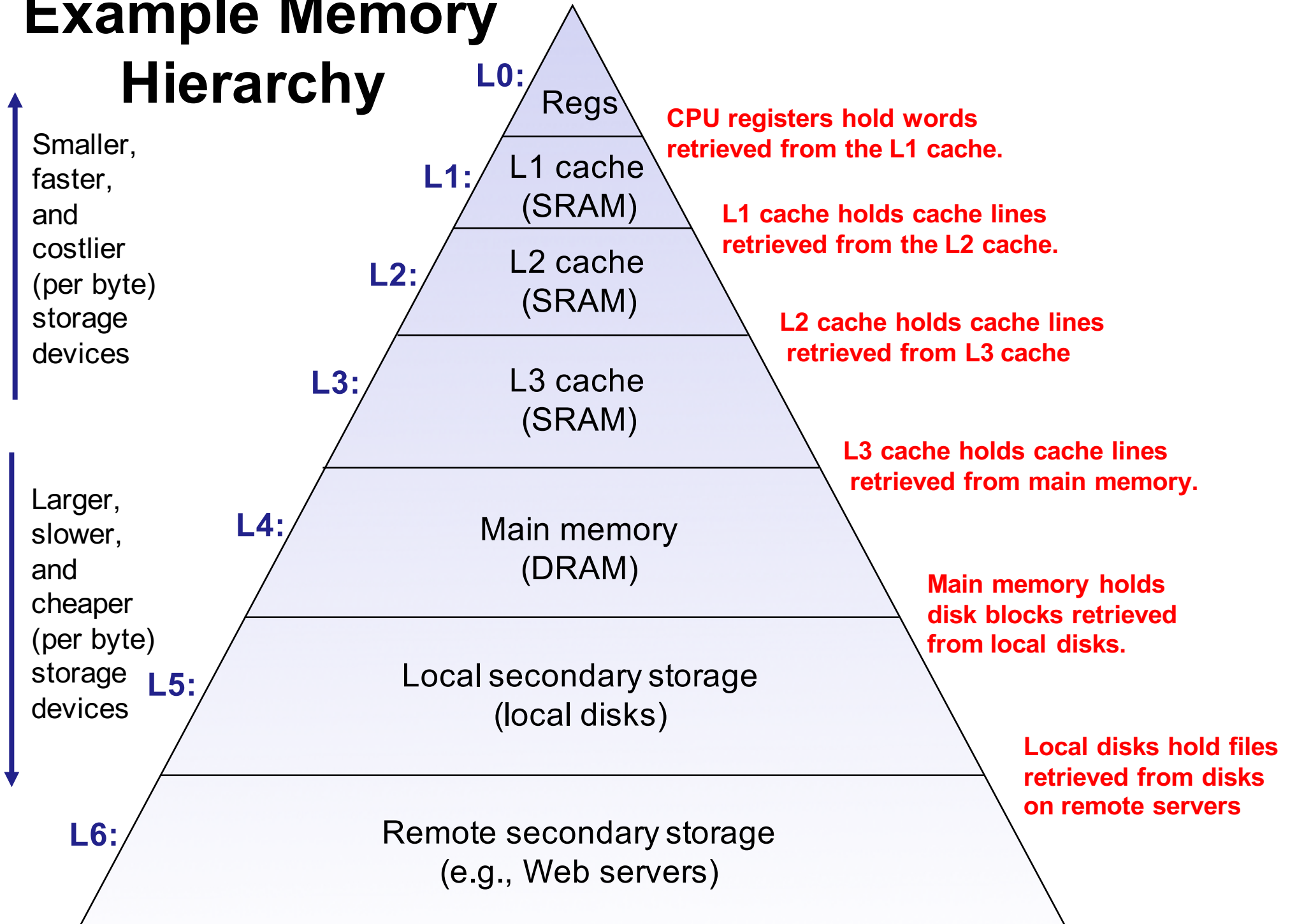# Cache Memories

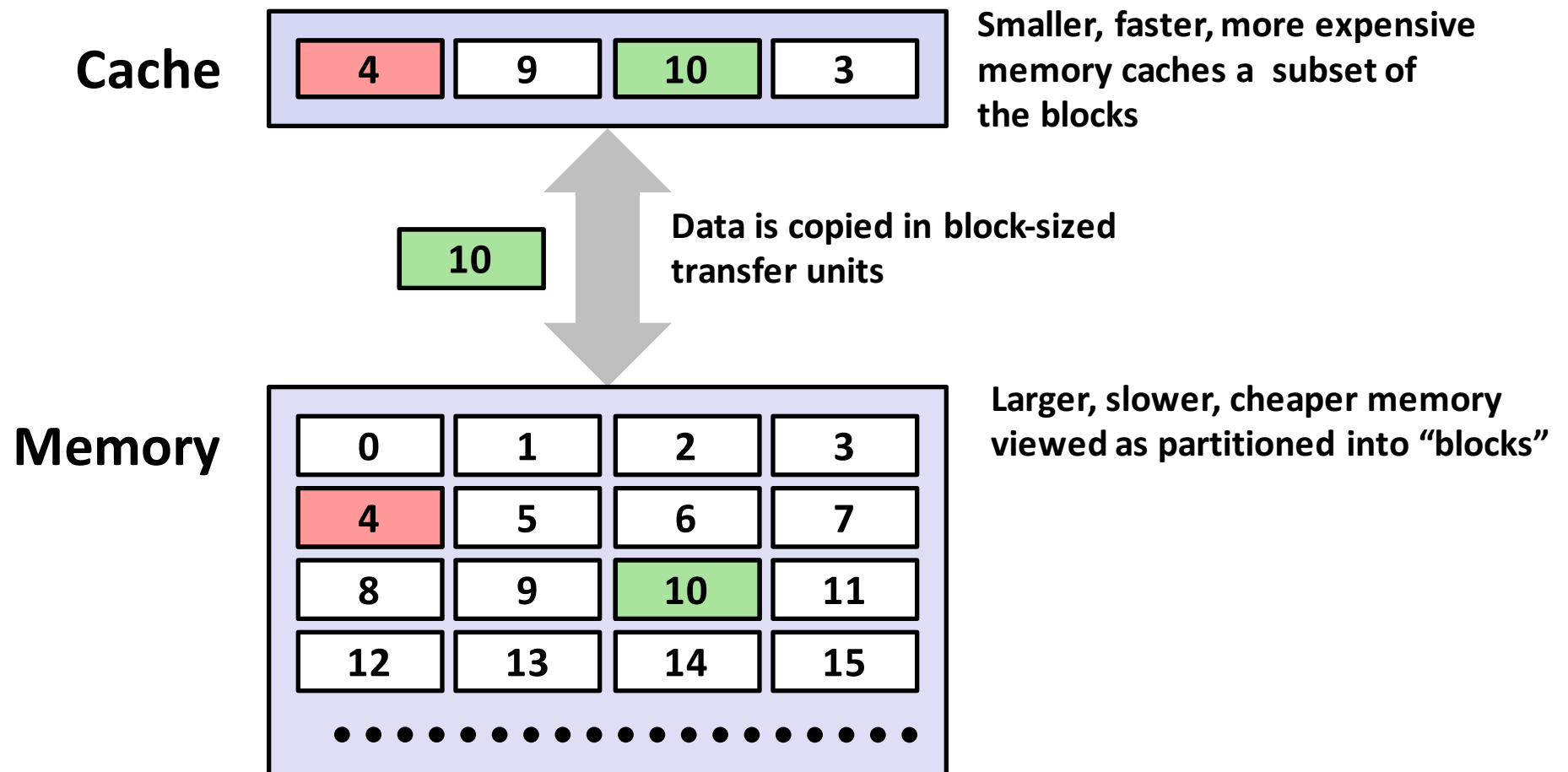From Bryant and O'Hallaron, *Computer Systems. A Programmer's Perspective*. Chapter 6.

# Today

- **Cache memory organization and operation**
