

Optimizing Cache Performance

Today's Lecture

- Use larger cache blocks to take advantage of **spatial locality**
- Use set associativity to resolve “hashing collisions”

For a byte-addressable machine with 16-bit addresses

Which picture best represents a cache that is direct-mapped, each block holds one byte, and has eight cache blocks

A

Block Index	2-bit Tag	8-bit data
000		
001		
010		
011		
100		
101		
110		
111		

B

Block Index	8-bit Tag	8-bit data
000		
001		
010		
011		
100		
101		
110		
111		

C

$16 - 3 = 13$

3 bits

Block Index	13-bit Tag	8-bit data
000		
001		
010		
011		
100		
101		
110		
111		

For a byte-addressable machine with 16-bit addresses

A cache has the following characteristics:

- It is **direct-mapped** (as discussed last time)
- Each block holds **one byte** — 8 bits data
- The cache index is the **four** least significant bits

4-bits index

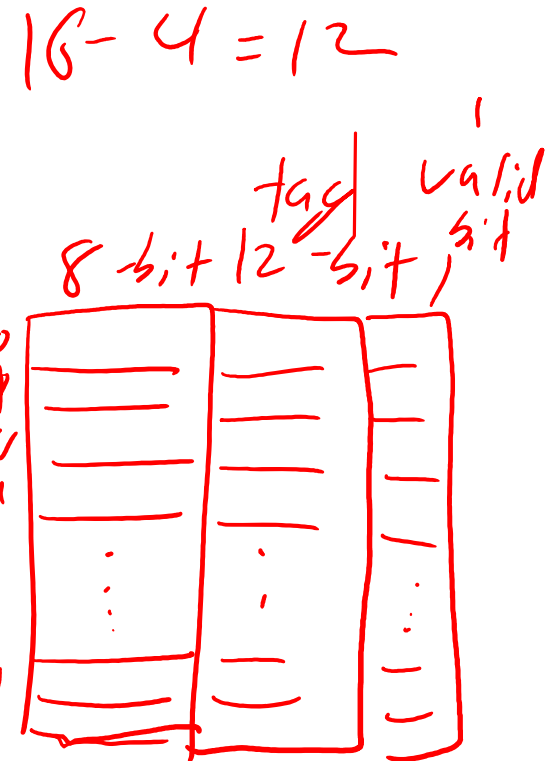
Two questions:

- How many blocks does the cache hold?

- a) 1
- b) 2
- c) 4
- d) 8
- e) 16

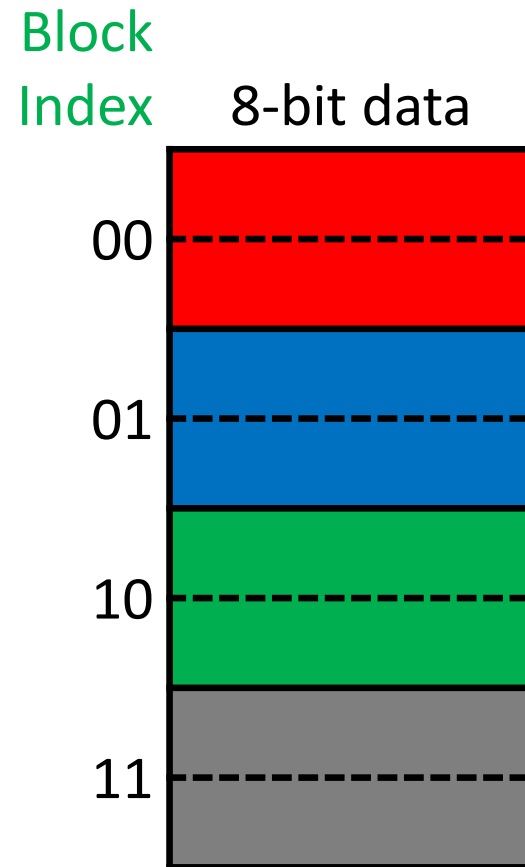
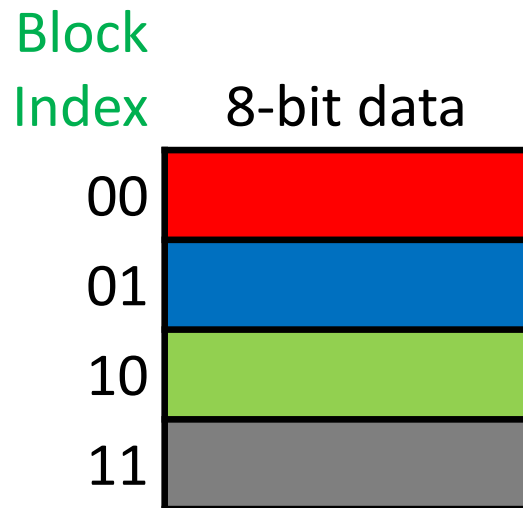
- How many bits are stored at each cache block (e.g., for the data array, tags, etc.)?

