

Memory & Caches III

CSE 351 Spring 2024

Instructor:

Elba Garza

Teaching Assistants:

Ellis Haker

Adithi Raghavan

Aman Mohammed

Brenden Page

Celestine Buendia

Chloe Fong

Claire Wang

Hamsa Shankar

Maggie Jiang

Malak Zaki

Naama Amiel

Nikolas McNamee

Shananda Dokka

Stephen Ying

Will Robertson



Playlist: [CSE 351 24Sp Lecture Tunes!](#)

Relevant Course Information

- ❖ HW 15 due tonight! HW16 due Monday
- ❖ HW 17/18 due following Friday (10 May)
 - Covers the major cache mechanics—big homework, start soon!
- ❖ Take-home Midterm, May 6th to May 7th
 - 48 hours, but should take 1-3 hours to complete
 - No in-person lecture on Monday the 6th—I will post a new recording instead
- ❖ Mid-Course Canvas Survey due May 6th by 11:59 PM
- ❖ Lab 3 due Wednesday, May 8th
- ❖ Lab 4 releasing soon afterward!
 - Can do Part 1 after today; will need Lecture 19 to do Part 2

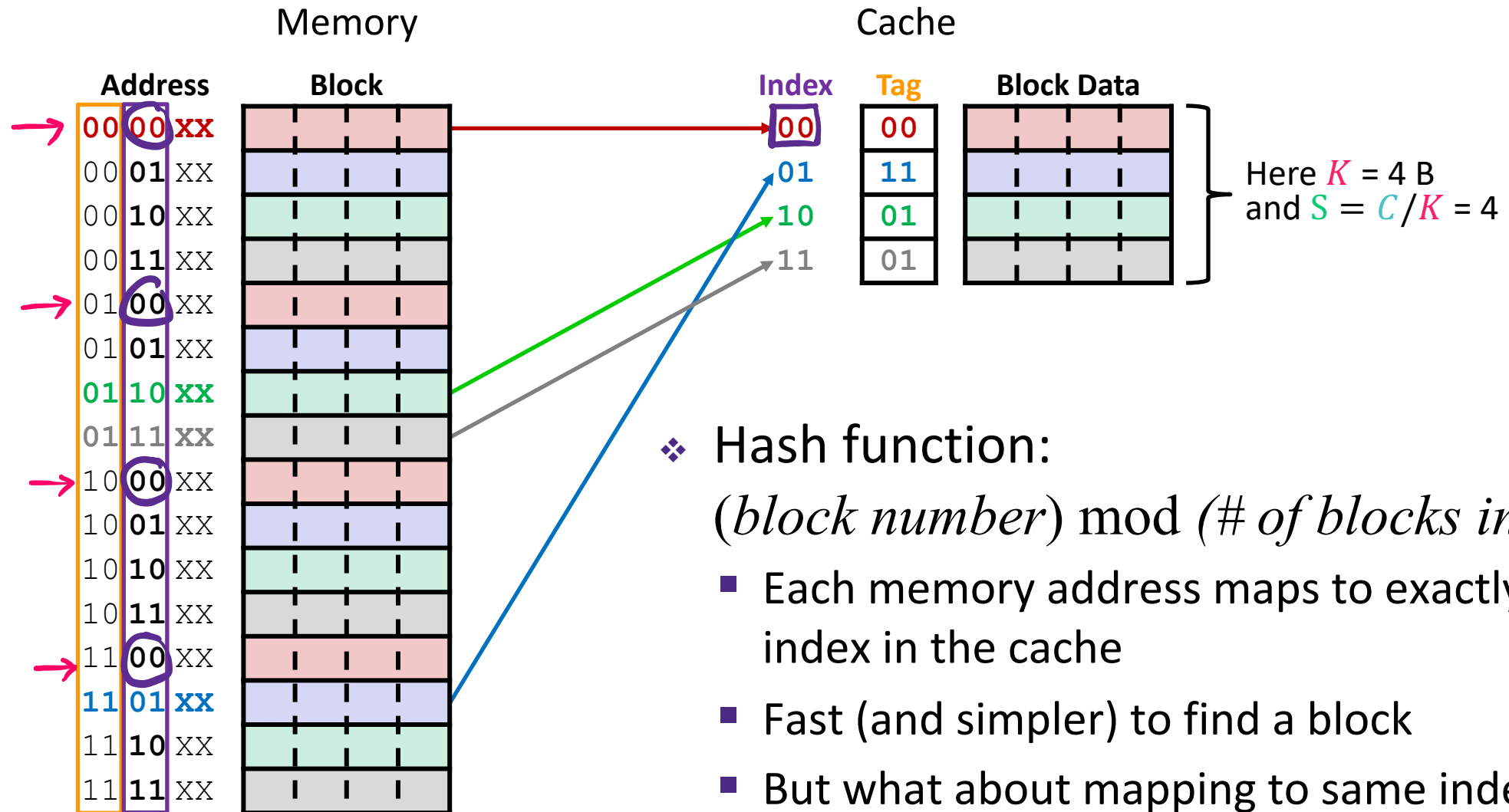
Making memory accesses fast!

