# The Memory Hierarchy

# Random-Access Memory (RAM)

- **Key features**
  - RAM is traditionally packaged as a chip.
  - Basic storage unit is normally a cell (one bit per cell).
  - Multiple RAM chips form a memory.

- **RAM comes in two varieties:**
  - SRAM (Static RAM)
  - DRAM (Dynamic RAM)

# SRAM vs DRAM Summary

| | Trans. per bit | Access time | Needs refresh? | Needs EDC? | Cost | Applications |
|---|---|---|---|---|---|---|
| SRAM | 4 or 6 | 1X | No | Maybe | 100x | Cache memories |
| DRAM | 1 | 10X | Yes | Yes | 1X | Main memories, frame buffers |

# Nonvolatile Memories

- **DRAM and SRAM are volatile memories**
    - Lose information if powered off.
- **Nonvolatile memories retain value even if powered off**
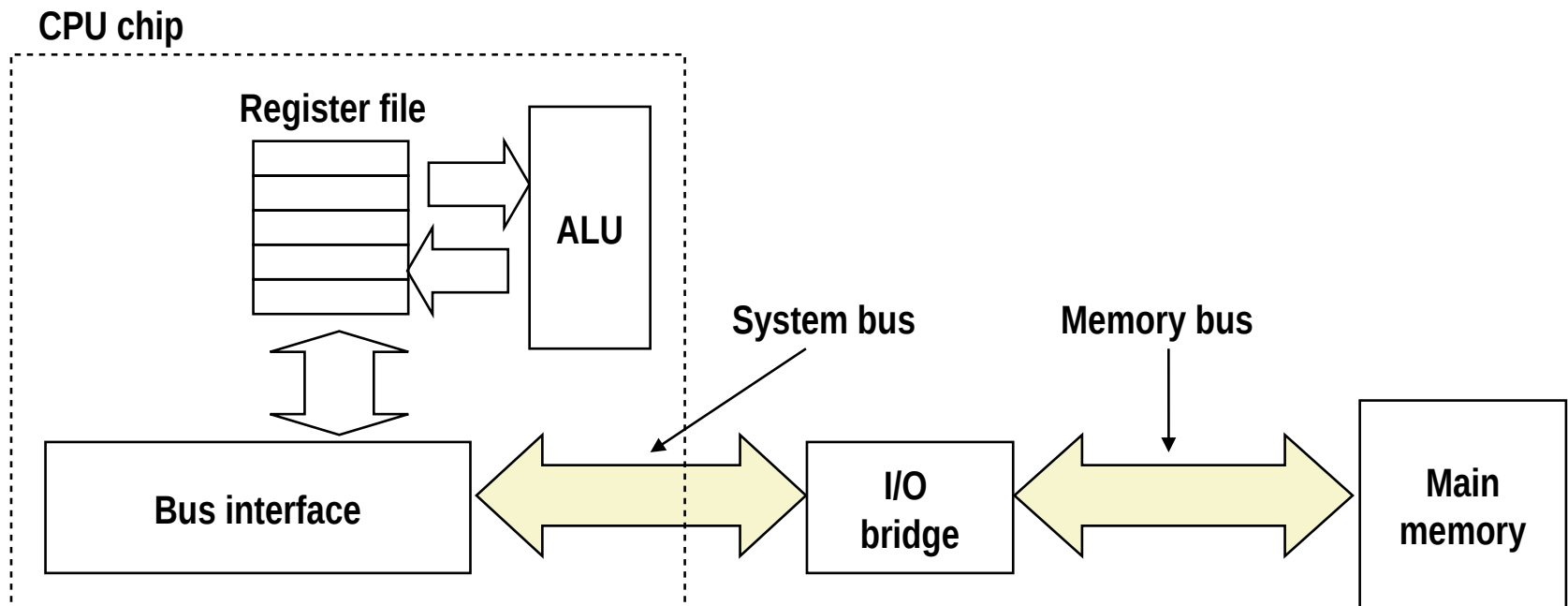    - Read-only memory (ROM): programmed during production
    - Programmable ROM (PROM): can be programmed once
    - Eraseable PROM (EPROM): can be bulk erased (UV, X-Ray)
    - Electrically eraseable PROM (EEPROM): electronic erase capability
    - Flash memory: EEPROMs. with partial (block-level) erase capability
        - Wears out after about 100,000 erasings
- **Uses for Nonvolatile Memories**
    - Firmware programs stored in a ROM (BIOS, controllers for disks, network cards, graphics accelerators, security subsystems,…)
    - Solid state disks (replace rotating disks in thumb drives, smart phones, mp3 players, tablets, laptops,…)
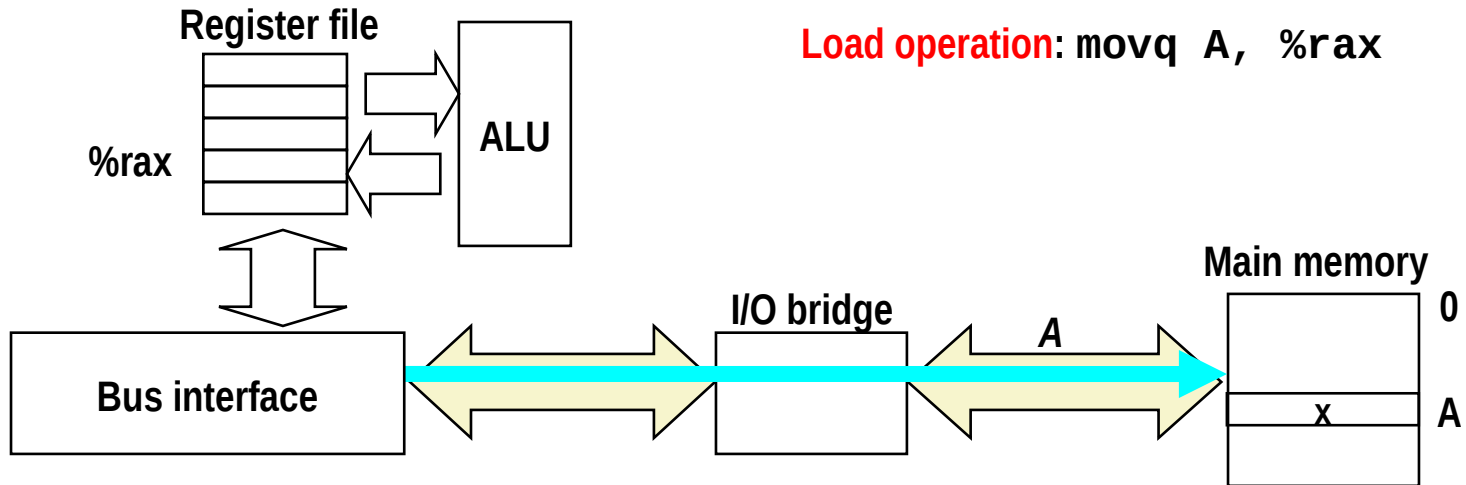    - Disk caches