- **Performance impact of caches**
  - The memory mountain
  - Rearranging loops to improve spatial locality
  - Using blocking to improve temporal locality
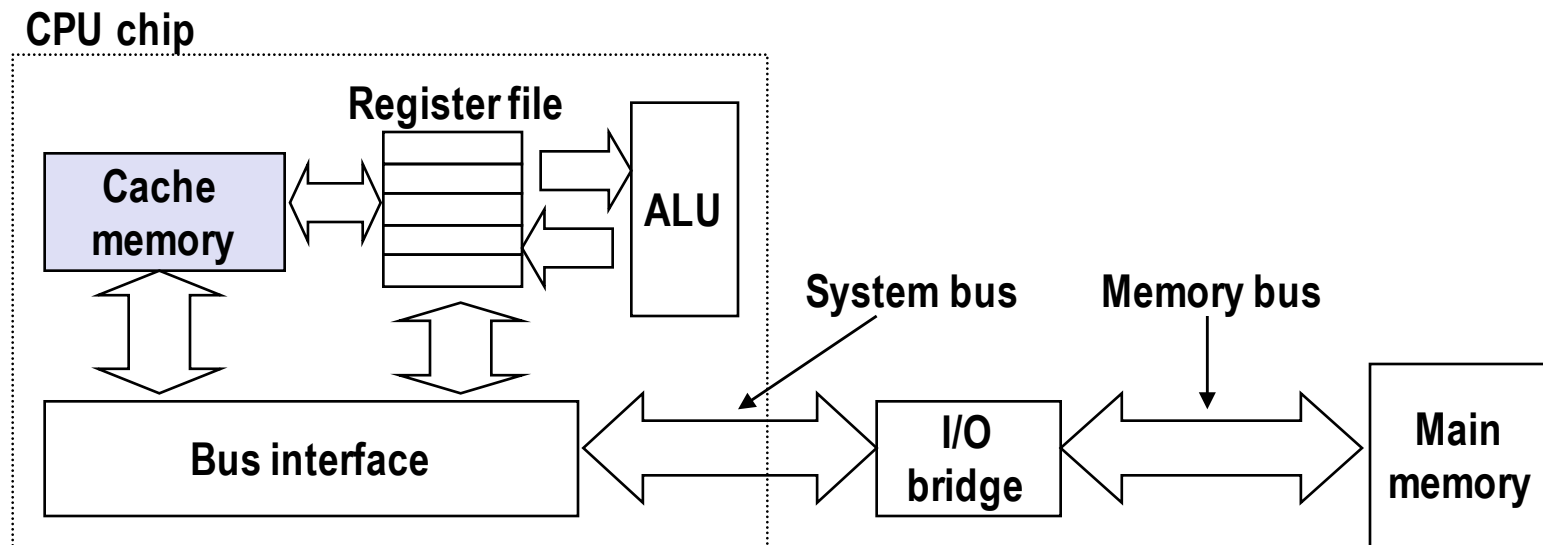
# Example Memory Hierarchy

