# Caches III
CSE 351 Spring 2019

**Instructor:**

Ruth Anderson

**Teaching Assistants:**

Gavin Cai
Jack Eggleston
John Feltrup
Britt Henderson
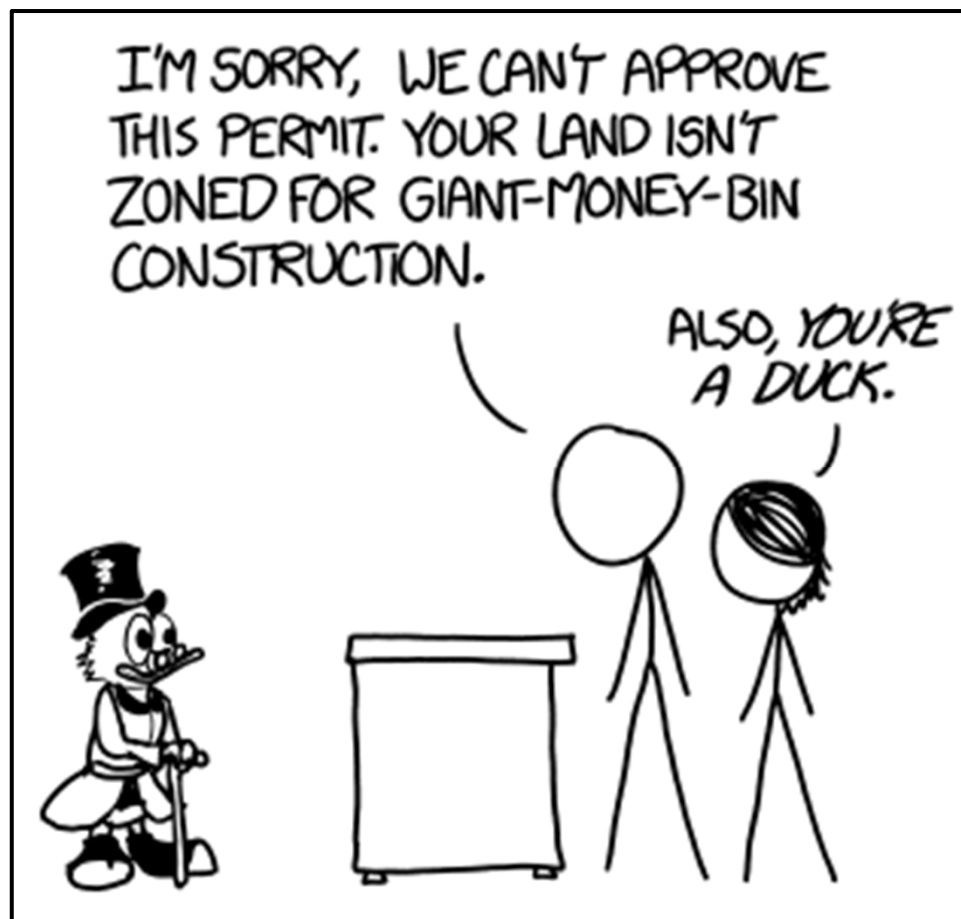Richard Jiang
Jack Skalitzky
Sophie Tian
Connie Wang
Sam Wolfson
Casey Xing
Chin Yeoh