- a) 8 b) 9 c) 12 d) 20 e) 21

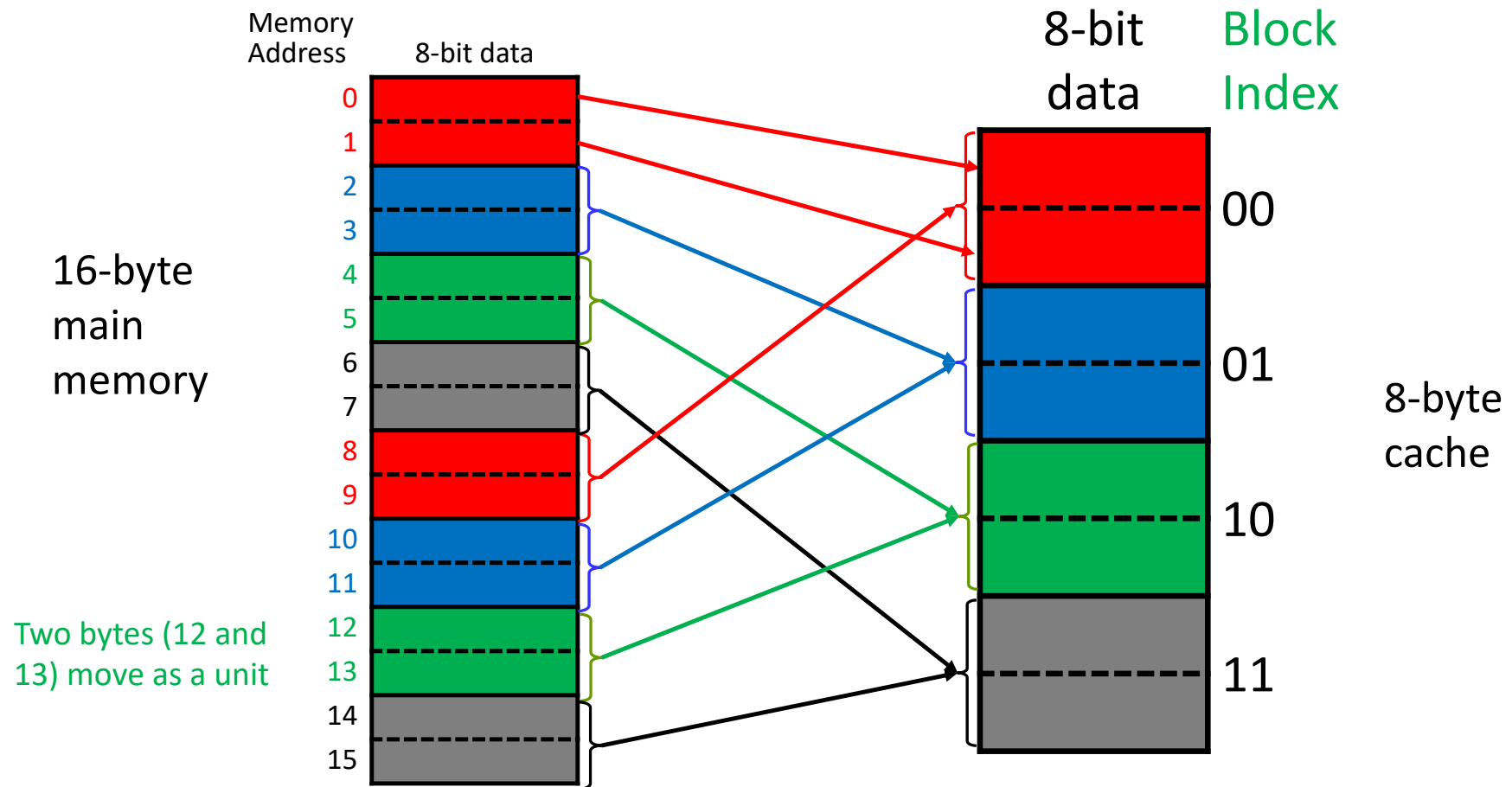


1-byte cache blocks do not take advantage of spatial locality

Create larger cache blocks



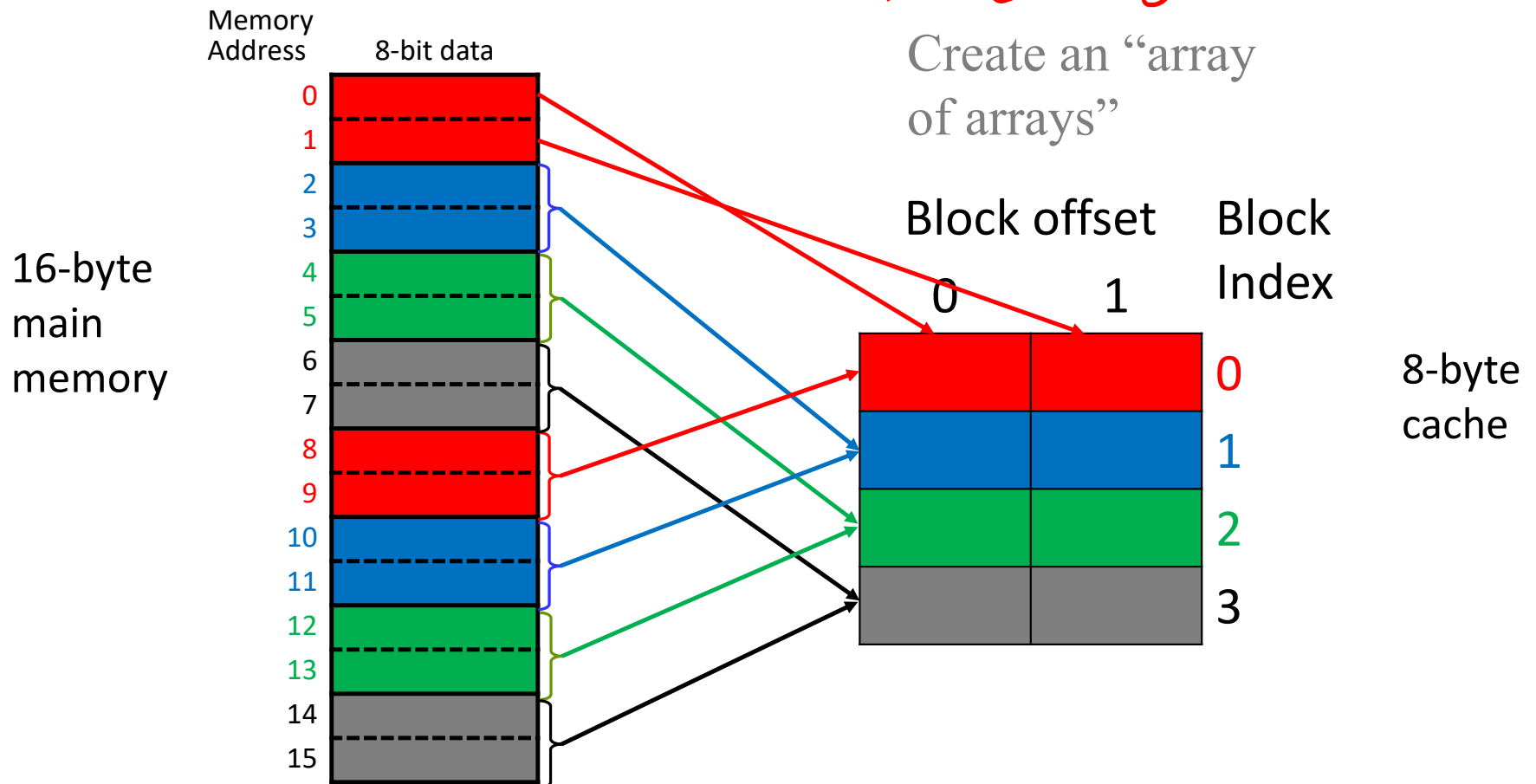
Creating larger cache blocks moves adjacent blocks as a unit



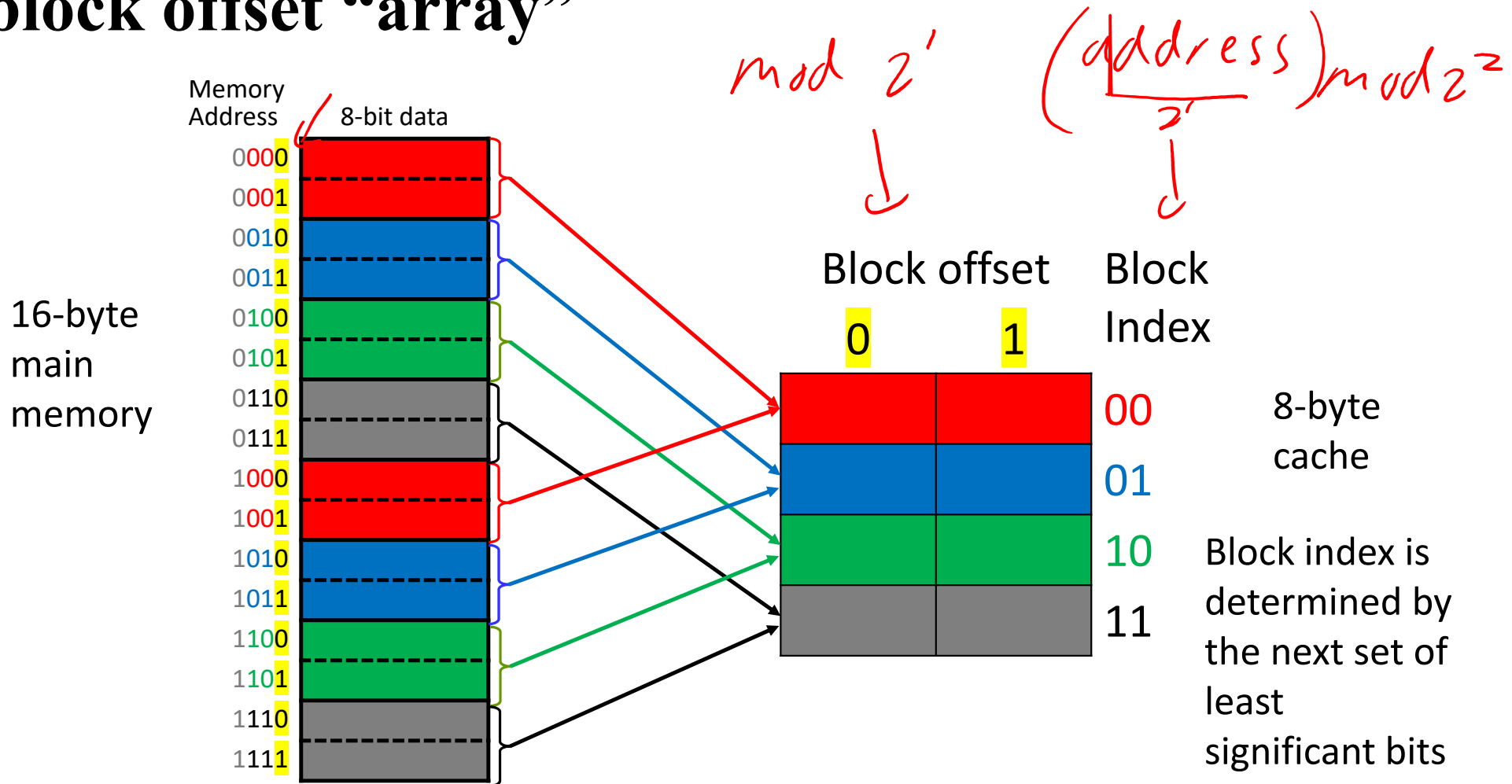
Each cache block is an array indexed by a block offset

