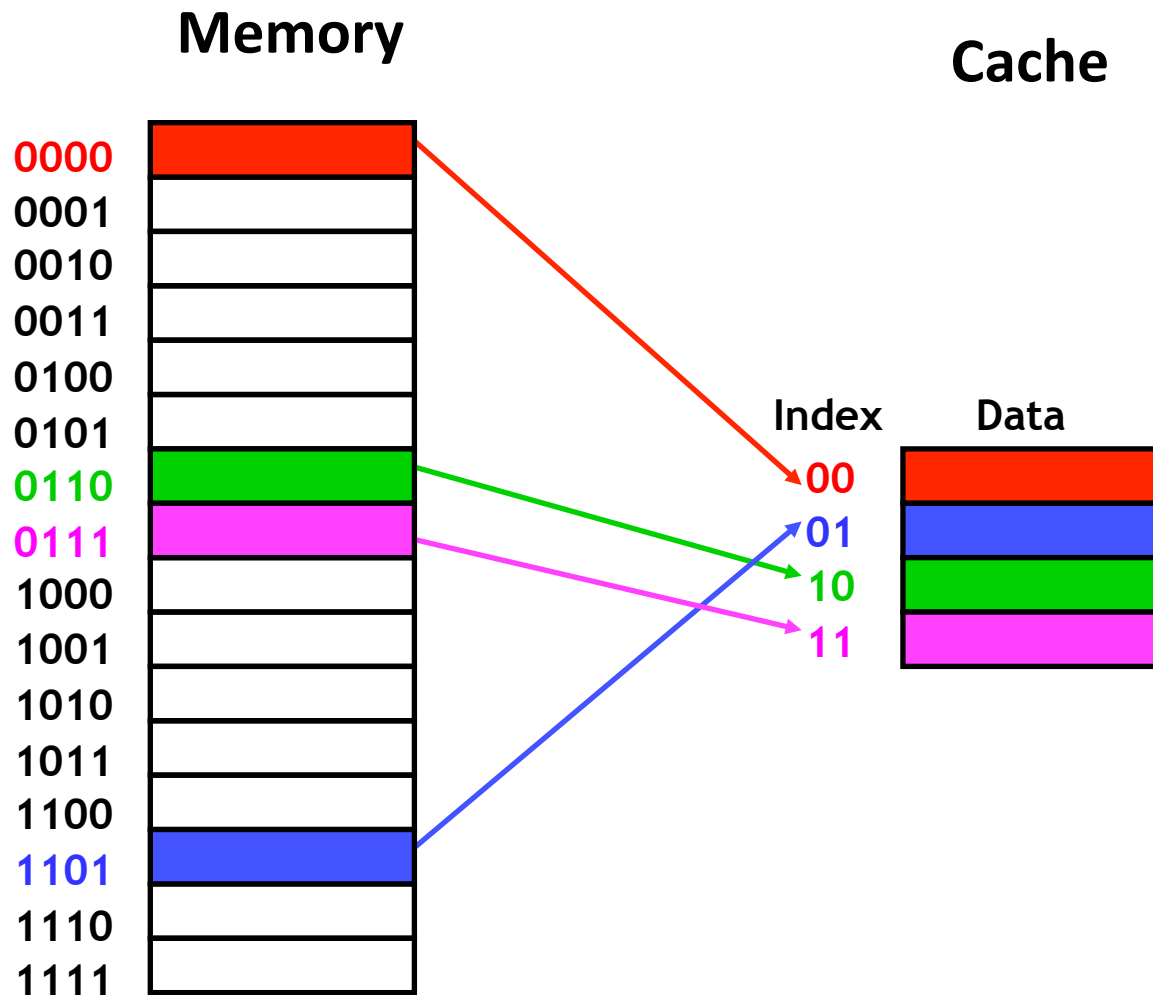


Section 7: Memory and Caches

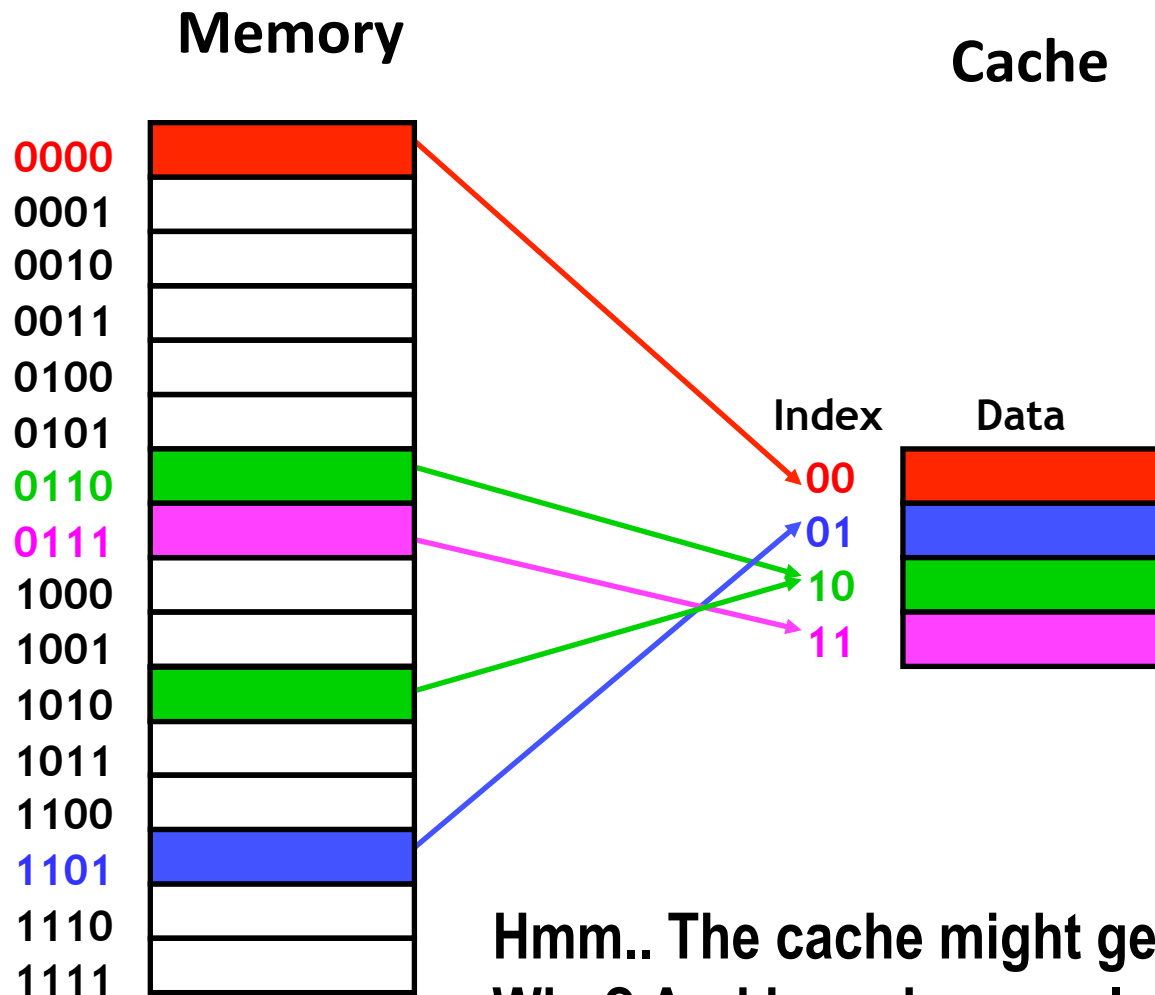
- ~~Cache basics~~
- ~~Principle of locality~~
- ~~Memory hierarchies~~
- Cache organization
- Program optimizations that consider caches

Where should we put data in the cache?



