# Caches III
## CSE 351 Summer 2020

**Instructor:**

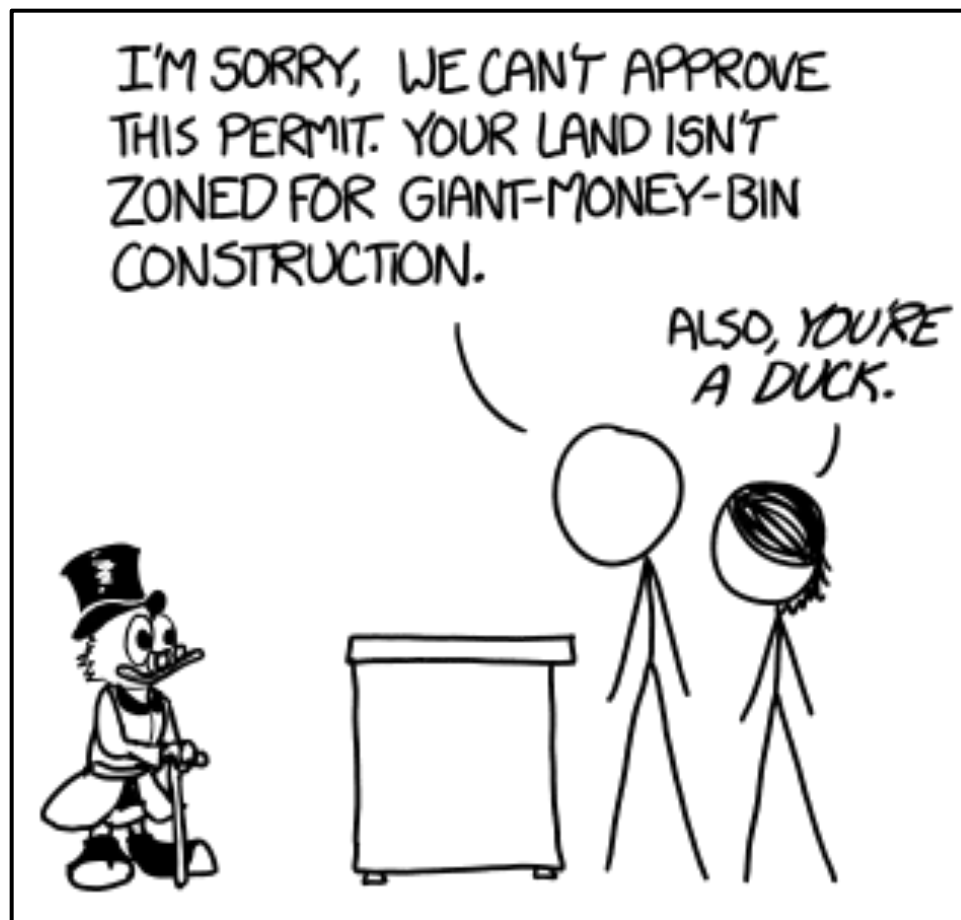Porter Jones

**Teaching Assistants:**

Amy Xu

Callum Walker

Sam Wolfson

Tim Mandzyuk



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

# Administrivia

❖ Questions doc: https://tinyurl.com/CSE351-7-31

❖ hw16 due Wednesday (8/5) – 10:30am

*long, start early!*

❖ Lab 3 due Tonight (7/31) – 11:59pm
  ▪ You get to write some buffer overflow exploits!

❖ Lab 4 released later today
  ▪ All about caches!

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

*Course Staff 1-on-1s – see Piazza*
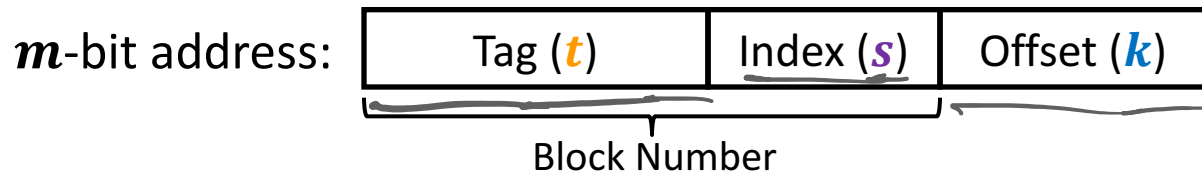
# Making memory accesses fast!

❖ 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

# Review: Cache Parameters

❖ **Block size (K):** basic unit of transfer between memory and the cache, given in bytes (e.g. 64 B).

❖ **Cache size (C):** Total amount of data that can be stored in the cache, given in bytes (e.g. 32 KiB).

- Must be multiple of block size
- Number of blocks in cache is calculated by C/K

❖ **Associativity (E):** Number of ways blocks can be stored in a cache set, or how many blocks in each set

❖ **Number of sets (S):** Number of unique sets that blocks can be placed into in a cache (calculated as C/K/E).