$M[b_i][b_o]$

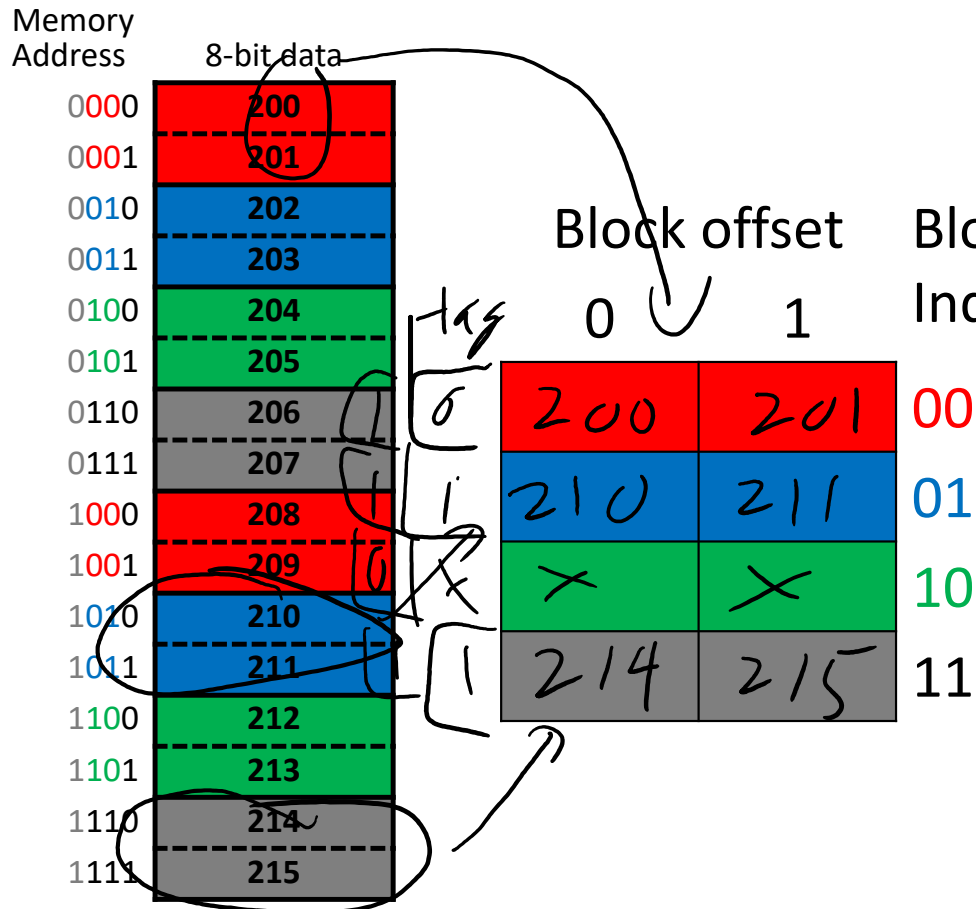
Create an “array of arrays”



The least-significant bit(s) is used to index the block offset “array”



Data gets moved in cache blocks and not bytes



Consider the following set of memory loads

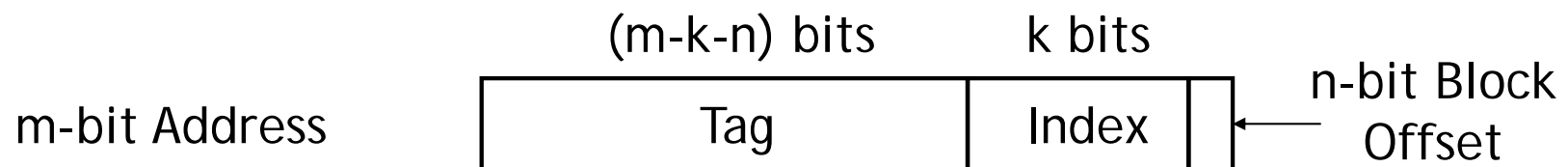
Load 0
Load 15
Load 10
Load 1

What is stored in block index 01, offset 1?

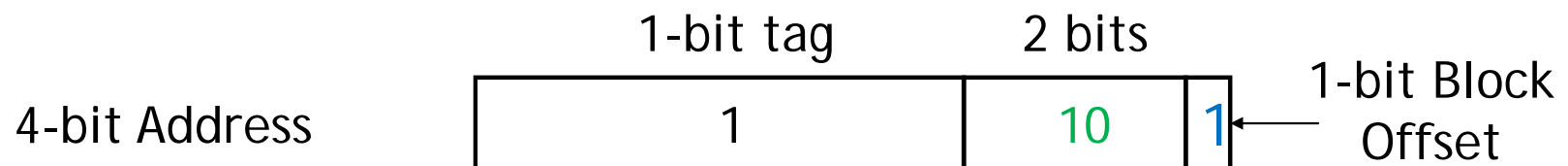
- a) 201
- b) 210
- c) 211
- d) 215
- e) invalid

To increase cache block size, subdivide an m-bit address into tag, index, and block offset