Smaller, faster, and costlier (per byte) storage devices

Larger, slower, and cheaper (per byte) storage devices

L0: Regs

L1: L1 cache (SRAM)

L2: L2 cache (SRAM)

L3: L3 cache (SRAM)

L4: Main memory (DRAM)

L5: Local secondary storage (local disks)

L6: Remote secondary storage (e.g., Web servers)

CPU registers hold words retrieved from the L1 cache.

L1 cache holds cache lines retrieved from the L2 cache.

L2 cache holds cache lines retrieved from L3 cache

L3 cache holds cache lines retrieved from main memory.

Main memory holds disk blocks retrieved from local disks.

Local disks hold files retrieved from disks on remote servers

# General Cache Concept

**Cache**

| 4 | 9 | 10 | 3 |
|---|---|----|---|

Smaller, faster, more expensive memory caches a subset of the blocks

| 10 |
|----|

Data is copied in block-sized transfer units

**Memory**

| 0 | 1 | 2 | 3 |
|---|---|---|---|
| 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 |

Larger, slower, cheaper memory viewed as partitioned into "blocks"
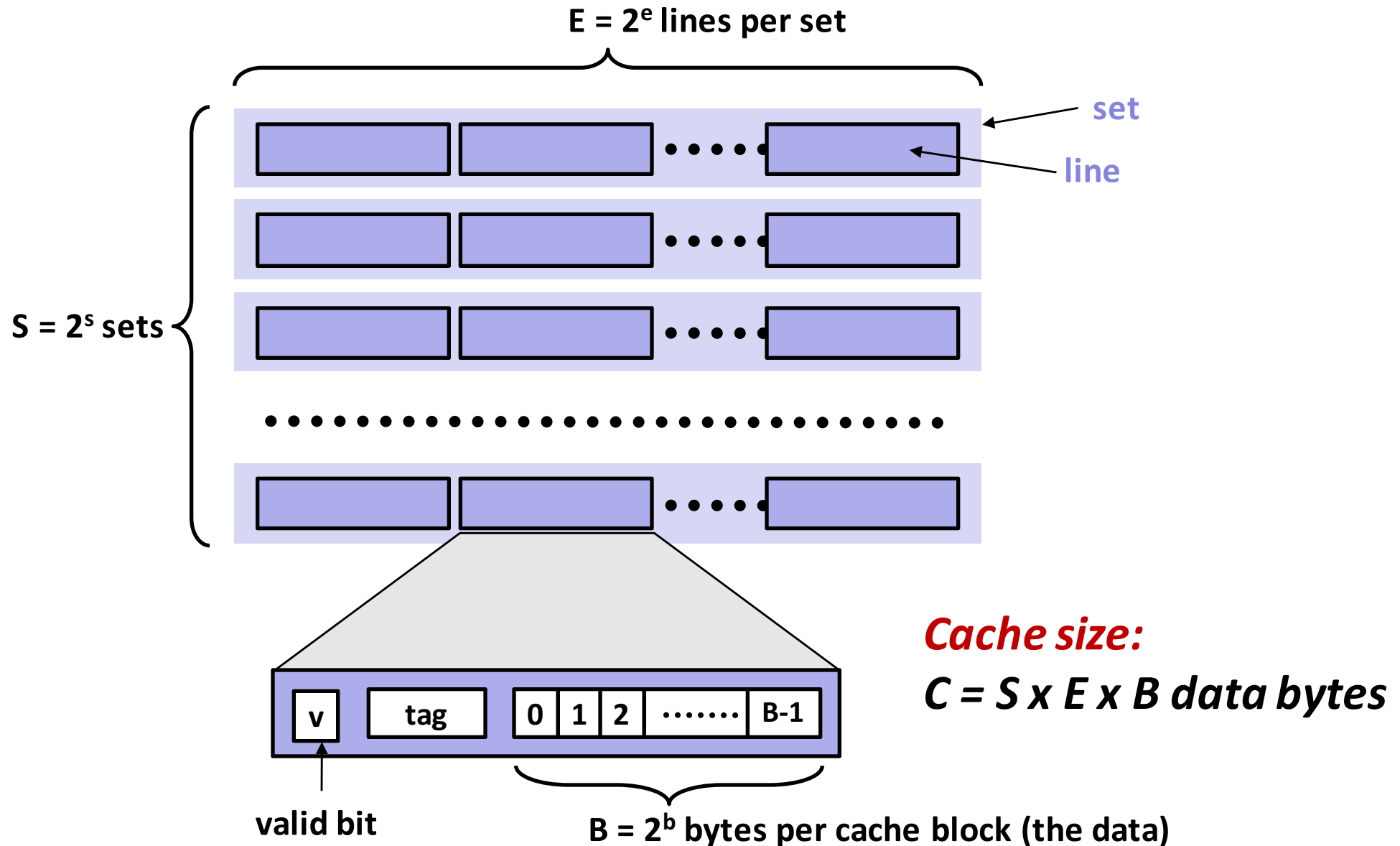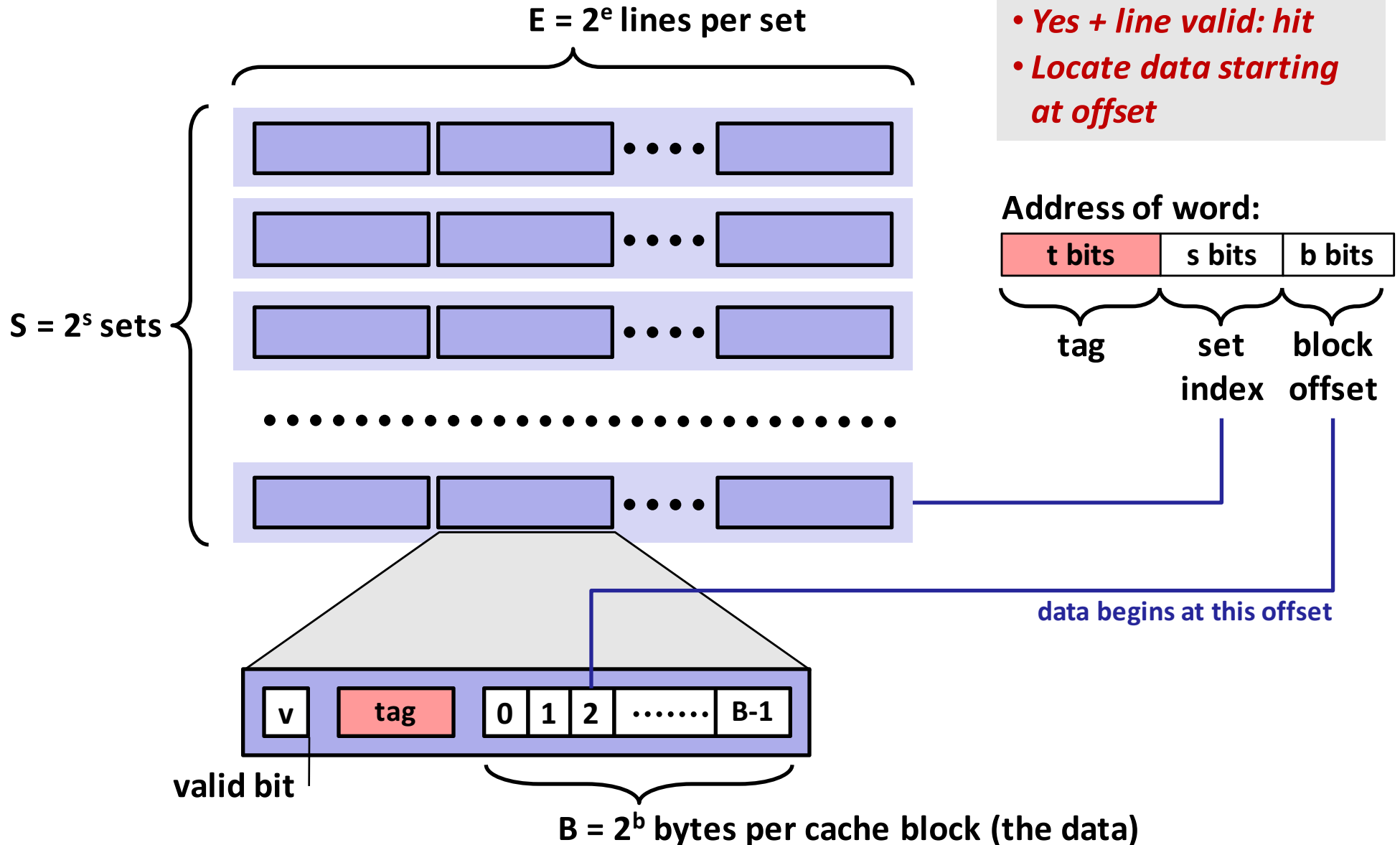
# Cache Memories

