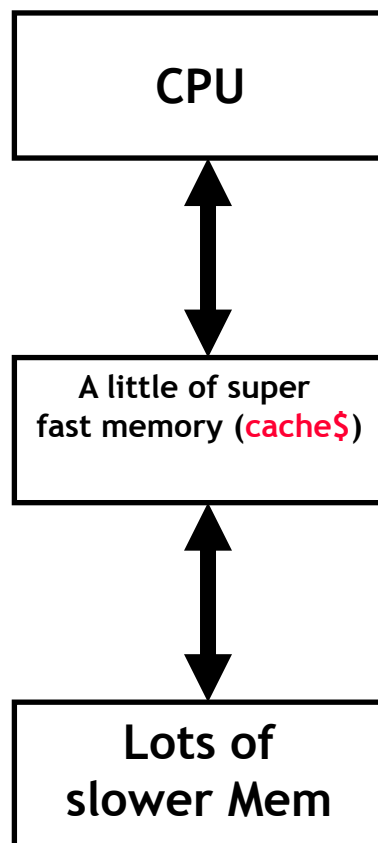


# **The Hardware/Software Interface**

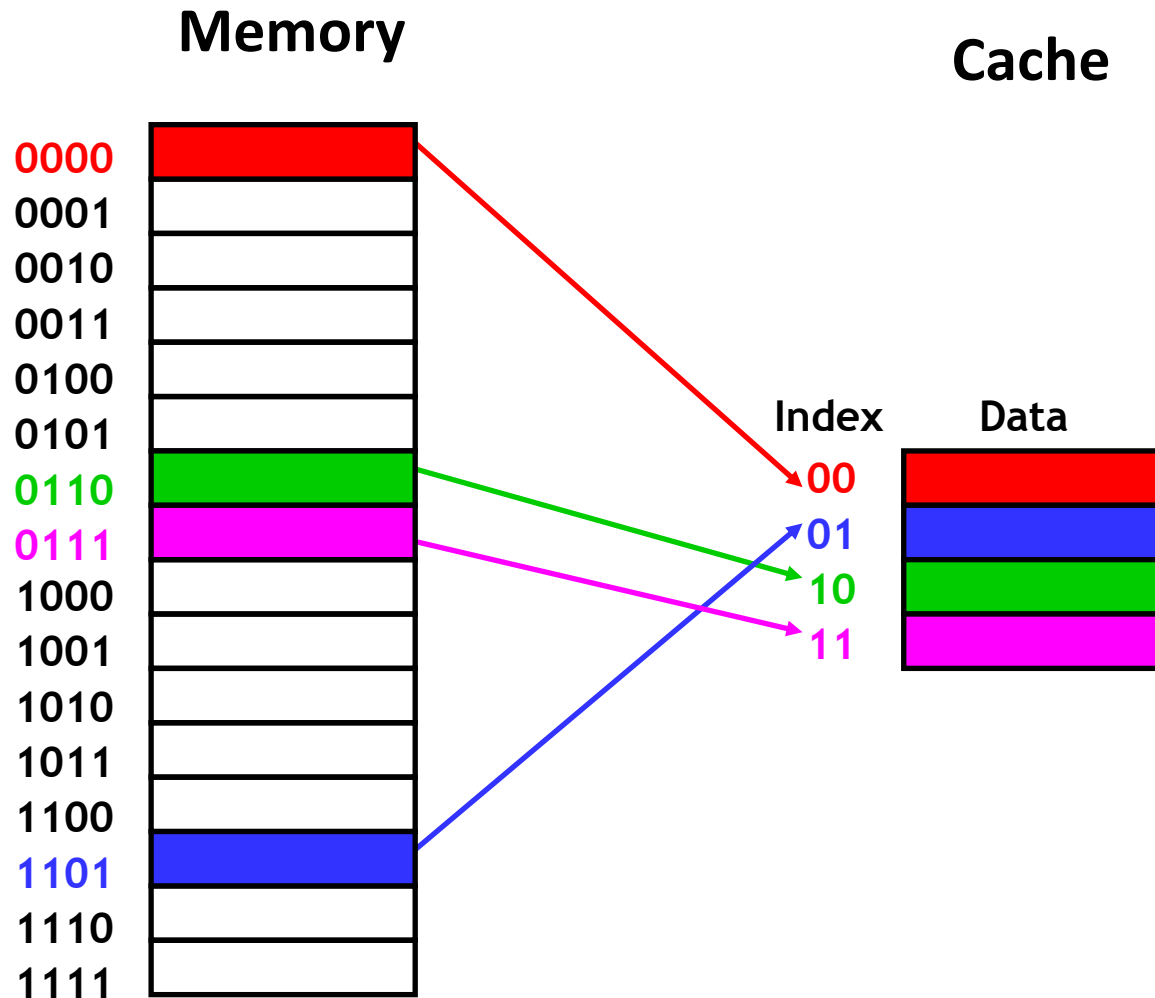
CSE351 Spring 2013

## **Memory and Caches II**

# Not to forget...

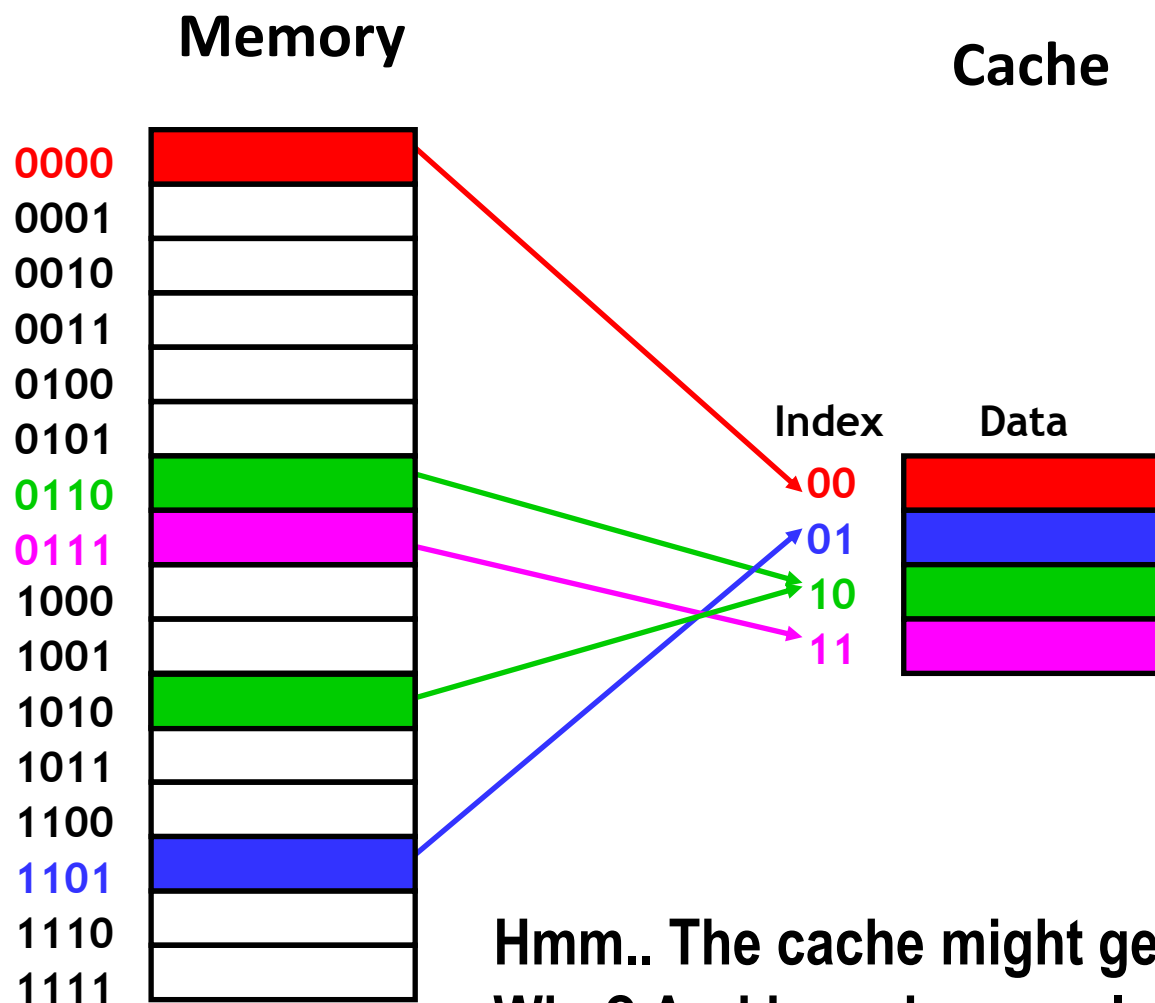


# Where should we put data in the cache?



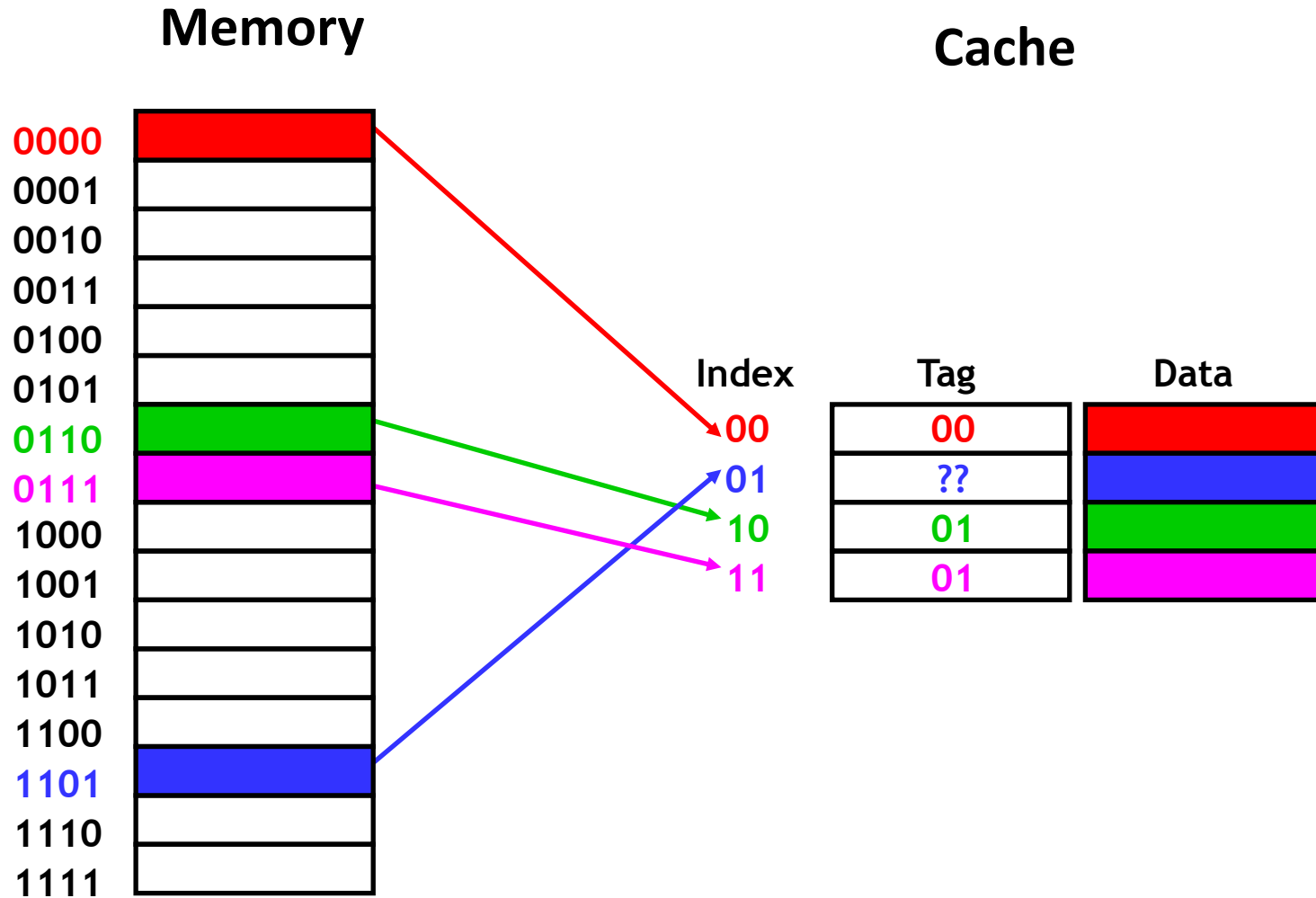
- How can we compute this mapping?