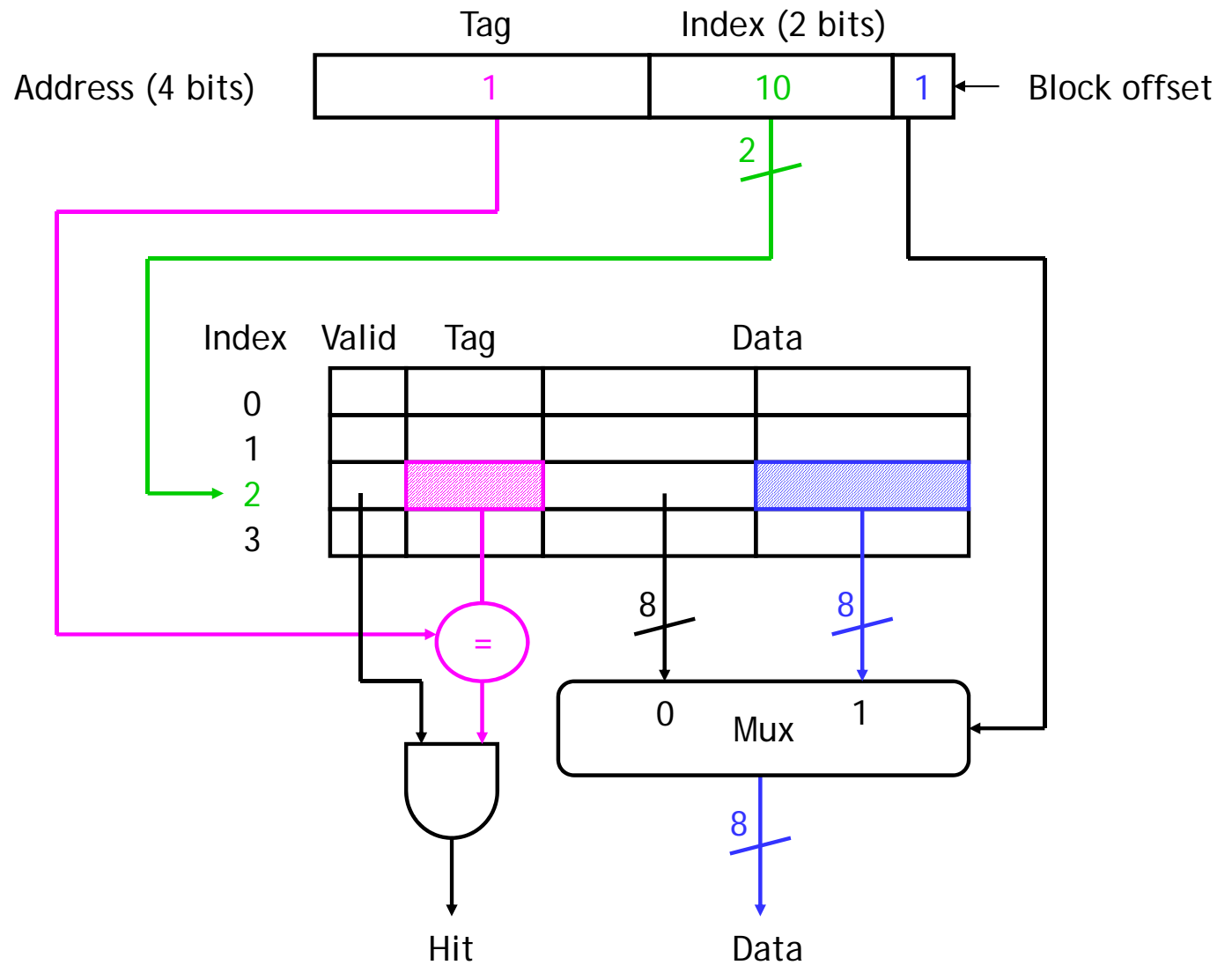
- Suppose we have a cache with 2^k blocks, each containing 2^n bytes.
 - Lowest n bits are the **block offset** that decides which of the 2^n bytes in the cache block will store the data.
 - Next k bits of the address select one of the 2^k cache blocks.



- Example: 2^2 -block cache with 2^1 bytes per block. Memory address 13 (1101) would be stored in offset **1** of cache block **2**.

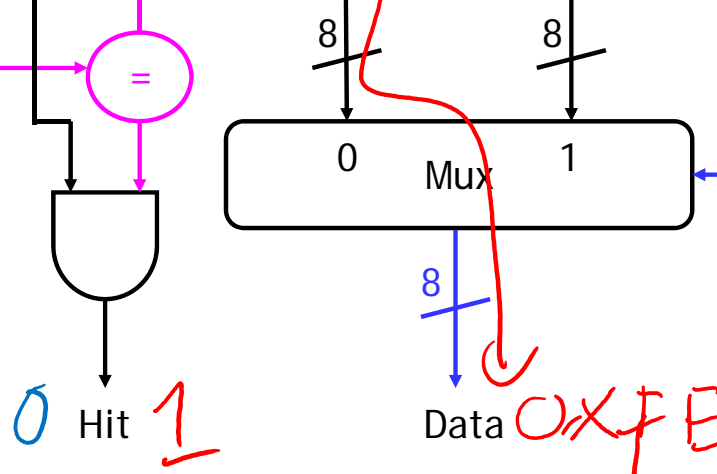


Implement block offset with a multiplexer





Index	Valid	Tag	Data	
0	1	0	0xCA	0xEF
1	1	1	0xFE	0xAD
2	1	0	0xBE	0xDE
3	0	1	0xFE	0xED



For the addresses below, what byte is read from the cache (or is there a miss)?

- 1010 - hit, 0xFE
- 1110 - miss, invalid
- 0001
- 1101

- a) Miss, tag mismatch
- b) Miss, invalid
- c) Hit, 0xDE
- d) Hit, 0xEF
- e) Hit, 0xFE

For a byte-addressable machine with 16-bit addresses

A cache has the following characteristics:

- It is direct-mapped
- Each block holds four bytes 2-bits offset
- The cache has 32 cache blocks 5-bits offset

How many bits are used for the tag?

- a) 2
- b) 4
- c) 5
- d) 7
- e) 9

$$16 - 5 - 2 = 9$$

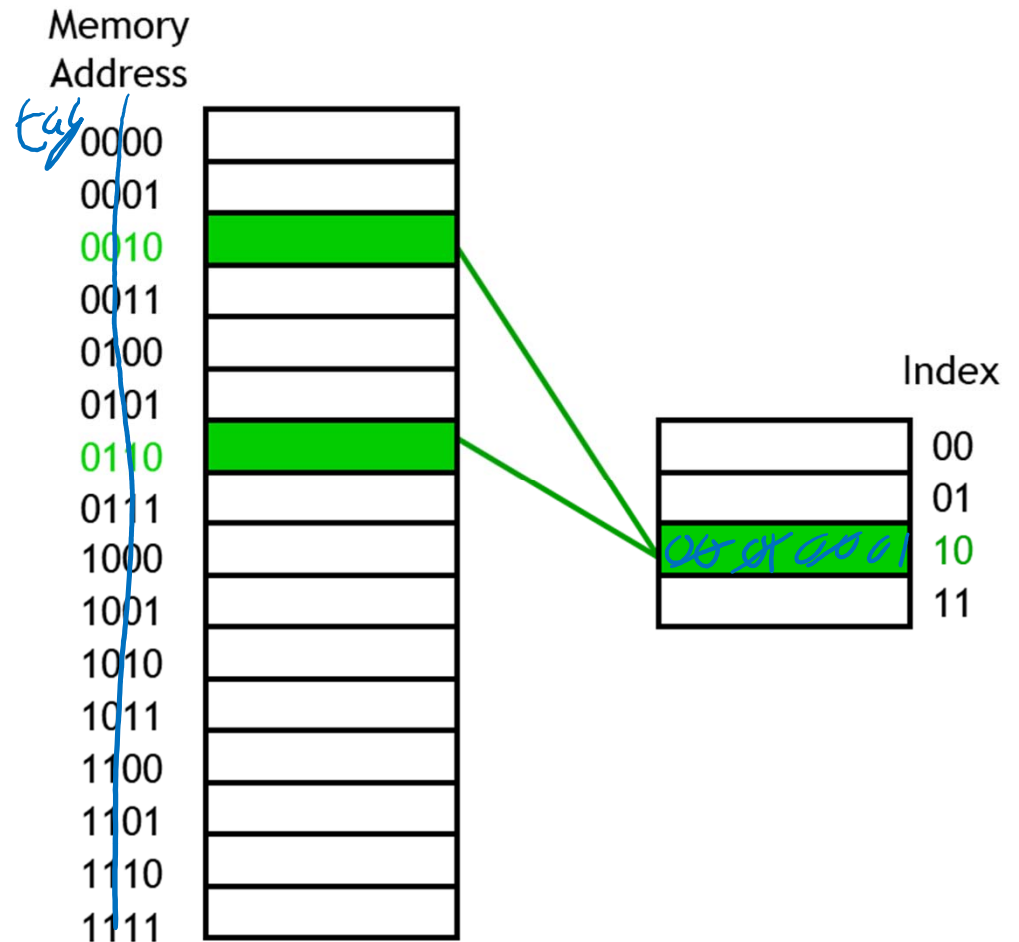
Direct-mapped caches fall short when addresses collide