- **Cache memories are small, fast SRAM-based memories managed automatically in hardware**
    - Hold frequently accessed blocks of main memory
- **CPU looks first for data in cache**
- **Typical system structure:**
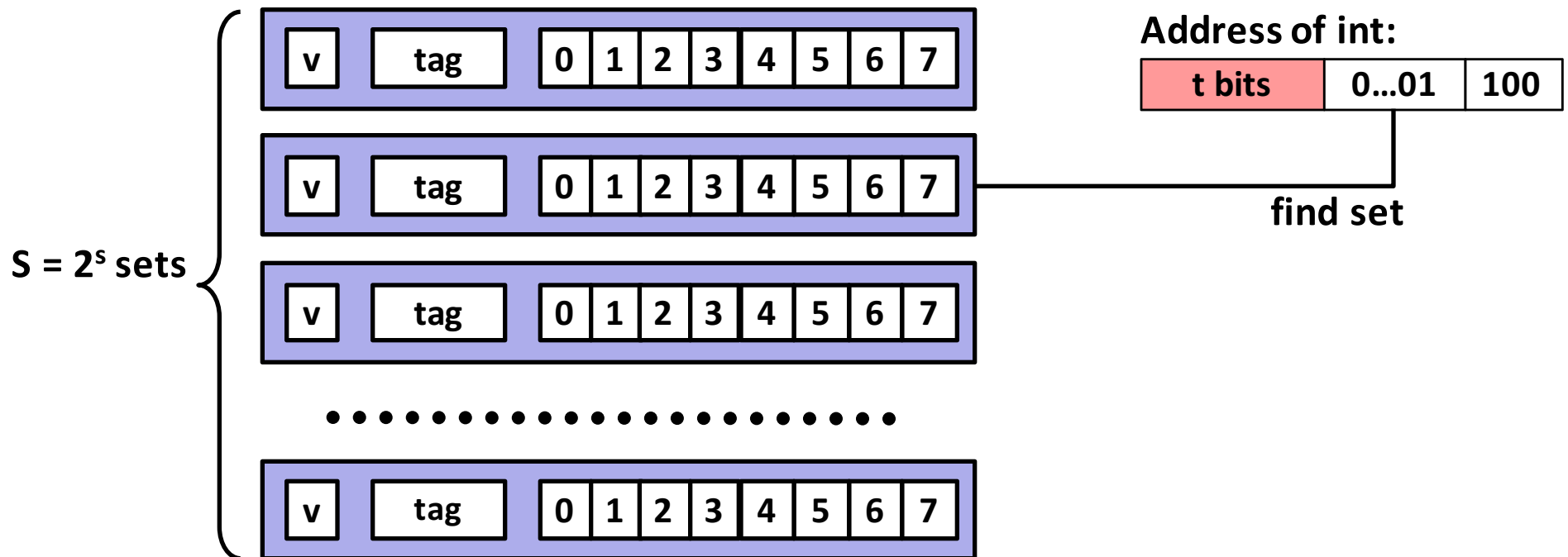
# General Cache Organization (S, E, B)

$E = 2^e$ lines per set

set

line

$S = 2^s$ sets

$$C = S \times E \times B \text{ data bytes}$$

Cache size:

| v | tag | 0 | 1 | 2 | ....... | B-1 |

valid bit

$B = 2^b$ bytes per cache block (the data)

# Cache Read

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

$E = 2^e$ lines per set

$S = 2^s$ sets

**Address of word:**

| t bits | s bits | b bits |
|--------|--------|--------|

tag      set index      block offset

data begins at this offset

| v | tag | 0 | 1 | 2 | ....... | B-1 |
|---|-----|---|---|---|---------|-----|

**valid bit**

$B = 2^b$ bytes per cache block (the data)

# Example: Direct Mapped Cache (E = 1)

**Direct mapped: One line per set**
**Assume: cache block size 8 bytes**

| v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
| v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

