# The Hardware/Software Interface

Memory & Caches II

#### **Instructors:**

Amber Hu, Justin Hsia

#### **Teaching Assistants:**

Anthony Mangus

**Grace Zhou** 

Jiuyang Lyu

Kurt Gu

Mendel Carroll

Naama Amiel

Rose Maresh

**Violet Monserate** 

Divya Ramu

Jessie Sun

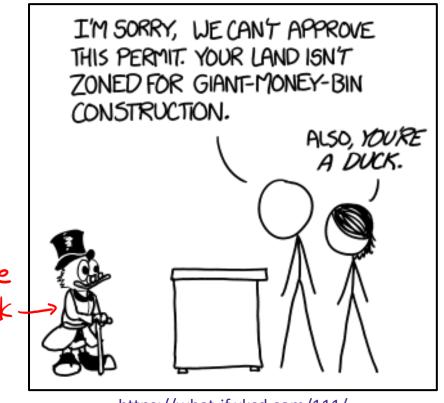
Kanishka Singh

Liander Rainbolt

Ming Yan

Pollux Chen

Soham Bhosale



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

#### **Relevant Course Information**

- HW 16 due Wednesday (11/5)
- HW 17 due Friday (11/7)
  - Don't wait too long, this is a big HW (includes this lecture)
- Lab 3 due Friday (11/7)
  - Late days open until Monday (11/10)
  - Monday is Veteran's Day
  - No lecture, but some office hours (see Ed)
- Midterm grades will be released when we can
  - Regrade requests will be available afterward

#### Lecture Outline (1/2)

- Blocks and Addresses
- Direct-Mapped Caches

#### **Cache Sizes (Review)**

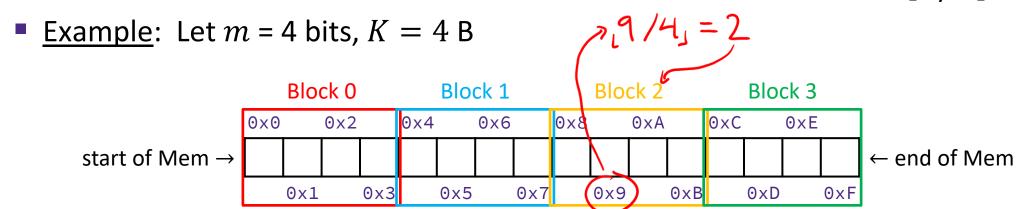
**Note:** The textbook uses "B" for block size

- $\bullet$  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 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)
  - Example: C = 32 KiB using 64-B blocks means that it can hold 512 blocks  $2^{5} \times 2^{10} = 2^{15} \times \frac{16}{2^6} \times \frac{$

#### Memory and Blocks (Review)

**Note:** The textbook uses "B" for block size

- $\bullet$  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!
- The address space is divided into adjacent, numbered blocks
  - For address A of width m bits, can determine its block number via [A/K]



### **Addresses and Blocks (Review)**

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

- $\bullet$  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!

W UNIVERSITY of WASHINGTON

$$\times$$
 %  $2^n = \text{value of the lowest } n \text{ bits}$ 

Ob ...  $\frac{1}{2^{n+1}} = \frac{1}{2^{n-1}} = \frac{1}{2^{n-1}} = \frac{1}{2^{n-1}}$ 

- Remaining address bits tell you ordered position of byte within the block
  - (address) modulo (# of bytes in a block) = A % K
  - Offset field is the low-order  $log_2(K) = k$  bits of address
    - $mod 2^n = n lowest bits$



m-k bits k bits m-bit address: Block Number Block Offset (refers to byte in memory)

#### **Addresses and Blocks Example**

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

- $\bullet$  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!