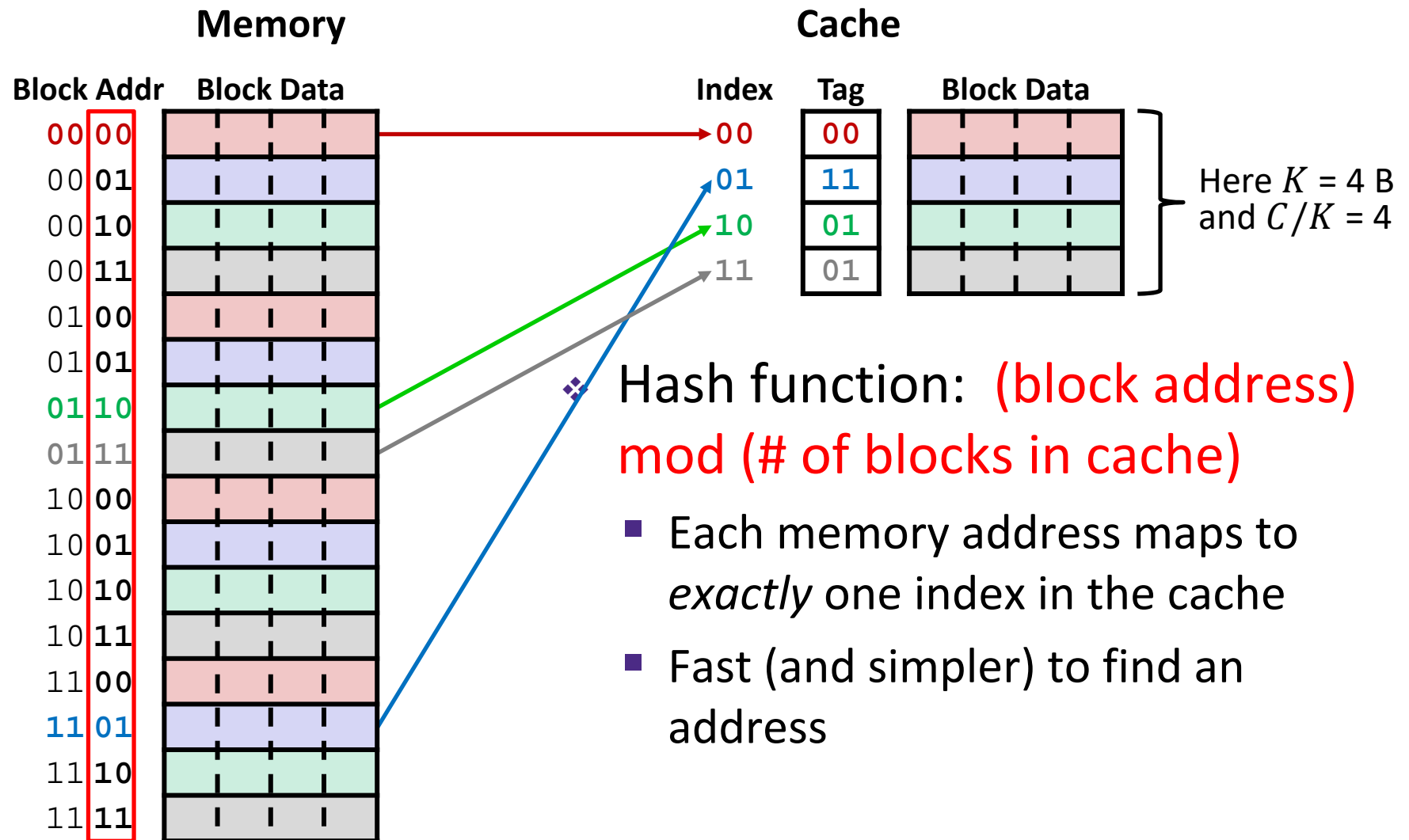https://what-if.xkcd.com/111/

# Administrivia

- ❖ Lab 3, due Wednesday (5/15)

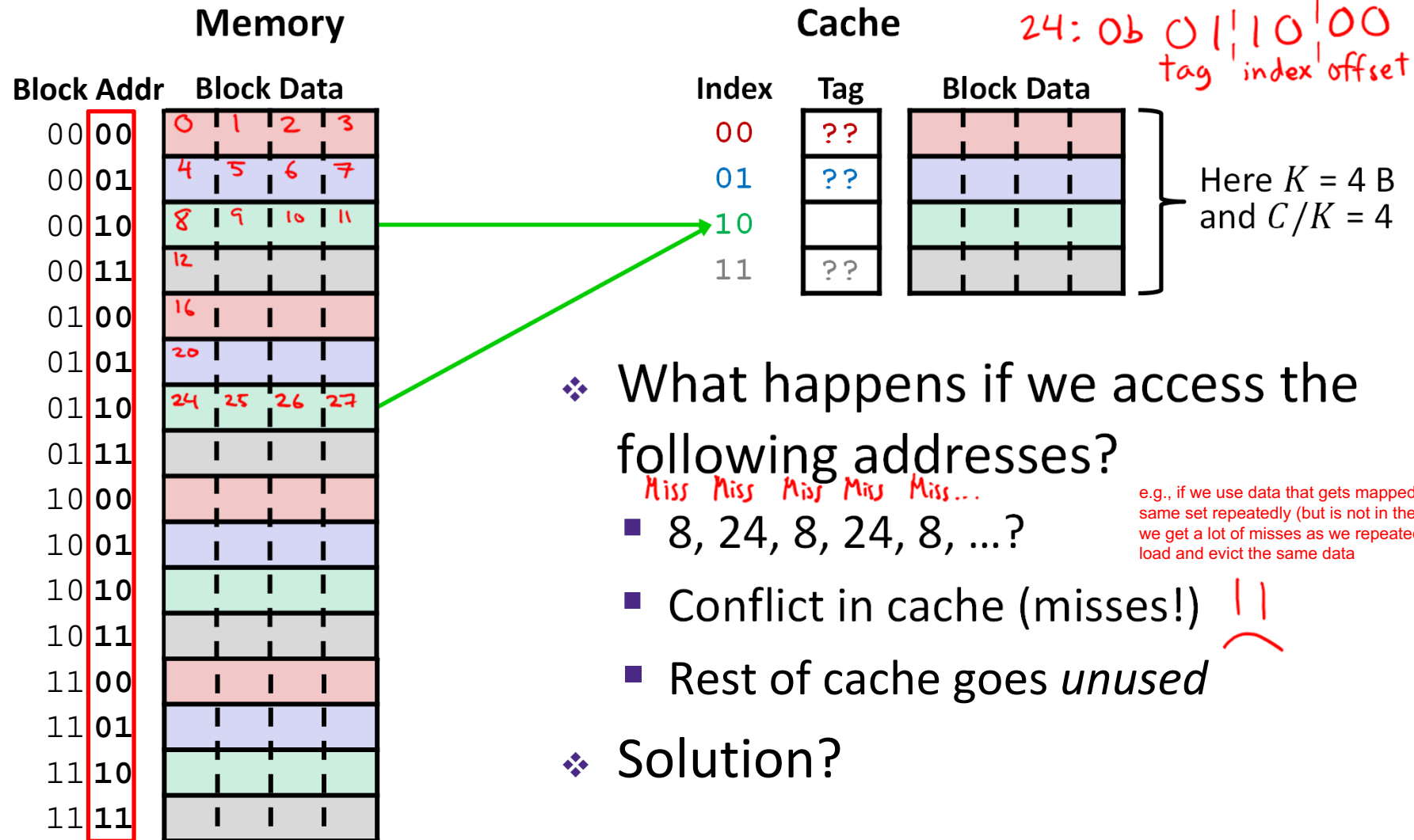- ❖ Homework 4 , due Wed (5/22) (Structs, Caches)

# 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

# Direct-Mapped Cache

**Memory**

**Cache**

**Block Addr**      **Block Data**

**Index**      **Tag**      **Block Data**

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

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

Hash function: (block address) mod (# of blocks in cache)

- Each memory address maps to *exactly* one index in the cache
- Fast (and simpler) to find an address

4

# Direct-Mapped Cache Problem

8: 0b 00 10 00
24: 0b 01 10 00
  tag  index  offset

**Memory**

**Block Addr    Block Data**

| 00 | 00 | 0 1 2 3 |
| 00 | 01 | 4 5 6 7 |
| 00 | 10 | 8 9 10 11 |
| 00 | 11 | 12 |
| 01 | 00 | 16 |
| 01 | 01 | 20 |
| 01 | 10 | 24 25 26 27 |
| 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 B
and $C/K$ = 4

❖ What happens if we access the following addresses?