$S = 2^s$ sets

| v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

• • • • • • • • • • • • • • • • • • • • • • • •

| v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

**Address of int:**

| t bits | 0...01 | 100 |

**find set**

# Example: Direct Mapped Cache (E = 1)

**Direct mapped: One line per set**
**Assume: cache block size 8 bytes**

valid?   +   match: assume yes = hit

Address of int:

| t bits | 0...01 | 100 |

| v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

block offset

# Example: Direct Mapped Cache (E = 1)

Direct mapped:  One line per set
Assume: cache block size 8 bytes

valid?   +  match: assume yes = hit

Address of int:

| t bits | 0...01 | 100 |

| v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

block offset

int (4 Bytes) is here

**If tag doesn't match:** old line is evicted and replaced

# Direct-Mapped Cache Simulation

| t=1 | s=2 | b=1 |
|-----|-----|-----|
| x | xx | x |

M=16 bytes (4-bit addresses),  B=2 bytes/block,
S=4 sets, E=1 Blocks/set

Address trace (reads, one byte per read):

| 0 | [$0\underline{00}0_2$], | miss |
|---|---|---|
| 1 | [$0\underline{00}1_2$], | hit |
| 7 | [$0\underline{11}1_2$], | miss |
| 8 | [$1\underline{00}0_2$], | miss |
| 0 | [$0\underline{00}0_2$] | miss |

|       | v | Tag | Block |
|-------|---|-----|-------|
| Set 0 | 1 | 0   | M[0-1] |
| Set 1 |   |     |       |
| Set 2 |   |     |       |
| Set 3 | 1 | 0   | M[6-7] |

# 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 |
|---|---|---|



find set

# 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 |
|--------|--------|-----|

compare both

valid? + match: yes = hit

| v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |    | v | tag | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

block offset

# 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 |

compare both

valid?  +  match: yes = hit

| v | tag | 0 1 2 3 4 5 6 7 | | v | tag | 0 1 2 3 4 5 6 7 |

short int (2 Bytes) is here

block offset

**No match:**
- **One line in set is selected for eviction and replacement**
- **Replacement policies: random, least recently used (LRU), …**

# 2-Way Set Associative Cache Simulation

t=2    s=1    b=1

| xx | x | x |
|----|---|---|

M=16 byte addresses, B=2 bytes/block,
S=2 sets, E=2 blocks/set

Address trace (reads, one byte per read):

| 0 | [$0000_2$], | miss |
|---|------------|------|
| 1 | [$0001_2$], | hit |
| 7 | [$0111_2$], | miss |
| 8 | [$1000_2$], | miss |
| 0 | [$0000_2$] | hit |

| | v | Tag | Block |
|--------|---|-----|--------|
| Set 0 | 1 | 00 | M[0-1] |
| | 1 | 10 | M[8-9] |
| Set 1 | 1 | 01 | M[6-7] |
| | 0 | | |

# What about writes?

- **Multiple copies of data exist:**
  - L1, L2, L3, Main Memory, Disk

- **What to do on a write-hit?**
  - Write-through (write immediately to memory)
  - Write-back (defer write to memory until replacement of line)
    - Need a dirty bit (line different from memory or not)

- **What to do on a write-miss?**
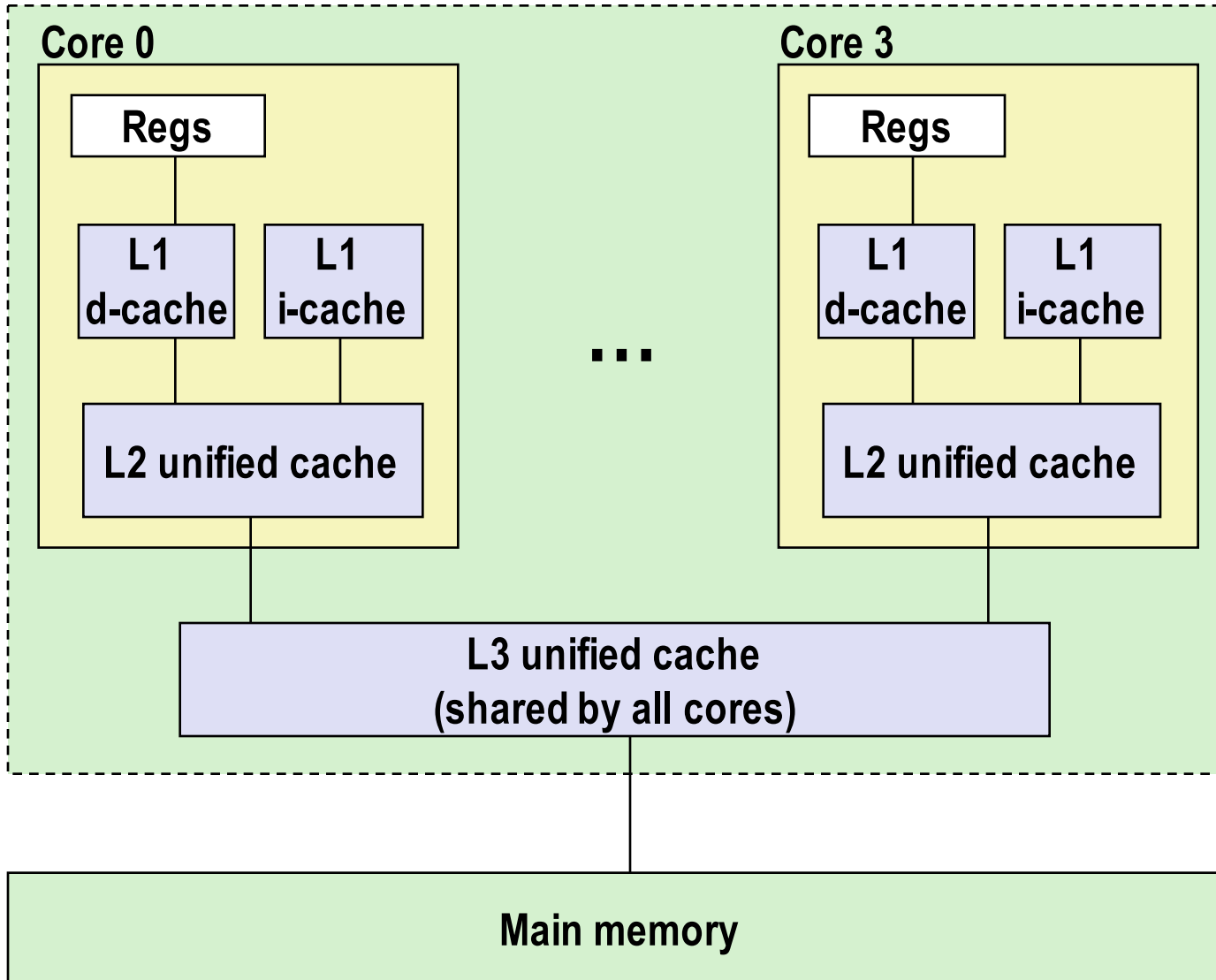  - Write-allocate (load into cache, update line in cache)
    - Good if more writes to the location follow
  - No-write-allocate (writes straight to memory, does not load into cache)