# Where should we put data in the cache?

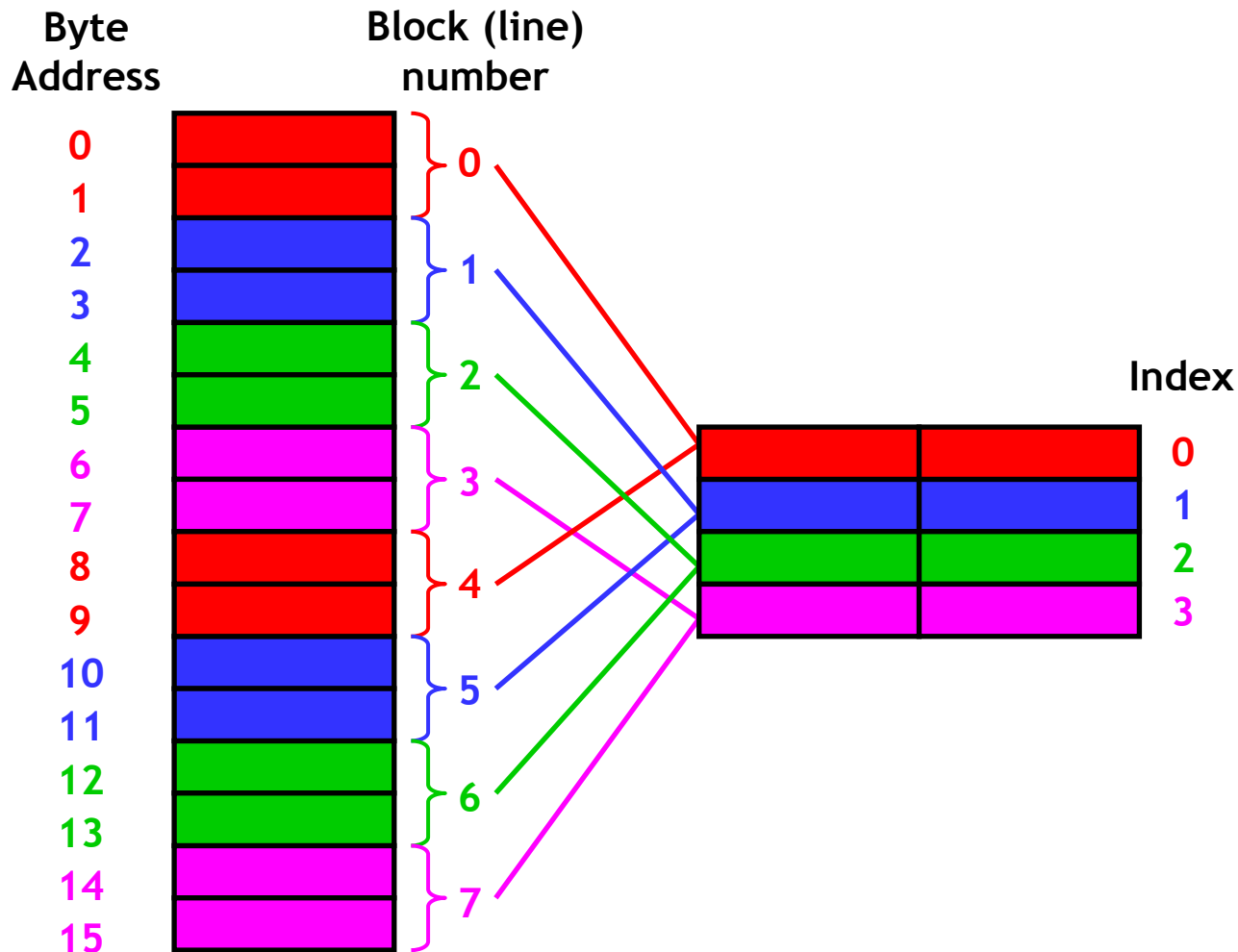


Hmm.. The cache might get confused later!  
Why? And how do we solve that?

# Use tags!



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

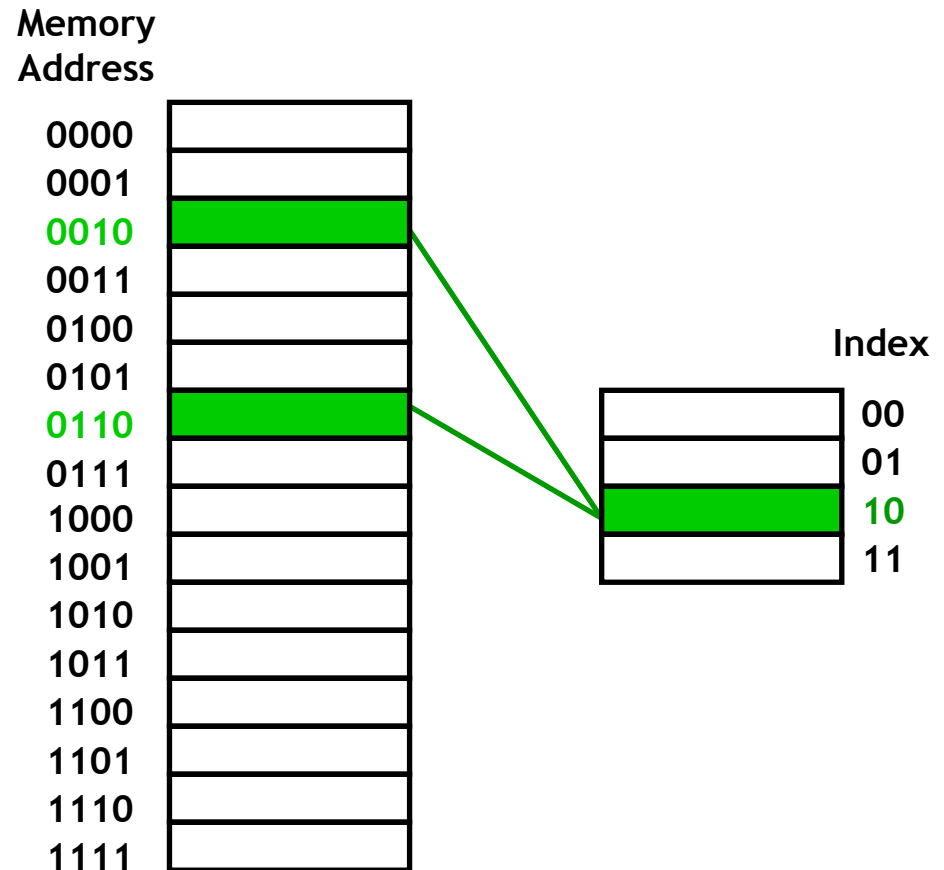


# A puzzle.

- 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

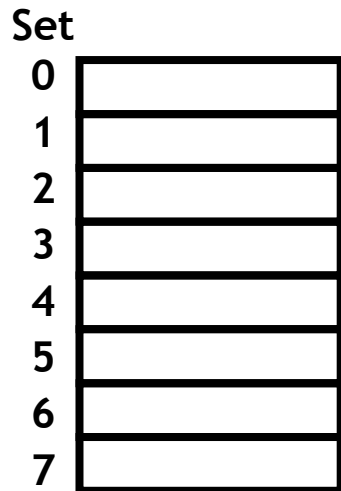
- What if we could store data in *any* place in the cache?

