# The Hardware/Software Interface

L18: Memory & Caches III

Memory & Caches III

### **Guest Instructor:**

Naama Amiel

## **Teaching Assistants:**

**Anthony Mangus** 

**Grace Zhou** 

Jiuyang Lyu

Kurt Gu

Mendel Carroll

Naama Amiel

Rose Maresh

**Violet Monserate** 

Divya Ramu

Jessie Sun

Kanishka Singh

Liander Rainbolt

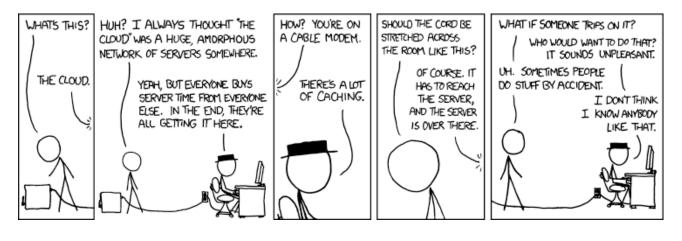
Ming Yan

Pollux Chen

Soham Bhosale

## **Instructors:**

Justin Hsia, Amber Hu

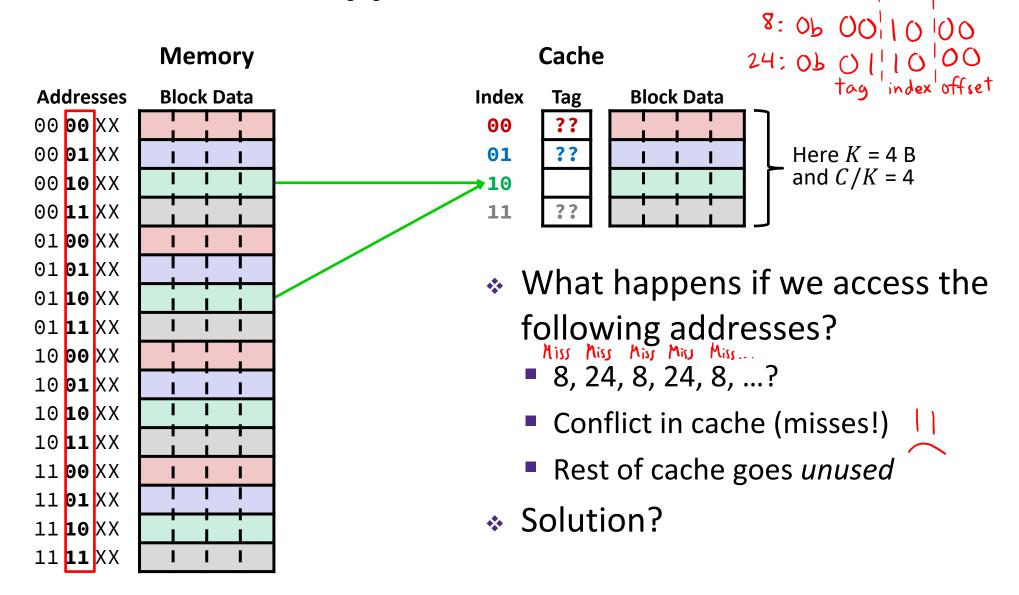


http://xkcd.com/908/

## Lecture Outline (1/3)

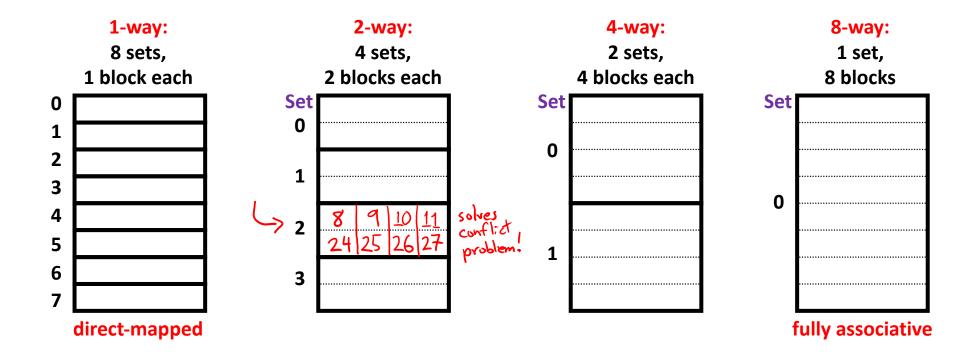
- Associativity and Replacement
- Cache Organization
- Example Cache Problems

