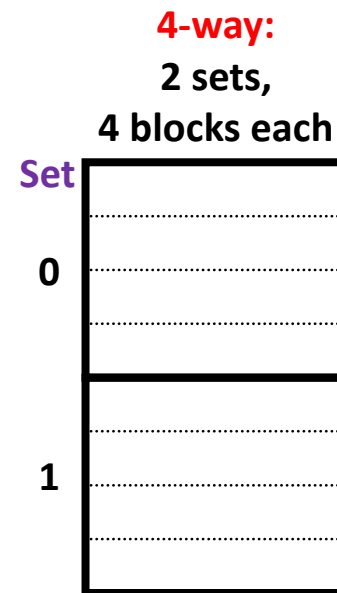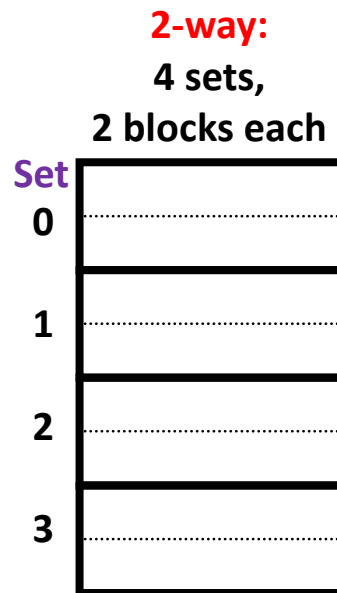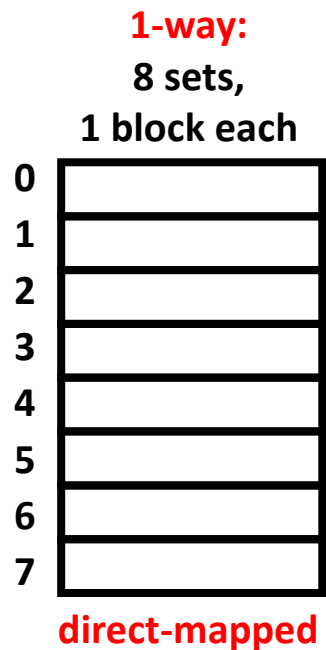# Cache Puzzle [Cache II–b] Vote at http://pollev.com/pbjones

❖ 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)

A. **4 bytes**

B. **8 bytes**

C. **16 bytes**

D. **32 bytes**

E. **We're lost…**

# Direct-Mapped Cache Problem

**Memory**

**Block Num**    **Block Data**

| 00 | 00 |
| 00 | 01 |
| 00 | 10 |
| 00 | 11 |
| 01 | 00 |
| 01 | 01 |
| 01 | 10 |
| 01 | 11 |
| 10 | 00 |
| 10 | 01 |
| 10 | 10 |
| 10 | 11 |
| 11 | 00 |
| 11 | 01 |
| 11 | 10 |
| 11 | 11 |

**Cache**

| Index | Tag | Block Data |
|-------|-----|------------|
| 00 | ?? | |
| 01 | ?? | |
| 10 | | |
| 11 | ?? | |

**Here $K$ = 4**
**and $C/K$ =**

❖ What happens if we access the following addresses?

 ▪ 8, 24, 8, 24, 8, …?

 ▪ Conflict in cache (misses!)

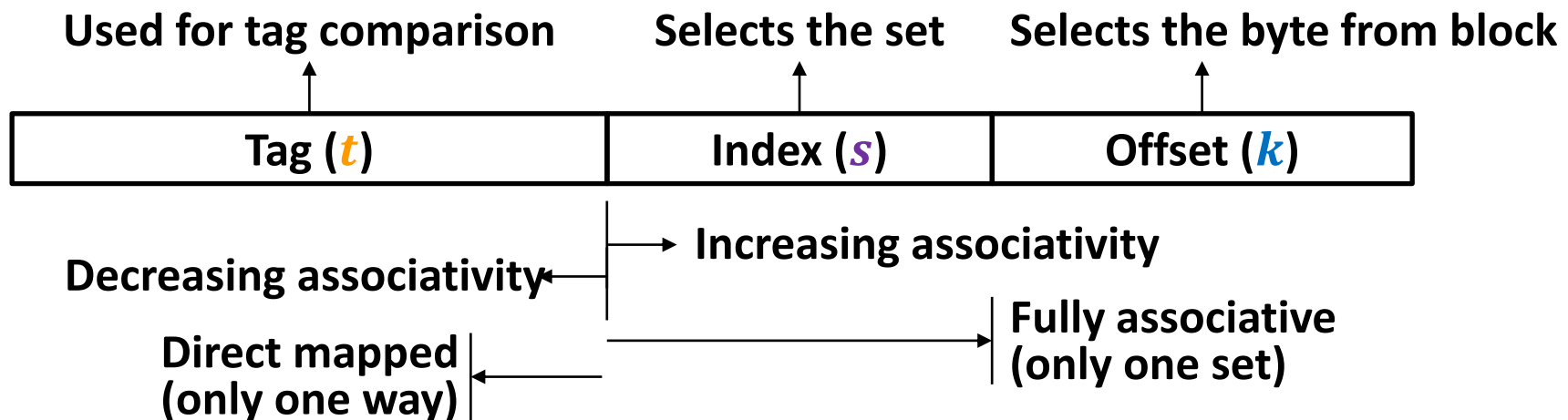 ▪ Rest of cache goes *unused*

❖ 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**



**1-way:**
8 sets,
1 block each

**2-way:**
4 sets,
2 blocks each

**4-way:**
2 sets,
4 blocks each

**8-way:**
1 set,
8 blocks

**direct-mapped**

**fully associative**

23

# 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 $S = C/K/E$
- 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

| Used for tag comparison | Selects the set | Selects the byte from block |
|:---:|:---:|:---:|
| Tag ($t$) | Index ($s$) | Offset ($k$) |

Increasing associativity

Decreasing associativity

Direct mapped (only one way)

Fully associative (only one set)

# Example Placement

| block size: | 16 B |
|---|---|
| capacity: | 8 blocks |
| address: | 16 bits |

❖ Where would data from address `0x1833` 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 ($k$) |
|---|---|---|

$s$ = ?
**Direct-mapped**

| Set | Tag | Data |
|---|---|---|
| 0 | | |
| 1 | | |
| 2 | | |
| 3 | | |
| 4 | | |
| 5 | | |
| 6 | | |
| 7 | | |

$s$ = ?
**2-way set associative**

| Set | Tag | Data |
|---|---|---|
| 0 | | |
| 1 | | |
| 2 | | |
| 3 | | |

$s$ = ?
**4-way set associative**

| Set | Tag | Data |
|---|---|---|
| 0 | | |
| 1 | | |