#### Example:

• If we have 6-bit addresses and block size K = 4 B, which block and byte does 0x19 refer to?

address: Ob 
$$O \perp 1$$
  $O / O \perp 1$  block number 6:

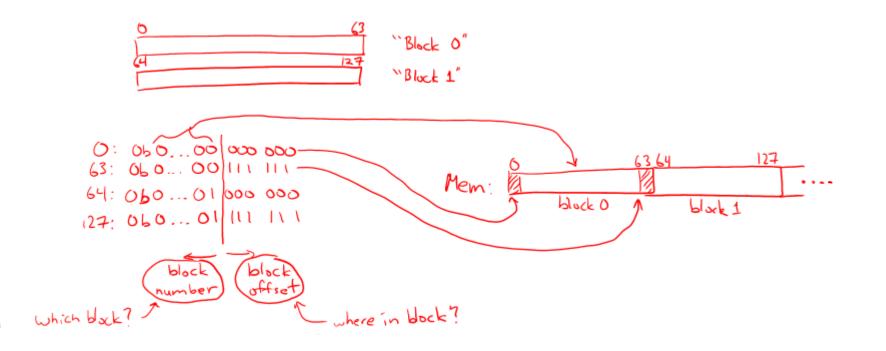
walne 6) (value 1) offset: 00 01 10 11

offset width =  $log_2(K) = log_2(4) = 2 log_2$ 

#### Lab 1a Callback

**Note:** The textbook uses "B" for block size

- The within\_same\_block function asks you to determine if the addresses stored by two pointers lie within the same block of 64-byte aligned memory
  - Solution: Right shift by 6 and compare



# Polling Questions (1/2)

- We have a cache with the following parameters:
  - Block size of 8 bytes  $K = 2^3 B$
  - Cache size of 4 KiB  $C = 2^{12} B$ 22 1 2210
- \* How many blocks can the cache hold?  $C/K = 2^{12.3} = 2^9 = \sqrt{5/2}$  blocks
- ♦ How many bits wide is the block offset field? 

  k=log<sub>2</sub>(k)=3 bits
- Which of the following addresses would fall under block number 3?

B. 
$$0x1F$$
 C.  $0x30$  D.  $0x38$ 
 $31/8 = 3$   $48/8 = 6$   $56/8 = 7$ 

#### Lecture Outline (2/2)

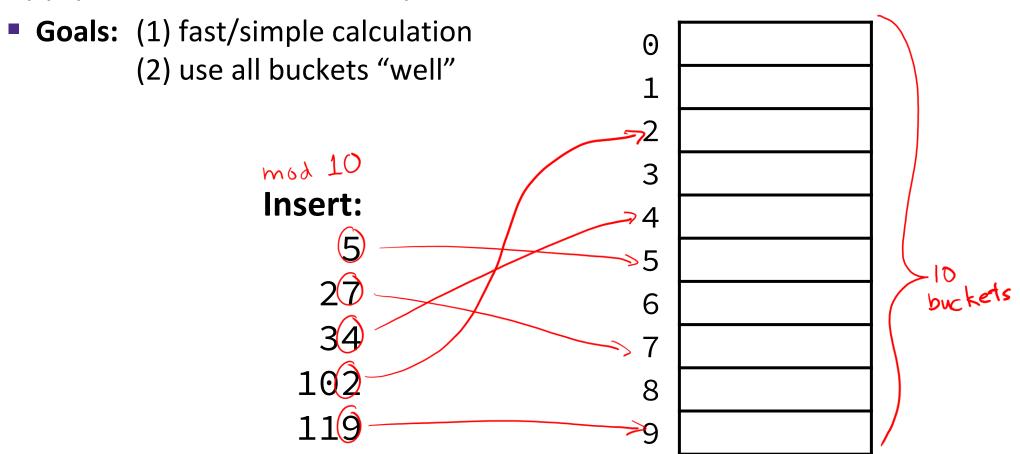
- Blocks and Addresses
- Direct-Mapped Caches

#### **Block Placement (Review)**

- 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!

#### **Hash Tables for Fast Lookup**

Apply hash function to map data to "buckets"