- Example: what happens if a program uses addresses

2, 6, 2, 6, 2, ...?

M M M M

thrashing



A fully-associative cache permits data to be stored in *any* cache block

- When data is fetched from memory, it can be placed in *any* unused block of the cache.
- Maximize temporal locality by keeping the most recently used data in the cache, replace the **least-recently used (LRU)** when cache is full

A fully associative cache can map addresses anywhere, eliminating thrashing

- Example: what happens if a program uses addresses

2, 6, 2, 6, 2, ...?

M M 1 1 1 1

Memory Address

0000

0001

0010

0011

0100

0101

0110

0111

1000

1001

1010

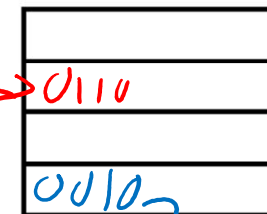
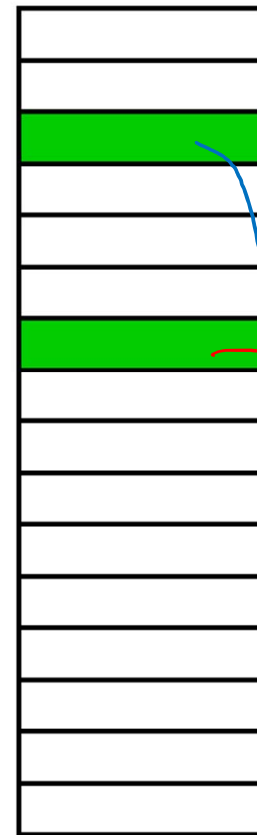
1011

1100

1101

1110

1111



Index

00

01

10

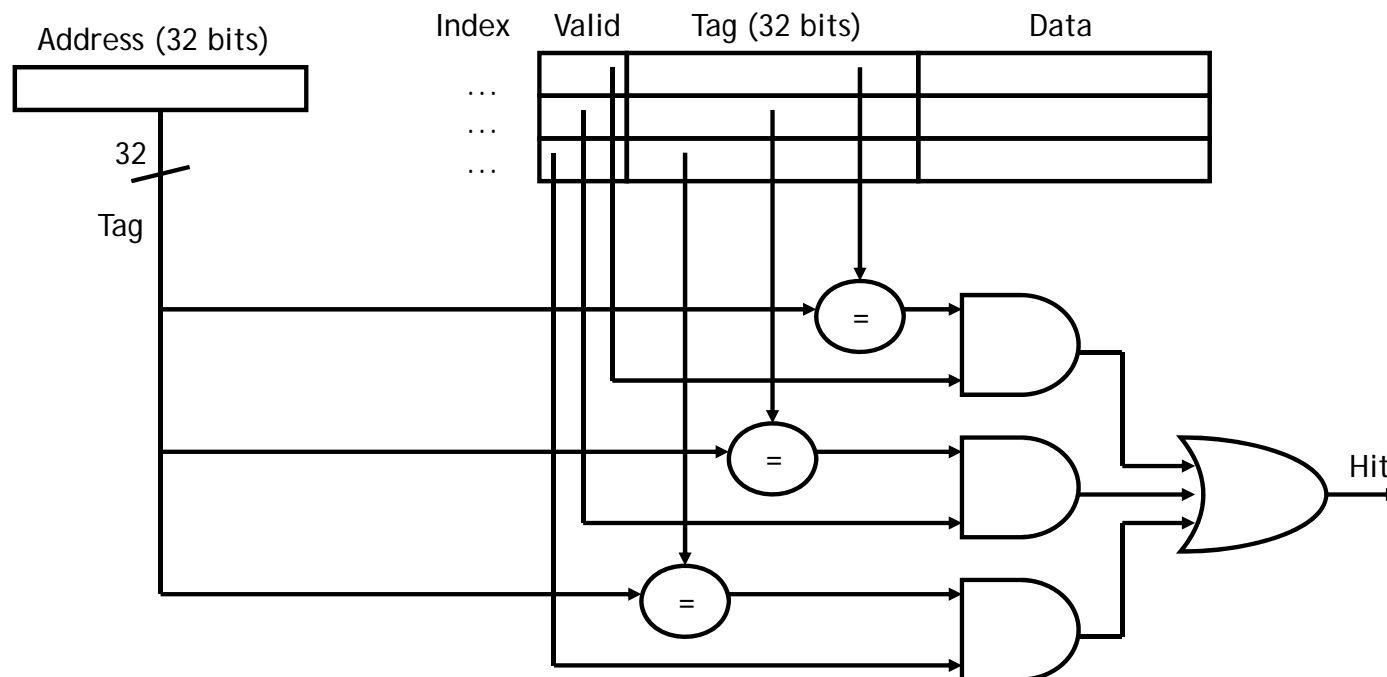
11

LRU

*↑
LRU*

A fully-associative cache is expensive because we need to store the entire tag!

- Data could be anywhere in the cache, so we must check the tag of *every* cache block. That's a lot of comparators!



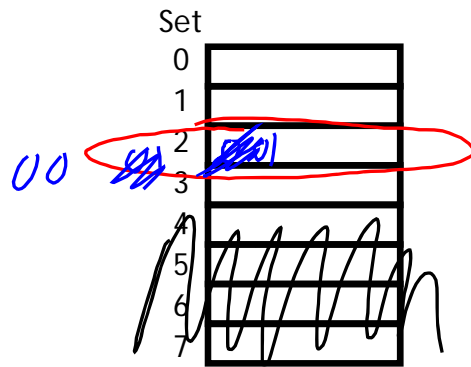
A **set-associative** cache organizes cache blocks into groups called **sets**

- Each memory address maps to exactly one set in the cache, but data may be placed in any block within that set.