# Review: TIO address breakdown
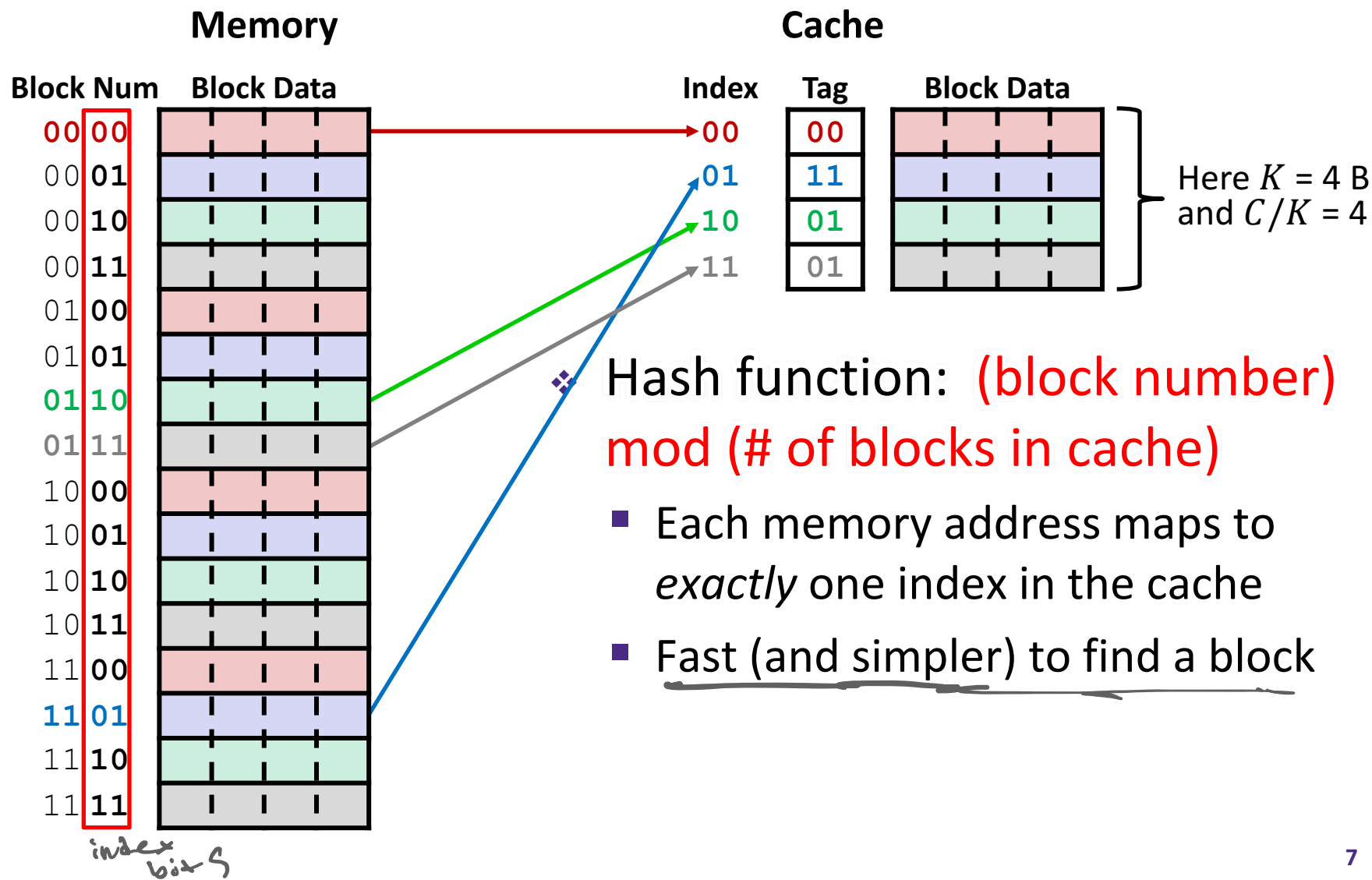
❖ TIO address breakdown:

$m$-bit address:

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

Block Number

- **Index (s)** field tells you where to look in cache *which set*
  - Number of bits is determined by number of sets ($\log_2(C/K/E)$) $== \log_2(s)$
  - Need enough bits to reference every set in the cache
- **Tag (t)** field lets you check that data is the block you want
  - Rest of the bits not used for index and offset ($m - s - k$)
- **Offset (k)** field selects specified start byte within block
  - Number of bits is determined by block size ($\log_2(K)$)
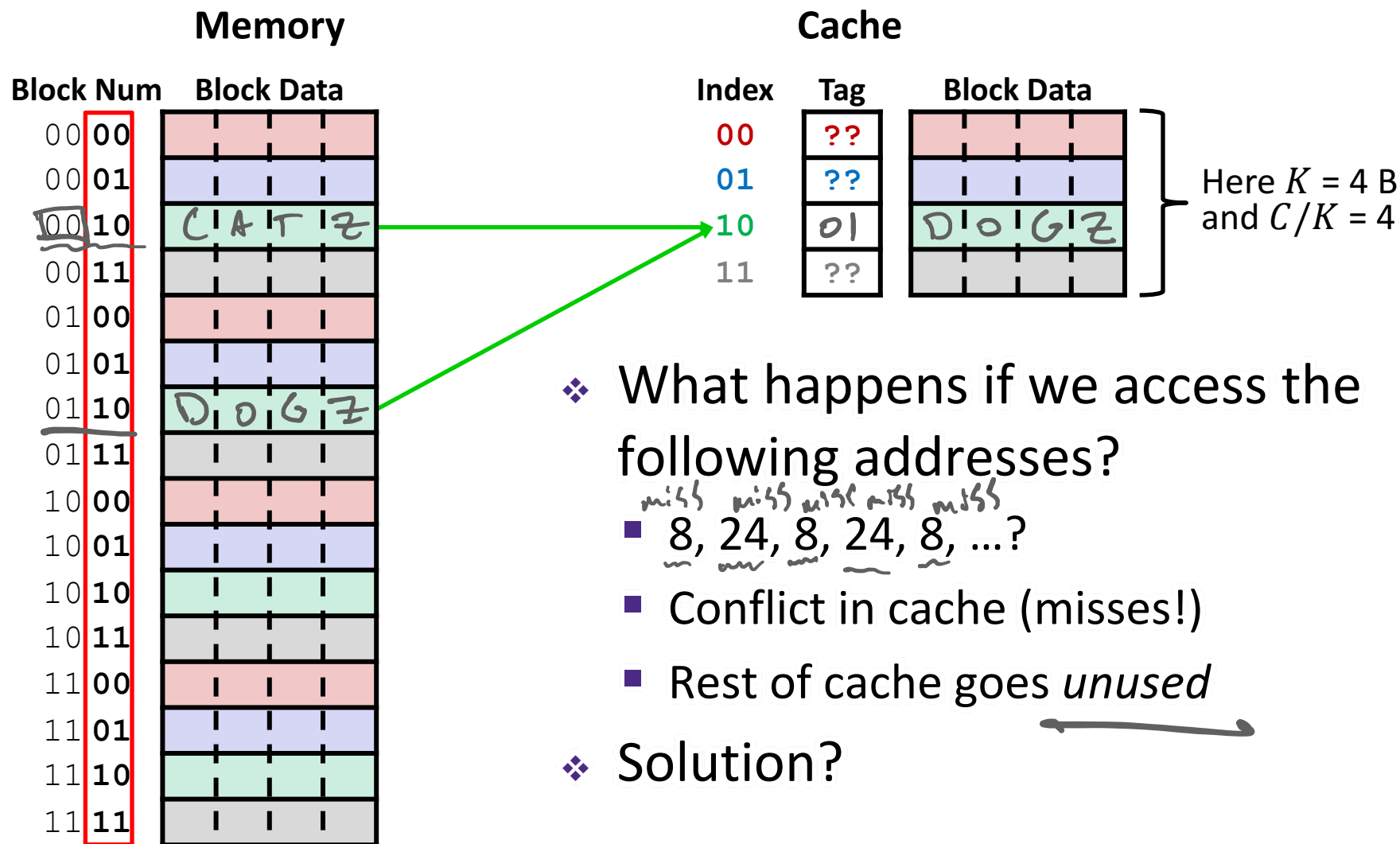  - Need enough bits to reference every byte in a block