0110

0111

1000

1001

1010

1011

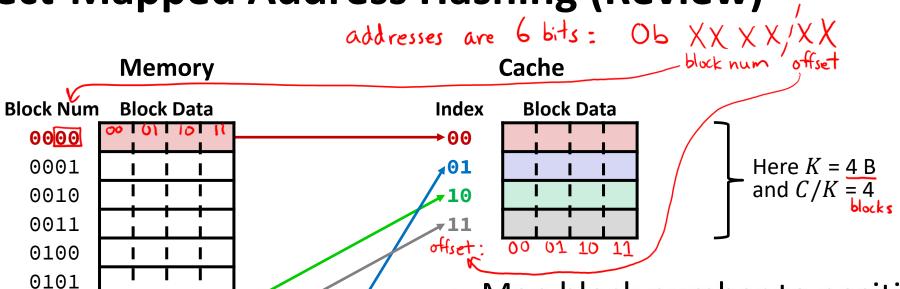
1100

1101

1110

1111

**Direct-Mapped Address Hashing (Review)** 



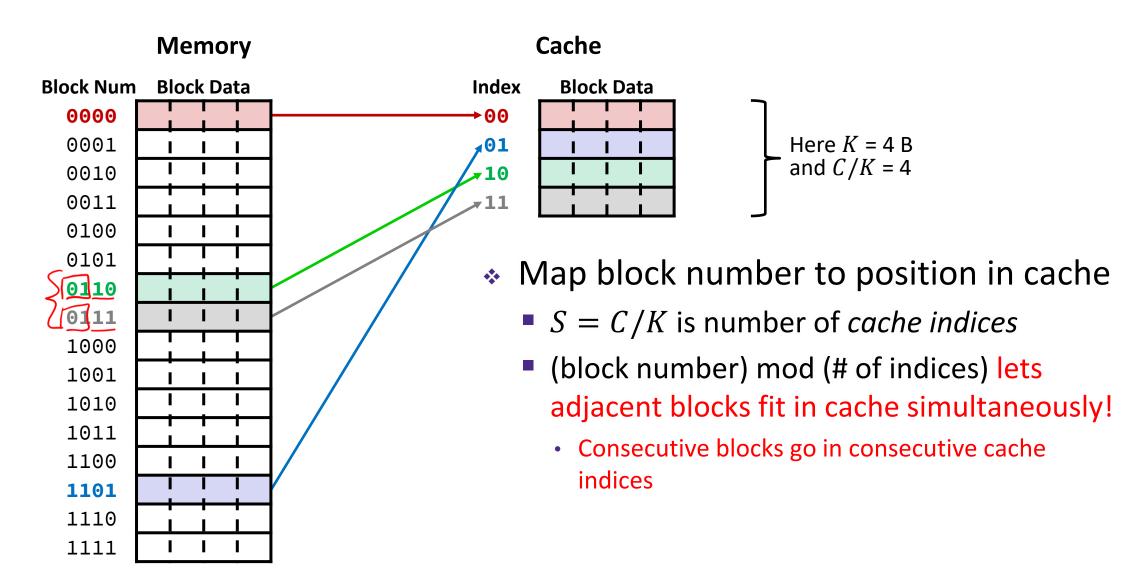
- Map block number to position in cache
  - S = C/K is number of cache indices
  - (block number) mod (# of indices)

• Index field is the low-order  $\log_2(S) = s$  bits

of the block number

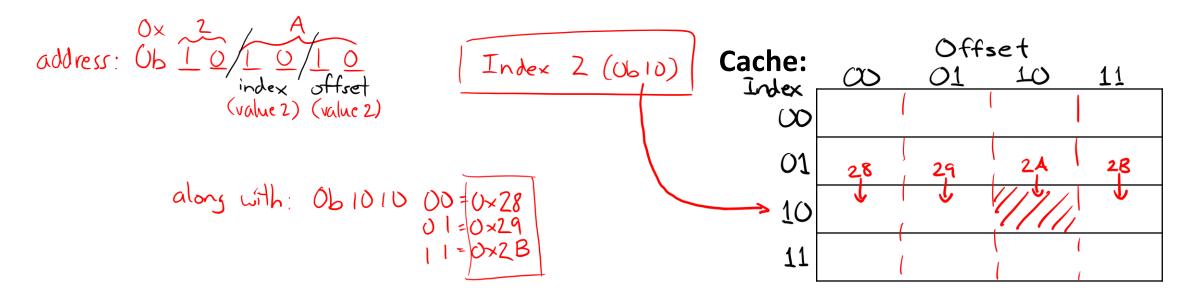
How many bits do I ) need to specify every set/index in my cache?