# Traditional Bus Structure Connecting CPU and Memory

- A **bus** is a collection of parallel wires that carry address, data, and control signals.
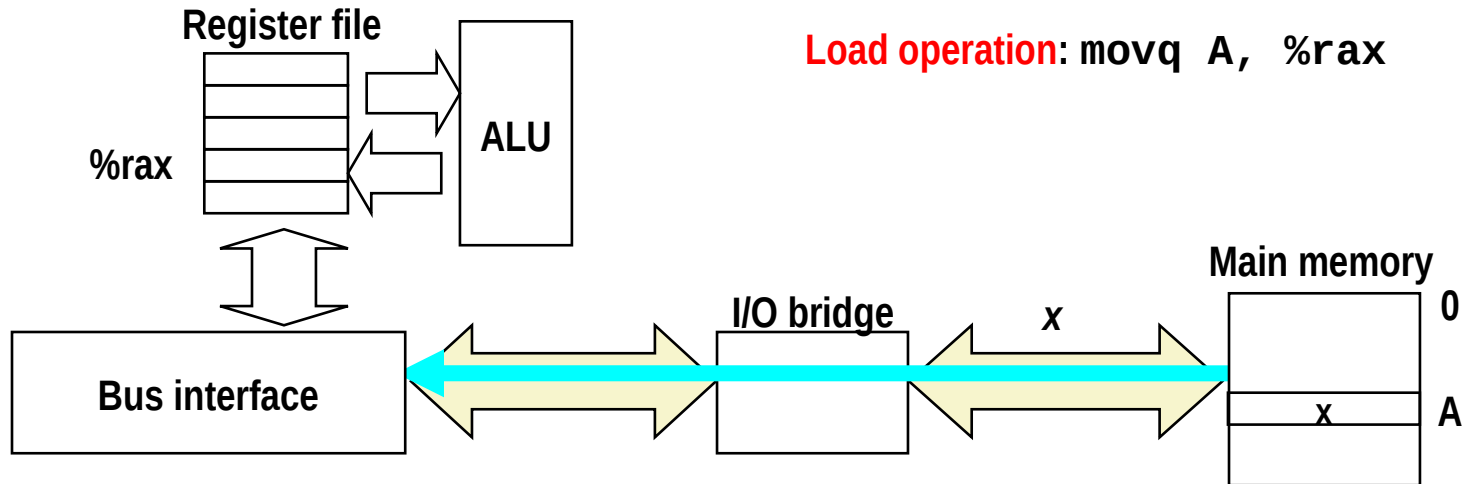- Buses are typically shared by multiple devices.

# Memory Read Transaction (1)

- **CPU places address A on the memory bus.**

**Register file**

**%rax**

**ALU**

**Load operation**: `movq A, %rax`
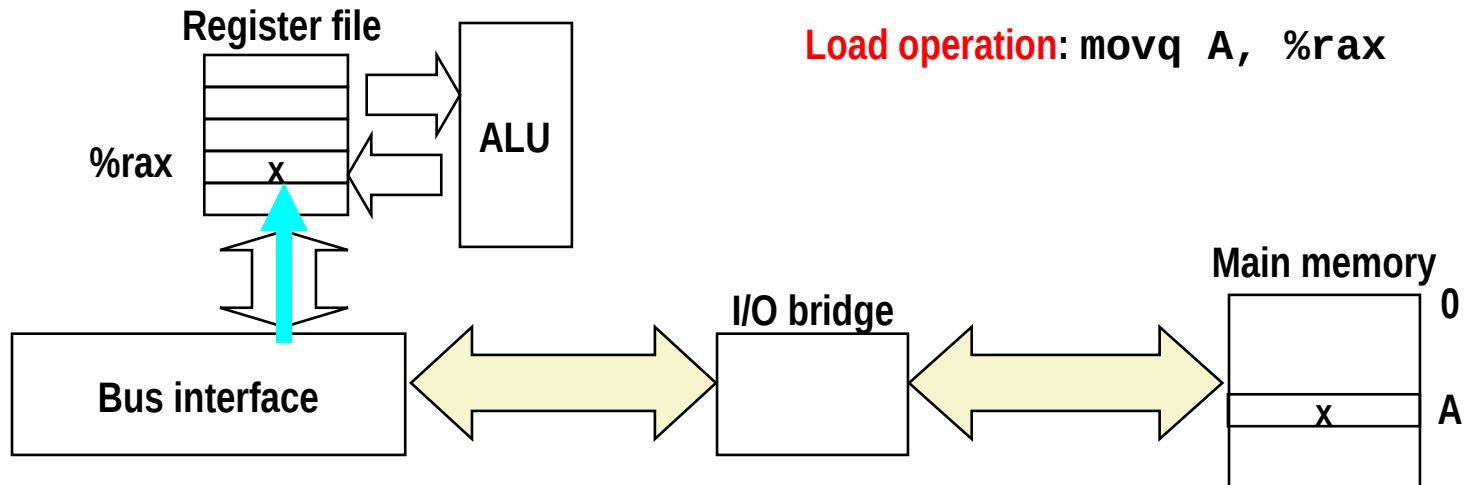
**Bus interface**

**I/O bridge**

**A**

**Main memory**

0

x    A

# Memory Read Transaction (2)

■ **Main memory reads A from the memory bus, retrieves word x, and places it on the bus.**