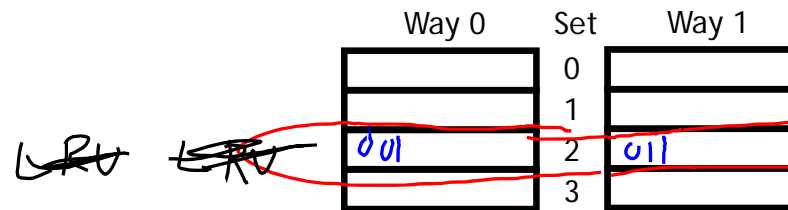
- If each set has 2^x blocks, the cache is a **2^x -way associative cache**.

00 | 10 01 | 10
tag index tag index

direct mapped
8 "sets", 1 block each

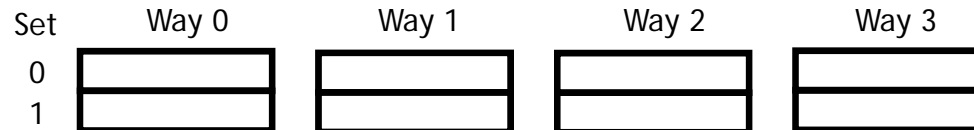


2-way associativity
4 sets, 2 blocks each



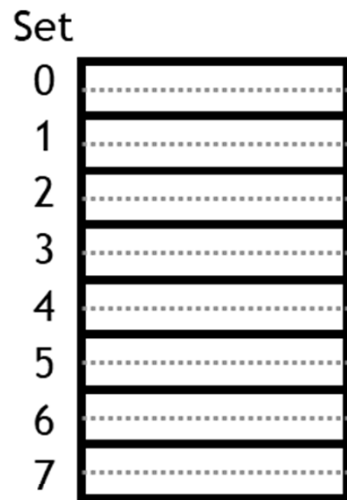
LD 2
MLD 6
1+LD 2
LRU

4-way associativity
2 sets, 4 blocks each

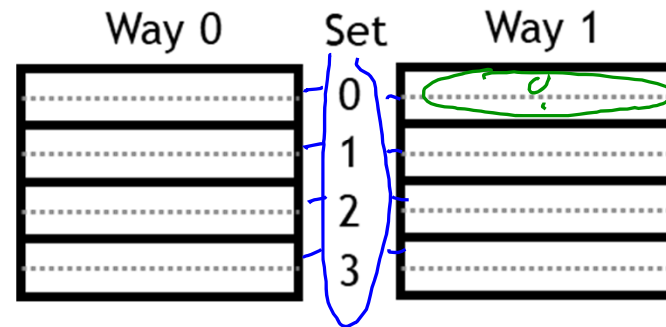


Blocks can still have multiple bytes

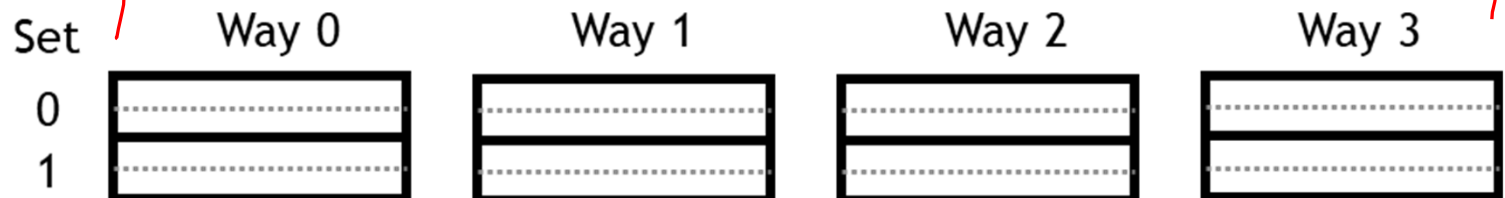
direct mapped
8 "sets", 1 block each



2-way associativity
4 sets, 2 blocks each

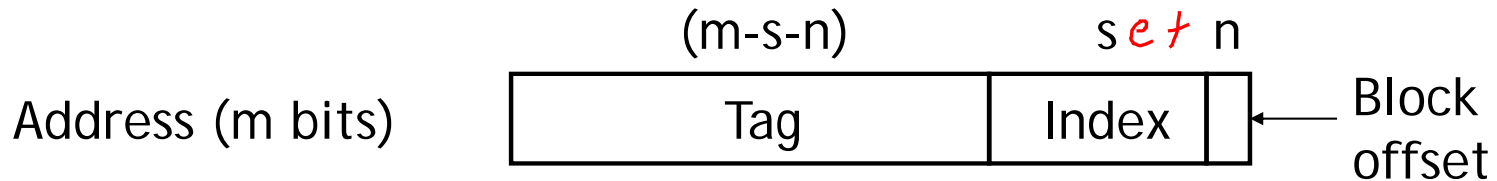


4-way associativity
2 sets, 4 blocks each



To find data, subdivide the address into tag, index, and block offset as before

- Memory has m-bit address
- Cache has 2^s sets and each block has 2^n bytes



- Our arithmetic computations now compute a **set index**, to select a *set* within the cache instead of an individual block.

$$\text{Block Offset} = \text{Memory Address} \bmod 2^n$$

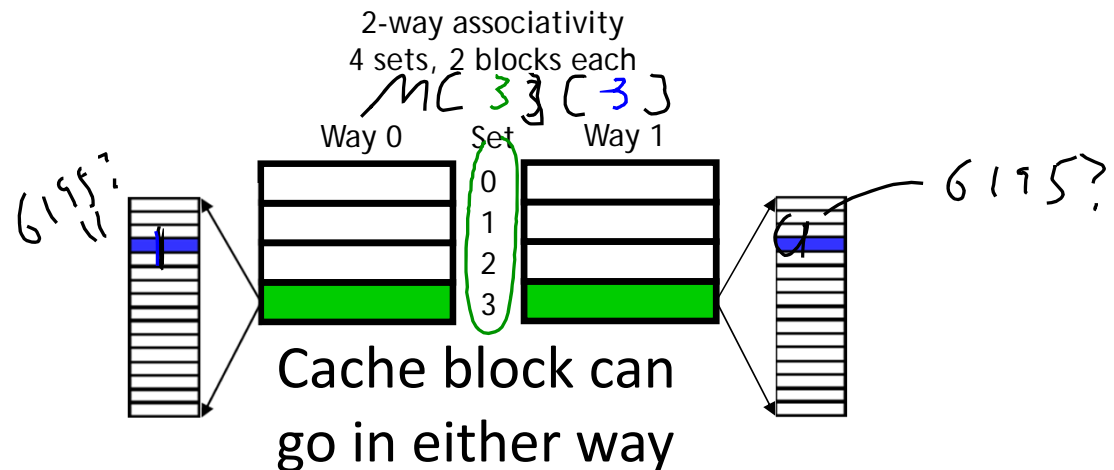
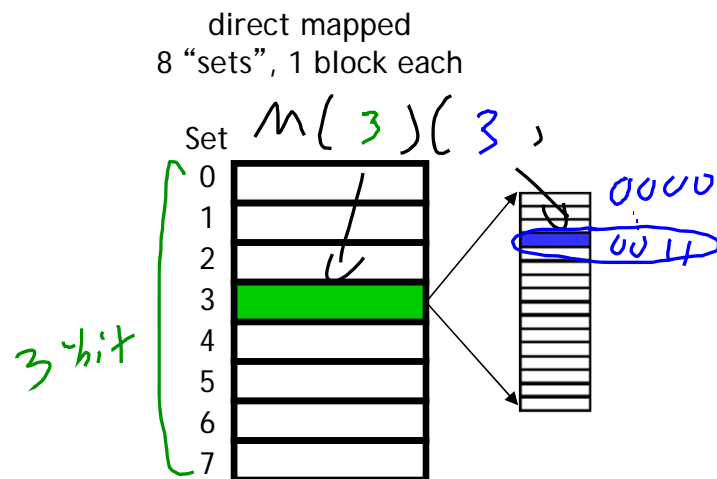
$$\text{Set Index} = (\text{Memory Address} / 2^n) \bmod 2^s$$

Example: Placement of address 6195 in caches with 8, 16-byte cache blocks

- 6195 in binary is 00...0110000|0|11|0011.
- Each block has 16 bytes, so the lowest 4 bits are the block offset.