- ❖ 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

Reading Review

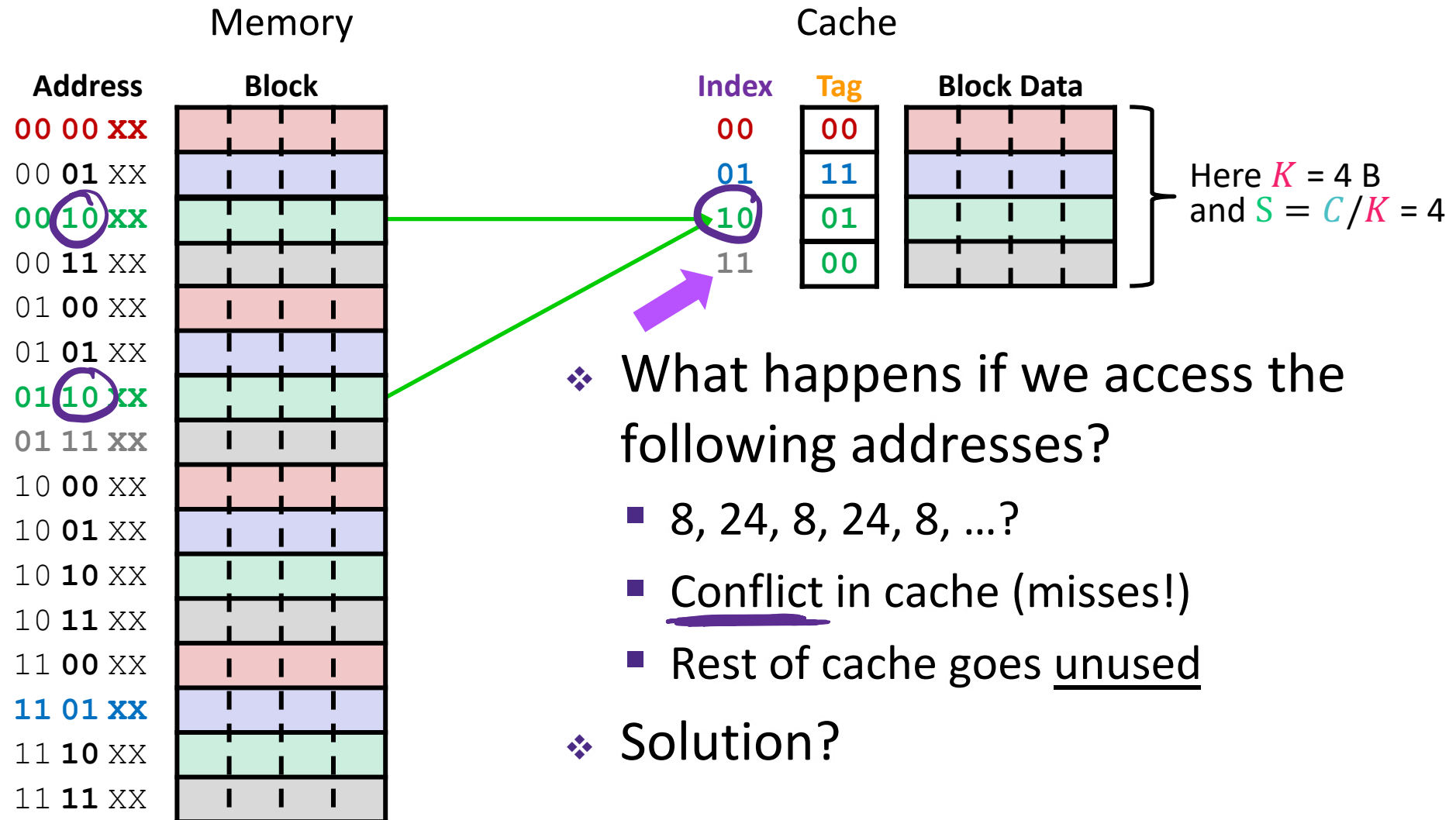
- ❖ Terminology:
 - Associativity: sets, fully-associative cache
 - Replacement policies: least recently used (LRU)
 - Cache line: cache block + management bits (valid, tag)
 - Cache misses: compulsory, conflict, capacity

Review: Direct-Mapped Cache



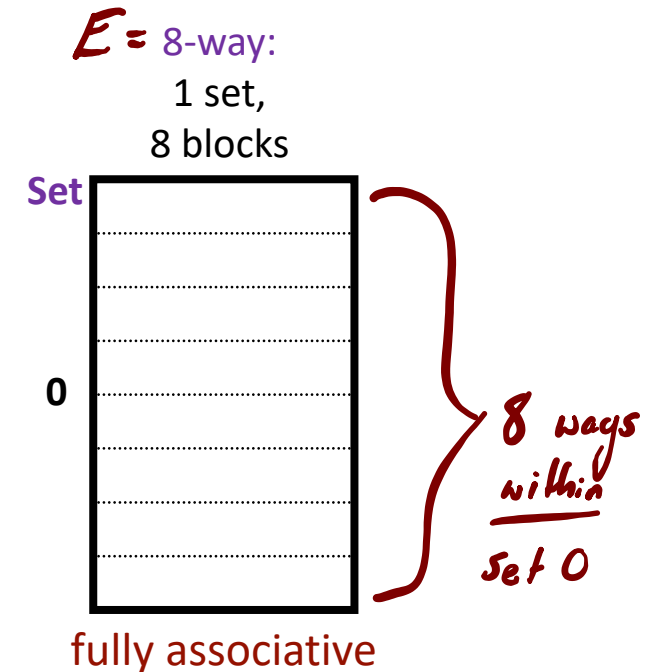
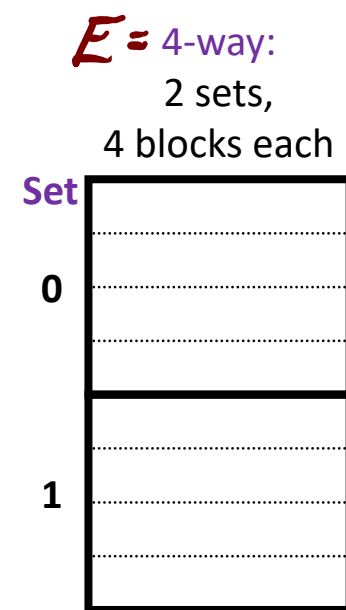
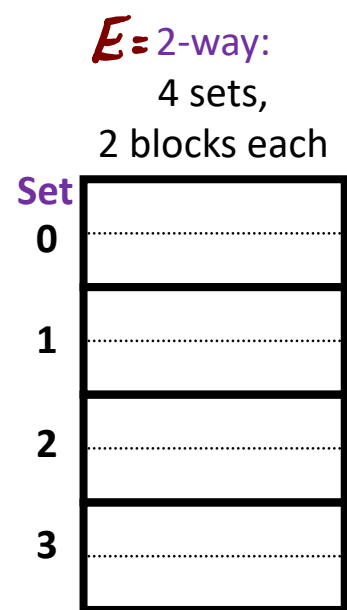
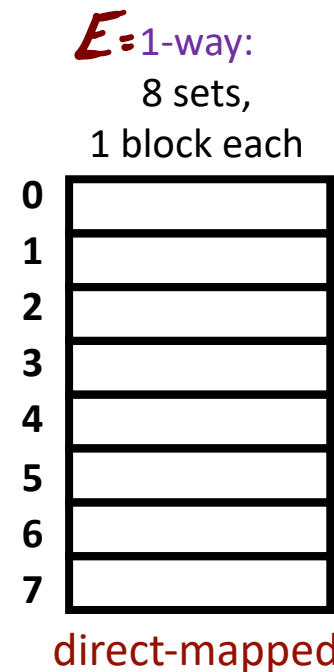
- ❖ Hash function: $(block\ number) \bmod (\#\ of\ blocks\ in\ cache)$
 - Each memory address maps to exactly one index in the cache
 - Fast (and simpler) to find a block
 - But what about mapping to same index?...