## **Review: Direct-Mapped Cache Problem**



## **Associativity Motivation**

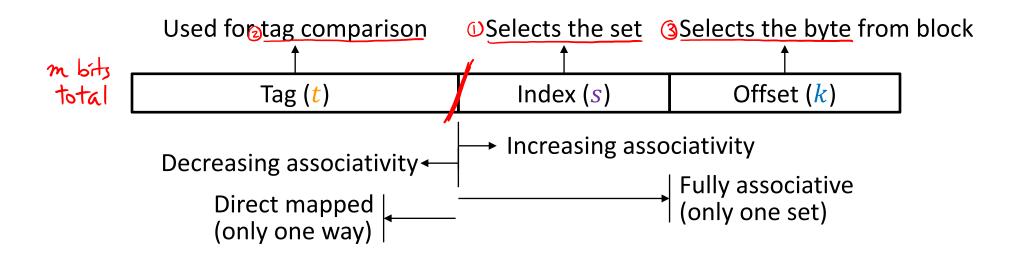
- What if we could store data in any place in the cache?
  - More complicated hardware = more power consumed, slower
- So we combine the two ideas:
  - Each address maps to exactly one set, each set can store block in some # of ways



## **Associativity (Review)**

**Note:** The textbook uses "b" for offset bits

- $\star$  Associativity (E): # of ways for each set
  - Such a cache is called an "E-way set associative cache"
  - We index into cache *sets*, of which there are S = (C/K)/E
  - Use lowest  $log_2(S) = s$  bits of block address
    - <u>Direct-mapped</u>: E = 1, so  $s = \log_2(C/K)$  as we saw previously
    - Fully associative: E = C/K, so s = 0 bits



## **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  $s = \log_2(C/K/E)$   $k = \log_2(K)$ m-bit address: Tag (t) Index (s) Offset (k)

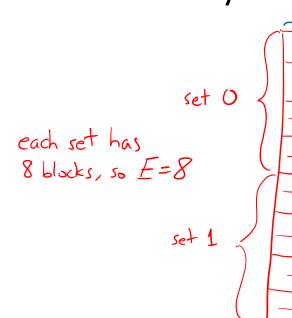
s = ?3 bits s = ?2 bAss = ?1 bit Direct-mapped 2-way set associative 4-way set associative Set Tag Set Tag **Data** Set Tag **Data Data**  $(\omega\omega)$  0 (OO)(001) 1 (0)0 (010) 2 (OI)1 (UII) 3 (<del>\)</del>သ) 4 (0)2 (101) 5 (1)1 (110) 6 (11)3 (III) 7

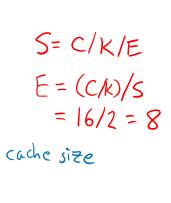
## **Polling Questions**

 $PC = 2^{11}$ 

- K=27B
- ❖ We have a cache of size 2 KiB with block size of 128 B.
  If our cache has 2 sets, what is its associativity?
  - cache holds C/K=211-7=24=16 blocks

- **A.** 2
- **B.** 4
- C. 8
- D. 16





m=16

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

$$k = log_2(K) = 7 \text{ bits}, s = log_2(S) = 1 \text{ bits}$$

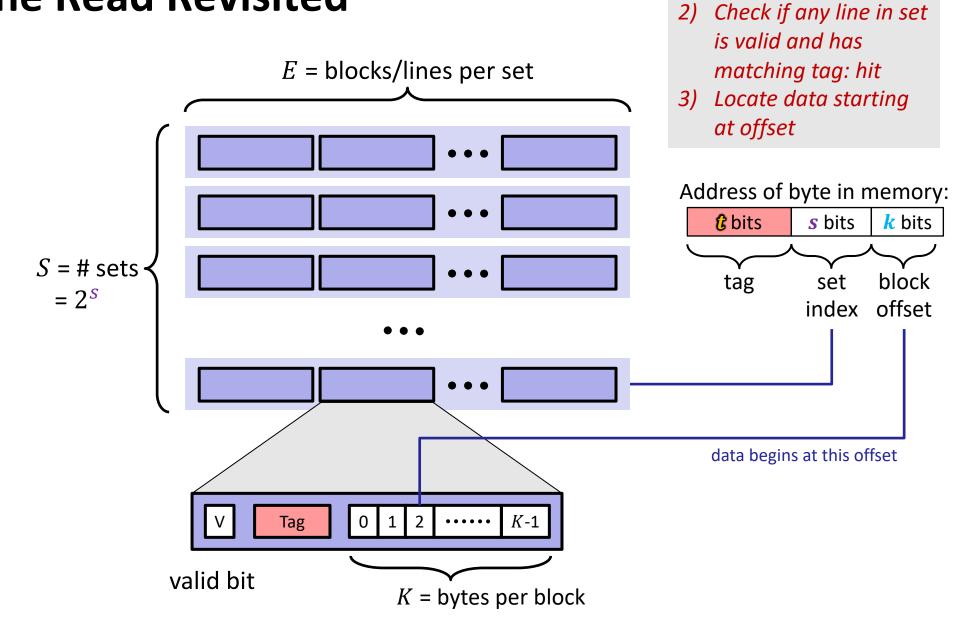
## **Block Placement and Replacement (Review)**