The next three bits (011) are the set index.

The next two bits (11) are the set index.



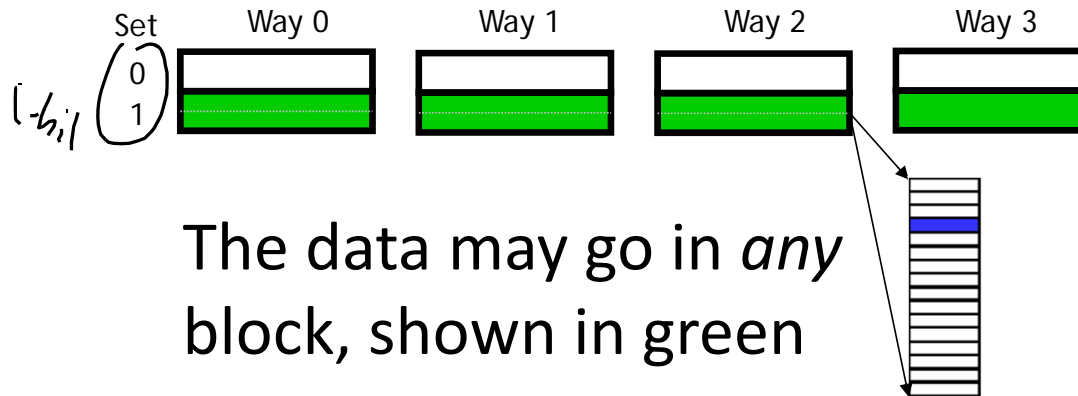
Example: Placement of address 6195 in caches with 8, 16-byte cache blocks

- 6195 in binary is 00...011000001|1|0011.
- Each block has 16 bytes, so the lowest 4 bits are the block offset.

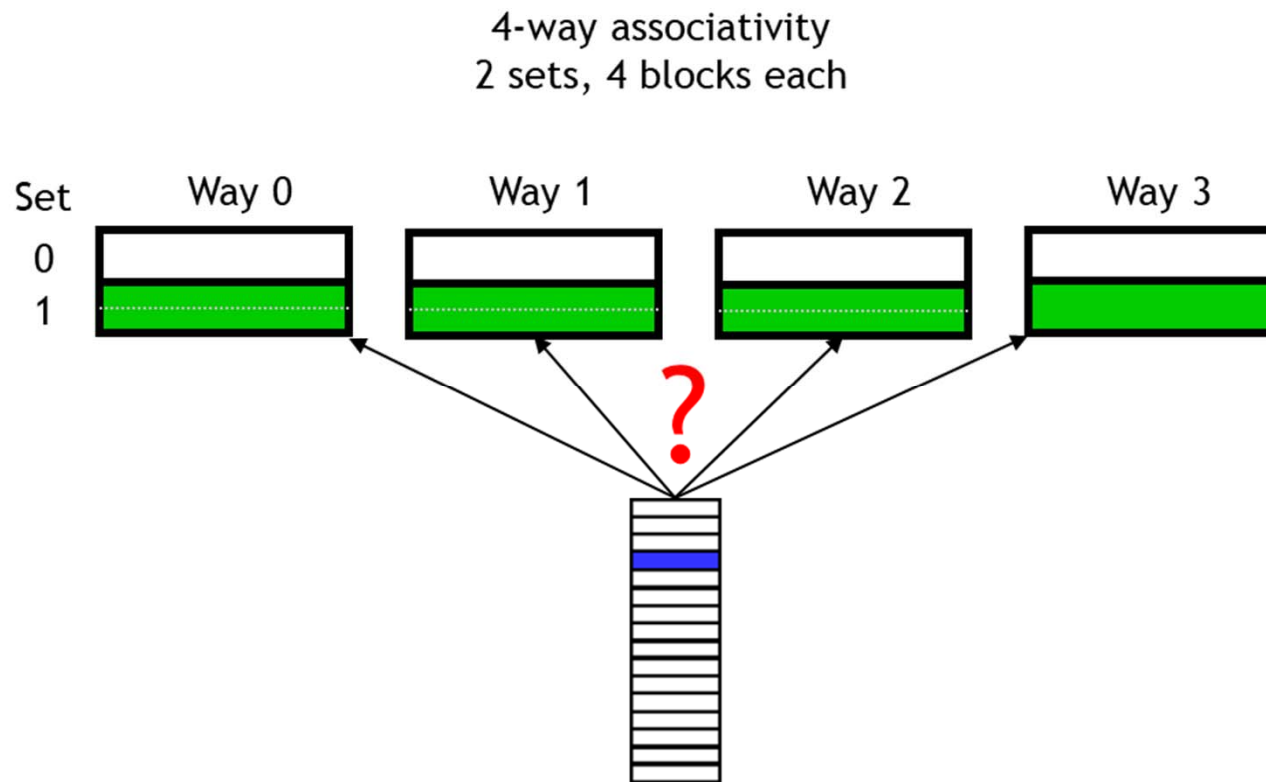
index

the next one bit (1) is
the set index

4-way associativity
2 sets, 4 blocks each



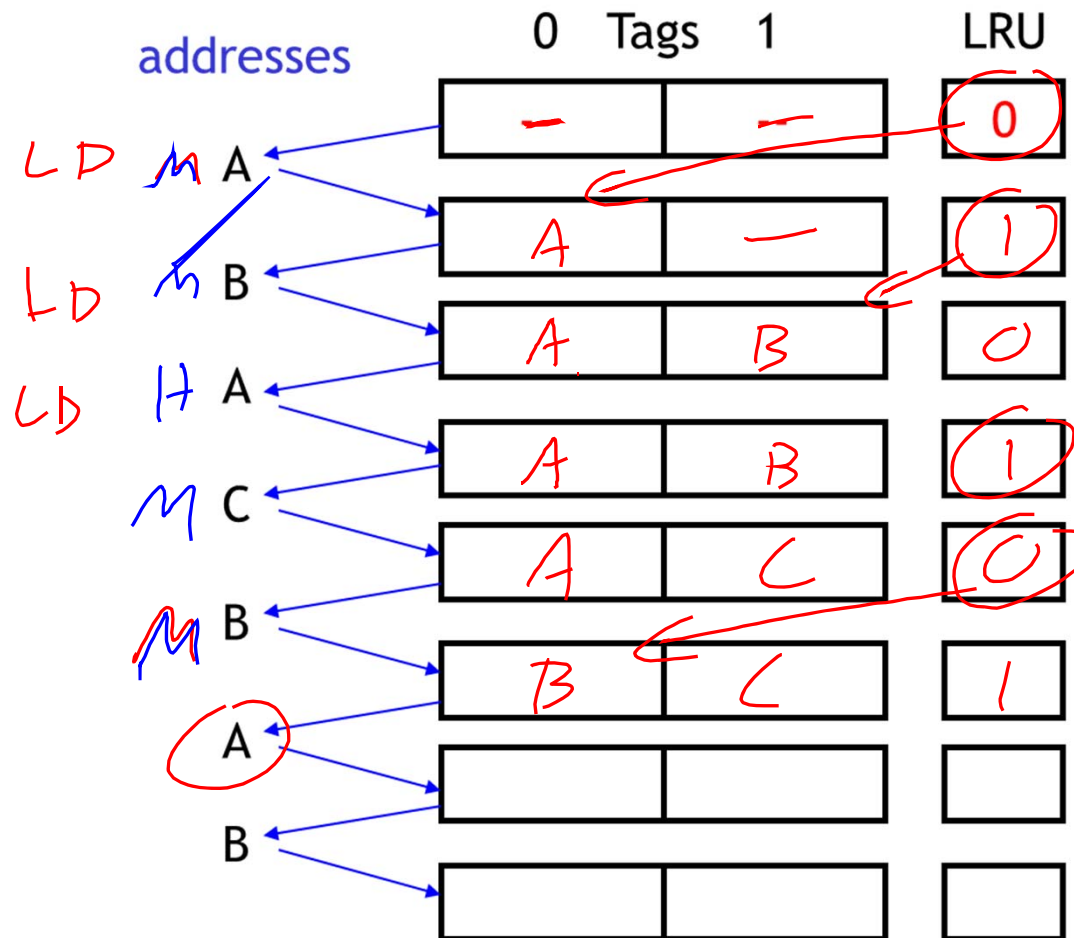
Place new data in empty cache blocks if possible, else, replace the least-recently used