Direct-Mapped: A Problem!



Associativity: A Solution!


- ❖ What if we could store **any** data in **any** place in the cache? 💡
 - But: requires more complicated hardware \Rightarrow more power consumed, slower
- ❖ Let's combine the two ideas:
 - Each address maps to exactly one **set**, but each set can store block in more than one **way** in the set!

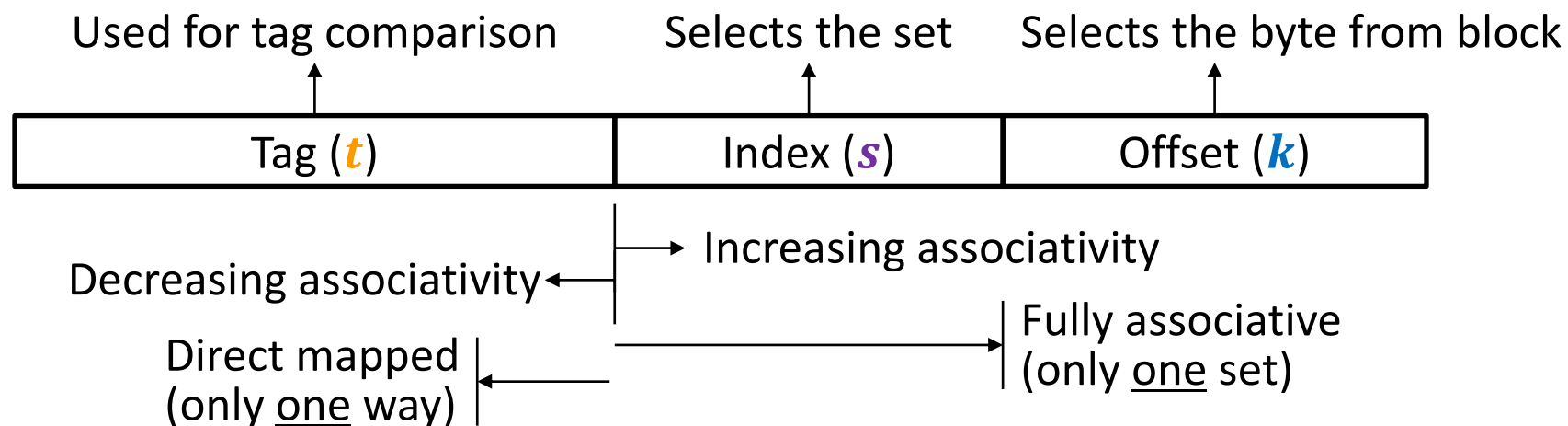


Cache Organization (3)

Note: The textbook uses "b" for offset bits

- ❖ **Associativity (E):** number of **ways** to store in 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

New variable? Kinda.
Direct-mapped means $E=1$, so equation was still correct


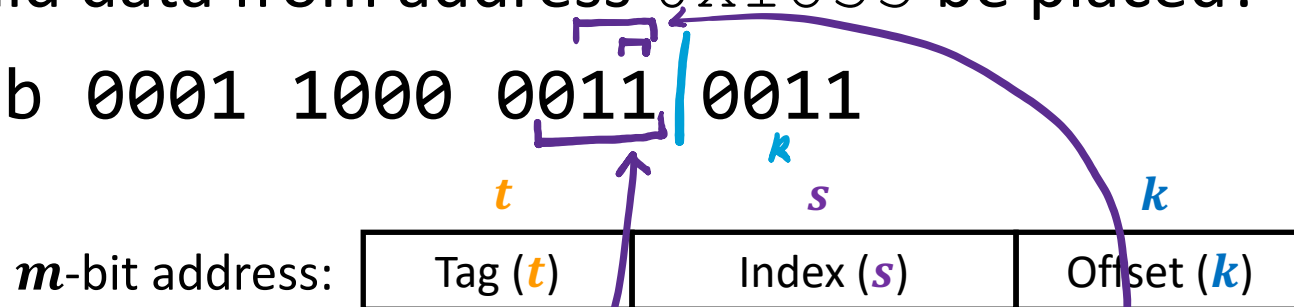


Example Placement

block size K :	16 B
Capacity C/K :	8 blocks *
Address m :	16 bits

❖ Where would data from address $0_{\times}1833$ be placed?