# Review: Cache Lookup Process

❖ CPU requests data at a given address
❖ Cache breaks down address into different bit fields
 ▪ Determines offset, index, and tag bits
❖ Cache checks to see if block containing address is already in the cache
 ▪ Uses index bits to find which set to look in
 ▪ Uses tag bits to make sure the block in the set matches

❖ If block is in the cache, it's a **cache hit** *fast!*
 ▪ Data is returned to CPU starting at byte offset
❖ If block is not in the cache, it's a **cache miss** *slow!*
 ▪ Block is loaded from memory into the cache, evicting other blocks from the cache if necessary
 ▪ Data is returned to CPU starting at byte offset

# Review: Direct-Mapped Cache

**Memory**

**Cache**

**Block Num   Block Data**

**Index   Tag   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 |

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

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

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

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

index bits

# Direct-Mapped Cache Problem

**Memory**

| Block Num | Block Data |
|---|---|
| 00 **00** | |
| 00 **01** | |
| 00 **10** | C A T Z |
| 00 **11** | |
| 01 **00** | |
| 01 **01** | |
| 01 **10** | D O G Z |
| 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 | 01 | D O G Z |
| 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, …?

■ Conflict in cache (misses!)

■ Rest of cache goes *unused*

❖ Solution?
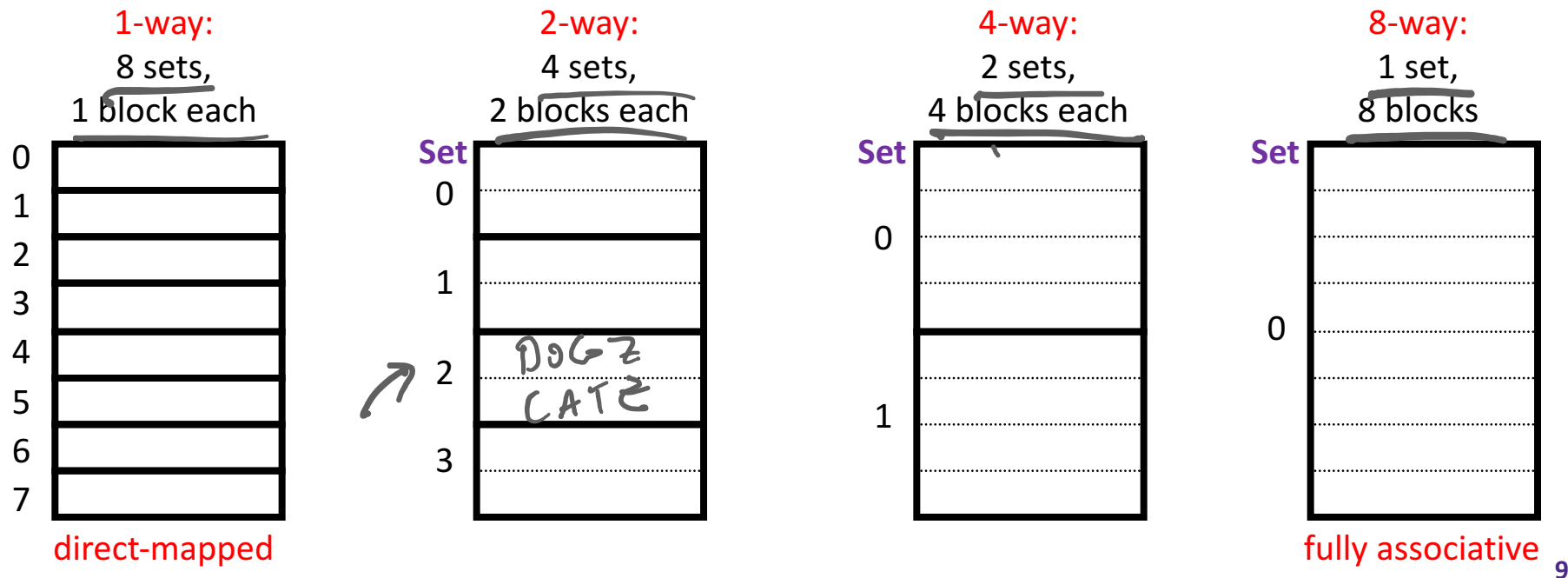
# Associativity

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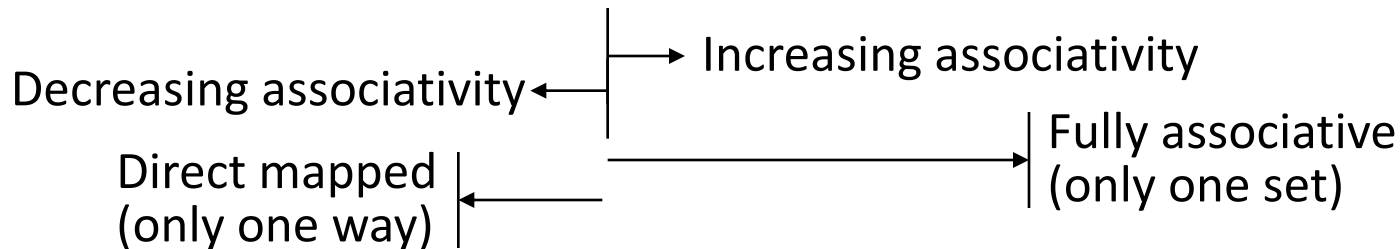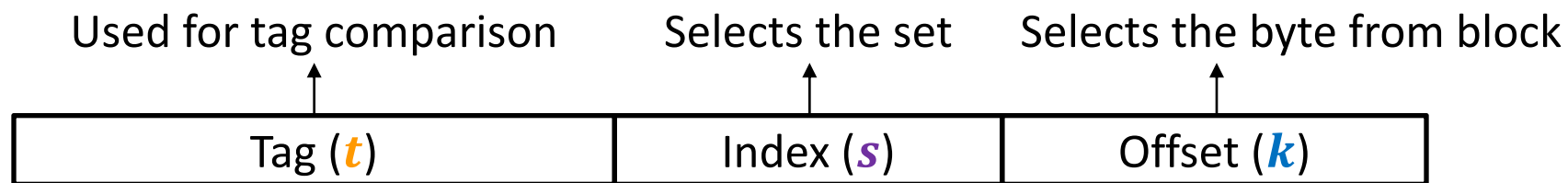  - ▪ Each set can store block in more than one **way**