- **Typical**
  - Write-through + No-write-allocate
  - **Write-back + Write-allocate**

# 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:  10 cycles

**L3 unified cache:**
   8 MB, 16-way,
   Access:  40-75 cycles

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

# Cache Performance Metrics

- **Miss Rate**
  - Fraction of memory references not found in cache (misses / accesses) = 1 – hit rate
  - Typical numbers (in percentages):
    - 3-10% for L1
    - can be quite small (e.g., < 1%) for L2, depending on size, etc.
- **Hit Time**
  - Time to deliver a line in the cache to the processor
    - includes time to determine whether the line is in the cache
  - Typical numbers:
    - 4 clock cycle for L1
    - 10 clock cycles for L2
- **Miss Penalty**
  - Additional time required because of a miss
    - typically 50-200 cycles for main memory (Trend: increasing!)

# Let's think about those numbers

- **Huge difference between a hit and a miss**
  - Could be 100x, if just L1 and main memory

- **Would you believe 99% hits is twice as good as 97%?**
  - Consider:
    cache hit time of 1 cycle
    miss penalty of 100 cycles

  - Average access time:
    97% hits:  1 cycle + 0.03 * 100 cycles = **4 cycles**
    99% hits:  1 cycle + 0.01 * 100 cycles = **2 cycles**

- **This is why "miss rate" is used instead of "hit rate"**

# Writing Cache Friendly Code

- **Make the common case go fast**
  - Focus on the inner loops of the core functions

- **Minimize the misses in the inner loops**
  - Repeated references to variables are good (temporal locality)
  - Stride-1 reference patterns are good (spatial locality)

**Key idea: Our qualitative notion of locality is quantified through our understanding of cache memories**

# The Memory Mountain

- **Read throughput (read bandwidth)**
  - Number of bytes read from memory per second (MB/s)

- **Memory mountain: Measured read throughput as a function of spatial and temporal locality.**
  - Compact way to characterize memory system performance.

# Memory Mountain Test Function

```c
long data[MAXELEMS];  /* Global array to traverse */

/* test - Iterate over first "elems" elements of
 *        array "data" with stride of "stride", using
 *        using 4x4 loop unrolling.
 */
int test(int elems, int stride) {
    long i, sx2=stride*2, sx3=stride*3, sx4=stride*4;
    long acc0 = 0, acc1 = 0, acc2 = 0, acc3 = 0;
    long length = elems, limit = length - sx4;

    /* Combine 4 elements at a time */
    for (i = 0; i < limit; i += sx4) {
        acc0 = acc0 + data[i];
        acc1 = acc1 + data[i+stride];
        acc2 = acc2 + data[i+sx2];
        acc3 = acc3 + data[i+sx3];
    }

    /* Finish any remaining elements */
    for (; i < length; i++) {
        acc0 = acc0 + data[i];
    }
    return ((acc0 + acc1) + (acc2 + acc3));
}
```