#### **Direct-Mapped Address Hashing Consequence**

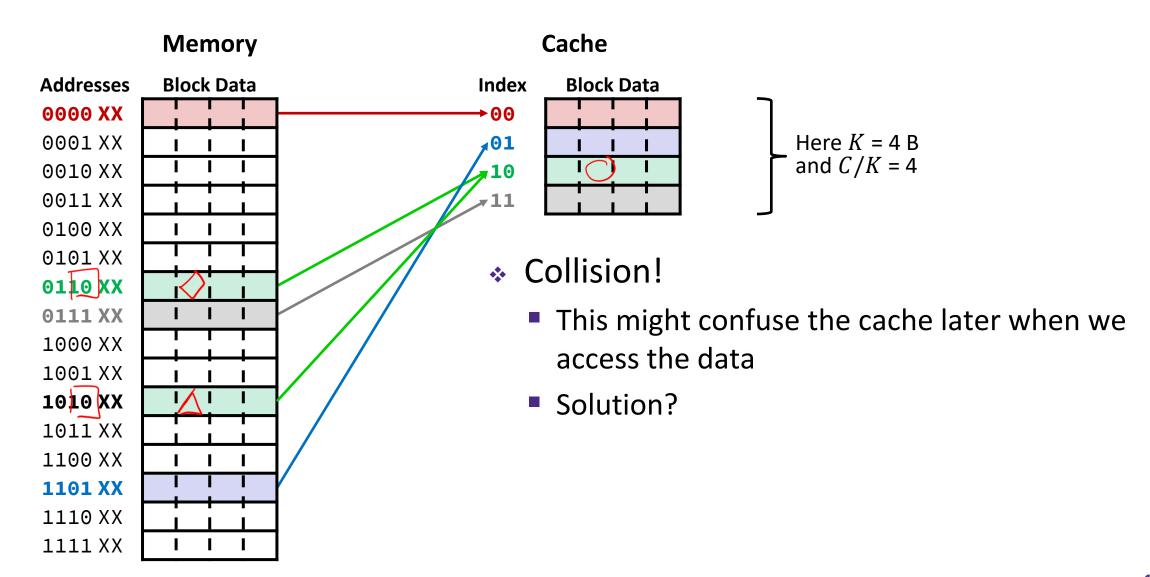


#### **Block Placement Example**

- \* 6-bit addresses, block size K = 4 B, and our cache holds S = 4 blocks. C/K ,  $S = log_2(4) = 2$  bits
- A request for address 0x2A results in a cache miss. Which index does this block get loaded into and which 3 other addresses are loaded along with it?



#### Place Data in Cache by Hashing Address



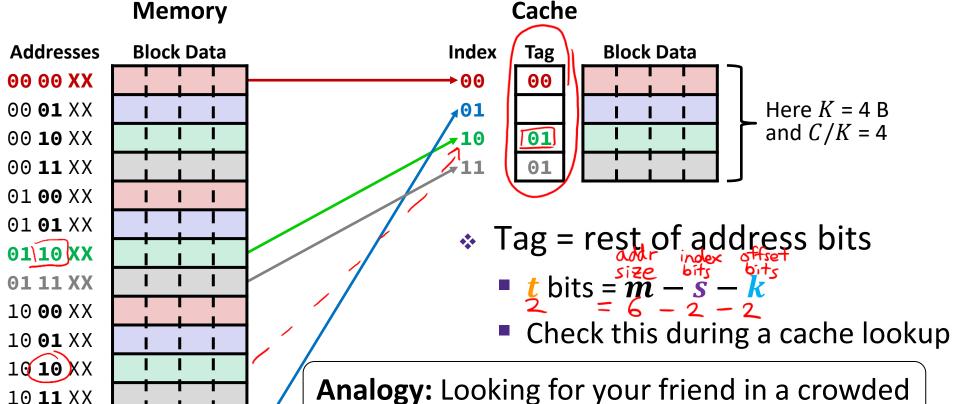
11 **00** XX

11 01 XX

11 **10** XX

11 **11** XX

#### **Tags Differentiate Blocks in Same Index**



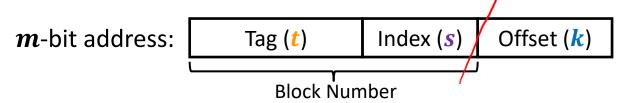
**Analogy:** Looking for your friend in a crowded room, so you yell their first name (index)

Several people respond "yeah?", so you get more specific: first name + last name (tag)