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

**9**

# 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 $S = C/K/E$

  *num blocks*    *blocks per set*

- 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   *one set!*
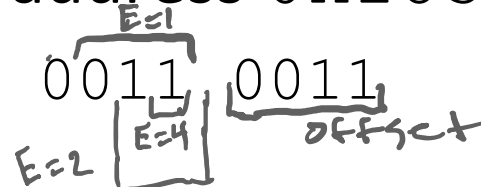
Used for tag comparison     Selects the set     Selects the byte from block

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

Decreasing associativity ←   → Increasing 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`

E=1

E=2   E=4   offset

$$t = m-s-k \quad s = \log_2(C/K/E) \quad k = \log_2(K) = 4$$

$m$-bit address:

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

E=1, S=C/K/E = 8
$s = ?\ \log_2(8) = 3$

E=2  S=C/K/E=4
$s = ?\ \log_2(4) = 2$

E=4  S=C/K/E=2
$s = ?\ \log_2(2) = 1$

Direct-mapped

2-way set associative

4-way set associative

| Set | Tag | Data |
|---|---|---|
| 0 000 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 0 | | |
| 1 1 | | ✓ ✓ ✓ ✓ |

11

# Block Replacement

*valid ≡ 0*

❖ *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    *hard to implement in hardware*

  ▪ Caches typically use something close to ***least recently used (LRU)***
    (hardware usually implements "*not most recently used*")

| 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 | | |

# Polling Question [Cache III]

$$C = 2 * 2^{10} B \qquad\qquad K = 2^7 B$$

❖ We have a cache of size 2 KiB with block size of 128 B. If our cache has 2 sets, what is its associativity?
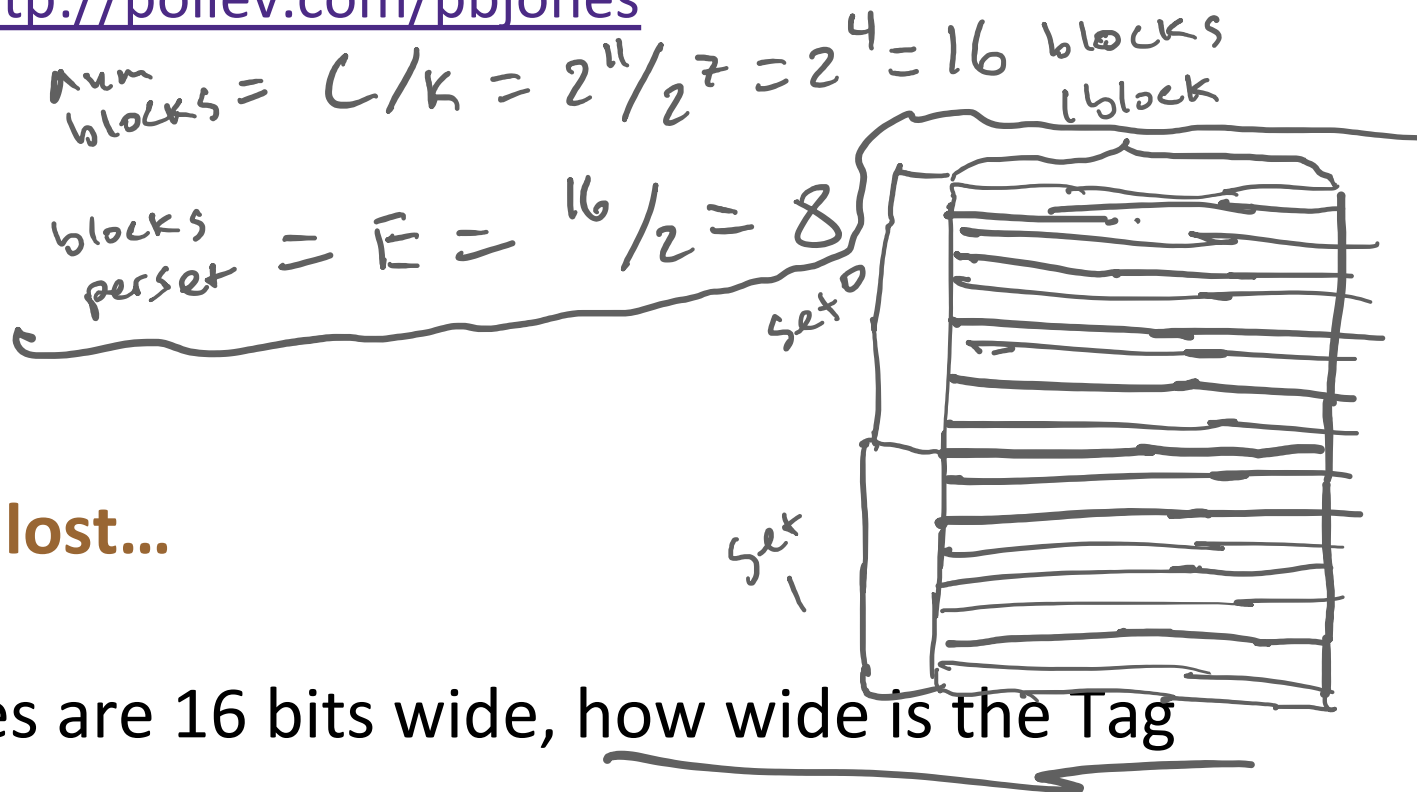
$S = 2$