Miss Miss Miss Miss Miss...

▪ 8, 24, 8, 24, 8, …?

e.g., if we use data that gets mapped to the same set repeatedly (but is not in the same block we get a lot of misses as we repeatedly load and evict the same data

▪ Conflict in cache (misses!)

▪ Rest of cache goes *unused*

❖ Solution?
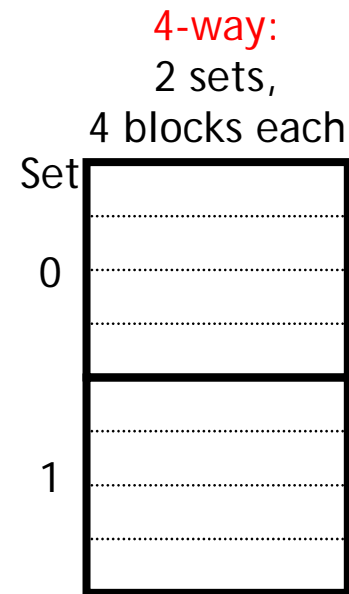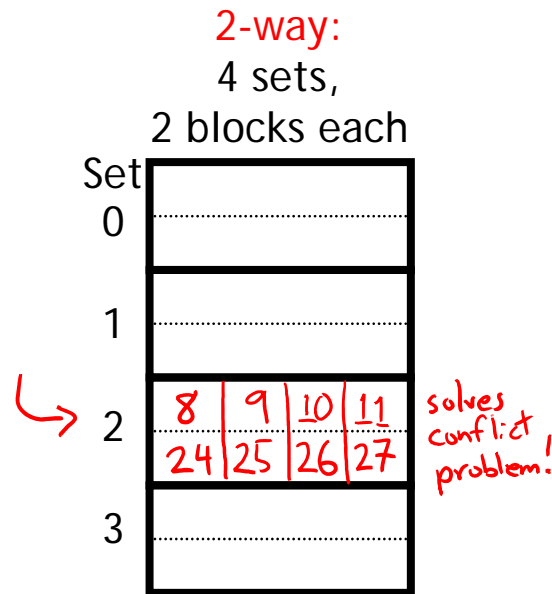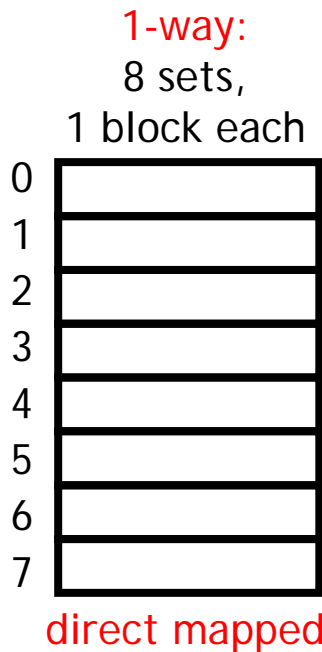
# Associativity

❖ 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**
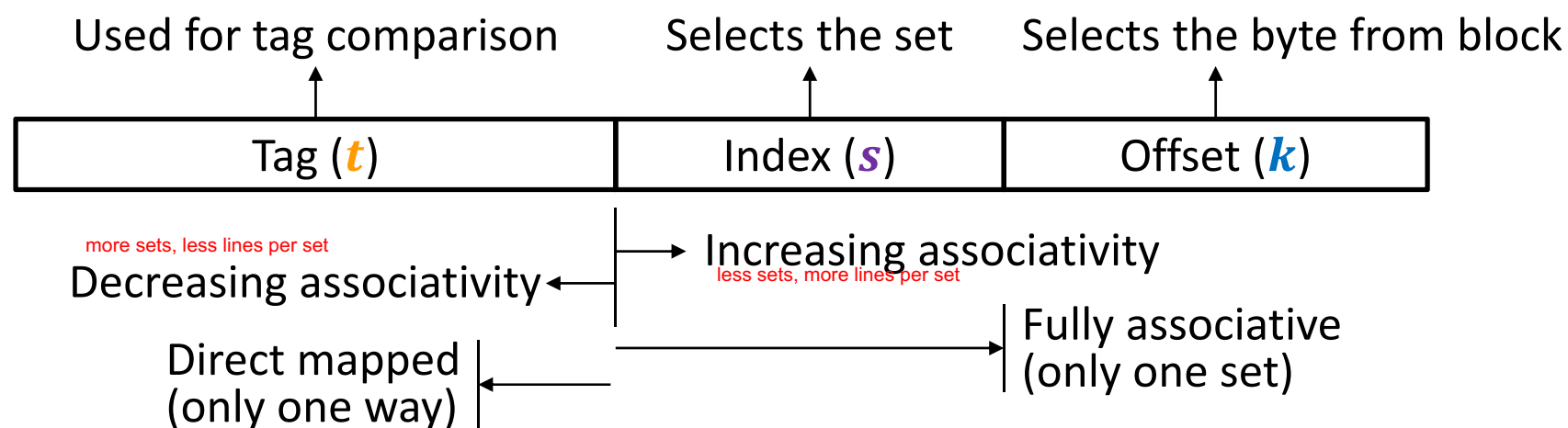  ▪ Each set can store block in more than one **way**

| 1-way:<br>8 sets,<br>1 block each | 2-way:<br>4 sets,<br>2 blocks each | 4-way:<br>2 sets,<br>4 blocks each | 8-way:<br>1 set,<br>8 blocks |
|---|---|---|---|



direct mapped                                                    fully associative

# Cache Organization (3)

❖ 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

    • <u>Direct-mapped</u>: $E$ = 1, so $s$ = $\log_2(C/K)$ as we saw previously

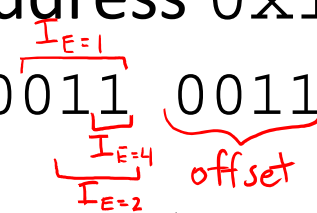    • <u>Fully associative</u>: $E = C/K$, so $s$ = 0 bits

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

more sets, less lines per set

Decreasing associativity ← → Increasing associativity
less sets, more lines per set

Direct mapped (only one way) ← → Fully associative (only one set)

# Example Placement

| block size: | 16 B | $K$ |
| capacity: | 8 blocks | $C/K$ |
| address: | 16 bits | $m$ |

❖ Where would data from address `0x1833` be placed?