■ Binary: 0b 0001 1000 0011 | 0011



$t = m - s - k$

$s = \log_2(C/K/E)$

$k = \log_2(K)$

$s = \log_2(C/K/E)$
8

(see top right*)

$E = 1$
 $s = \log_2(8/1) = 3 \text{ bits}$

Direct-mapped

Set	Tag	Data
(000) 0		
1		
2		
(011) 3		
4		
5		
6		
(11) 7		

$E = 2$
 $s = \log_2(8/2) = 2 \text{ bits}$

2-way set associative

Set	Tag	Data
0 (00)		
1 (01)		
2 (10)		
3 (11)		

$E = 4$
 $s = \log_2(8/4) = 1 \text{ bit}$

4-way set associative

Set	Tag	Data
0 (0)		
1 (1)		

Block Placement and Replacement

- ❖ Any empty block in the correct set may be used to store block
 - Valid bit for each cache block indicates if valid (1) or mystery (0) data
- ❖ If there are no empty blocks, which one should we replace? i.e. replacement policy
 - No choice for direct-mapped caches—gotta replace what’s there. Super easy.
 - Otherwise, caches typically use something close to **least recently used (LRU)** (hardware usually implements “*not most recently used*”)

↓ Direct-mapped

Set	V	Tag	Data
0			
1			
2			
3			
4			
5			
6			
7			

↓ 2-way set associative

Set	V	Tag	Data
0			
1			
2			
3			

↓ 4-way set associative

Set	V	Tag	Data
0			
1			

Polling Questions

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

$$K = 2^7 B \Rightarrow k = 7 \text{ bits}$$

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

- A. 2
- B. 4
- C. 8
- D. 16
- E. We're lost...

$$S = 2^1 \Rightarrow s = 1 \text{ bit}$$

$$E = ???$$

$$(C/K) = 2^{11} / 2^7 = 2^4$$

$$S = C/K/E$$

$$\Rightarrow 2 = 2^4 / E = 16 / E \quad \therefore E = 8$$

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

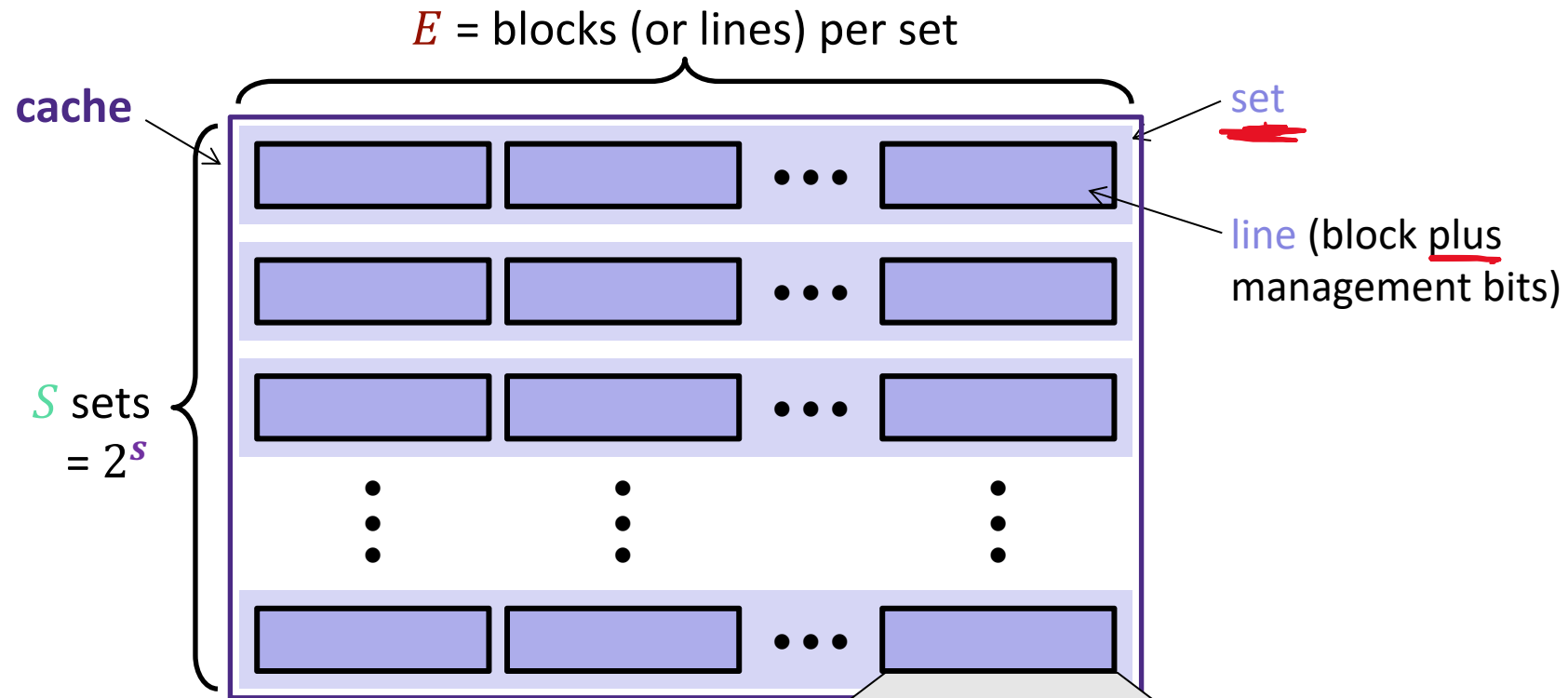
$$K = 2^7 \quad S = 2^1 \quad m = 16$$

$$k = 7 \quad s = 1$$

$$t = m - k - s$$

$$\boxed{8} = 16 - 7 - 1$$

General Cache Organization (S, E, K)