# 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**

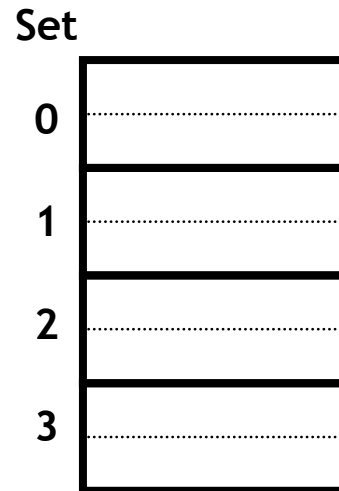
8 sets,  
1 block each



**direct mapped**

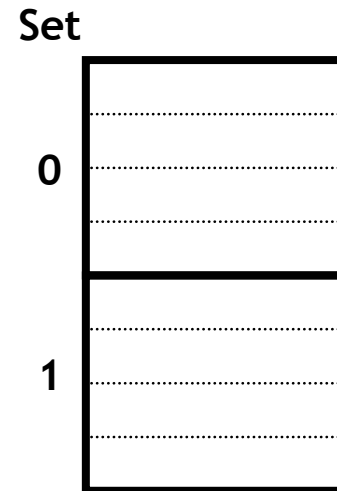
**2-way**

4 sets,  
2 blocks each



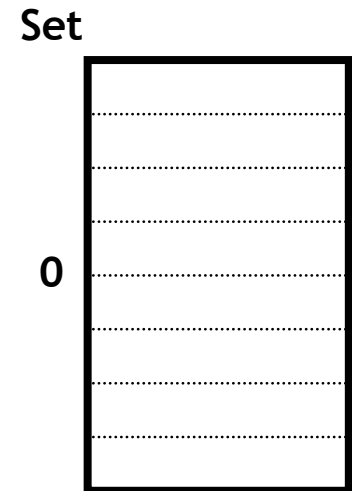
**4-way**

2 sets,  
4 blocks each



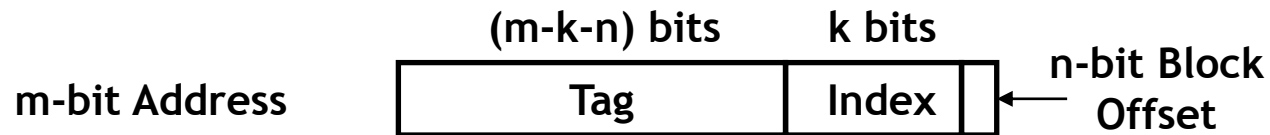
**8-way**

1 set,  
8 blocks

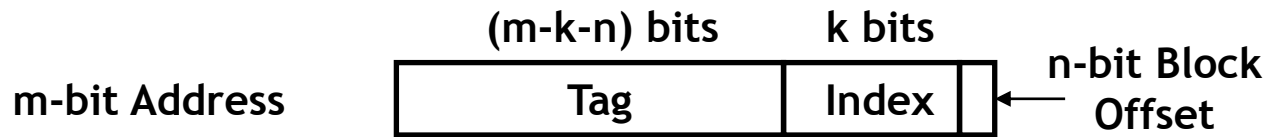


**fully associative**

# Now how do I know where data goes?



# But now how do I know where data goes?

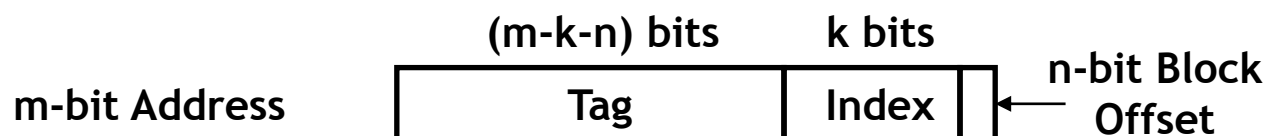


Our example used a  $2^2$ -block cache with  $2^1$  bytes per block. Where would 13 (1101) be stored?



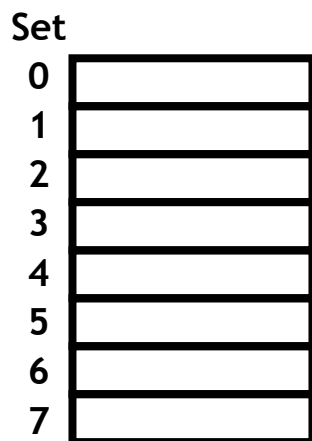
# Example placement in set-associative caches

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



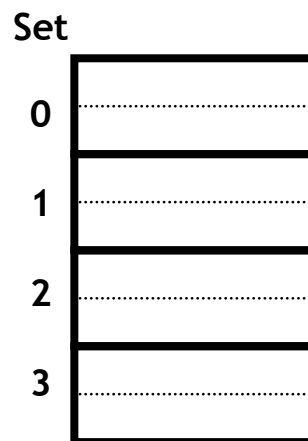
k = ?

1-way associativity  
8 sets, 1 block each



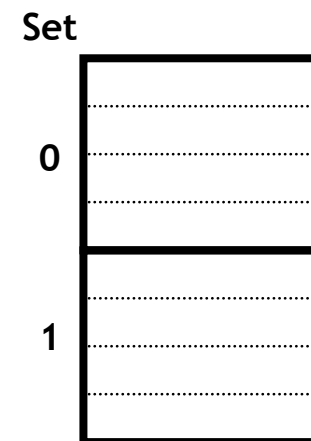
k = ?

2-way associativity  
4 sets, 2 blocks each



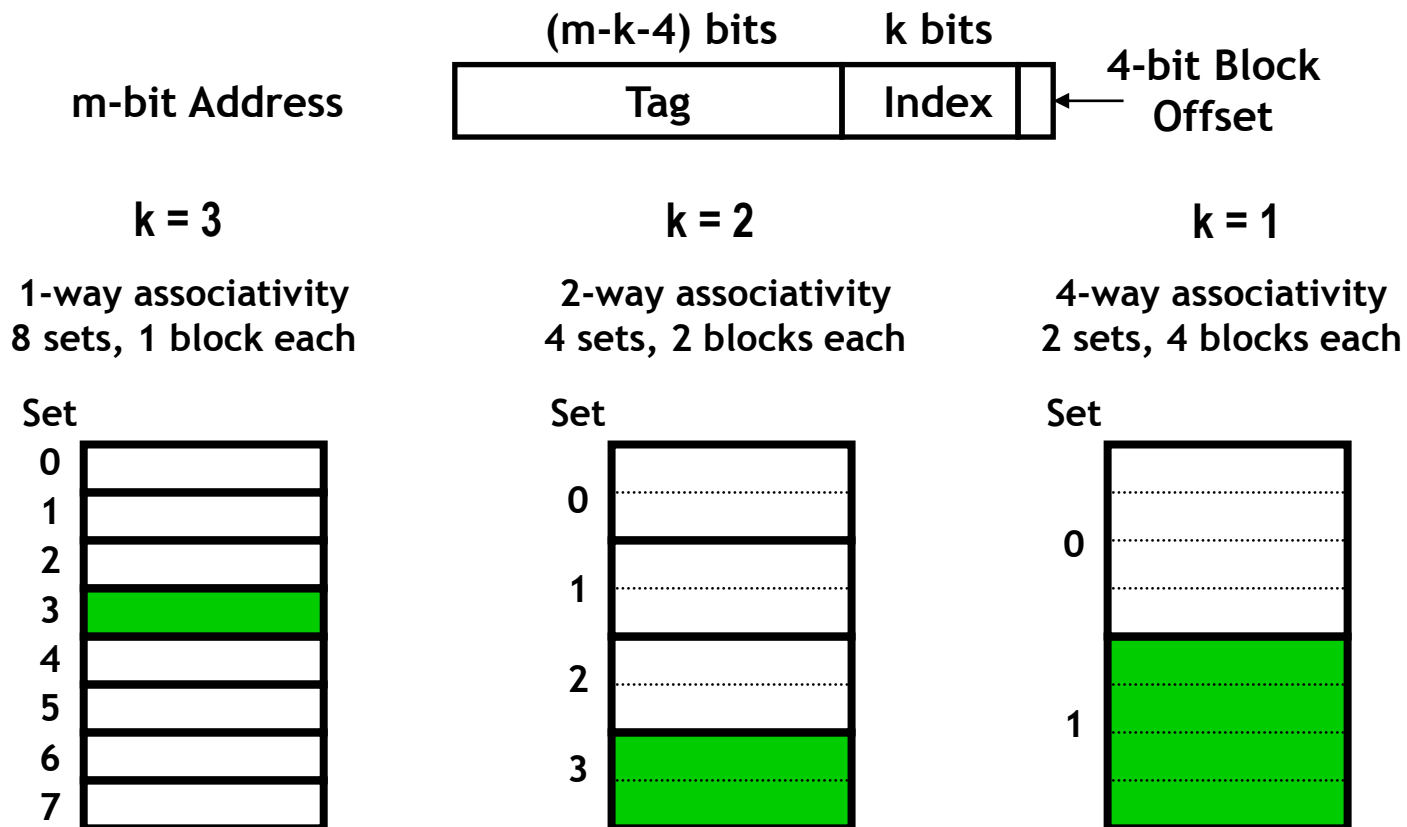
k = ?

4-way associativity  
2 sets, 4 blocks each



# Example placement in set-associative caches

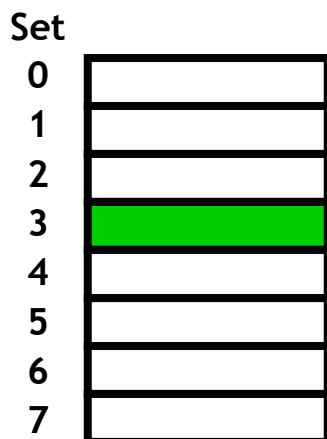
- Where would data from address 0x1833 be placed?
  - Block size is 16 bytes.
- 0x1833 in binary is 00...0110000 **011** **0011**.



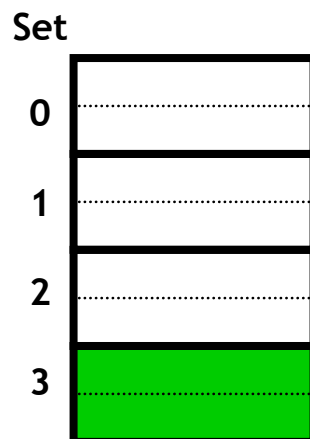
# Block replacement

- Any empty block in the correct set may be used for storing data.
- If there are no empty blocks, which one should we replace?

1-way associativity  
8 sets, 1 block each



2-way associativity  
4 sets, 2 blocks each



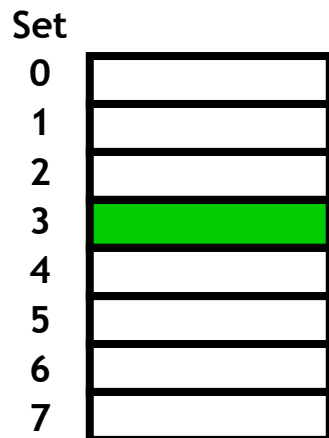
4-way associativity  
2 sets, 4 blocks each



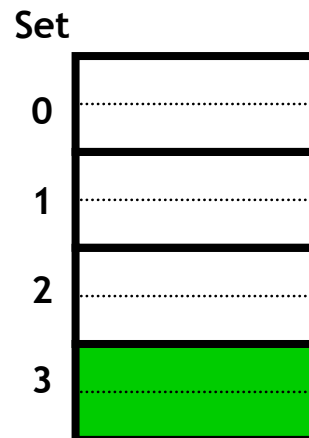
# Block replacement

- Replace something, of course, but what?