▪ Vote at http://pollev.com/pbjones

A. **2**

B. **4**

C. **8**
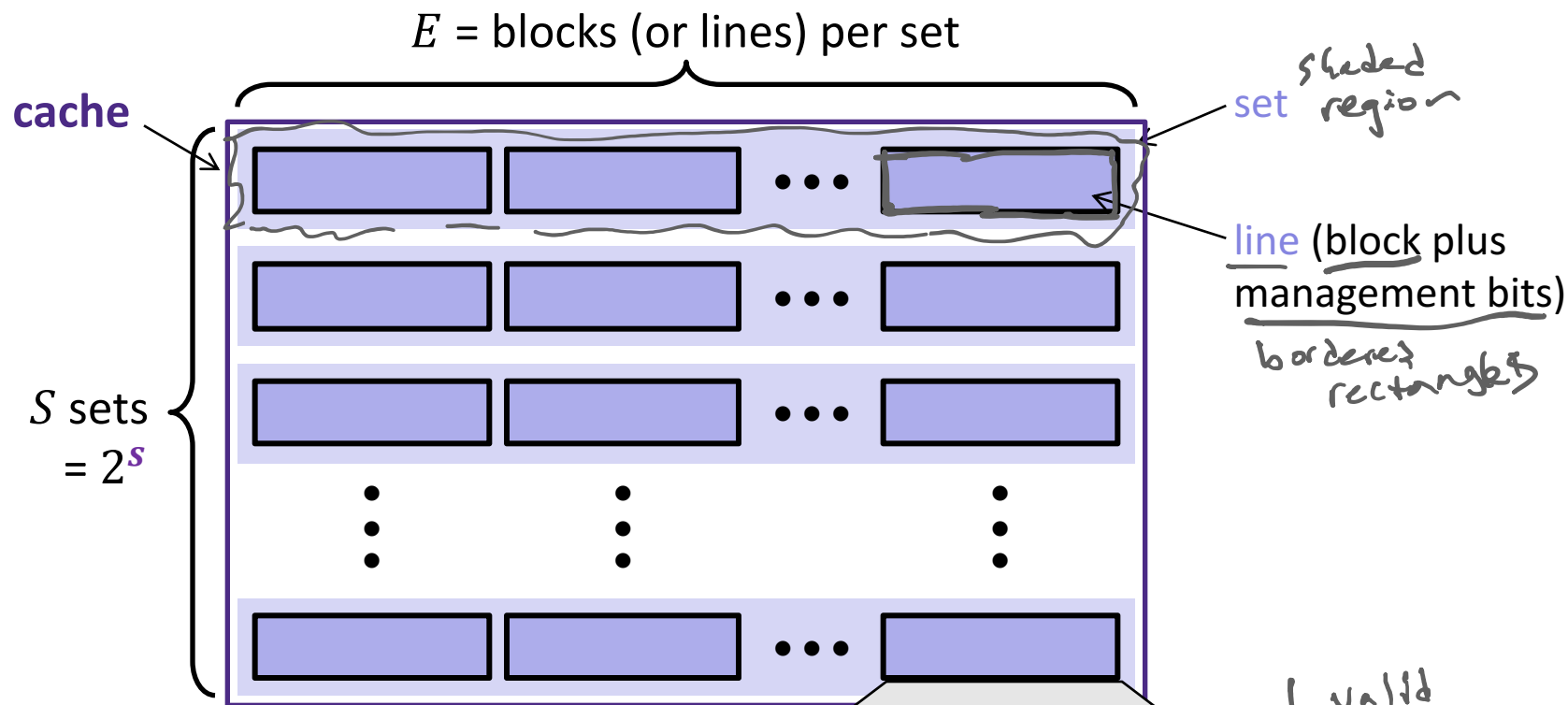
D. **16**

E. **We're lost...**

$$\text{num blocks} = C/K = 2^{11}/2^7 = 2^4 = 16 \text{ blocks}$$
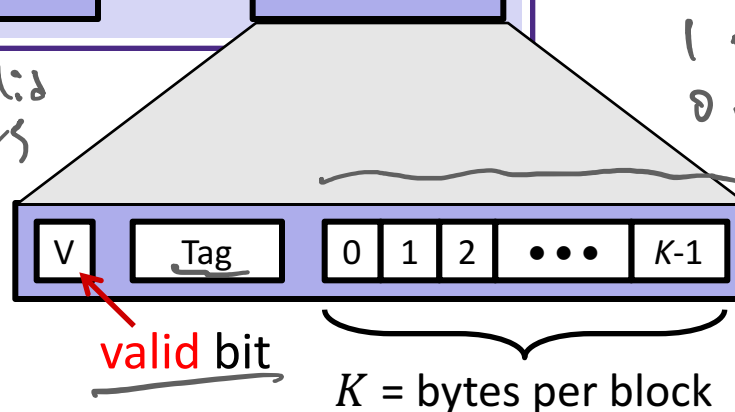
1 block

$$\text{blocks per set} = E = 16/2 = 8$$

set 0

set 1

❖ If addresses are 16 bits wide, how wide is the Tag field?