**Register file**

**%rax**

**ALU**

**Load operation**: `movq A, %rax`

**Bus interface**

**I/O bridge**

*x*

**Main memory**

0

x    A

# Memory Read Transaction (3)

- **CPU read word x from the bus and copies it into register %rax.**

Register file

Load operation: `movq A, %rax`

ALU

%rax | x

Main memory

0

Bus interface

I/O bridge

x | A

# Memory Write Transaction (1)

■ **CPU places address A on bus. Main memory reads it and waits for the corresponding data word to arrive.**

**Register file**

**%rax**

y

**ALU**

**Store operation:** `movq %rax, A`

**Bus interface**

**I/O bridge**

A

**Main memory**

0

A

# Memory Write Transaction (2)

■ **CPU places data word y on the bus.**

**Register file**

**Store operation:** `movq %rax, A`

**%rax** y

**ALU**

**Main memory**

**Bus interface**

**I/O bridge** y

0

A

# Memory Write Transaction (3)

- **Main memory reads data word y from the bus and stores it at address A.**

**Register file**

**ALU**

**%rax**  | y |

**Store operation**: `movq %rax, A`

**Bus interface**

**I/O bridge**
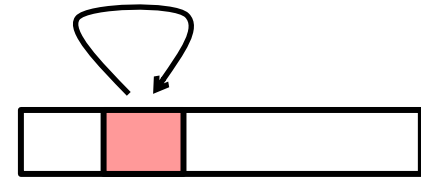
**main memory**

0

y | A

# Locality

- **Principle of Locality: Programs tend to use data and instructions with addresses near or equal to those they have used recently**


- **Temporal locality:**
  - Recently referenced items are likely
    to be referenced again in the near future


- **Spatial locality:**
  - Items with nearby addresses tend
    to be referenced close together in time

# Locality Example

```
sum = 0;
for (i = 0; i < n; i++)
    sum += a[i];
return sum;
```

- **Data references**
  - Reference array elements in succession (stride-1 reference pattern).
  - Reference variable `sum` each iteration.
- **Instruction references**
  - Reference instructions in sequence.
  - Cycle through loop repeatedly.