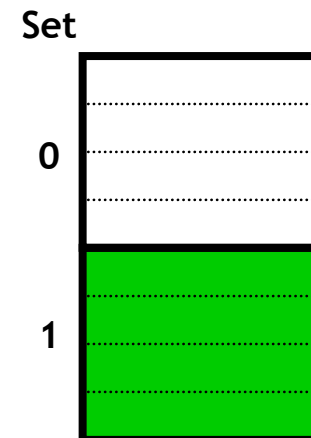
1-way associativity  
8 sets, 1 block each



2-way associativity  
4 sets, 2 blocks each



4-way associativity  
2 sets, 4 blocks each

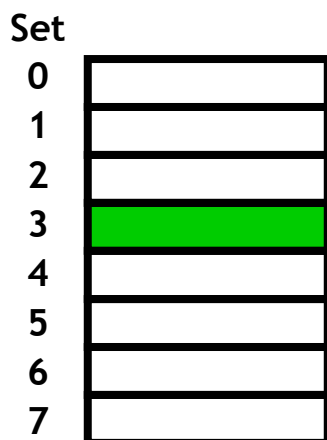




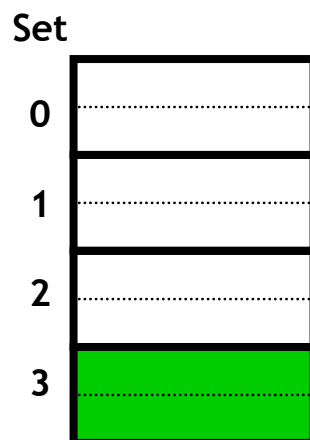
# Block replacement

- 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



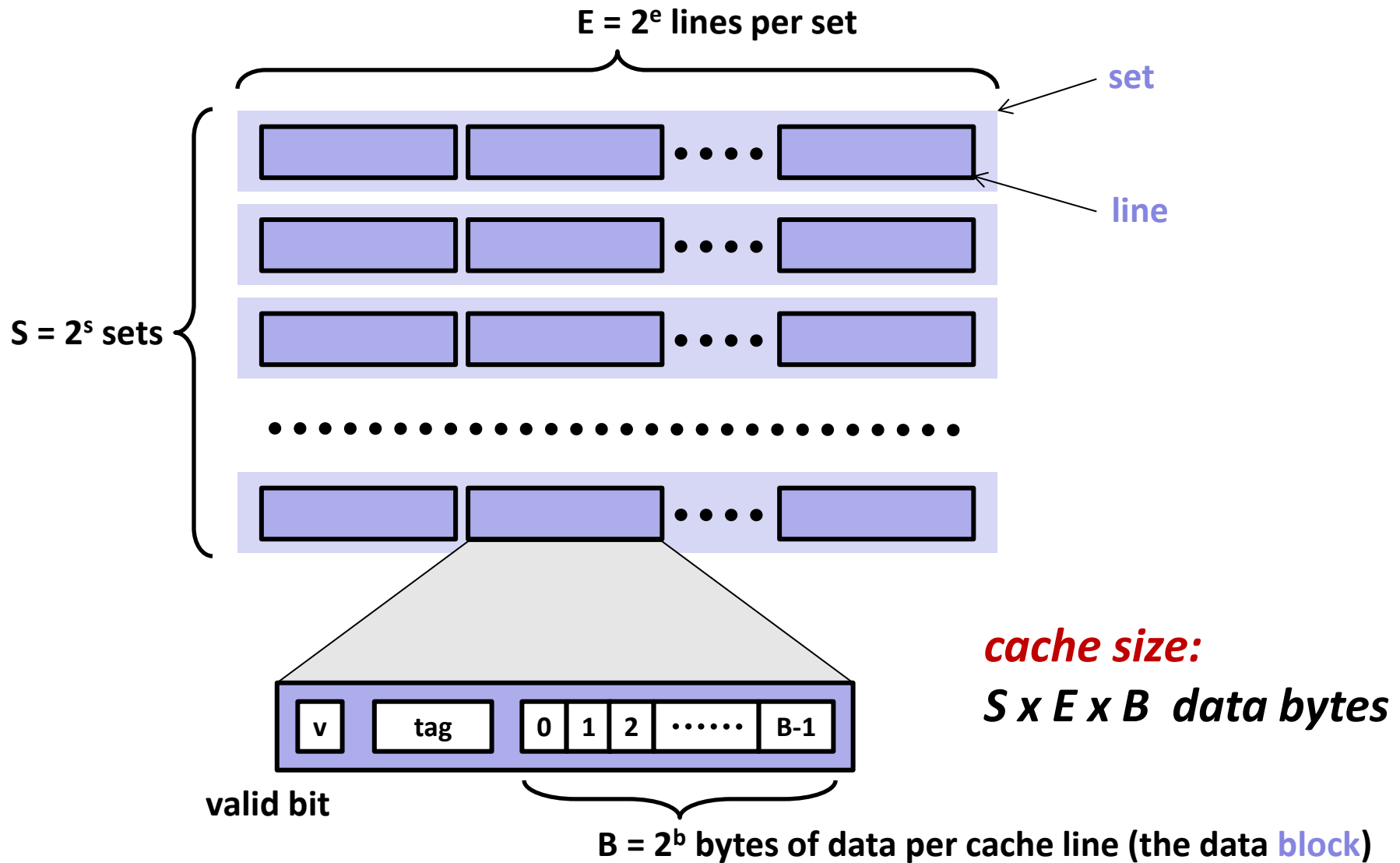
4-way associativity  
2 sets, 4 blocks each



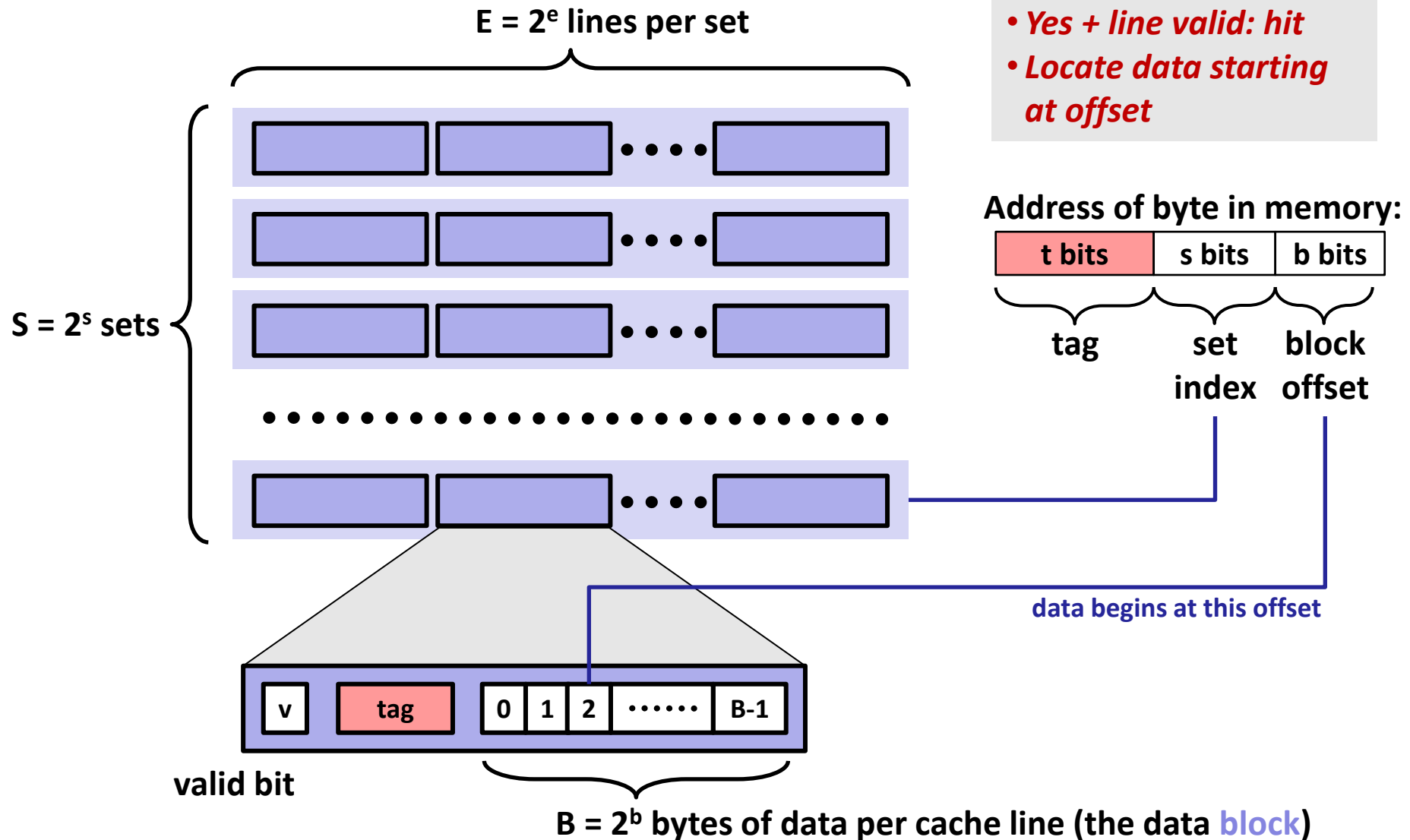
# Another puzzle.

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

# General Cache Organization (S, E, B)



# Cache Read

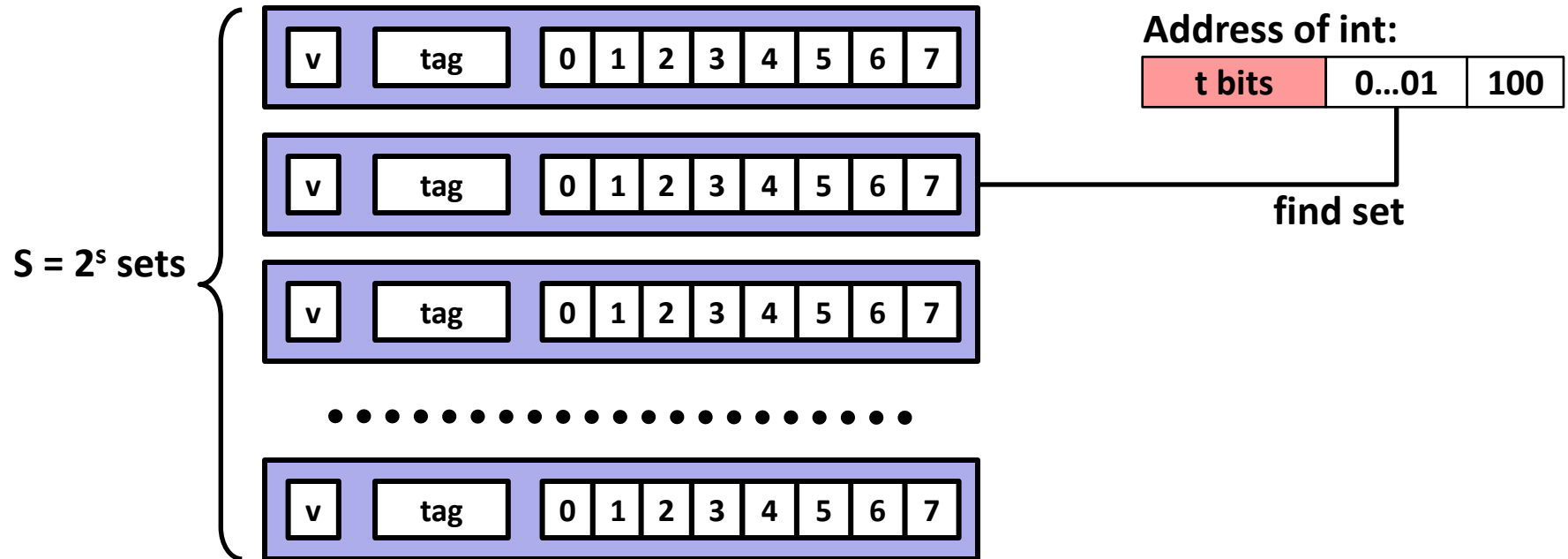


- *Locate set*
- *Check if any line in set has matching tag*
- *Yes + line valid: hit*
- *Locate data starting at offset*