- Binary: `0b 0001 1000 0011 0011`

  $I_{E=1}$    $I_{E=4}$   $I_{E=2}$   offset

$$t = m-s-k \qquad s = \log_2(C/K/E) \qquad k = \log_2(K) = 4 \text{ bits}$$

$m$-bit address:

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

$s = ? \quad \log_2(8/1) = 3 \text{ bits}$
Direct-mapped (E=1)

| Set | Tag | Data |
|-----|-----|------|
| 000  0 | | |
| 001  1 | | |
| 010  2 | | |
| 011  3 | — | ⬭ |
| 100  4 | | |
| 101  5 | | |
| 110  6 | | |
| 111  7 | | |

$s = ? \quad \log_2(8/2) = 2 \text{ bits}$
2-way set associative (E=2)

| Set | Tag | Data |
|-----|-----|------|
| 00  0 | | |
| 01  1 | | |
| 10  2 | | ⬭ |
| 11  3 | ~ | ⬭ |

$s = ? \quad \log_2(8/4) = 1 \text{ bit}$
4-way set associative (E=4)

| Set | Tag | Data |
|-----|-----|------|
| 0 | | |
|  | — | ⬭ |
|  | — | ⬭ |
| 1 | — | ⬭ |
|  | — | ⬭ |

# Block Replacement

❖ *Any* empty block in the correct set may be used to store block

❖ If there are no empty blocks, which one should we replace?

    ■ No choice for direct-mapped caches

    ■ Caches typically use something close to **least recently used (LRU)** <span style="color:red">Replace the line that has been used the fewest number of times in some given time frame,.</span>
(hardware usually implements "*not most recently used*")

| Direct-mapped | 2-way set associative | 4-way set associative |
|---|---|---|

**Direct-mapped**

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

**2-way set associative**

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

**4-way set associative**

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

# Peer Instruction Question

$K = 2^7$ B

$\rightarrow C = 2^{11}$ B

❖ We have a cache of size 2 KiB with block size of 128 B. If our cache has 2 sets, what is its associativity?
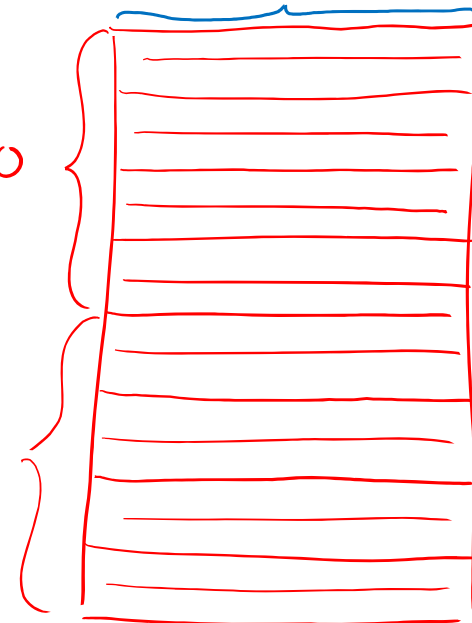
  ■ Vote at http://pollev.com/rea