**Spatial locality**

**Temporal locality**

**Spatial locality**

**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**

14

# Caches

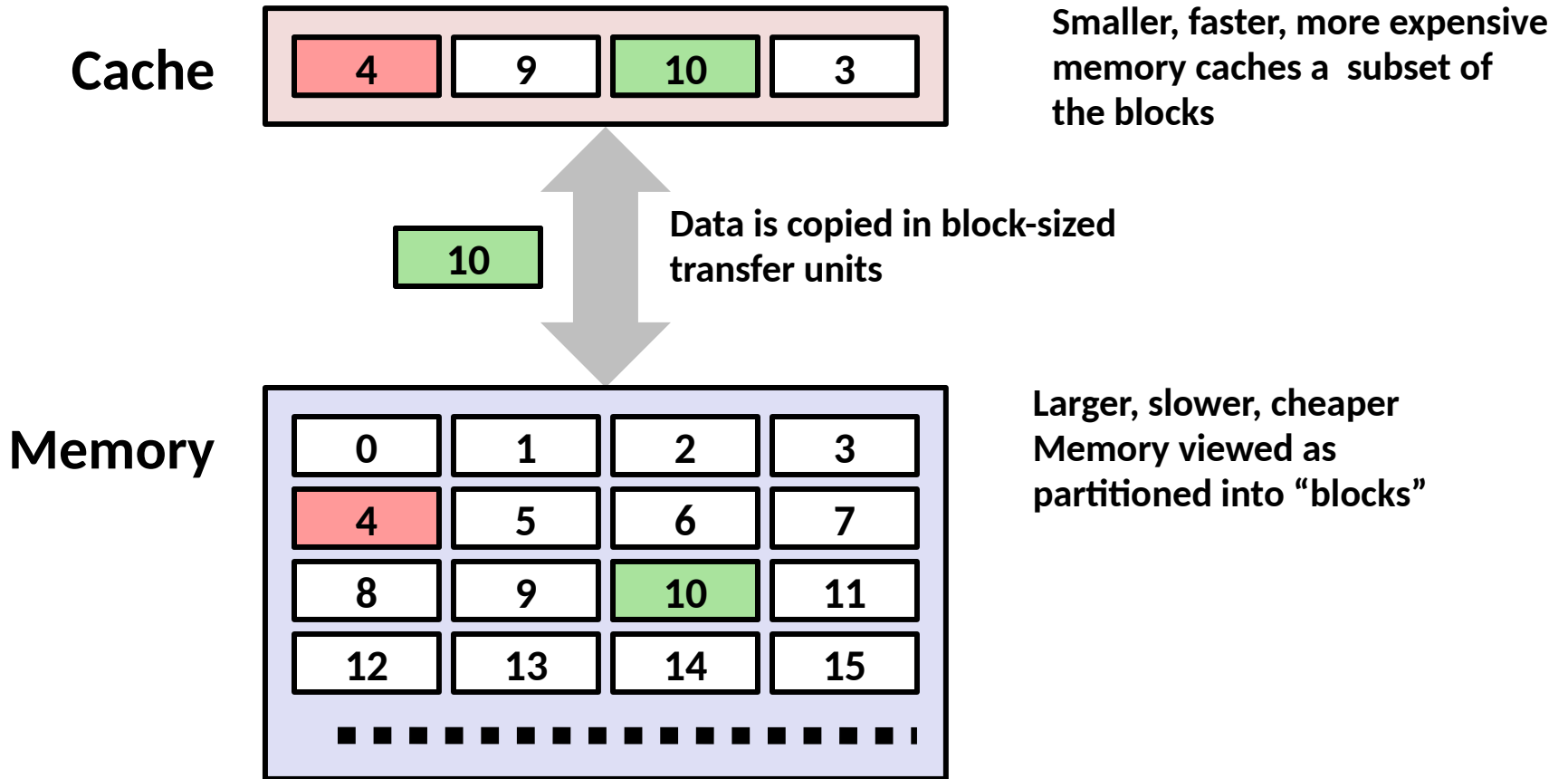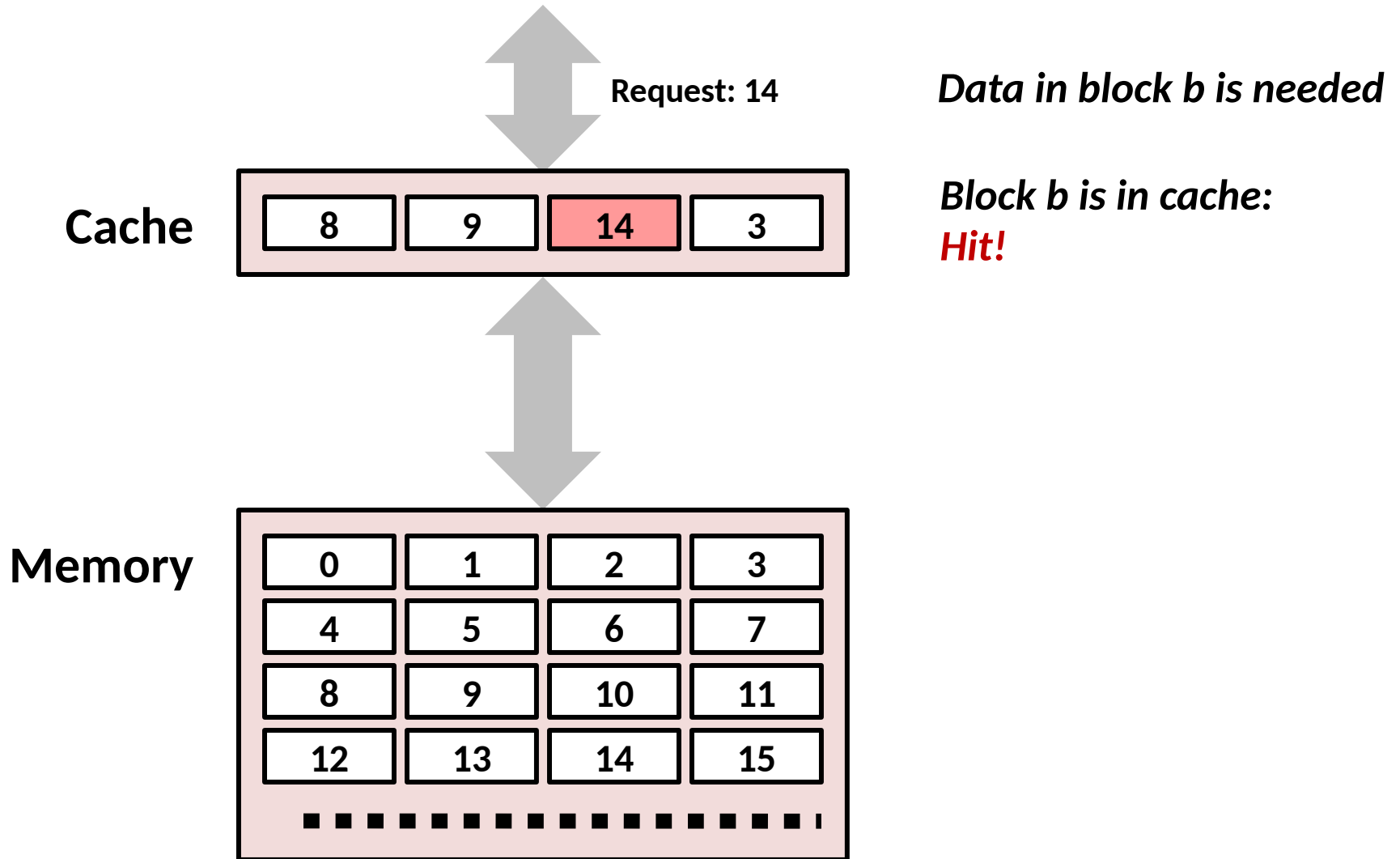- *Cache:* **A smaller, faster storage device that acts as a staging area for a subset of the data in a larger, slower device.**
- **Fundamental idea of a memory hierarchy:**
  - For each k, the faster, smaller device at level k serves as a cache for the larger, slower device at level k+1.
- **Why do memory hierarchies work?**
  - Because of locality, programs tend to access the data at level k more often than they access the data at level k+1.
  - Thus, the storage at level k+1 can be slower, and thus larger and cheaper per bit.
- *Big Idea:* **The memory hierarchy creates a large pool of storage that costs as much as the cheap storage near the bottom, but that serves data to programs at the rate of the fast storage near the top.**