#### **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:



- 1) Index field tells you where to look in cache
- 2) Tag field lets you check that data is the block you want
- 3) Offset field selects specified start byte within block
- Note: t and s sizes will change based on hash function

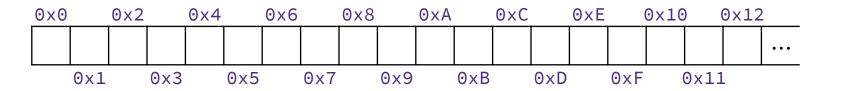
# Polling Questions (2/2)

- Based on the following behavior, which of the following block sizes is
  - NOT possible for our cache? hit: block with data already in \$ miss: data not in \$, pulls block containing data from Mem
  - Cache starts empty, also known as a cold cache
  - Access (addr: hit/miss) stream:
    - (0xE: miss), (0xF: hit), (0x10: miss)

      3 16 is in a different block

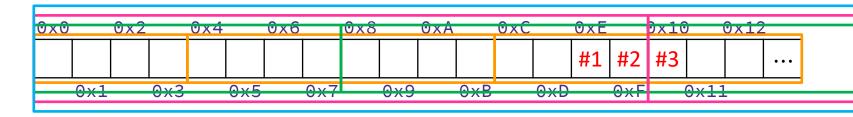
      2 14 \$ 15 are in the same block

      4 pulls block containg 14 into \$
  - A. 4 bytes
  - 8 bytes
  - C. 16 bytes
  - D. 32 bytes

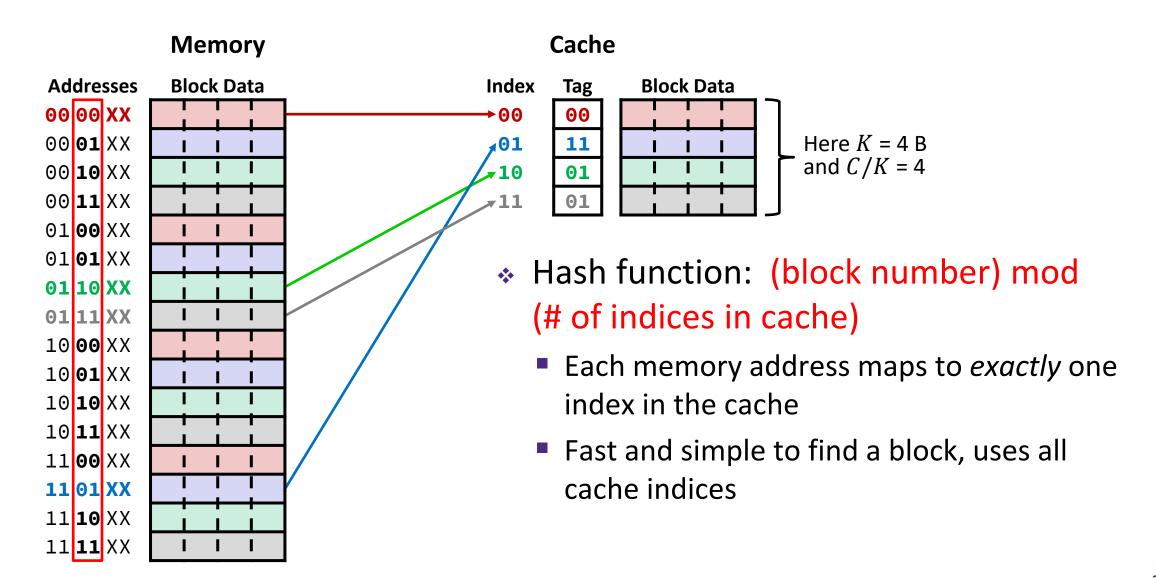


### Polling Questions Solution (2/2)

- 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:
    - (0xE: miss), (0xF: hit), (0x10: miss)
    - Need 0xE and 0xF in same block; 0x10 in different block
  - A. 4 bytes
  - B. 8 bytes
  - C. 16 bytes
  - D. 32 bytes