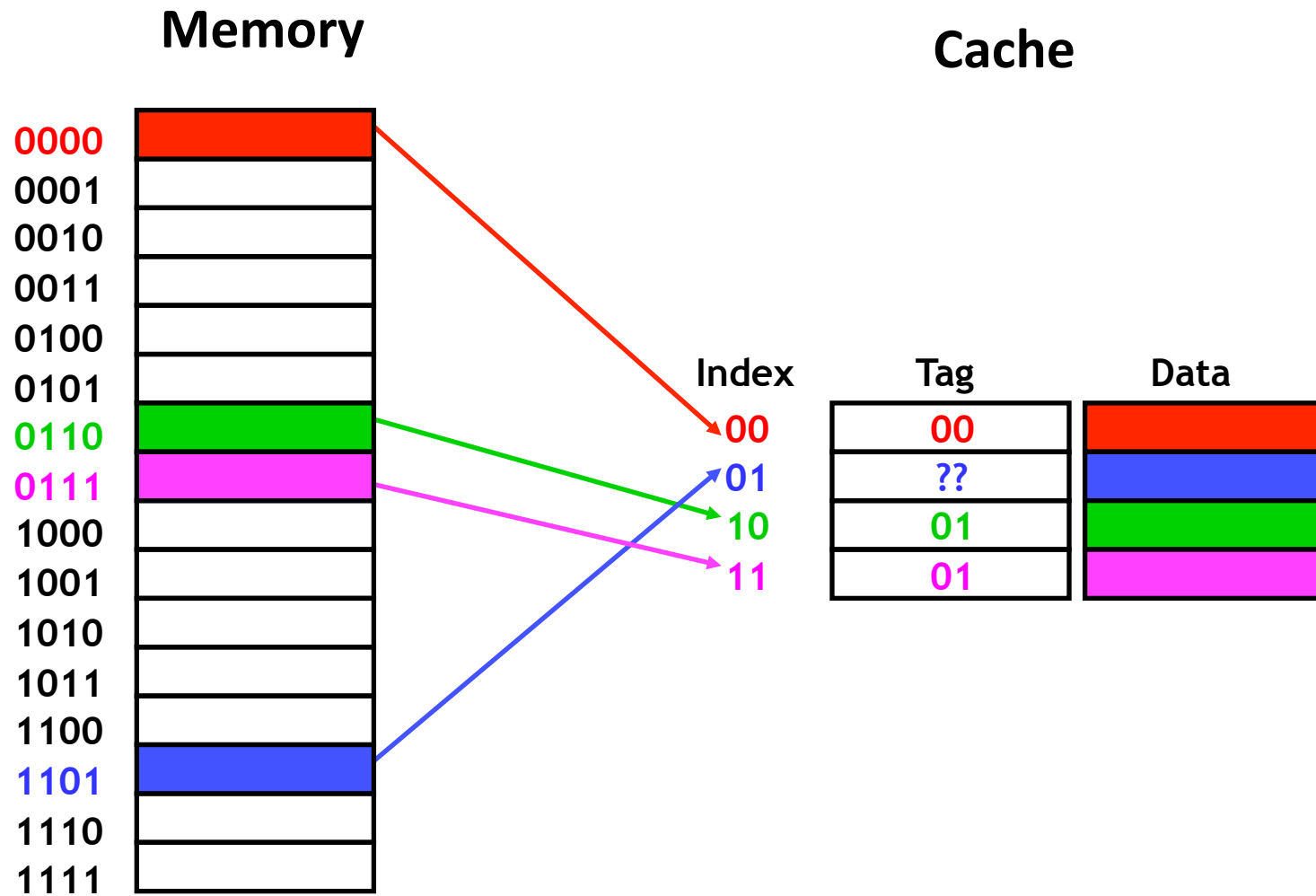
- How can we compute this mapping?

Where should we put data in the cache?