Cache size:

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

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

Example Cache Parameters Problem

❖ 1 KiB address space, 125 cycles to go to memory. → MP

Fill in the following table:

Cache Size C	64 B = 2^6
Block Size K	8 B = 2^3
Associativity E	2-way = 2^1
Hit Time	3 cycles
Miss Rate	20% = 0.2
Address width (m)	10 bits
Tag Bits (t)	$(m - s - k)$ 5 bits
Index Bits (s)	3 bits
Offset Bits (k)	3 bits
AMAT	

$M = 2^{10} B$

$S = 2^6 / 2^3 / 2^1$
 $= 2^2 = 4$

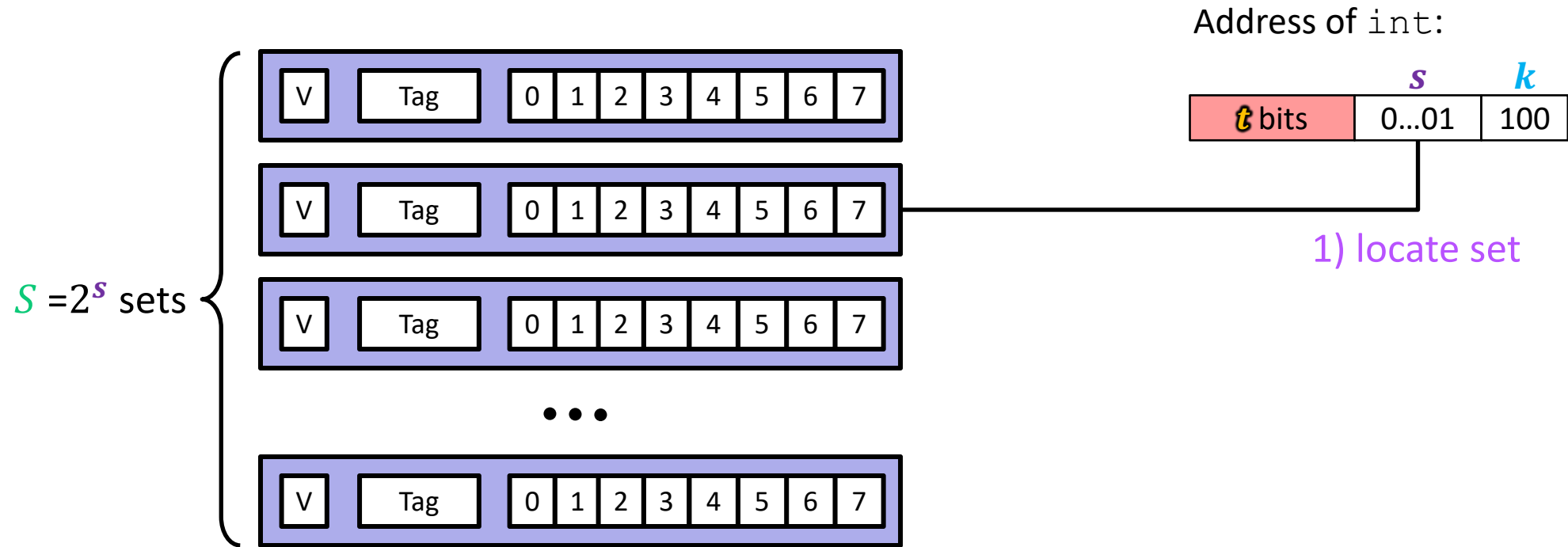
$m = \log_2(2^{10}) = 10$

$t = 10 - 2 - 3$
 $s = \log_2(4) = 3$
 $k = \log_2(2^3) = 3$

Read: Direct-Mapped Cache ($E = 1$)

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

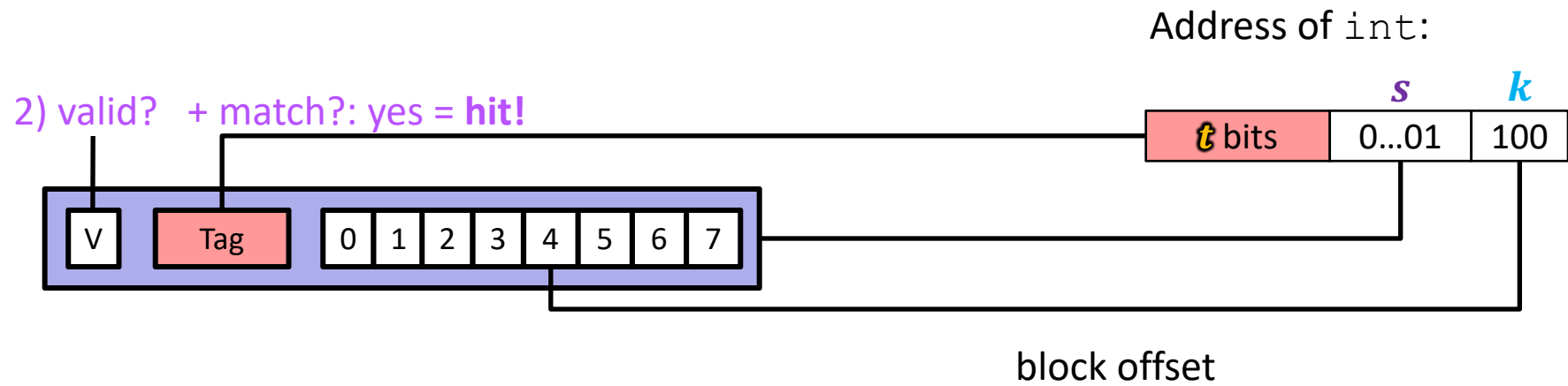
- 1) Locate set
- 2) Check if any line in set is valid and has matching tag: **hit!**
- 3) Locate data starting at offset