- 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 <u>replace</u>?
  - No choice for direct-mapped caches
  - Caches typically use something close to least recently used (LRU)
     (hardware usually implements "not most recently used")

	Direct-mapped		2-way set associative					4-way set associative			
Set	V	Tag	Data	Set	V	Tag	Data	Set	V	Tag	Data
0				0							
1				U				0			
2				1				U			
3					† <u> </u>						
4				2							
5				۷				1			
6				3							
7				3							

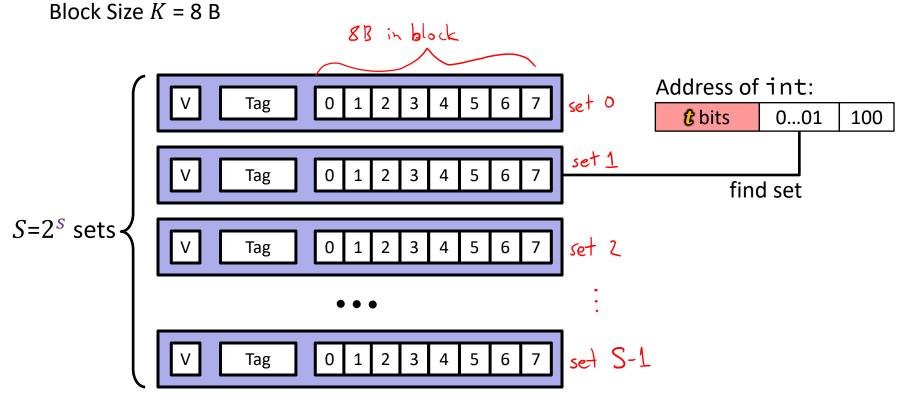
Locate set





# Example: Direct-Mapped (E = 1) Cache (1/3)

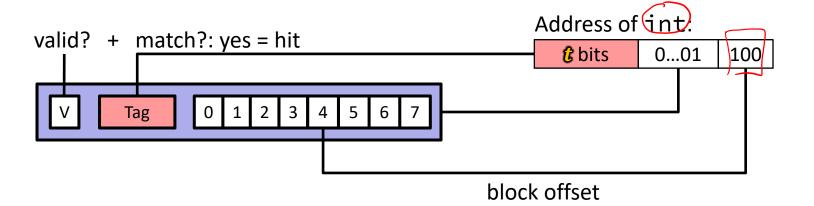