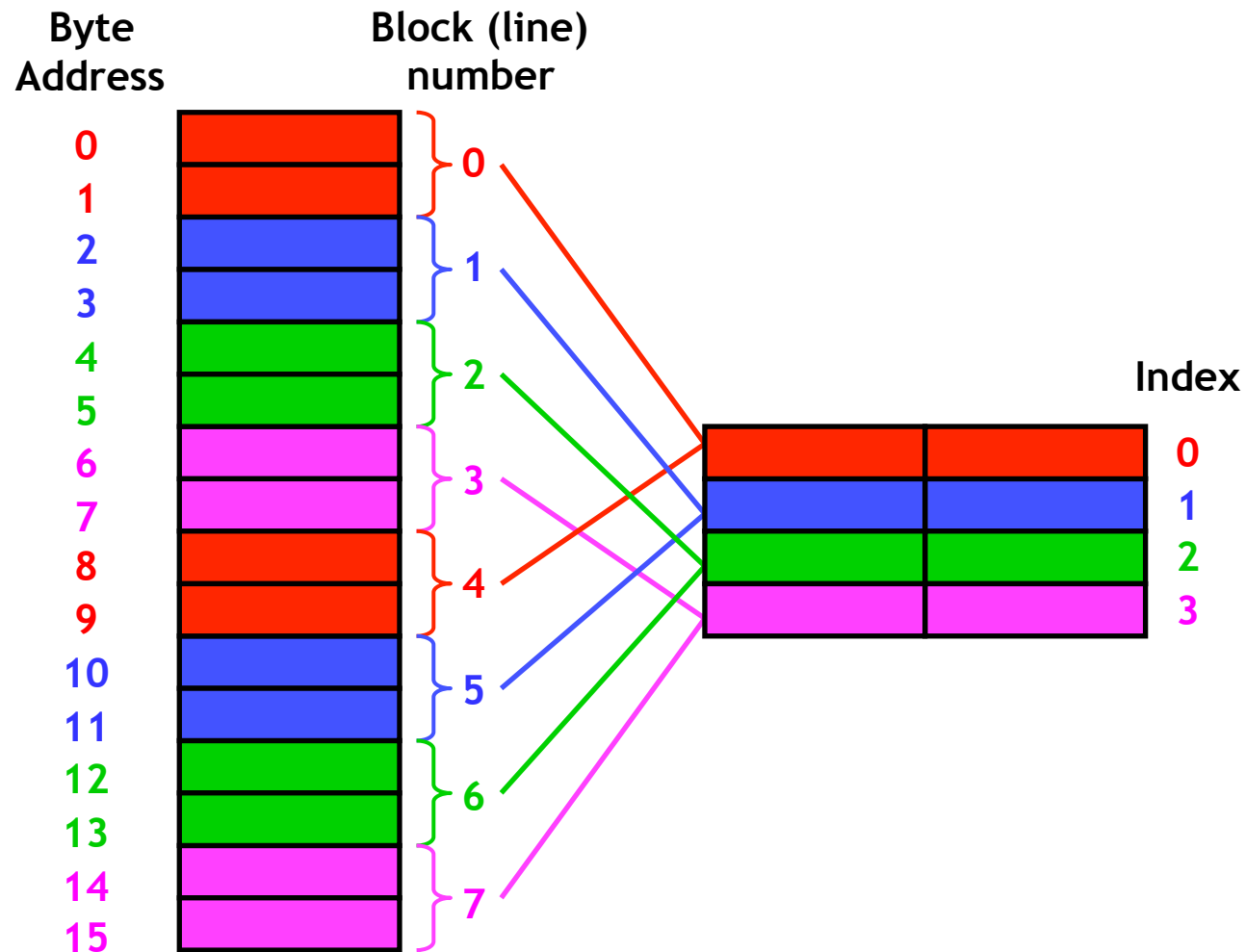


Hmm.. The cache might get confused later!
Why? And how do we solve that?

Use tags!



What's a cache block? (or *cache line*)