Read: Direct-Mapped Cache ($E = 1$)

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

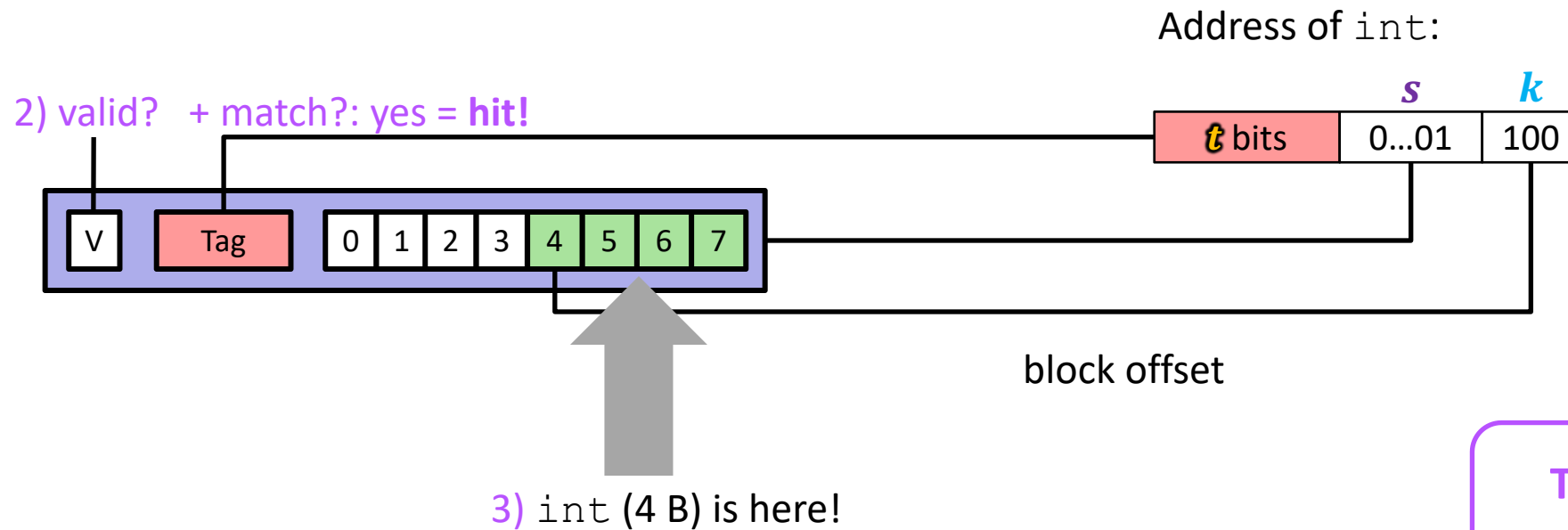
- 1) Locate set
- 2) Check if any line in set is valid and has matching tag: **hit!**
- 3) Locate data starting at offset



Read: Direct-Mapped Cache ($E = 1$)

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

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



This is why we want alignment!

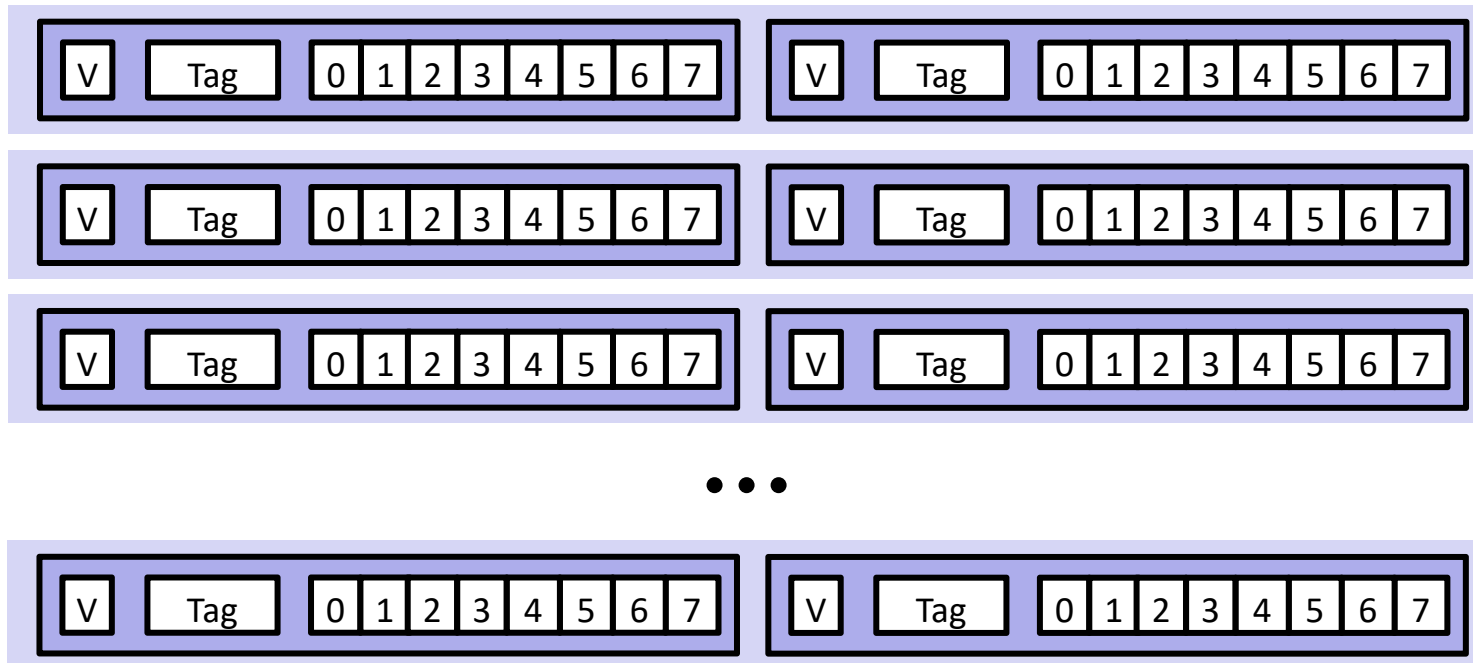
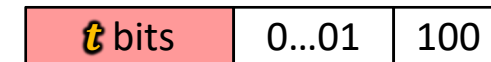
No match? Then old line/block gets evicted and replaced!