cache holds $C/K = 2^{11-7} = 2^4 = 16$ blocks

1 block

**A.  2**

**B.  4**

set 0

$S = C/K/E$

each set has 8 blocks, so $E = 8$

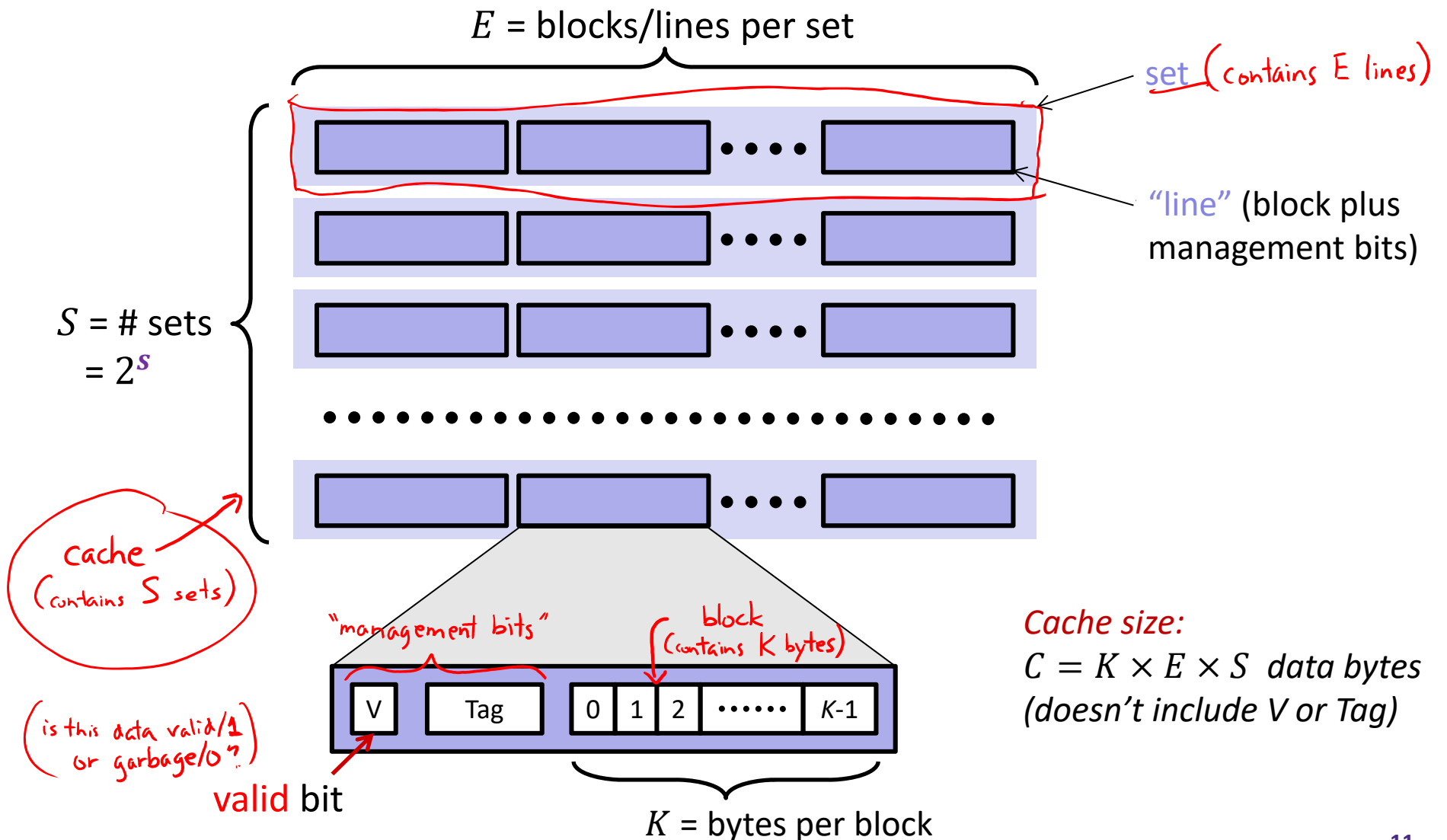**C.  8**

$E = (C/K)/S$
$= 16/2 = 8$

**D.  16**

cache size

set 1

**E.  We're lost…**

$m = 16$

❖ If addresses are 16 bits wide, how wide is the Tag field?   $k = \log_2(K) = 7$ bits,  $s = \log_2(S) = 1$ bit,  $t = m - s - k = \boxed{8 \text{ bits}}$

# General Cache Organization ($S$, $E$, $K$)

*associativity*

*sets*                *block size*

$E$ = blocks/lines per set

set (contains E lines)

"line" (block plus management bits)

$S$ = # sets
= $2^s$

cache (contains S sets)

(is this data valid/**1** or garbage/0 ?)

"management bits"

block (contains K bytes)

| V | Tag | 0 | 1 | 2 | ⋯⋯ | $K$-1 |
|---|-----|---|---|---|-----|-------|

valid bit

$K$ = bytes per block

*Cache size:*
$C = K \times E \times S$ *data bytes*
*(doesn't include V or Tag)*

# Notation Review

- ❖ We just introduced a lot of new variable names!
  - Please be mindful of block size notation when you look at past exam questions or are watching videos

blocks/lines per set

| Variable | This Quarter | Formulas |
|---|---|---|
| Block size | $K$ ($B$ in book) | |
| Cache size | $C$ | $M = 2^m \leftrightarrow m = \log_2 M$ |
| Associativity | $E$ | $S = 2^s \leftrightarrow s = \log_2 S$ |
| Number of Sets | $S$ | $K = 2^k \leftrightarrow k = \log_2 K$ |
| Address space | $M$ | |
| Address width | $m$ | $C = K \times E \times S$ |
| Tag field width | $t$ | $s = \log_2(C/K/E)$ |
| Index field width | $s$ | $m = t + s + k$ |
| Offset field width | $k$ ($b$ in book) | |

# Example Cache Parameters Problem
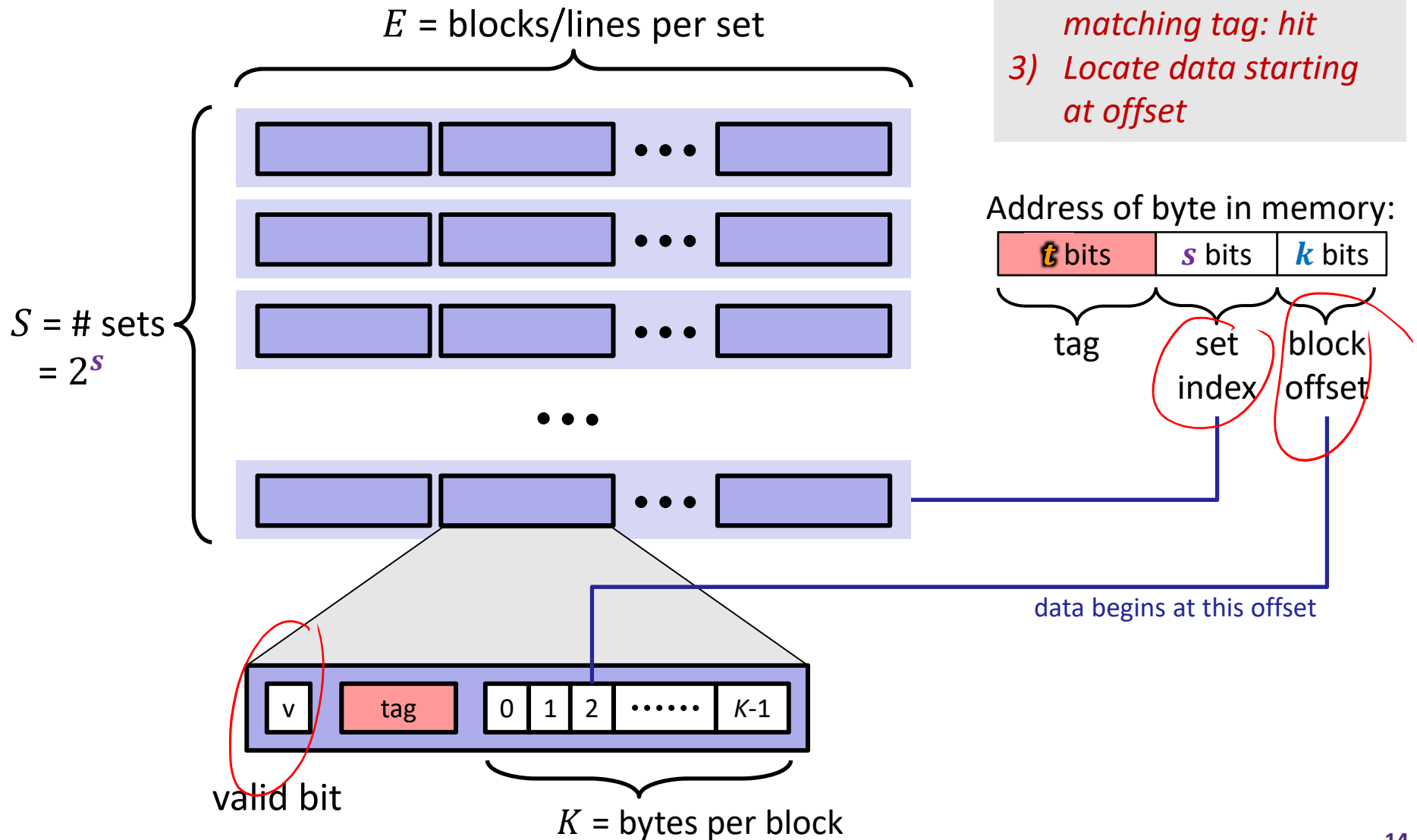
$\rightarrow 2^{12} B \iff m = 12 \text{ bits}$