Direct-mapped: One line per set



# Example: Direct-Mapped (E = 1) Cache (2/3)

Direct-mapped: One line per set

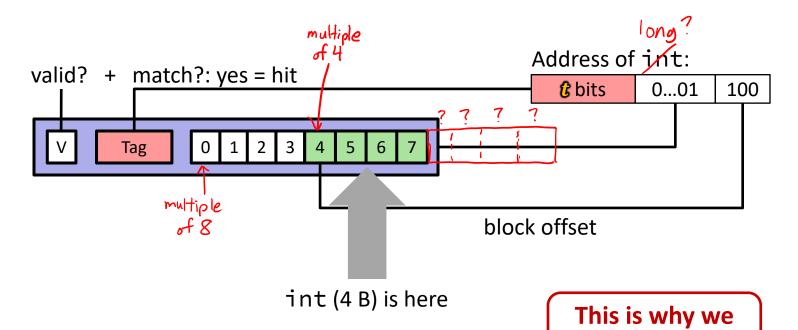
Block Size K = 8 B



# Example: Direct-Mapped (E = 1) Cache (3/3)

Direct-mapped: One line per set

Block Size K = 8 B

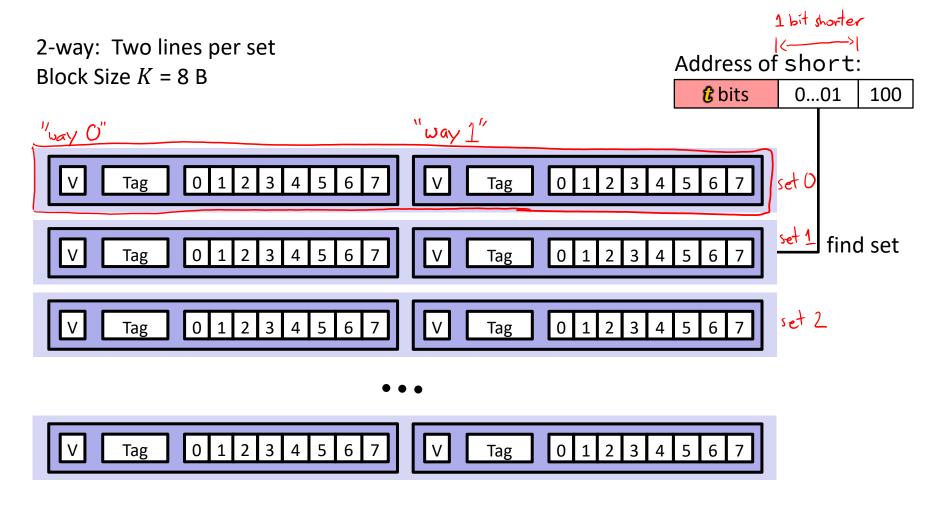


No match? Then old line gets evicted and replaced

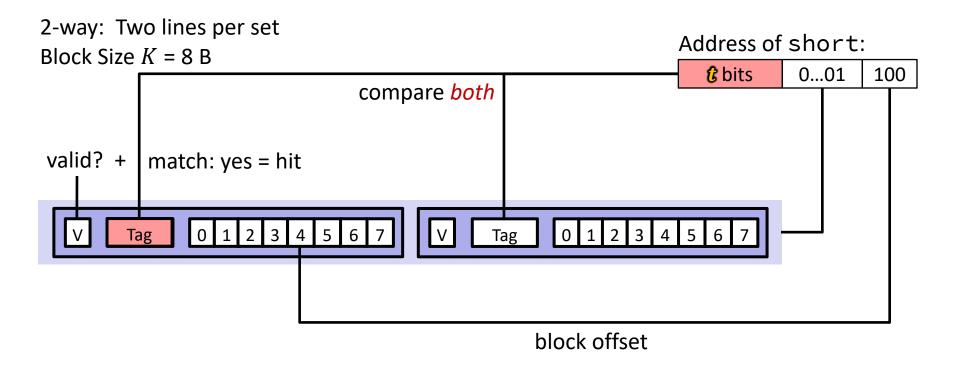
no unnecessary extra cache accesses across block boundaries

want alignment!

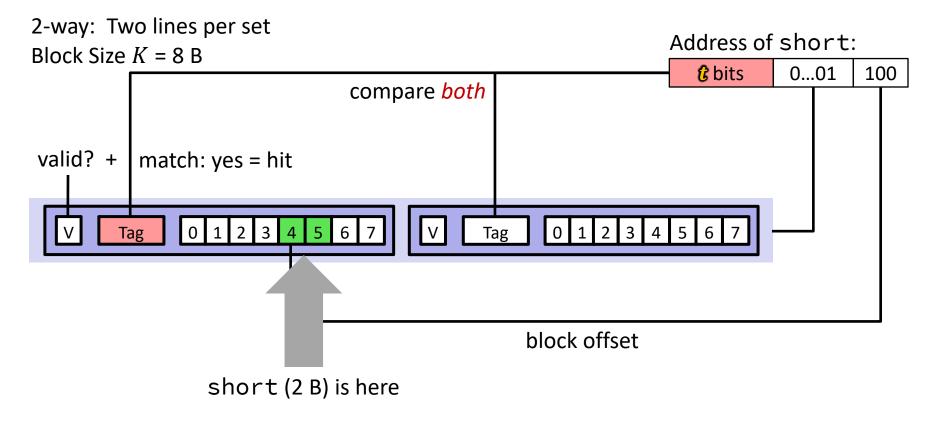
## Example: Set-Associative (E = 2) Cache (1/?)