Approximate
LRU for high
associativity

Given a fully-associative cache with two blocks, which memory accesses miss in the cache?

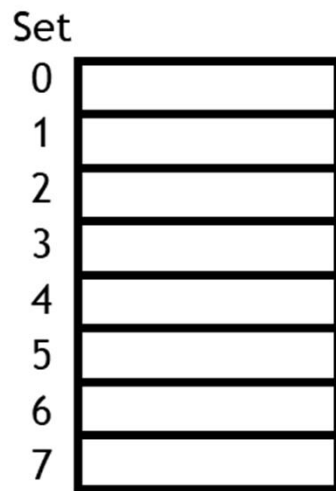
Assume distinct addresses go to distinct blocks



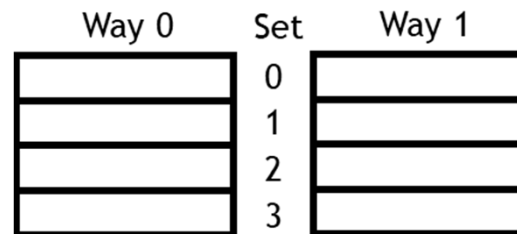
a) Hit
b) Miss

Direct-mapped and fully-associative caches are types of set-associative caches

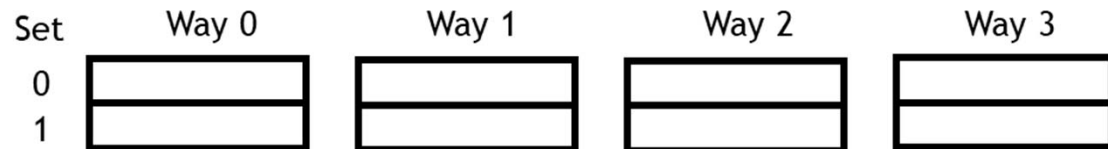
1-way, **direct mapped**
8 “sets”, 1 block each



2-way associativity
4 sets, 2 blocks each



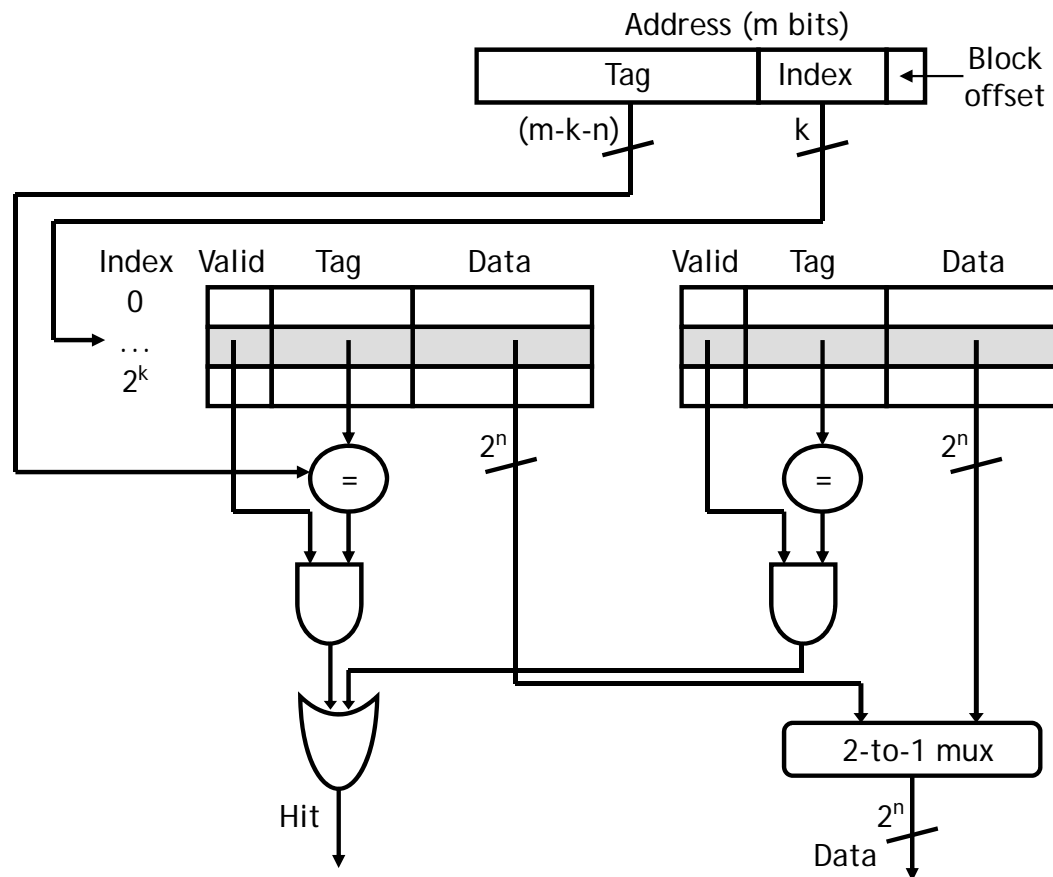
4-way associativity, 2 sets, 4 blocks each



8-way associativity (**fully associative**), 1 set, 8 blocks each



Set associativity is implemented by using a multiplexer to choose from two identical caches



Summary

- Larger **block** sizes can take advantage of **spatial locality** by loading data from not just one address, but also nearby addresses, into the cache.
- **Associative caches** assign each memory address to a particular set within the cache, but not to any specific block within that set.
 - Set sizes range from 1 (**direct-mapped**) to 2^k (**fully associative**).
 - Larger sets and higher associativity lead to fewer cache conflicts and lower miss rates, but they also increase the hardware cost.
 - In practice, 2-way through 16-way set-associative caches strike a good balance between lower miss rates and higher costs.
- Next time, we'll talk more about measuring cache performance, and also discuss the issue of *writing* data to a cache.