❖ 4 KiB address space, 125 cycles to go to memory. [MP]
  Fill in the following table:

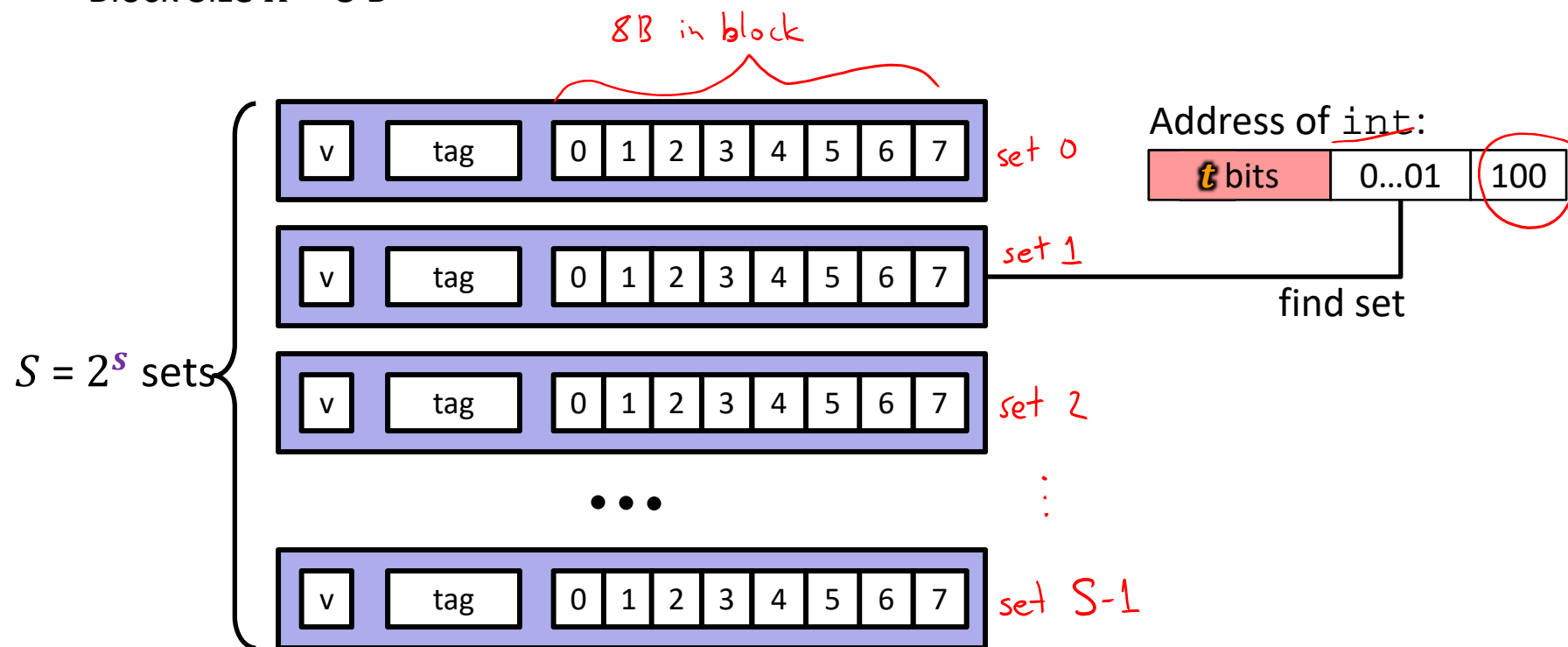| | | | |
|---|---|---|---|
| C | **Cache Size** | 256 B | $2^8$ |
| K | **Block Size** | 32 B | $2^5$ |
| E | **Associativity** | 2-way | $2^1$ |
| HT | **Hit Time** | 3 cycles | |
| MR | **Miss Rate** | 20% | |
| $t = m - s - k$ | **Tag Bits** | 5 | |
| $s = \log_2(C/K/E)$ | **Index Bits** | 2 | $2^8/2^5/2^1$ |
| $k = \log_2(K)$ | **Offset Bits** | 5 | |
| AMAT = HT + MR * MP | **AMAT** | $3 + 0.2(125) = 28$ clock cycles | |

# Cache Read

1) *Locate set*
2) *Check if any line in set is valid and has matching tag: hit*
3) *Locate data starting at offset*