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

Direct-mapped: One line per set

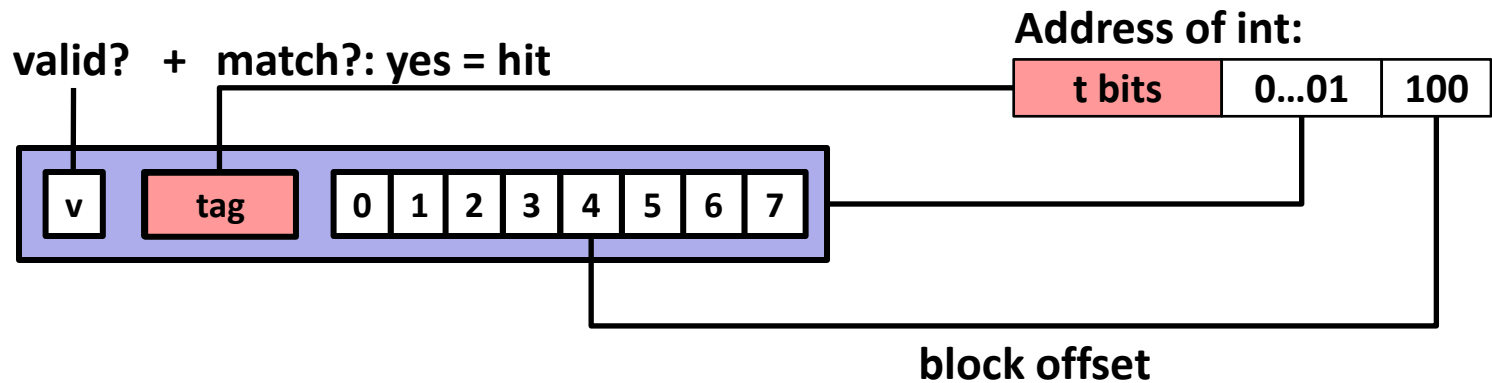
Assume: cache block size 8 bytes



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

Direct-mapped: One line per set

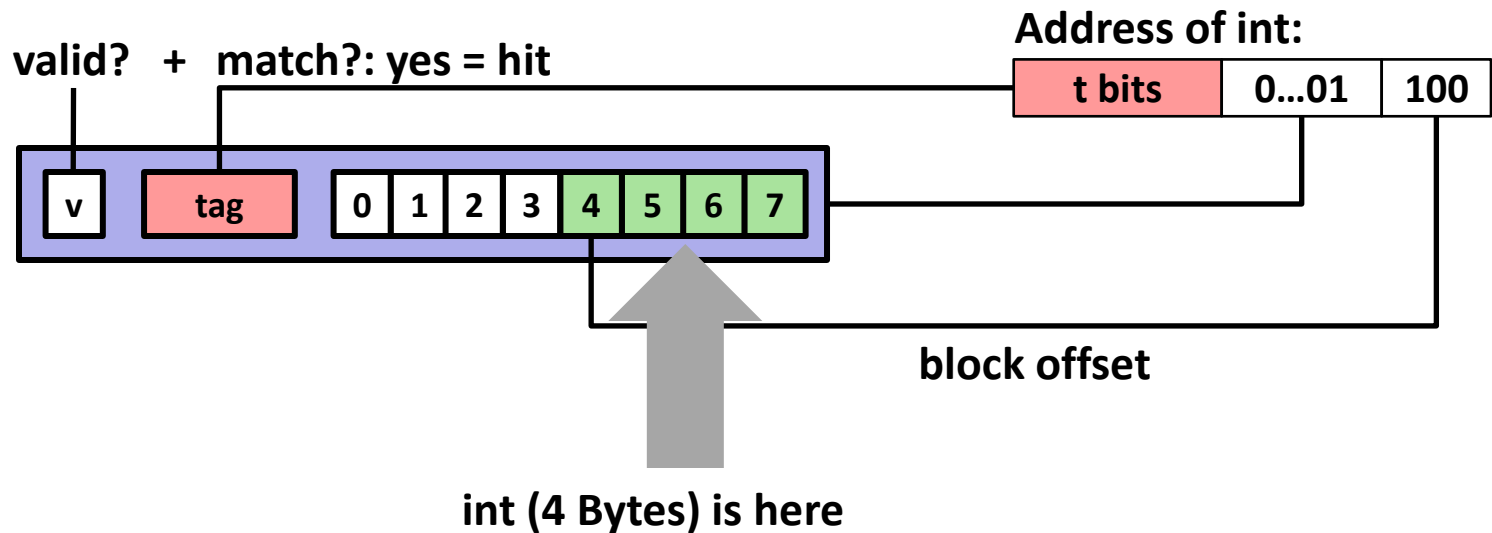
Assume: cache block size 8 bytes



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

Direct-mapped: One line per set

Assume: cache block size 8 bytes



**No match:** old line is evicted and replaced

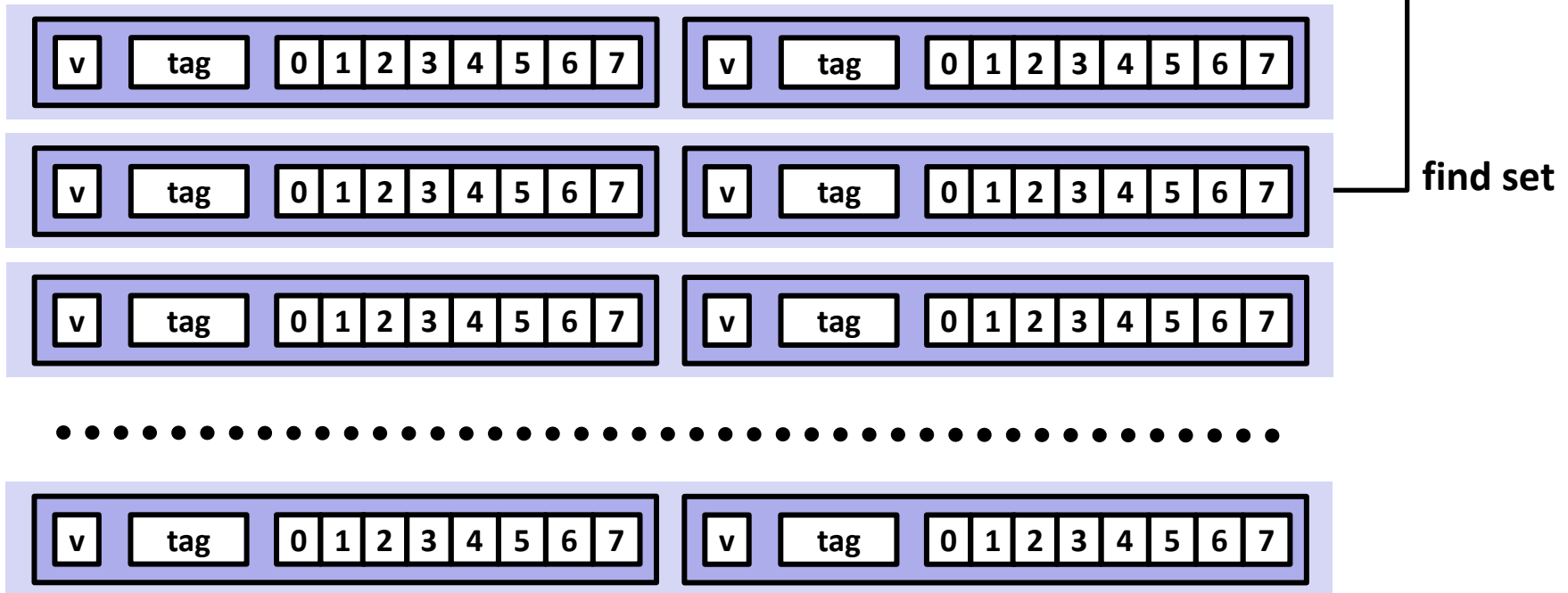
# E-way Set-Associative Cache (Here: E = 2)

E = 2: Two lines per set

Assume: cache block size 8 bytes

Address of short int:

t bits	0...01	100
--------	--------	-----

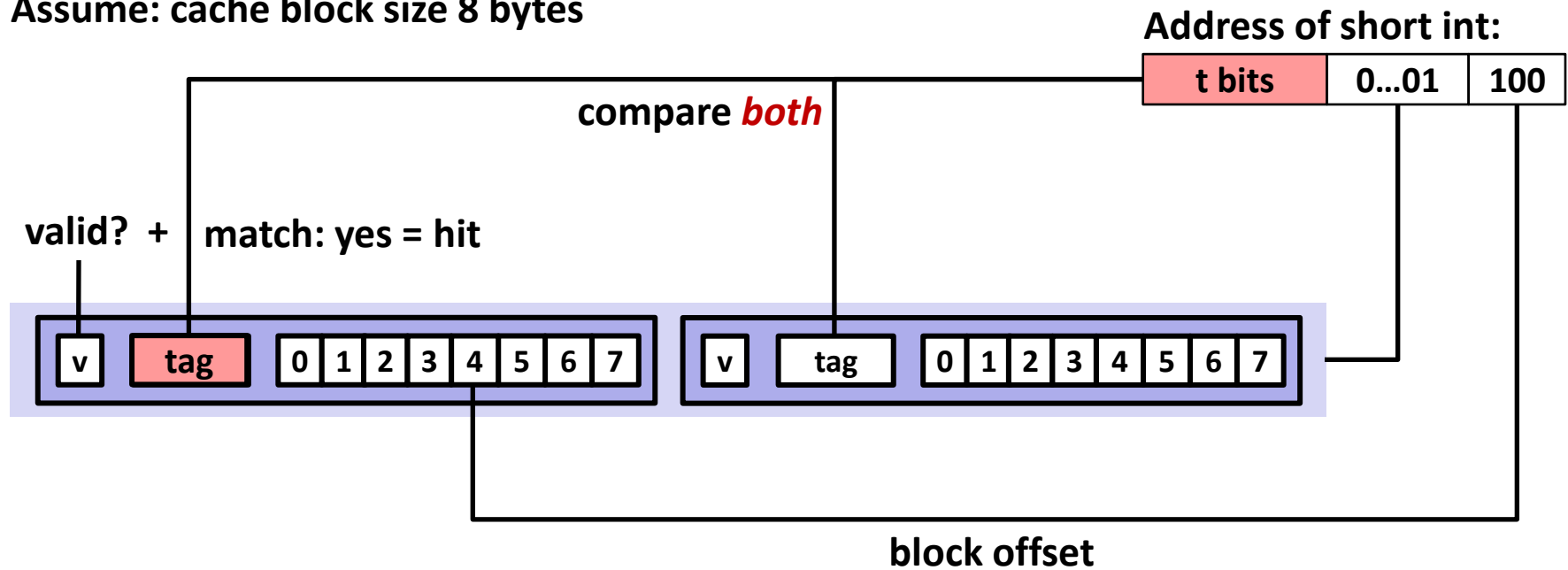




# E-way Set-Associative Cache (Here: E = 2)

E = 2: Two lines per set

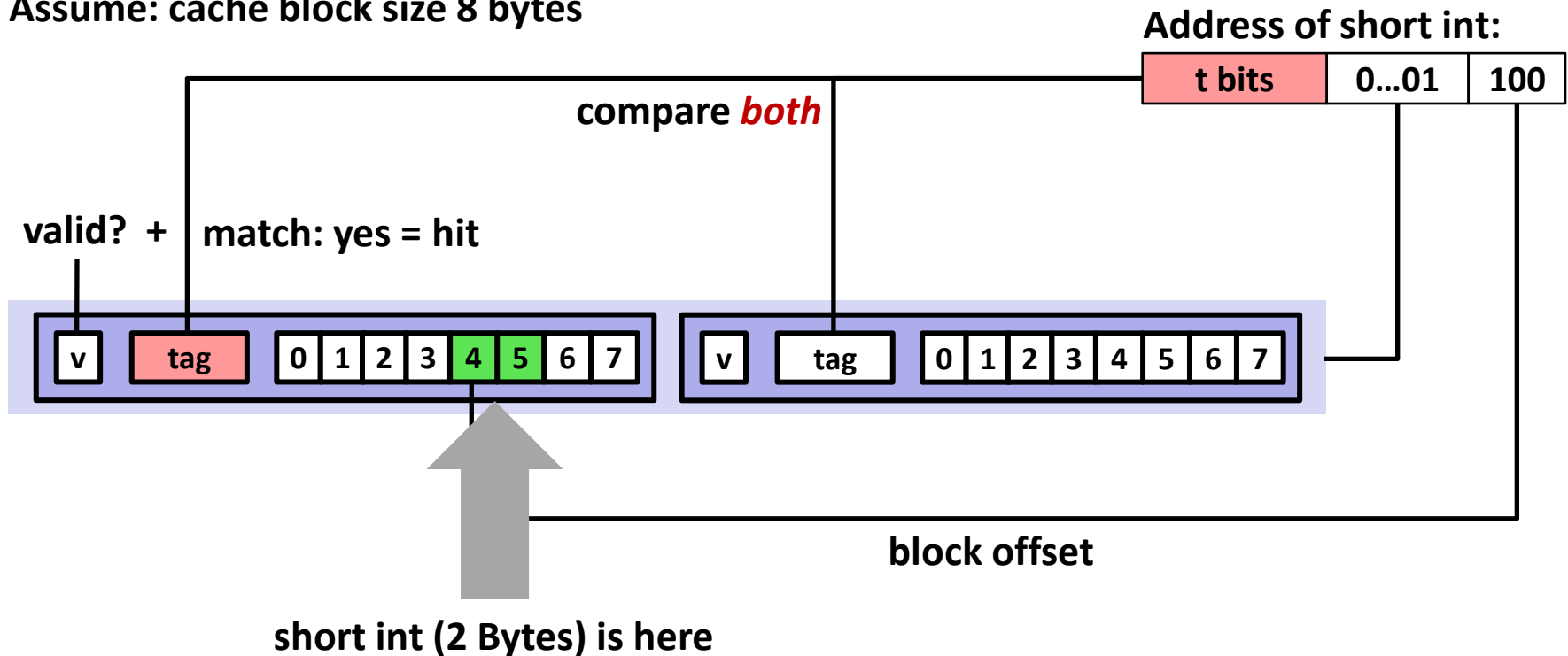
Assume: cache block size 8 bytes



# E-way Set-Associative Cache (Here: E = 2)

E = 2: Two lines per set

Assume: cache block size 8 bytes



## No match:

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

# Types of Cache Misses

- **Cold (compulsory) miss**
  - Occurs on first access to a block

# Types of Cache Misses

## ■ Cold (compulsory) miss

- Occurs on first access to a block

## ■ Conflict miss

- Most hardware caches limit blocks to a small subset (sometimes just one) of the available cache slots
  - if one (e.g., block  $i$  must be placed in slot  $(i \bmod \text{size})$ ), direct-mapped
  - if more than one,  $n$ -way set-associative (where  $n$  is a power of 2)
- 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, ... would miss every time=

# Types of Cache Misses

## ■ Cold (compulsory) miss

- Occurs on first access to a block

## ■ Conflict miss

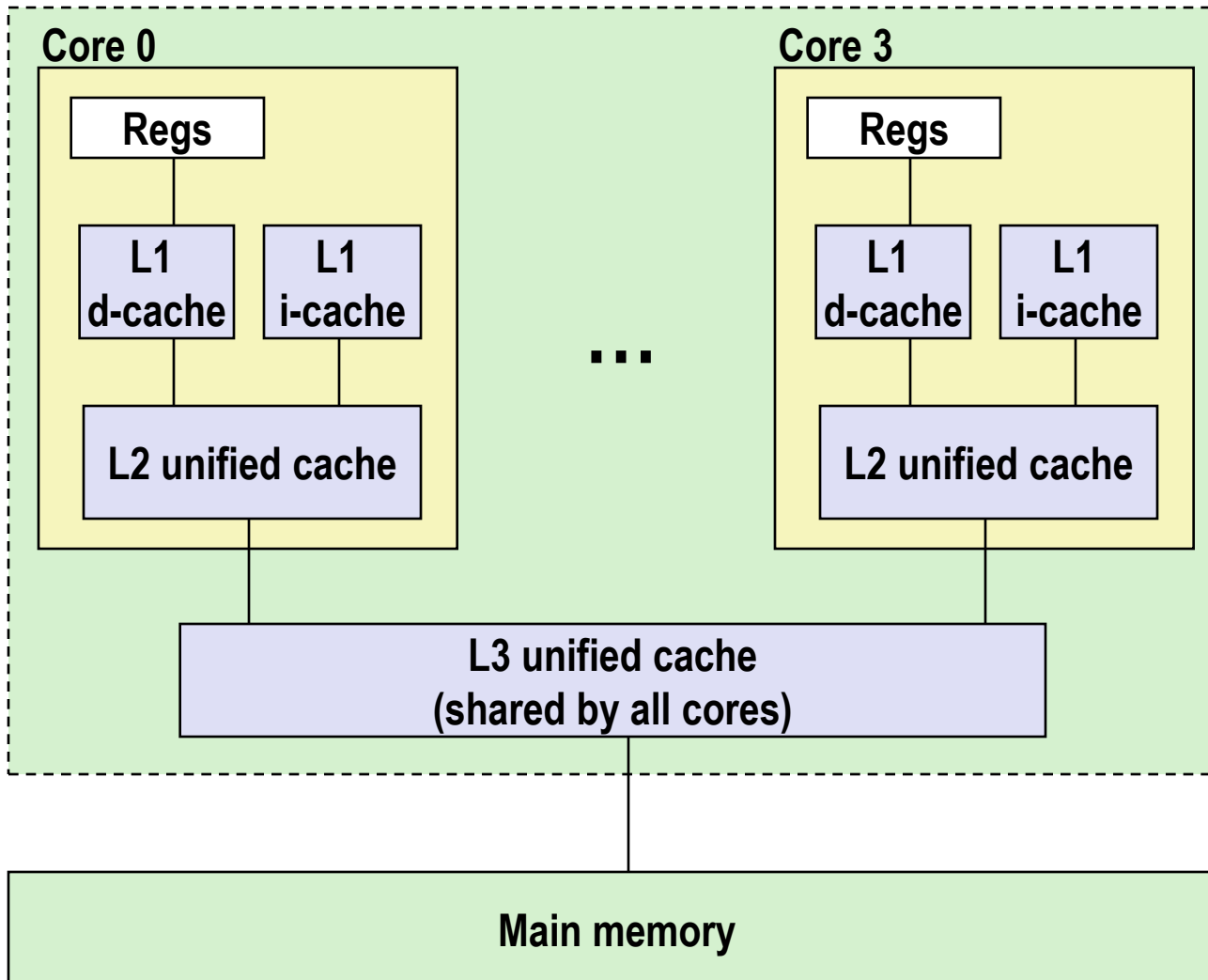
- Most hardware caches limit blocks to a small subset (sometimes just one) of the available cache slots
  - if one (e.g., block  $i$  must be placed in slot  $(i \bmod \text{size})$ ), direct-mapped
  - if more than one,  $n$ -way set-associative (where  $n$  is a power of 2)
- Conflict misses occur when the cache is large enough, but multiple data objects all map to the same slot
  - e.g., referencing blocks 0, 8, 0, 8, ... would miss every time

## ■ Capacity miss

- Occurs when the set of active cache blocks (the working set) is larger than the cache (just won't fit)

# Intel Core i7 Cache Hierarchy

Processor package



**L1 i-cache and d-cache:**

32 KB, 8-way,  
Access: 4 cycles

**L2 unified cache:**

256 KB, 8-way,  
Access: 11 cycles

**L3 unified cache:**

8 MB, 16-way,  
Access: 30-40 cycles

**Block size:** 64 bytes for  
all caches.

# What about writes?

- **Multiple copies of data exist:**
  - L1, L2, possibly L3, main memory
- **What is the main problem with that?**

# What about writes?

- **Multiple copies of data exist:**
  - L1, L2, possibly L3, main memory
- **What to do on a write-hit?**
  - **Write-through** (write immediately to memory)
  - **Write-back** (defer write to memory until line is evicted)
    - Need a *dirty bit* to indicate if line is different from memory or not
- **What to do on a write-miss?**
  - **Write-allocate** (load into cache, update line in cache)
    - Good if more writes to the location follow
  - **No-write-allocate** (just write immediately to memory)
- **Typical caches:**
  - Write-back + Write-allocate, usually
  - Write-through + No-write-allocate, occasionally



# Where else is caching used?

# Software Caches are More Flexible

## ■ Examples

- File system buffer caches, web browser caches, etc.