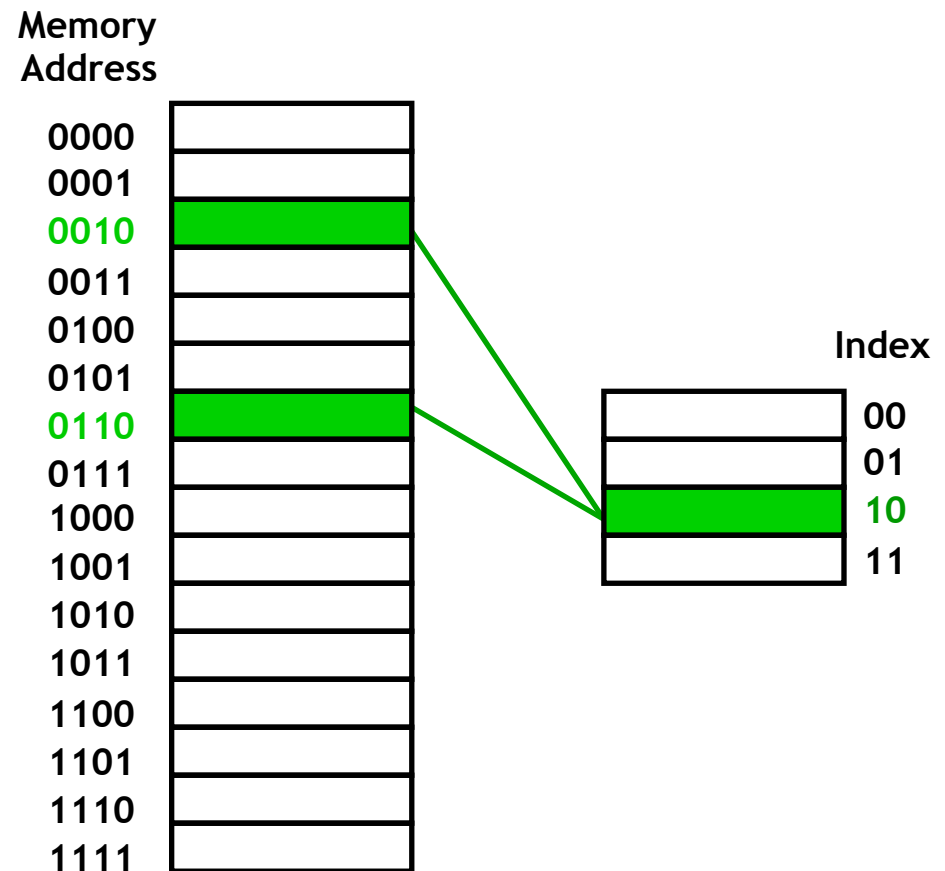


A puzzle.

- What can you infer from this:
- Cache starts *empty*
- Access (addr, hit/miss) stream
- (10, miss), (11, hit), (12, miss)

Problems with direct mapped caches?

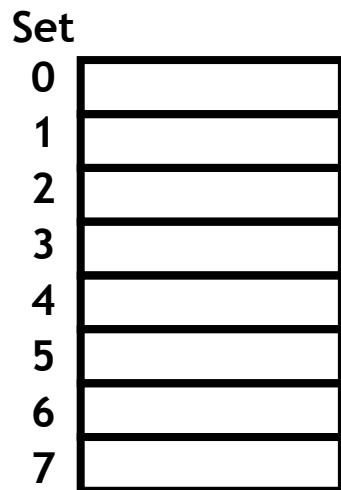
- What happens if a program uses addresses 2, 6, 2, 6, 2, ...?



Associativity

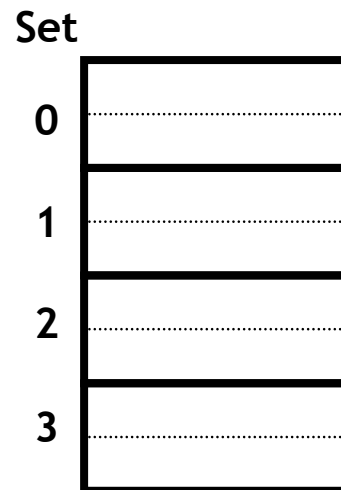
- What if we could store data in *any* place in the cache?
- But that might slow down caches... so we do something in between.

1-way
8 sets,
1 block each

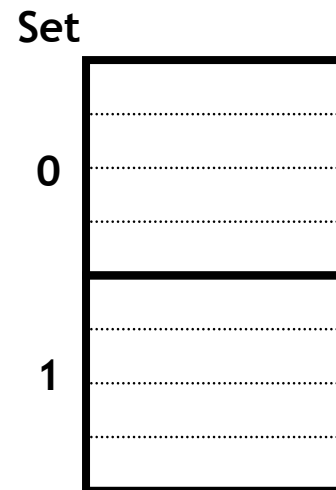


direct mapped

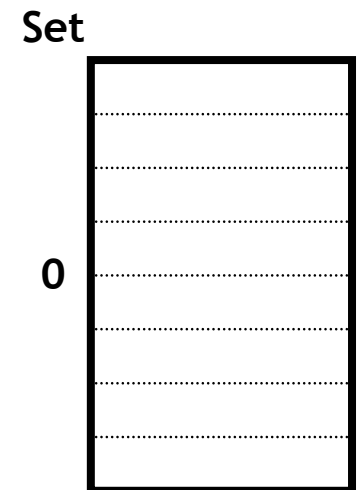
2-way
4 sets,
2 blocks each



4-way
2 sets,
4 blocks each

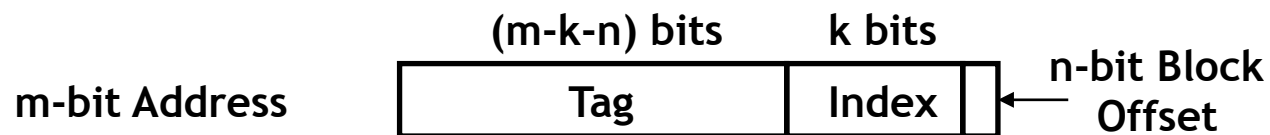


8-way
1 set,
8 blocks



fully associative

But now how do I know where data goes?