$E$ = blocks/lines per set

$S$ = # sets = $2^s$

Address of byte in memory:

| $t$ bits | $s$ bits | $k$ bits |
|---|---|---|
| tag | set index | block offset |

data begins at this offset

| v | tag | 0 | 1 | 2 | ...... | K-1 |

valid bit

$K$ = bytes per block

# Example:  Direct-Mapped Cache ($E$ = 1)

Direct-mapped:  One line per set
Block Size $K$ = 8 B

8 B in block



set 0

Address of int:

| $t$ bits | 0...01 | 100 |

find set

$S = 2^s$ sets

set 1

set 2

set S-1

# Example:  Direct-Mapped Cache ($E$ = 1)

Direct-mapped:  One line per set
Block Size $K$ = 8 B

valid?  +  match?: yes = hit

Address of `int`:

| v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

| $t$ bits | 0…01 | 100 |

block offset

# Example: Direct-Mapped Cache ($E$ = 1)

Direct-mapped: One line per set
Block Size $K$ = 8 B

valid? + match?: yes = hit

multiple
of 4

long

Address of `int`!

| v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

??

| $t$ bits | 0...01 | 100 |

multiple
of 8

block offset

`int` (4 B) is here

**This is why we want alignment!**

No match? Then old line gets evicted and replaced

no unnecessary extra
cache accesses across
block boundaries

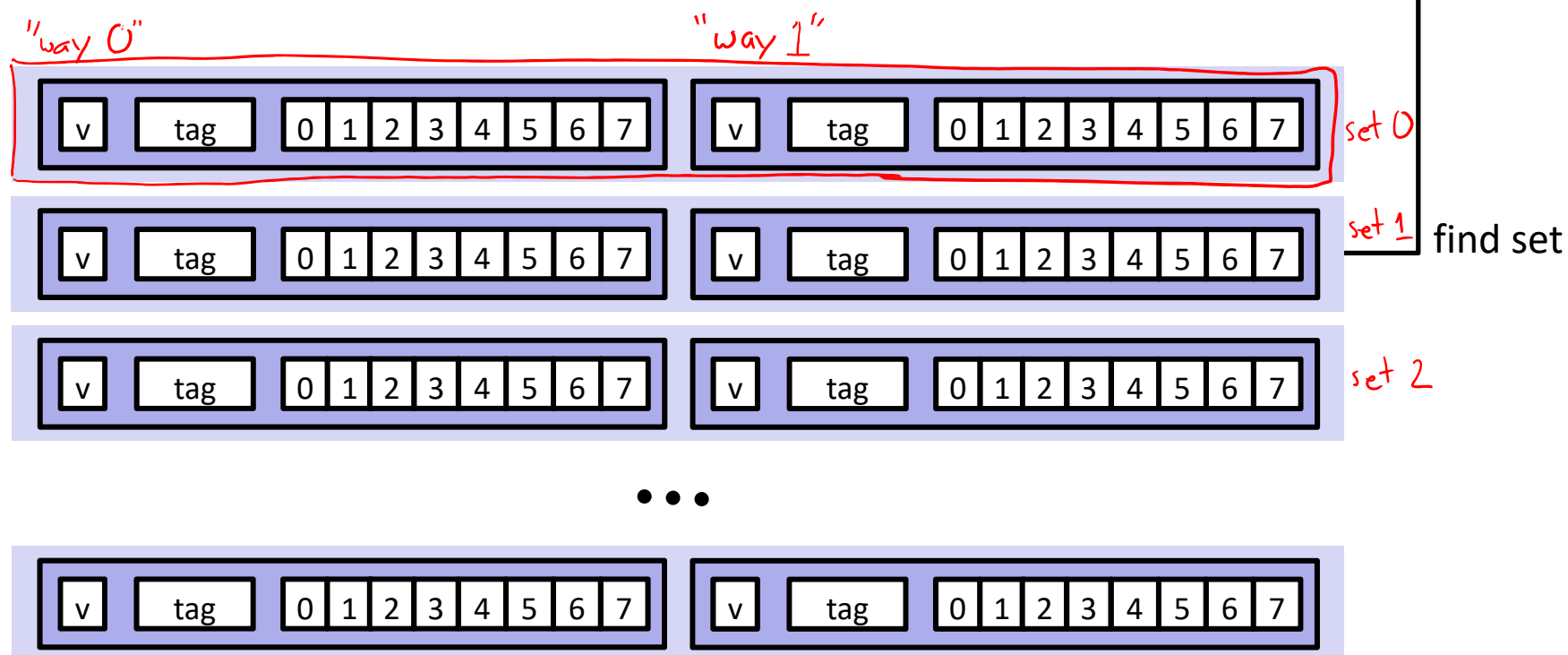e.g. data is always stored in one block?