Read: Set-Associative Cache ($E = 2$)

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

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

Address of short int:

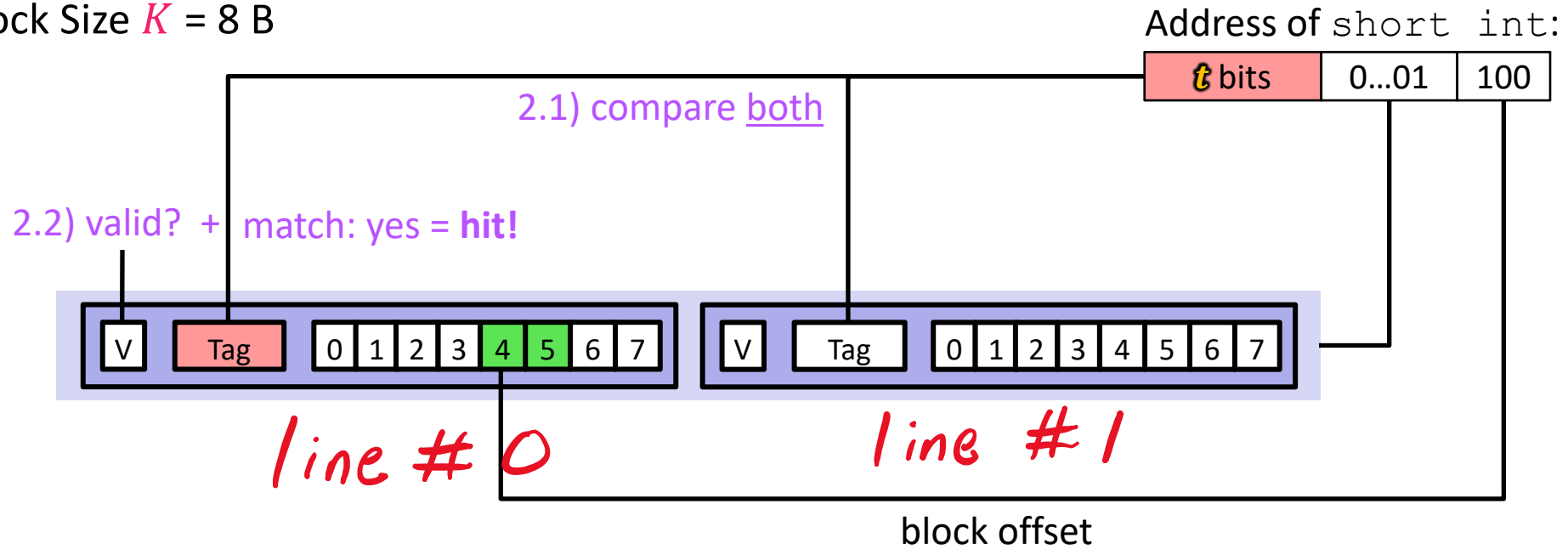


1) locate set

Read: Set-Associative Cache ($E = 2$)