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

$E$ = blocks (or lines) per set



set — *shaded region*

line (block plus management bits)

*bordered rectangles*

**cache**

$S$ sets
= $2^s$

*Only hit on valid entries*

1 valid
0 not valid

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

| V | Tag | 0 | 1 | 2 | ••• | $K$-1 |

valid bit

$K$ = bytes per block

# 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

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

# Example Cache Parameters Problem

$2^2 * 2^{10} = 2^{12}$

MP

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

$C/k = \frac{2^8}{2^5} = 2^3 = 8$

$m = \log_2(M)$
$= \log_2(2^2 * 2^{10})$
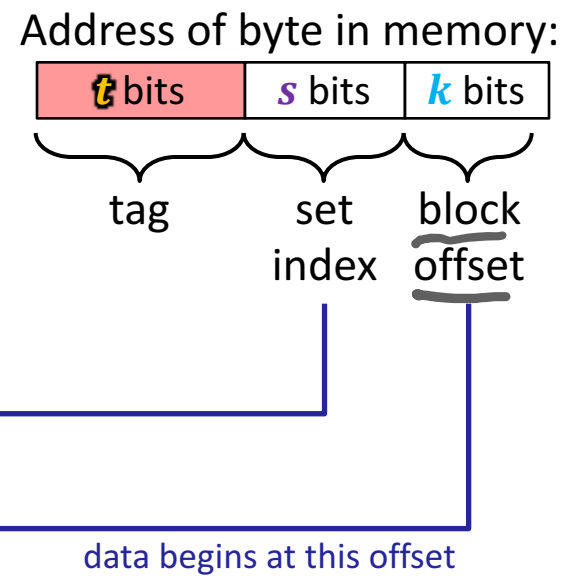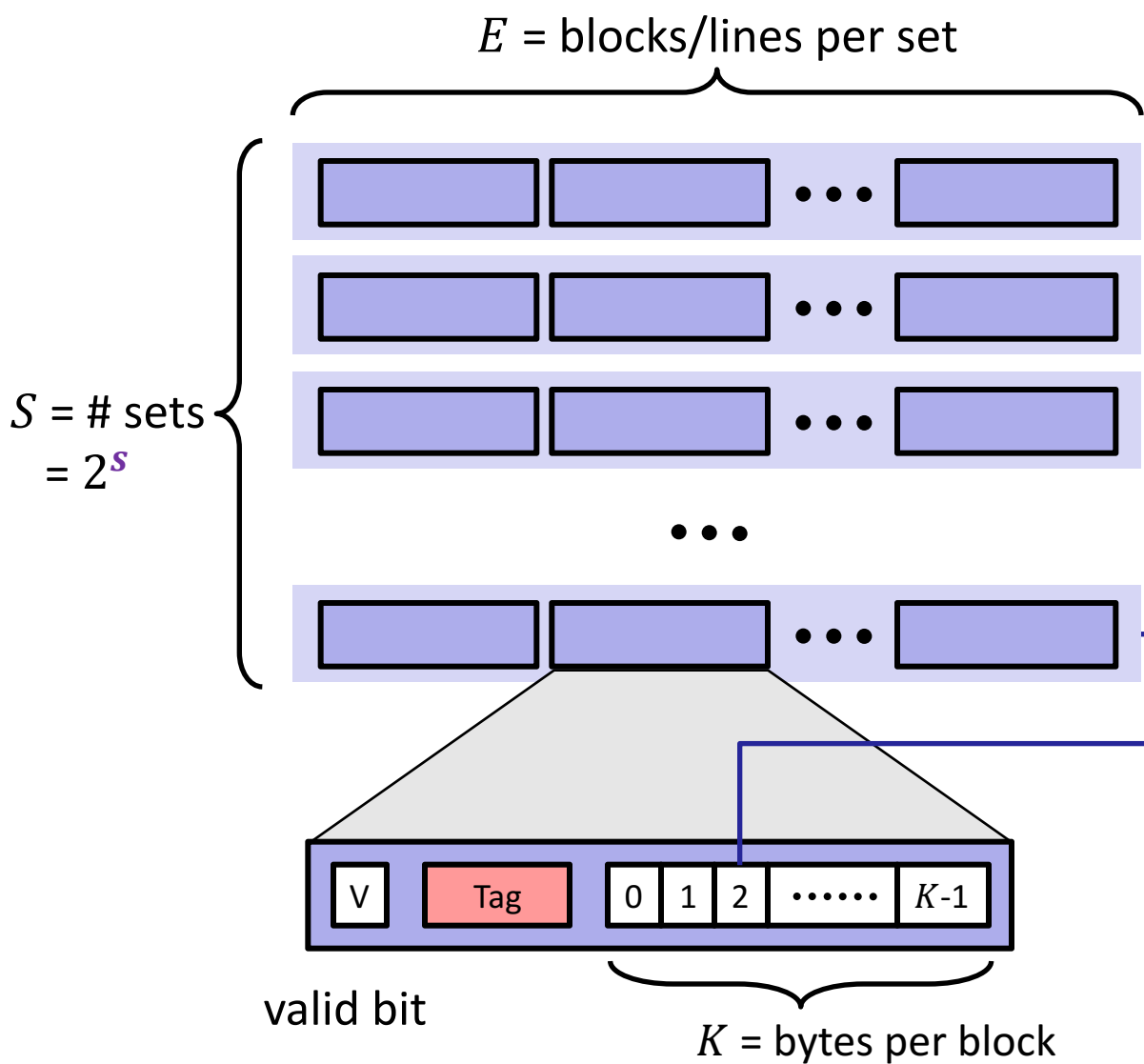$= \log_2(2^{12})$
$= 12 \text{ bits}$