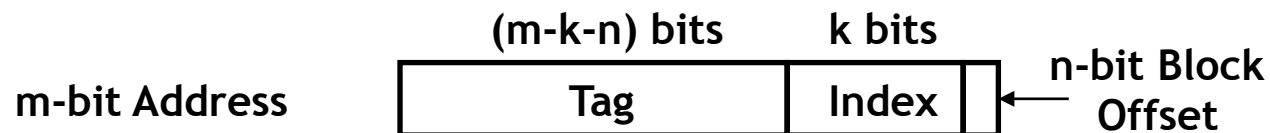


Our example used a 2^2 -block cache with 2^1 bytes per block. Where would 13 (1101) be stored?



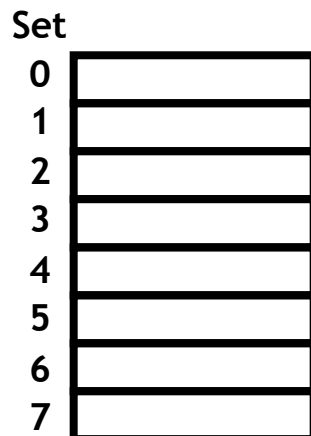
Example placement in set-associative caches

- Where would data from address 0x1833 be placed?
 - Block size is 16 bytes.
- 0x1833 in binary is 00...0110000 **011** **0011**.



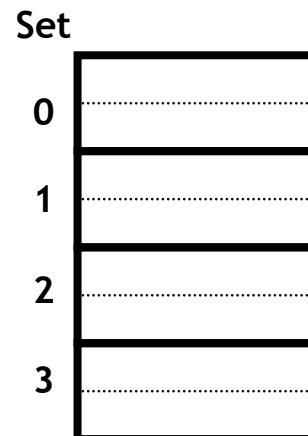
k = ?

1-way associativity
8 sets, 1 block each



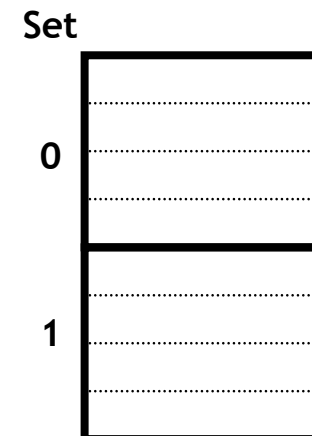
k = ?