# Example: Set-Associative (E = 2) Cache (2/3)



# Example: Set-Associative (E = 2) Cache (3/3)



### No match?

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

## Lecture Outline (2/3)

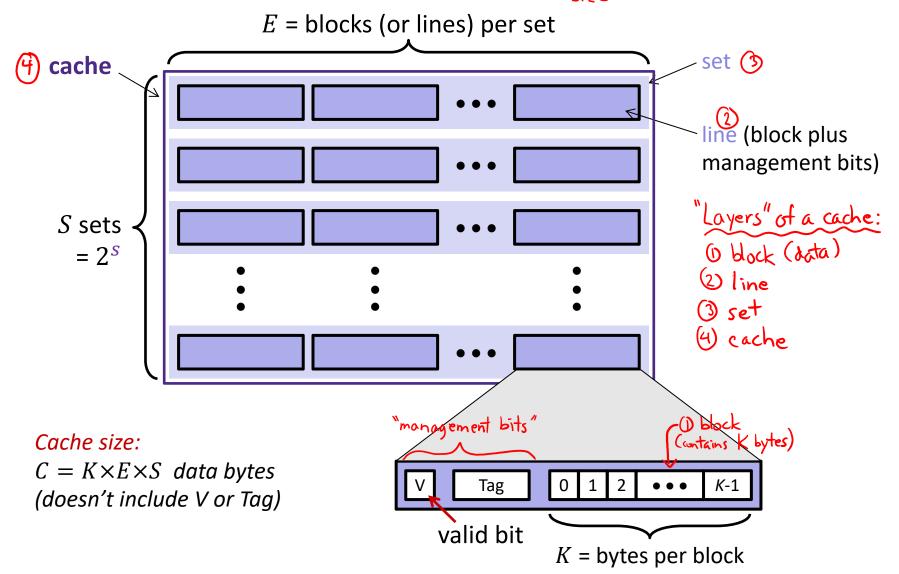
- Associativity and Replacement
- Cache Organization
- Example Cache Problems

## **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

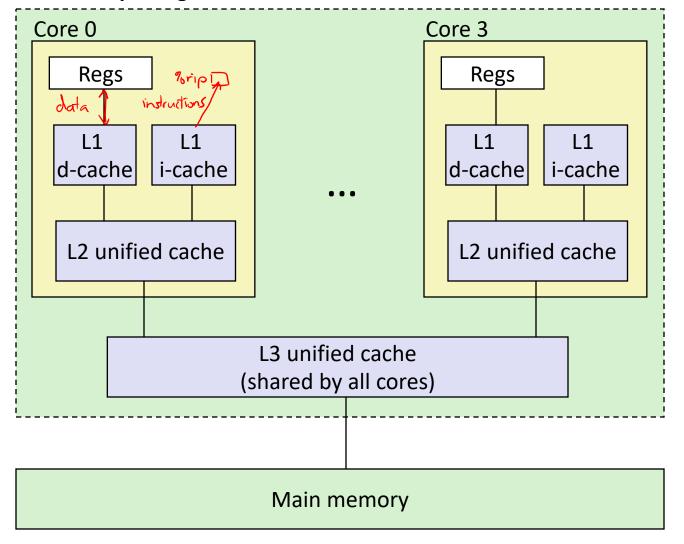
Parameter	Variable	Formulas
Block size	K (B in book)	
Cache size	С	$M = 2^m \leftrightarrow m = \log_2 M$
Associativity	E	$S = 2^{s} \leftrightarrow s = \log_{2} S$
Number of Sets	S	$K = 2^{k} \leftrightarrow k = \log_2 K$
Address space	M	$C = K \times E \times S$
Address width	m	$\mathbf{s} = \log_2(C/K/E)$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

# General Cache Organization (S, E, K)



## **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

## **Learning About Your Machine**

### \* Linux:

- Iscpu
- Is /sys/devices/system/cpu/cpu0/cache/index0/
  - <u>Example</u>: cat /sys/devices/system/cpu/cpu0/cache/index\*/size

## Windows:

- wmic memcache get <query> (all values in KB)
- Example: wmic memcache get MaxCacheSize
- Modern processor specs: <a href="http://www.7-cpu.com/">http://www.7-cpu.com/</a>

## Lecture Outline (3/3)

- Associativity and Replacement
- Cache Organization
- **\* Example Cache Problems**

## **Example Cache Parameters Problem**

Fill in the following table:

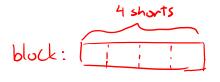
		_
Cache Size	64 B	26
Block Size	8 B	23
Associativity	2-way	2'
Hit Time	3 cycles	
Miss Rate	20%	
Tag Bits	5	
<b>Index Bits</b>	2	2 /23
Offset Bits	3	
AMAT	3+().2(125)= 28 clock cycles	
	Block Size Associativity Hit Time Miss Rate Tag Bits Index Bits Offset Bits	Block Size 8 B  Associativity 2-way  Hit Time 3 cycles  Miss Rate 20%  Tag Bits 5  Index Bits 2  Offset Bits 3

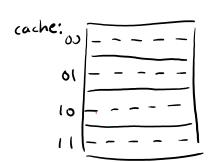
	8 B
set	<u>~</u>
0	
1	
2	
1	

## **Example Code Analysis Problem**

- \* Assuming the cache starts <u>cold</u> (all blocks invalid) and sum, i, and j are stored in registers, calculate the **miss rate**:
  - m = 10 bits, C = 64 B, K = 8 B, E = 2 to bits, S = 2 bits, K = 3 bits

```
#define SIZE 8
short ar[SIZE][SIZE], sum = 0; // &ar=0x200
for (int i = 0; i < SIZE; i++)
for (int j = 0; j < SIZE; j++)
sum += ar[j][i];
```





```
ar [0] [0] address Ob 10 0000 0000 -> M (14 way)

ar [1] [0] -> Ob 10 0000 0000 -> M (14 way)

ar [3] [0] -> Ob 10 0000 0000 -> M (2nd way)

ar [3] [0] -> Ob 10 0000 0000 -> M (replacement)

ar [0] [1] -> Ob 10 0000 0000 -> M (conflict!)

jumping by a

row skips two block
```

## **Cache Simulator**

- https://courses.cs.washington.edu/courses/cse351/cachesim/
  - From course website: Simulators → Cache Simulator
  - Allows you to play around with the effects of cache parameters and policies
  - Lots of neat features like highlighting, hover text, ability to rewind and replay accesses, and copy-and-paste access patterns

## Ways to use:

- Take advantage of "explain mode" and navigable history to test your own hypotheses and answer your own questions
- A <u>tutorial Ed lesson</u> is available
- Will be used in HW19 Lab 4 Preparation

## **Cache Simulator Demo**

- https://courses.cs.washington.edu/courses/cse351/cachesim/
  - From course website: Simulators → Cache Simulator
  - Allows you to play around with the effects of cache parameters and policies
  - Lots of neat features like highlighting, hover text, ability to rewind and replay accesses, and copy-and-paste access patterns
- Let's simulate the example code analysis problem:

```
#define SIZE 8
short ar[SIZE][SIZE], sum = 0; // &ar=0x200
for (int i = 0; i < SIZE; i++)
   for (int j = 0; j < SIZE; j++)
     sum += ar[j][i];</pre>
```

# **Summary (1/2)**

- Associativity gives us flexibility in where to place blocks in the cache
  - Group E slots into sets, means there are E ways to place block within each set
    - Direct-mapped is E = 1, fully associative is E = # of slots in cache (i.e., S = 1)
    - Helps avoid conflicts in each set at the expense of slightly longer and more complex searching & placing in the cache
    - By default, we will replace the **least-recently used** block in a set
  - Index  $\rightarrow$  Set, S = (C/K)/E, still  $s = \log_2(S)$

**Direct-mapped** 

 Set
 Tag
 Data

 0
 ...

 1
 ...

 2
 ...

 3
 ...

 4
 ...

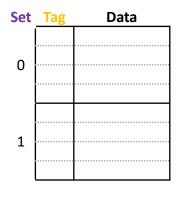
 5
 ...

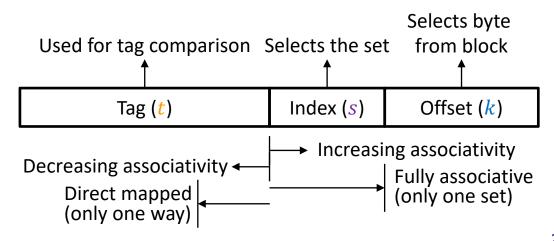
 6
 ...

 7
 ...

2-way set associative 4-way set associative

Set	Tag	Data	
0			
1			
2			
3			





CSE351, Autumn 2025

# **Summary (2/2)**

- Management bits
  - Information needed for proper management of the cache & its data, but not counted in the cache size
  - Valid bit for validity of data
  - Tag bits for identifying which block