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

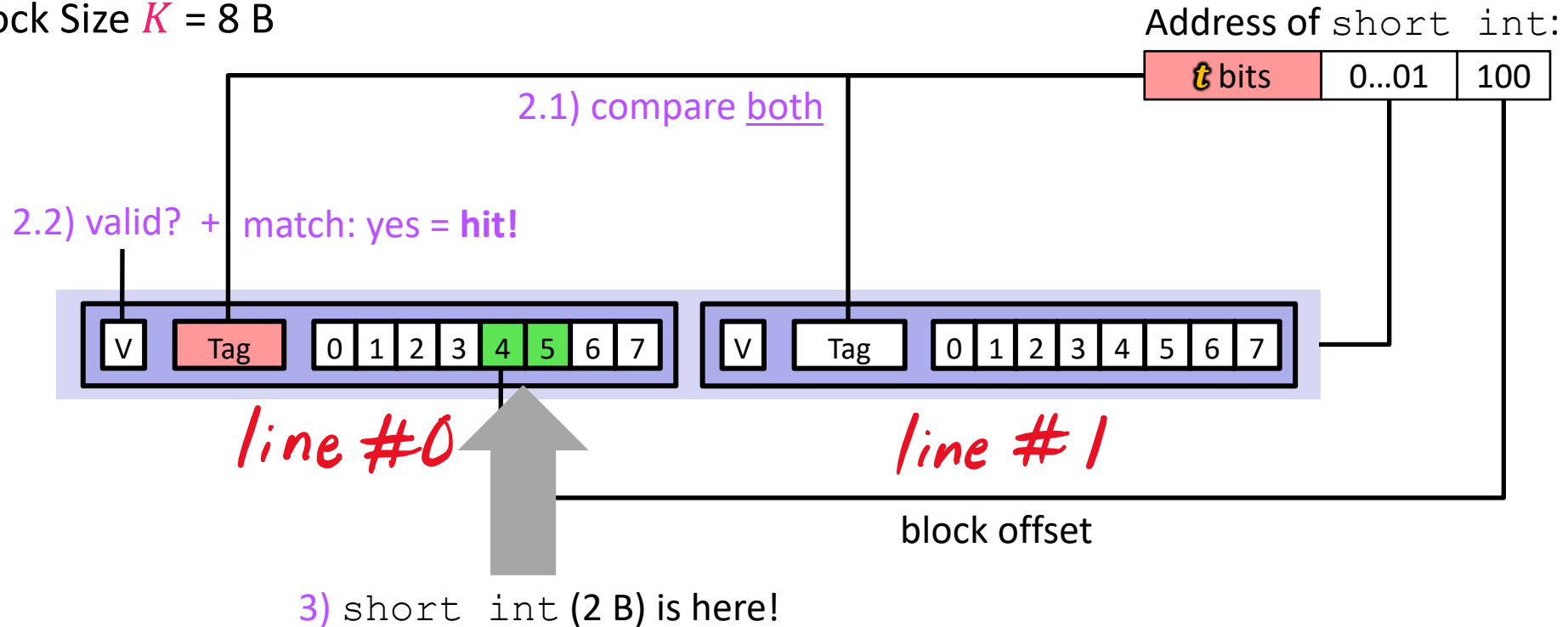
2-way: Two lines per set
 Block Size $K = 8$ B



Read: Set-Associative Cache ($E = 2$)

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

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

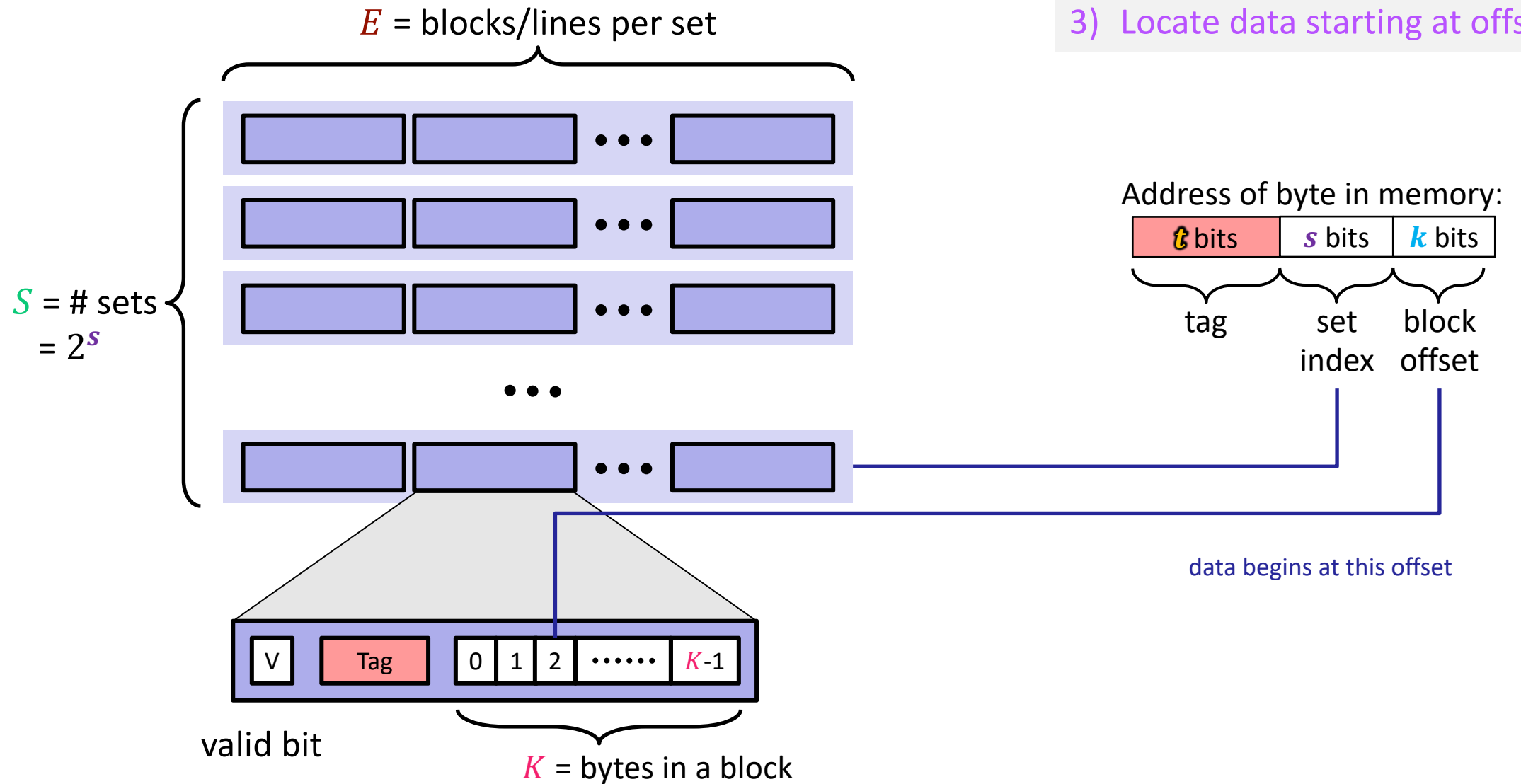


No match?

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

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



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