2-way associativity
4 sets, 2 blocks each



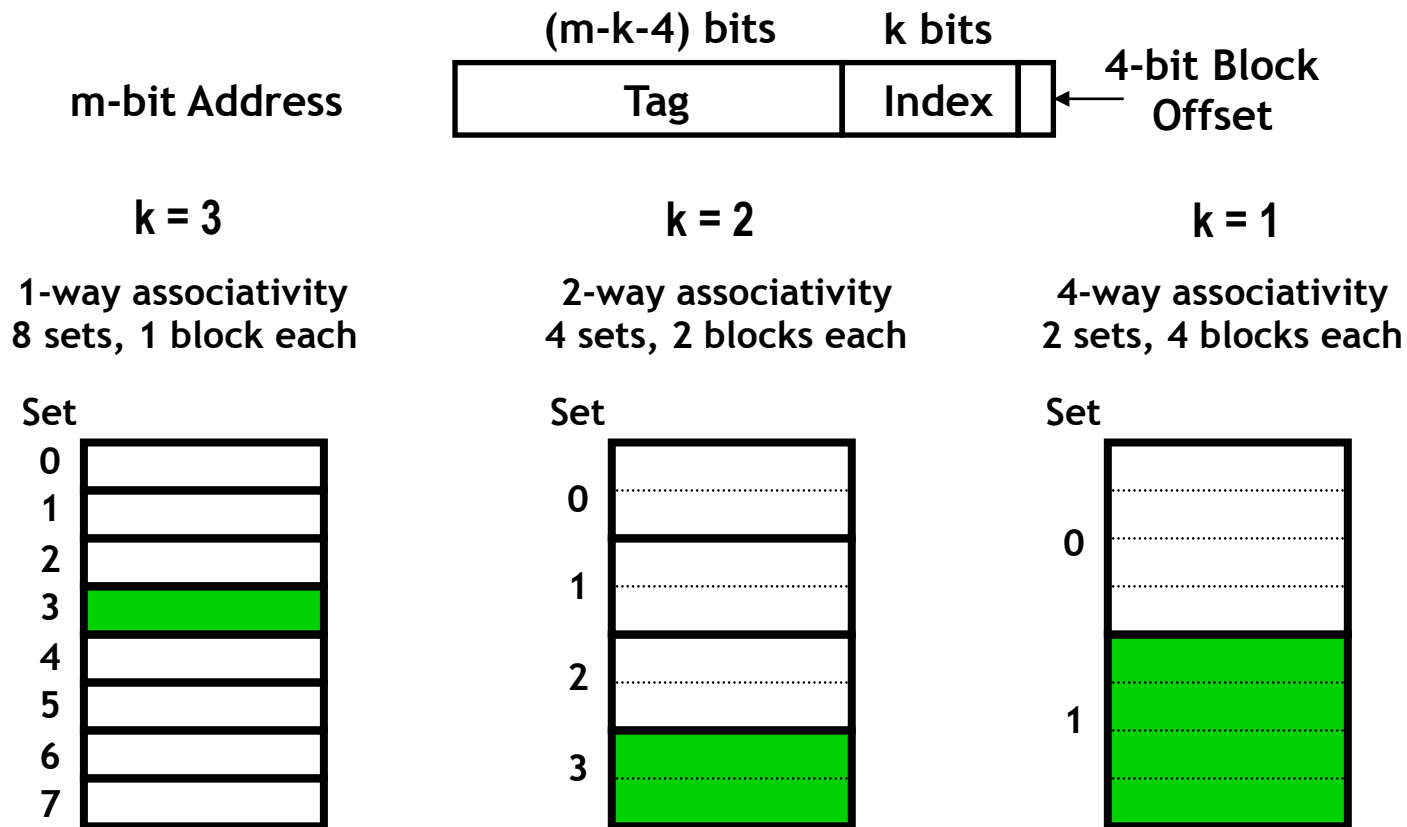
k = ?

4-way associativity
2 sets, 4 blocks each



Example placement in set-associative caches

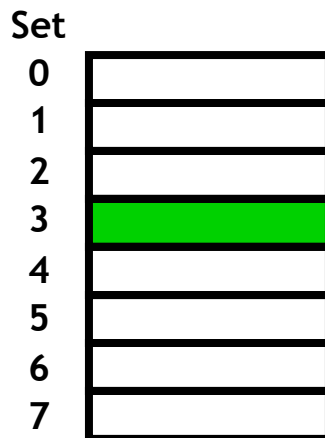
- Where would data from address 0x1833 be placed?
 - Block size is 16 bytes.
- 0x1833 in binary is 00...0110000 **011** **0011**.



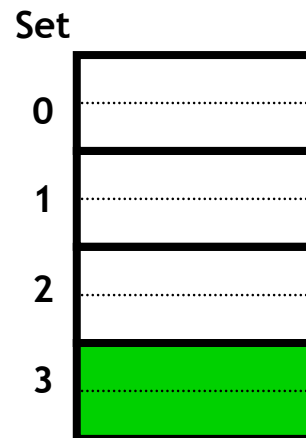
Block replacement

- Any empty block in the correct set may be used for storing data.
- If there are no empty blocks, which one should we replace?
- Replace something, of course, but what?
 - Caches typically use something close to least-recently-used

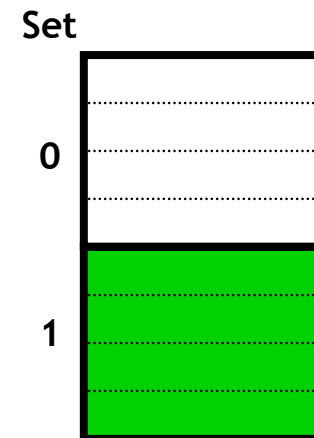
1-way associativity
8 sets, 1 block each



2-way associativity
4 sets, 2 blocks each



4-way associativity
2 sets, 4 blocks each



Another puzzle.

- What can you infer from this:
- Cache starts *empty*
- Access (addr, hit/miss) stream
- (10, miss); (12, miss); (10, miss)