$2^8$   $C/K/E = 8/2 = 4$

C

K   $2^5$

E   $2$

| | |
|---|---|
| **Cache Size** | 256 B |
| **Block Size** | 32 B |
| **Associativity** | 2-way |
| **Hit Time** | 3 cycles |
| **Miss Rate** | 20% |
| **Tag Bits** | $= 12 - 2 - 5 = 5 \text{ bits}$ |
| **Index Bits** | $= \log_2(4) = 2 \text{ bits}$ |
| **Offset Bits** | $= \log_2(32) = 5 \text{ bits}$ |
| **AMAT** | $= 3 + .2(125) \approx 28 \text{ cycles}$ |

HT

MR

$m - s - k$

$\begin{cases} \log_2(C/K/E) \\ \log_2(K) \end{cases}$

$HT + MR * MP$

# Cache Read

$E$ = blocks/lines per set

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

$S$ = # 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 ($E = 1$)

*one block per set*

Direct-mapped: One line per set

Block Size $K$ = 8 B

*8 bytes per block*



$S = 2^s$ sets

Set 0

Set 1

Set 2

Set $S-1$

Address of `int`:

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

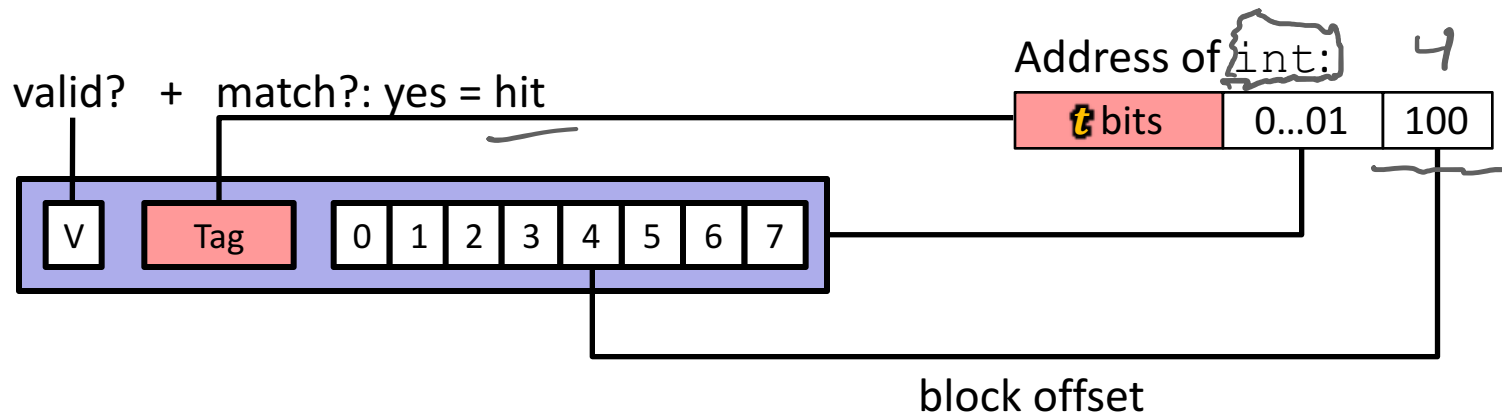find set

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

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

*Start reading at offset*
*read sizeof (int) bytes*

Address of int:   4

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

valid? + match?: yes = hit

| V | | Tag | | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

block offset

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

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

multiply of 2

Address of int:

valid?  +  match?: yes = hit

| $t$ bits | 0…01 | 100 |

multiply of 8

block offset

int (4 B) is here

Don't want to have to access 2 blocks for one int

**This is why we want alignment!**
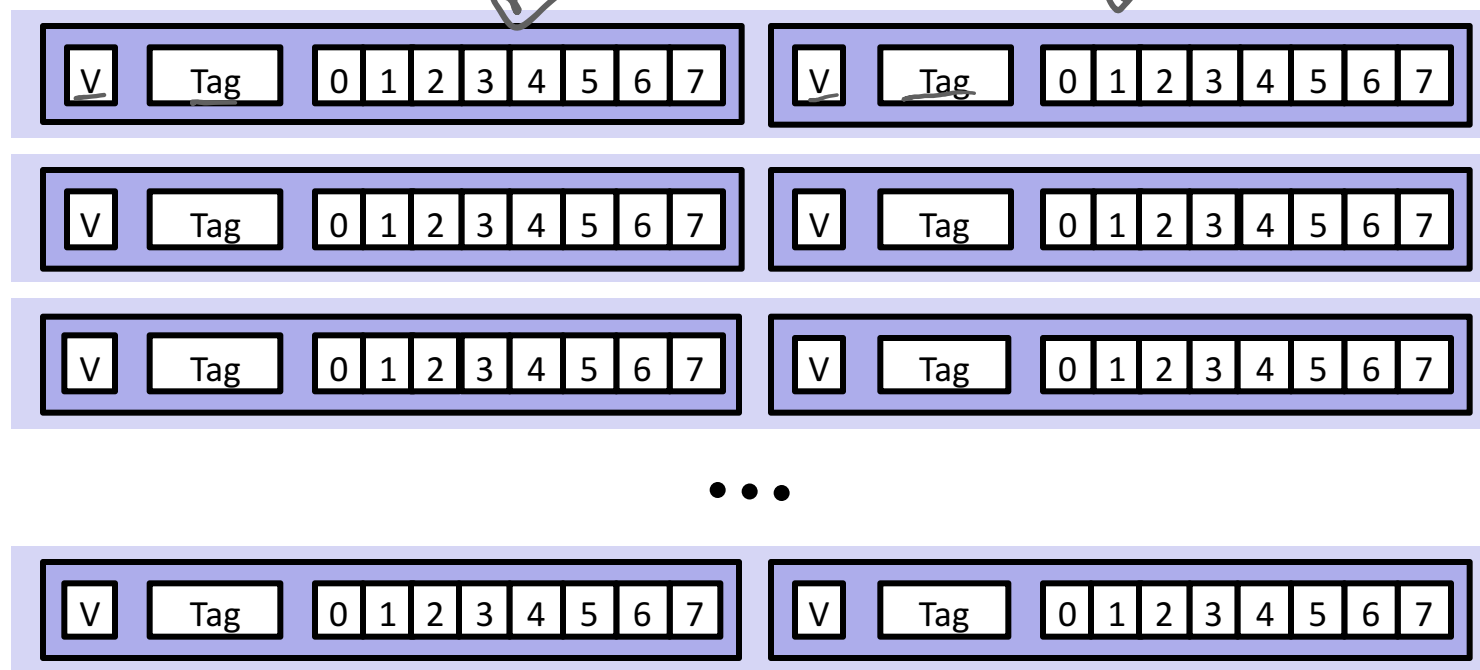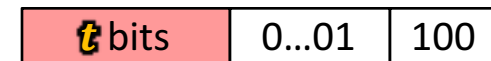
No match? Then old line gets evicted and replaced

by the new block from memory

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

2 blocks per set

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

Address of `short int`:

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

find set

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

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

Address of short int:
2 bytes

| $t$ 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 |

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 |
|---|---|---|---|---|---|---|---|---|---|

block offset

`short int` (2 B) is here

## No match?

- One line in set is selected for eviction and replacement
- Replacement policies: random, least recently used (LRU), …

*not most recently used*

# 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)

*to store data we are working with*

❖ **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*)
  
  *one set*
  
  ▪ **Note:** *Fully-associative* only has Compulsory and Capacity misses

# Example Code Analysis Problem

$K = \log_2(K) = 5$ bits
$S = \log_2(S) = 2$ bits
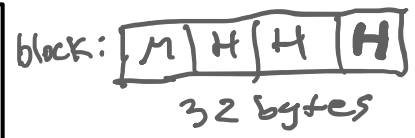$t = m - K - s = 5$ bits

❖ Assuming the cache starts <u>cold</u> (all blocks invalid) and `sum`, `i`, and `j` are stored in registers, calculate the **miss rate**:
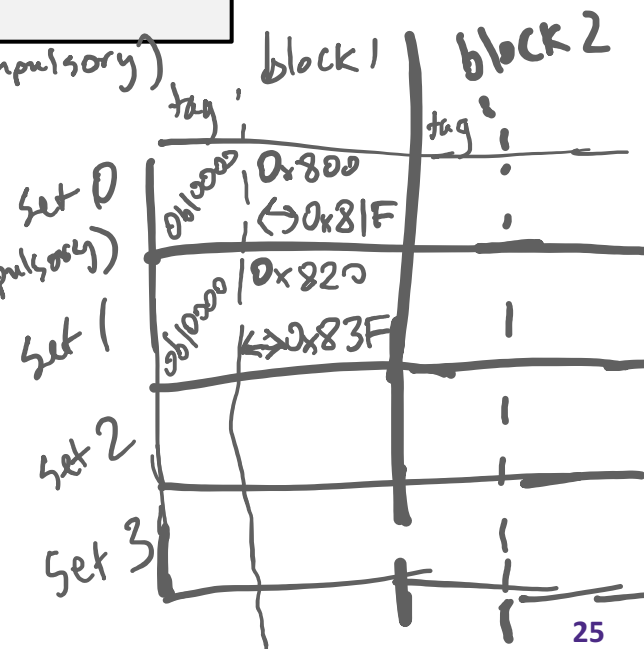
- $m$ = 12 bits, $C$ = 256 B, $K$ = 32 B, $E$ = 2    $S = C/K/E = 4$

4 longs

block: | M | H | H | H |

32 bytes

```
#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 byte elements

ar[0][0] →addr→ 0x800 = 0b1000 0000 0000 | miss! (compulsory)
  tag / index / offset

ar[0][1] → 0x808 = 0b1000 0000 1000  hit!
  tag

ar[0][2] → 0x810 = same block, hit!

ar[0][3] → 0x818 = same block, hit!

ar[0][4] → 0x820 = 0b1000 0010 0000  miss! (compulsory)
  tag / set / offset

ar[0][5] → 0x828 hit!   tag

Pattern continues...

4 longs/accesses per block,
one compulsory miss, then 3 hits

miss rate = 1/4

block 1   block 2
tag       tag

Set 0 | 0b10000 | 0x800
        |        | ↔ 0x81F
Set 1 | 0b10000 | 0x820
        |        | ↔ 0x83F
Set 2

Set 3

25

# Notes Diagrams

$E$ = blocks (or lines) per set



set

line

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

valid bit

| V | Tag | 0 | 1 | 2 | $\cdots$ | $K$-1 |

management
bits

$K$ = bytes per block

1-way:
8 sets,
1 block each

Set

0
1
2
3
4
5
6
7

direct-mapped

2-way:
4 sets,
2 blocks each

Set

0

1

2

3

4-way:
2 sets,
4 blocks each

Set

0

1

8-way:
1 set,
8 blocks

Set

0

fully associative