# Example: Set-Associative Cache ($E$ = 2)

2-way: Two lines per set

Block Size $K$ = 8 B
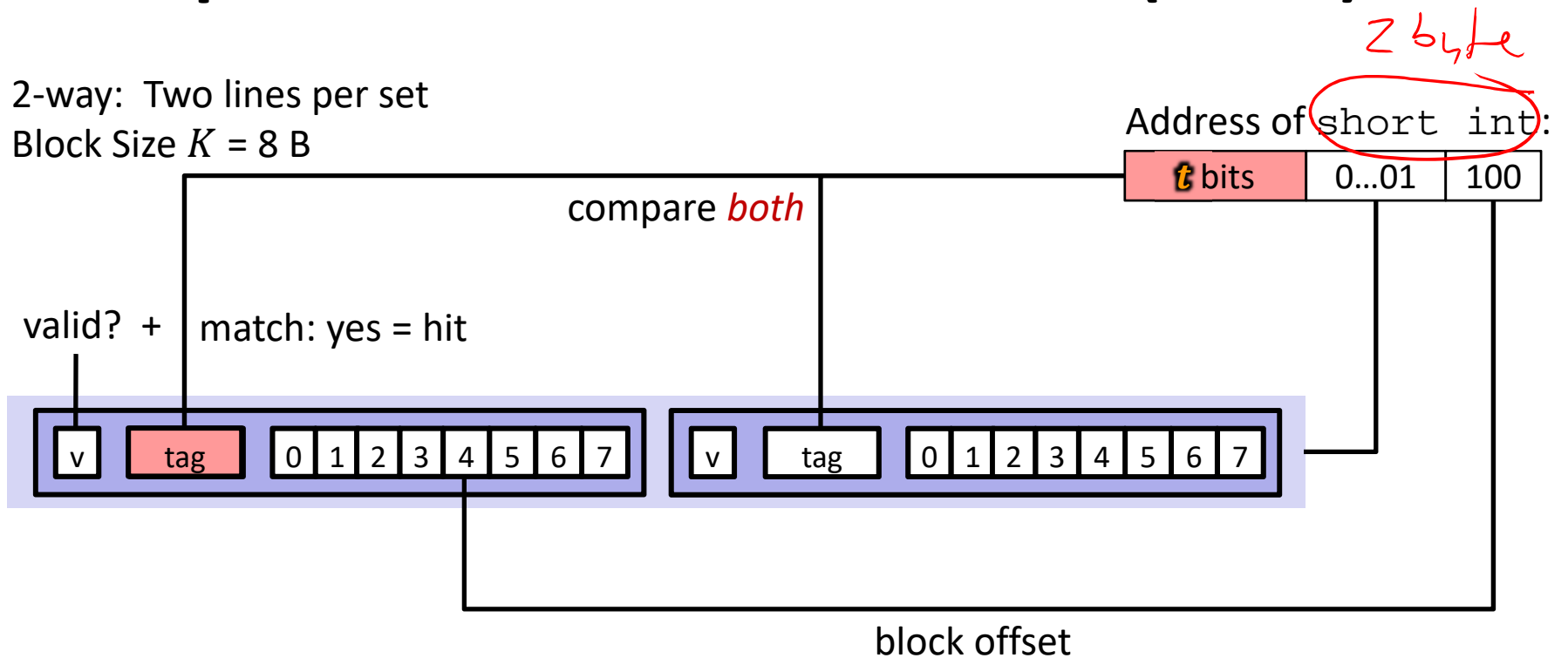
Address of `short int`:

| $t$ bits | 0...01 | 100 |
|---|---|---|

"way 0"  "way 1"

| v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | | v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | set 0 |

| v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | | v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | set 1 find set |

| v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | | v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | set 2 |

• • •

| v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | | v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

# Example: Set-Associative Cache ($E = 2$)

2-way: Two lines per set

Block Size $K$ = 8 B

Z byte

Address of short int:

compare *both*

| $t$ bits | 0...01 | 100 |

valid? + match: yes = hit

| v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | | v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

block offset

# Example: Set-Associative Cache ($E$ = 2)

2-way: Two lines per set
Block Size $K$ = 8 B

Address of `short int`:

| $t$ bits | 0...01 | 100 |

compare *both*

valid? +  match: yes = hit

| v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

| v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

`short int` (2 B) is here

block offset

**No match?**
- One line in set is selected for eviction and replacement
- Replacement policies: random, least recently used (LRU), ...

# Types of Cache Misses: 3 C's!

❖ Compulsory (cold) miss
- Occurs on first access to a block

❖ Conflict miss
- Conflict misses occur when the cache is large enough, but multiple data objects all map to the same slot
  - *e.g.* referencing blocks 0, 8, 0, 8, ... could miss every time
- Direct-mapped caches have more conflict misses than $E$-way set-associative (where $E > 1$)

❖ Capacity miss
- Occurs when the set of active cache blocks (the *working set*) is larger than the cache (just won't fit, even if cache was *fully-associative*)
- **Note:** *Fully-associative* only has Compulsory and Capacity misses

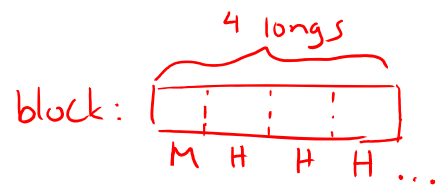# Example Code Analysis Problem

❖ Assuming the cache starts <u>cold</u> (all blocks invalid) and `sum` is stored in a register, calculate the **miss rate**: 1/4

- ▪ $m$ = 12 bits, $C$ = 256 B, $K$ = 32 B, $E$ = 2

  $t = 5$ bits, $s = 2$ bits, $k = 5$ bits

```
#define SIZE 8
long ar[SIZE][SIZE], sum = 0;   // &ar=0x800
for (int i = 0; i < SIZE; i++)
    for (int j = 0; j < SIZE; j++)
        sum += ar[i][j];
```

8 bytes per element

block: 4 longs, M H H H …

|  | tag | index | offset |  |
|---|---|---|---|---|
| ar[0][0] address → | 0b 1000 | 0 0 0 0 | 0000 → | M (compulsory) |
| ar[0][1] → | 0b 1000 | 0 0 0 0 | 1000 → | H |
| ar[0][2] → | 0b 1000 | 0 0 0 1 | 0000 → | H |
| ar[0][3] → | 0b 1000 | 0 0 0 1 | 1000 → | H |
| ar[0][4] → | 0b 1000 | 0 0 1 0 | 0000 → | M (compulsory) |

Challenge: what is the miss rate if we switch the ordering of the for-loops?

22

# What about writes?

❖ Multiple copies of data exist:
  - L1, L2, possibly L3, main memory

❖ What to do on a **write-hit**? *(block/data already in $)*
  - Write-through: write immediately to next level
  - Write-back: defer write to next level until line is evicted (replaced)
    - Must track which cache lines have been modified ("*dirty bit*") ← *extra management bit only for write-back cache*

❖ What to do on a **write-miss**? *(block/data not currently in $)*
  - Write-allocate: ("fetch on write") load into cache, update line in cache
    - Good if more writes or reads to the location follow
  - No-write-allocate: ("write around") just write immediately to ~~memory~~ *next level*

❖ Typical caches:
  - Write-back + Write-allocate, usually ⭐
  - Write-through + No-write-allocate, occasionally