*mountain/mountain.c*

Call `test()` with many combinations of `elems` and `stride`.
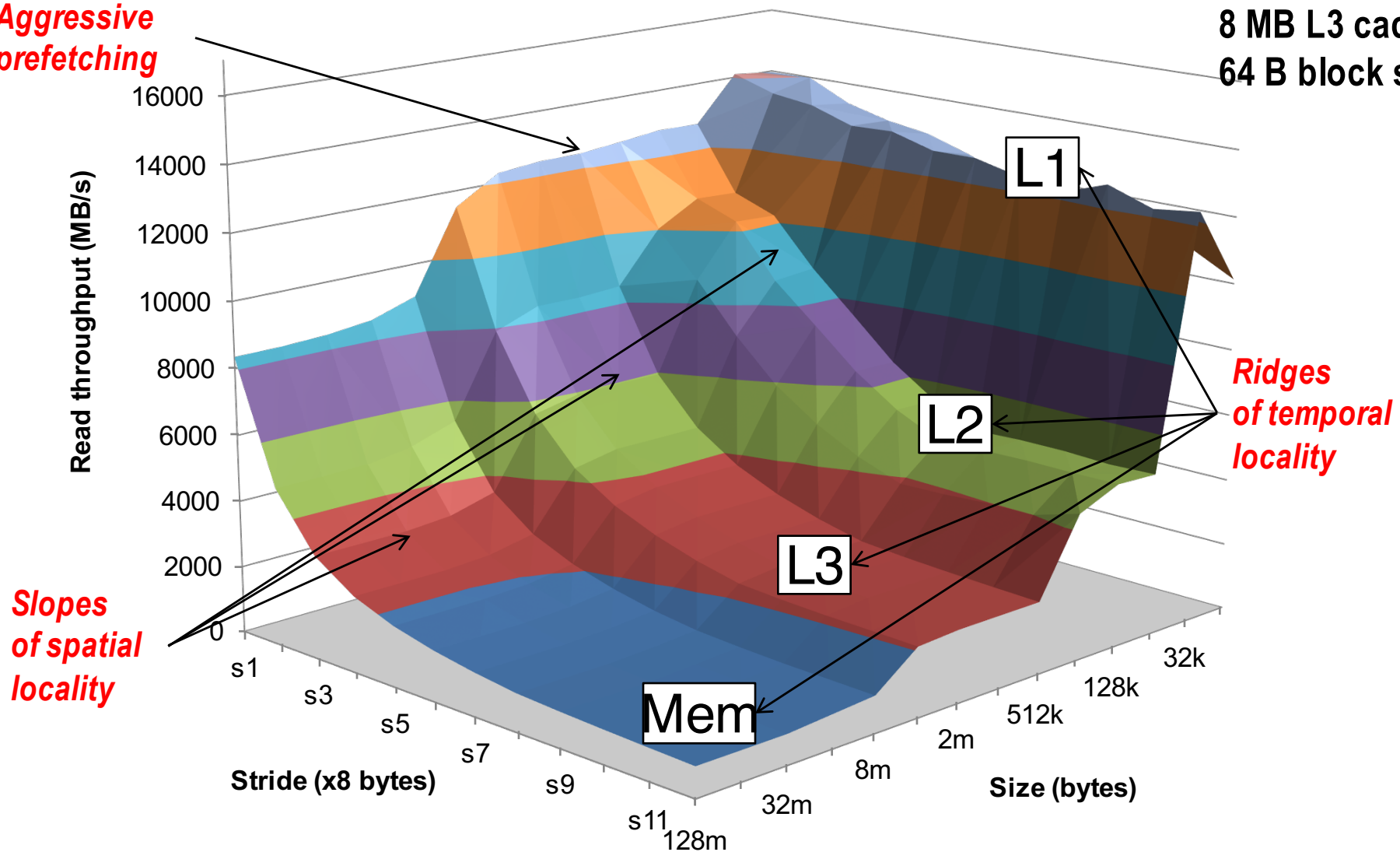
For each `elems` and `stride`:

1. Call `test()` once to warm up the caches.

2. Call `test()` again and measure the read throughput(MB/s)

# The Memory Mountain

Core i7 Haswell
2.1 GHz
32 KB L1 d-cache
256 KB L2 cache
8 MB L3 cache
64 B block size

**Aggressive prefetching**

**Slopes of spatial locality**

**Ridges of temporal locality**

L1

L2

L3

Mem

Read throughput (MB/s)

16000
14000
12000
10000
8000
6000
4000
2000
0

Stride (x8 bytes)

s1
s3
s5
s7
s9
s11

Size (bytes)

128m
32m
8m
2m
512k
128k
32k

# Writing Cache-Friendly Code

- **Cache memories can have significant performance impact**

- **You can write your programs to exploit this**
  - Focus on the inner loops, where bulk of computations and memory accesses occur.
  - Try to maximize spatial locality by reading data objects with sequentially with stride 1.
  - Try to maximize temporal locality by using a data object as often as possible once it's read from memory.