## ■ Some design differences

- Almost always fully-associative
  - so, no placement restrictions
  - index structures like hash tables are common (for placement)
- Often use complex replacement policies
  - misses are very expensive when disk or network involved
  - worth thousands of cycles to avoid them
- Not necessarily constrained to single “block” transfers
  - may fetch or write-back in larger units, opportunistically

# Optimizations for the Memory Hierarchy

## ■ Write code that has locality

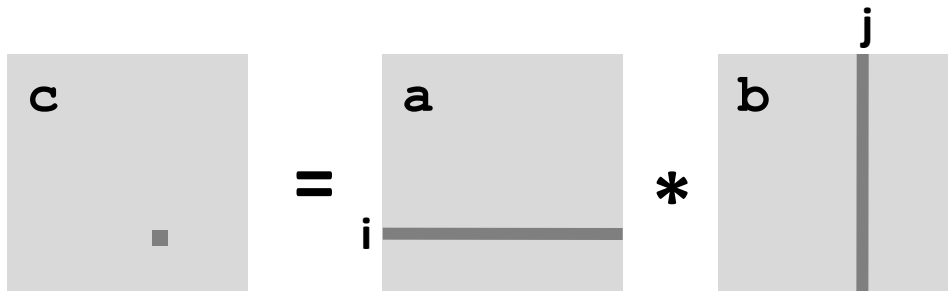
- Spatial: access data contiguously
- Temporal: make sure access to the same data is not too far apart in time

## ■ How to achieve?

- Proper choice of algorithm
- Loop transformations

# Example: Matrix Multiplication

```
c = (double *) calloc(sizeof(double), n*n);  
  
/* Multiply n x n matrices a and b */  
void mmm(double *a, double *b, double *c, int n) {  
    int i, j, k;  
    for (i = 0; i < n; i++)  
        for (j = 0; j < n; j++)  
            for (k = 0; k < n; k++)  
                c[i*n + j] += a[i*n + k]*b[k*n + j];  
}
```



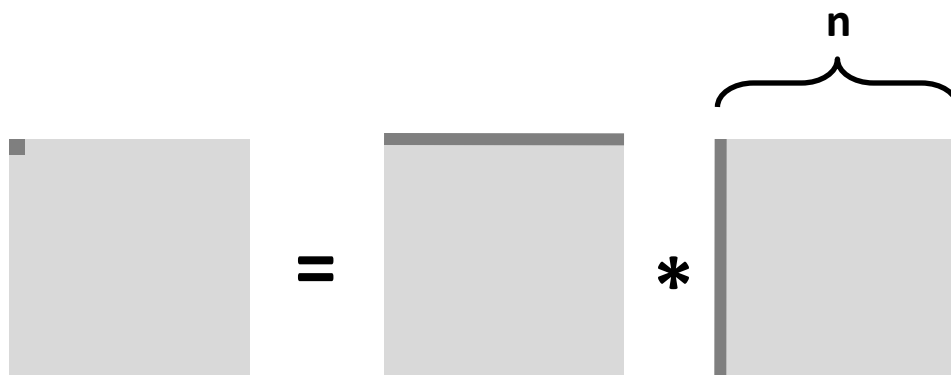
# Cache Miss Analysis

## ■ Assume:

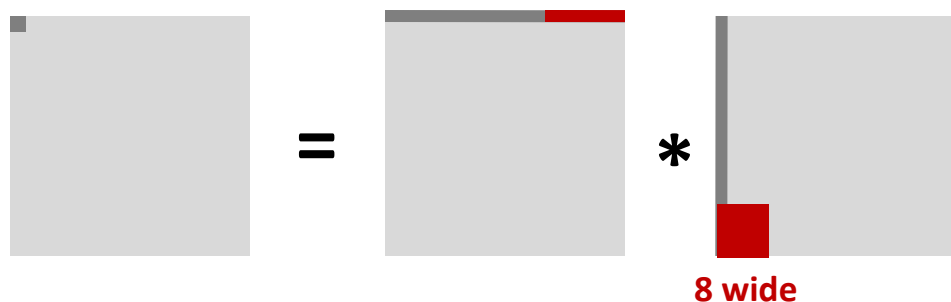
- Matrix elements are doubles
- Cache block = 64 bytes = 8 doubles
- Cache size  $C \ll n$  (much smaller than  $n$ )

## ■ First iteration:

- $n/8 + n = 9n/8$  misses  
(omitting matrix  $c$ )



- Afterwards **in cache:**  
(schematic)



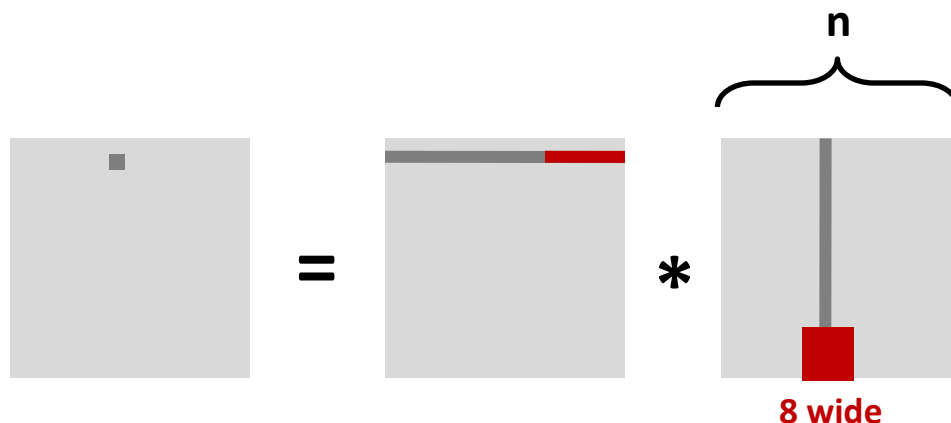
# Cache Miss Analysis

## ■ Assume:

- Matrix elements are doubles
- Cache block = 64 bytes = 8 doubles
- Cache size  $C \ll n$  (much smaller than  $n$ )

## ■ Other iterations:

- Again:  
 $n/8 + n = 9n/8$  misses  
 (omitting matrix  $c$ )



## ■ Total misses:

- $9n/8 * n^2 = (9/8) * n^3$

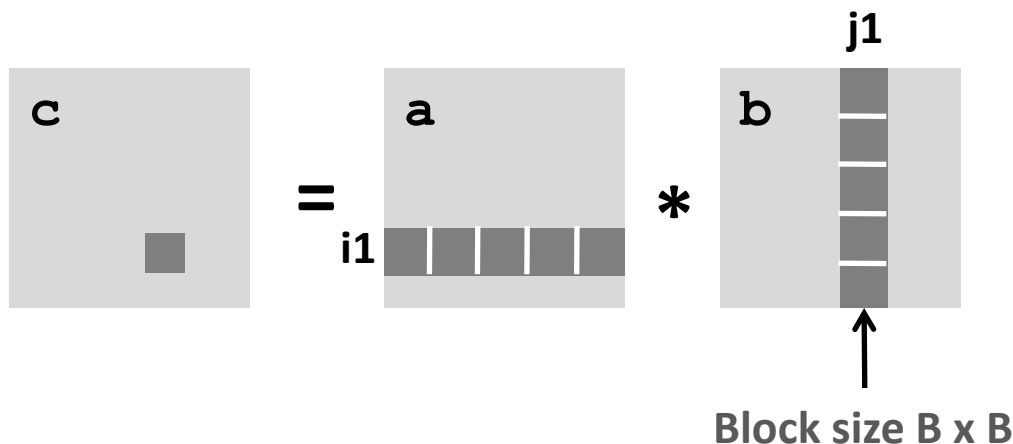
# Blocked Matrix Multiplication

```

c = (double *) calloc(sizeof(double), n*n);


/* Multiply n x n matrices a and b */
void mmm(double *a, double *b, double *c, int n) {
    int i, j, k;
    for (i = 0; i < n; i+=B)
        for (j = 0; j < n; j+=B)
            for (k = 0; k < n; k+=B)
                /* B x B mini matrix multiplications */
                for (i1 = i; i1 < i+B; i1++)
                    for (j1 = j; j1 < j+B; j1++)
                        for (k1 = k; k1 < k+B; k1++)
                            c[i1*n + j1] += a[i1*n + k1]*b[k1*n + j1];
}

```



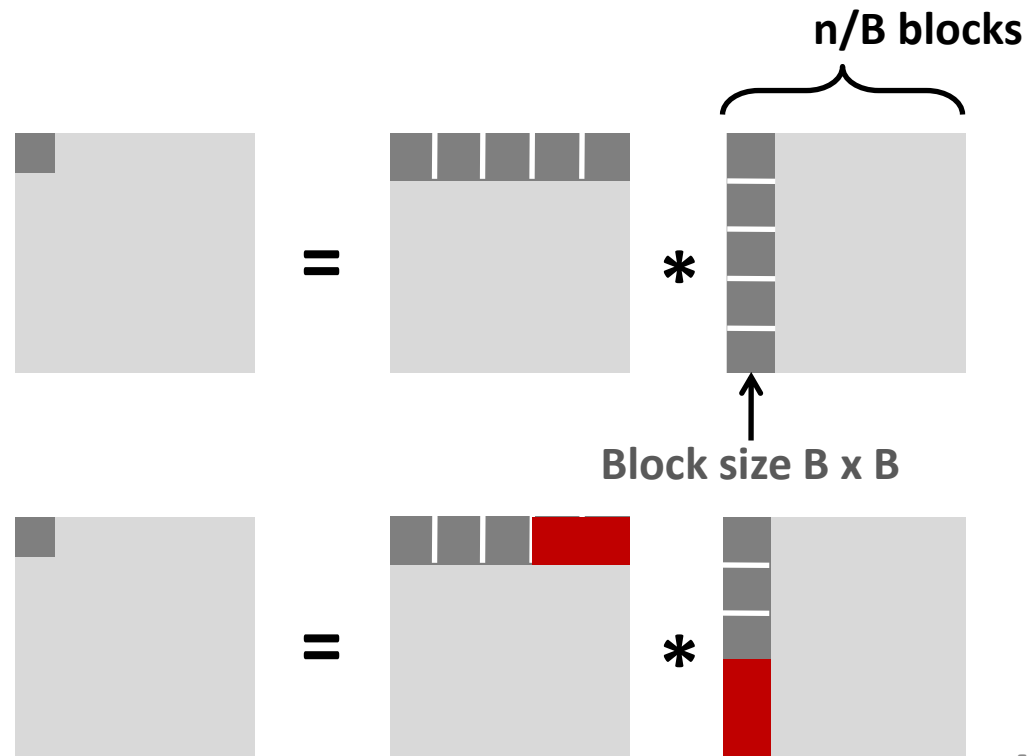
# Cache Miss Analysis

## ■ Assume:

- Cache block = 64 bytes = 8 doubles
- Cache size  $C \ll n$  (much smaller than  $n$ )
- Three blocks  fit into cache:  $3B^2 < C$

## ■ First (block) iteration:

- $B^2/8$  misses for each block
- $2n/B * B^2/8 = nB/4$   
(omitting matrix  $c$ )




- Afterwards in cache  
(schematic)