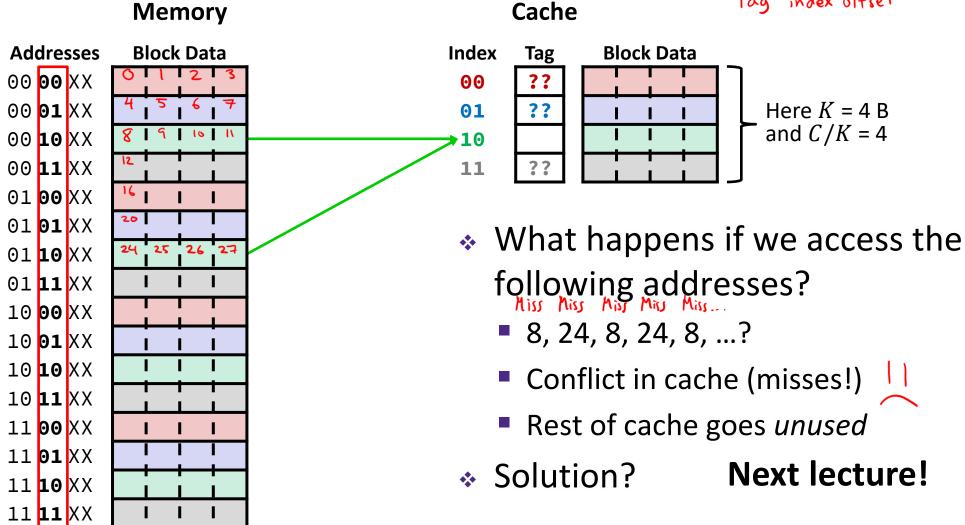


#### **Direct-Mapped Cache Summary**



#### **Direct-Mapped Cache Problem**





#### Homework Setup (1/2)

```
K struct WolfPos {
     float x;
                  struct WolfPos:
    float y;
                                   id
    float z; offset: 0
     int id;
 }; Kmax = 4
 struct WolfPos grid[16][16];
    Assume &grid = 0
  C is row-major:
                 16 columns
```

What are the addresses of the following pieces of data?

• & (grid[0][0].id) = 
$$\frac{12 = 0.00}{0}$$

• & (grid[3][4].x) = 
$$832 = 0x340$$
  
 $646+4)*16=832$   
 $620$   
 $646+4)*16=832$   
 $620$   
 $646+4)*16=832$ 

W UNIVERSITY of WASHINGTON

k grid [0][0].x = 
$$0x0$$
 =  $0b 0000 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0000 | 0$ 

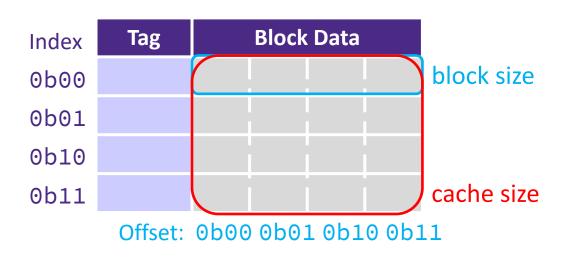
$$S = C/K = cache holds 64 blocks$$
  
 $s = log_2(C/K) = 6 bls$ 

\* Cold direct-mapped cache with C = 1024 B and K = 16 B k= 4 ks

• What happens if we access grid[0][0].x and then grid[4][0].x?

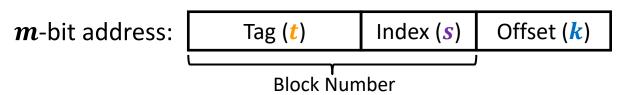
# **Summary (1/2)**

- Cache parameters define the cache geometry:
  - Block size is number of bytes per block
  - Cache size is number of bytes (or blocks) of data the cache can hold
- Finding a byte in the cache:
  - Offset refer to which byte in block
  - Index refers to which block in cache
- Example:
  - K = 4 B, C = 16 B = 4 blocks



# **Summary (2/2)**

- Direct-mapped cache: each block in cache is assigned a unique index
  - Uses hash function of (block number) mod (# of cache indices)
    - Deterministic placement of each block, with many blocks mapping into the same index
    - Tag bits stored in cache and used to distinguish between blocks that map to same index
- Accessing the cache:(TIO address breakdown)



- 1) Index field tells you where to look in cache (width  $s = \log_2 S$ )
- 2) Tag field lets you check that data is the block you want (width t = m s k)
- 3) Offset field selects specified start byte within block (width  $k = \log_2 K$ )