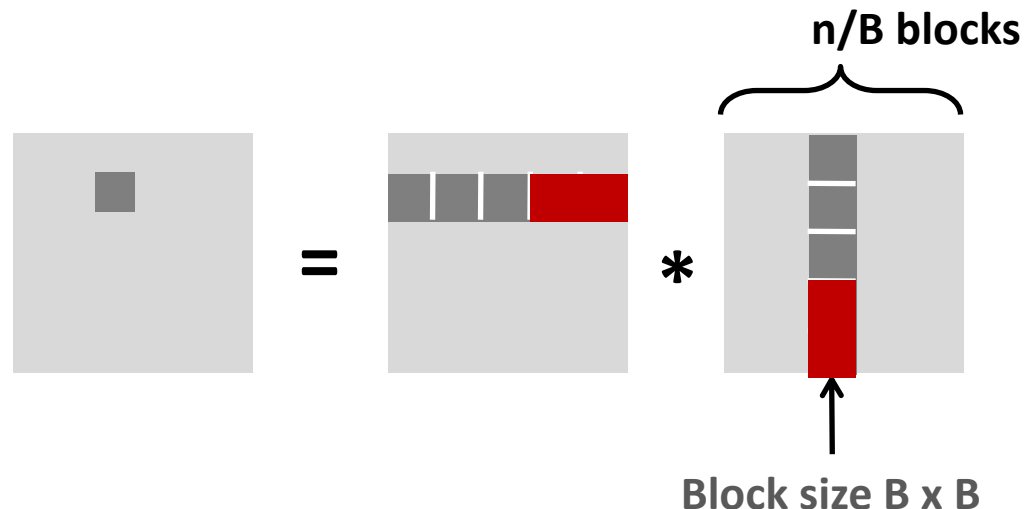
# Cache Miss Analysis

## ■ Assume:

- Cache block = 64 bytes = 8 doubles
- Cache size  $C \ll n$  (much smaller than  $n$ )
- Three blocks  fit into cache:  $3B^2 < C$

## ■ Other (block) iterations:

- Same as first iteration
- $2n/B * B^2/8 = nB/4$



## ■ Total misses:

- $nB/4 * (n/B)^2 = n^3/(4B)$

# Summary

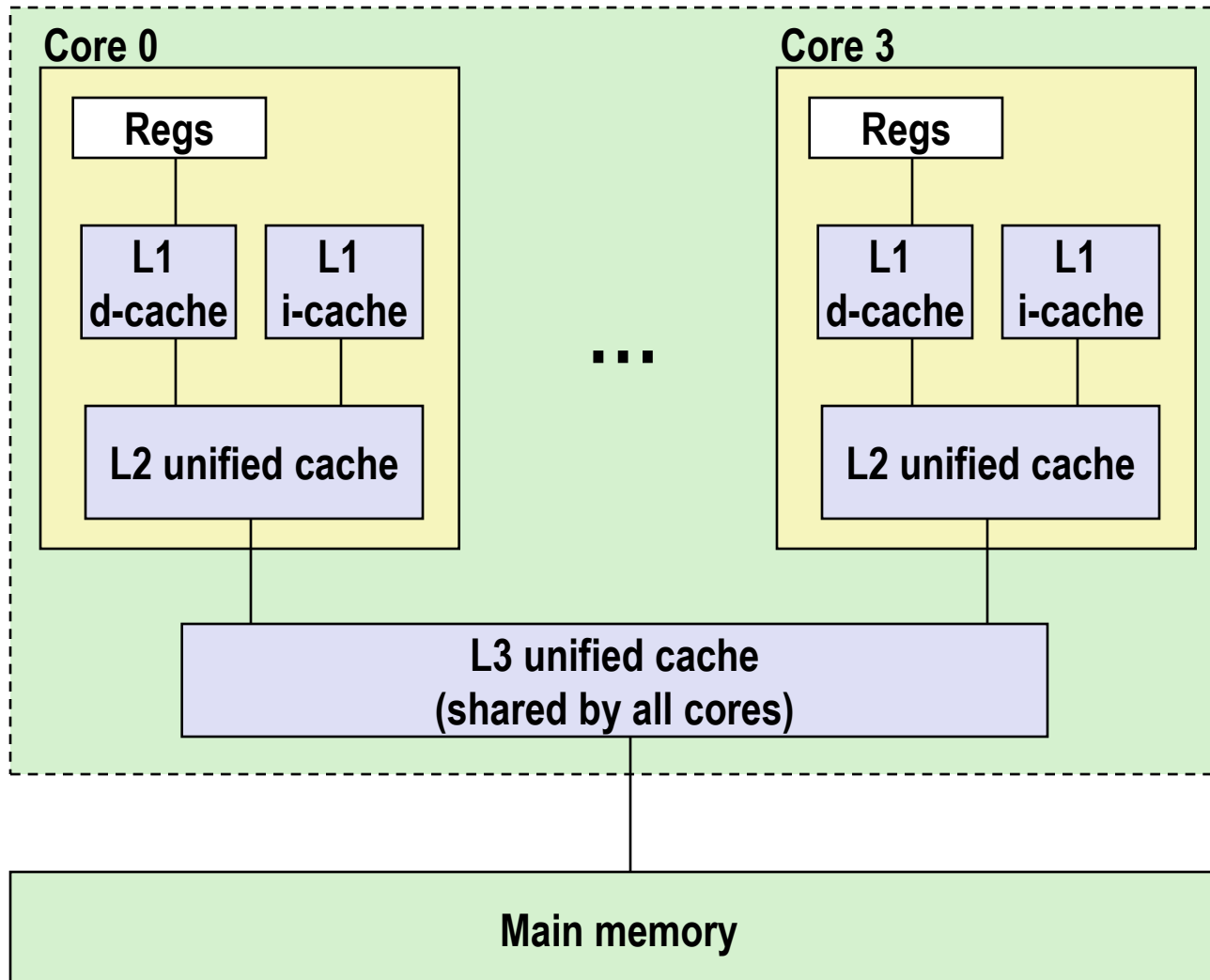
- No blocking:  $(9/8) * n^3$
- Blocking:  $1/(4B) * n^3$
- If  $B = 8$  difference is  $4 * 8 * 9 / 8 = 36x$
- If  $B = 16$  difference is  $4 * 16 * 9 / 8 = 72x$
- Suggests largest possible block size  $B$ , but limit  $3B^2 < C$ !
- Reason for dramatic difference:
  - Matrix multiplication has inherent temporal locality:
    - Input data:  $3n^2$ , computation  $2n^3$
    - Every array element used  $O(n)$  times!
  - But program has to be written properly

# Cache-Friendly Code

- **Programmer can optimize for cache performance**
  - How data structures are organized
  - How data are accessed
    - Nested loop structure
    - Blocking is a general technique
- **All systems favor “cache-friendly code”**
  - Getting absolute optimum performance is very platform specific
    - Cache sizes, line sizes, associativities, etc.
  - Can get most of the advantage with generic code
    - Keep working set reasonably small (temporal locality)
    - Use small strides (spatial locality)
    - Focus on inner loop code

# Intel Core i7 Cache Hierarchy

Processor package



**L1 i-cache and d-cache:**

32 KB, 8-way,  
Access: 4 cycles

**L2 unified cache:**

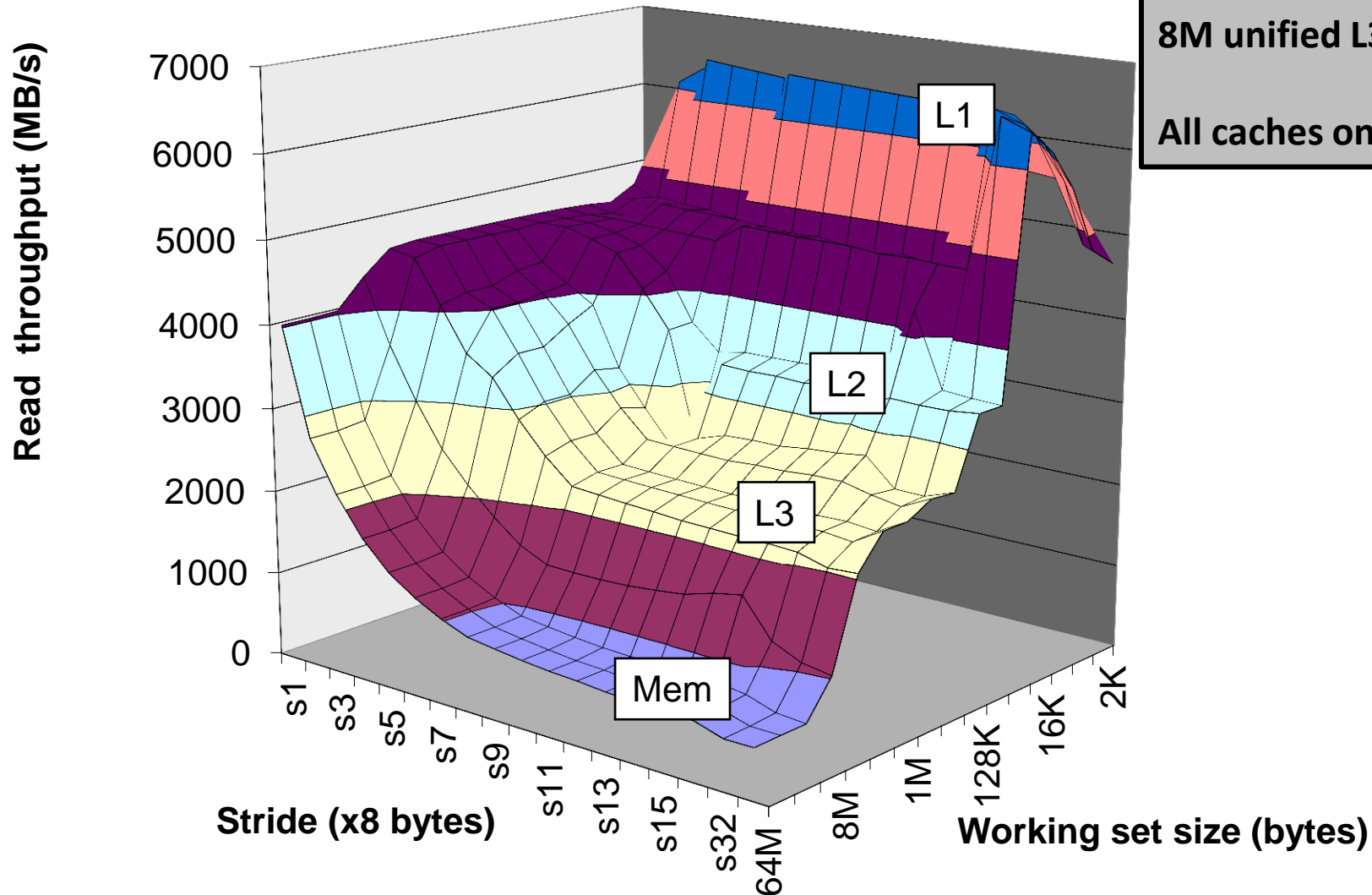
256 KB, 8-way,  
Access: 11 cycles

**L3 unified cache:**

8 MB, 16-way,  
Access: 30-40 cycles

**Block size:** 64 bytes for  
all caches.

# The Memory Mountain



Intel Core i7

32 KB L1 i-cache

32 KB L1 d-cache

256 KB unified L2 cache

8M unified L3 cache

All caches on-chip