# General Cache Concepts

**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"

# General Cache Concepts: Hit

**Request: 14**

*Data in block b is needed*

**Cache**

| 8 | 9 | 14 | 3 |

*Block b is in cache:*
*Hit!*

**Memory**

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

# General Cache Concepts: Miss

Request: 12

| 8 | 12 | 14 | 3 |

Cache

12

Request: 12

Memory

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

*Data in block b is needed*

*Block b is not in cache:*
*Miss!*

*Block b is fetched from memory*

*Block b is stored in cache*
- Placement policy: determines where b goes
- Replacement policy: determines which block gets evicted (victim)

# General Caching Concepts: Types of Cache Misses

- ## Cold (compulsory) miss
  - Cold misses occur because the cache is empty.
- ## Conflict miss
  - Most caches limit blocks at level k+1 to a small subset (sometimes a singleton) of the block positions at level k.
    - E.g. Block i at level k+1 must be placed in block (i mod 4) at level k.
  - Conflict misses occur when the level k cache is large enough, but multiple data objects all map to the same level k block.
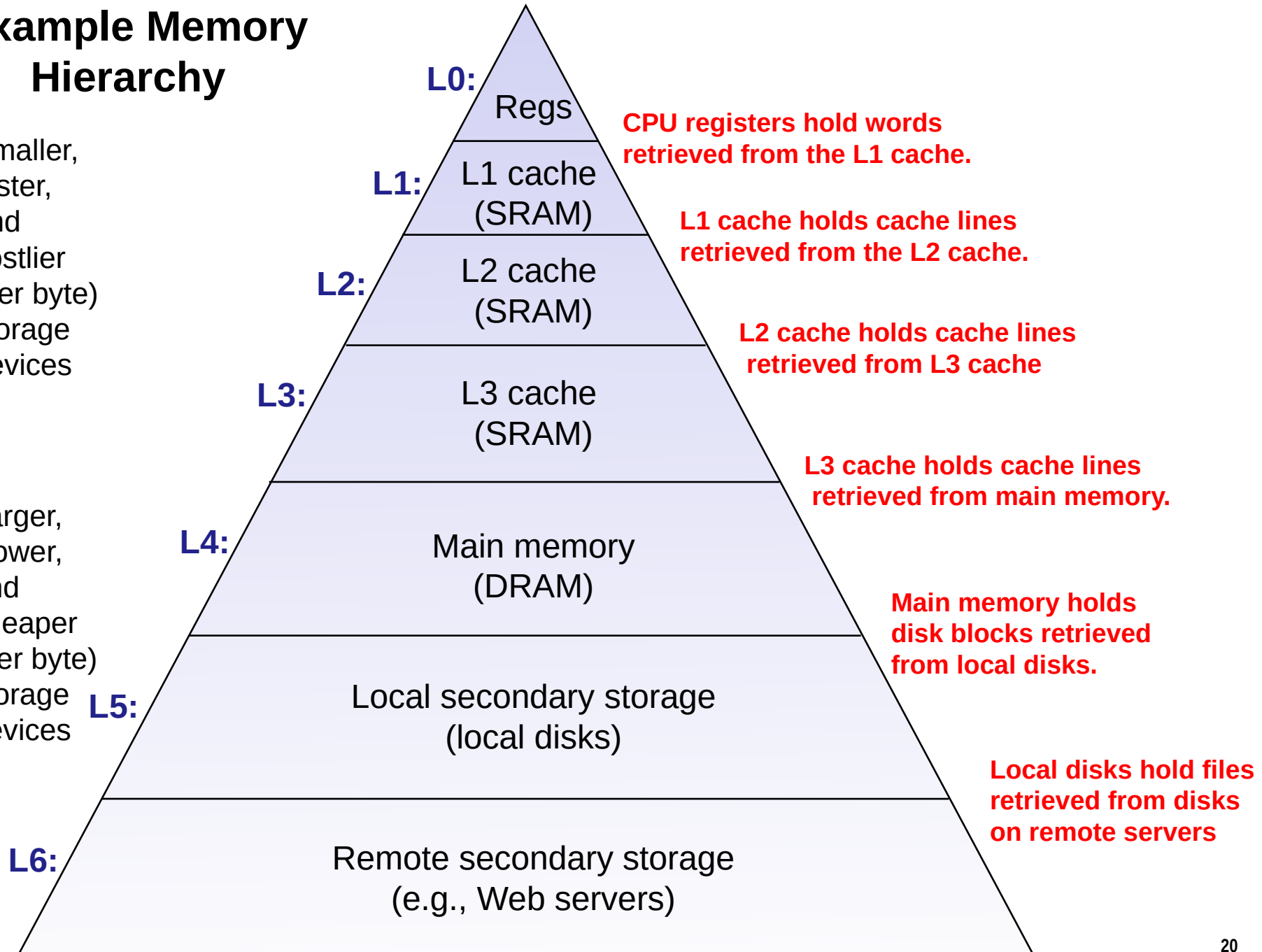    - E.g. Referencing blocks 0, 8, 0, 8, 0, 8, … would miss every time.
- ## Capacity miss
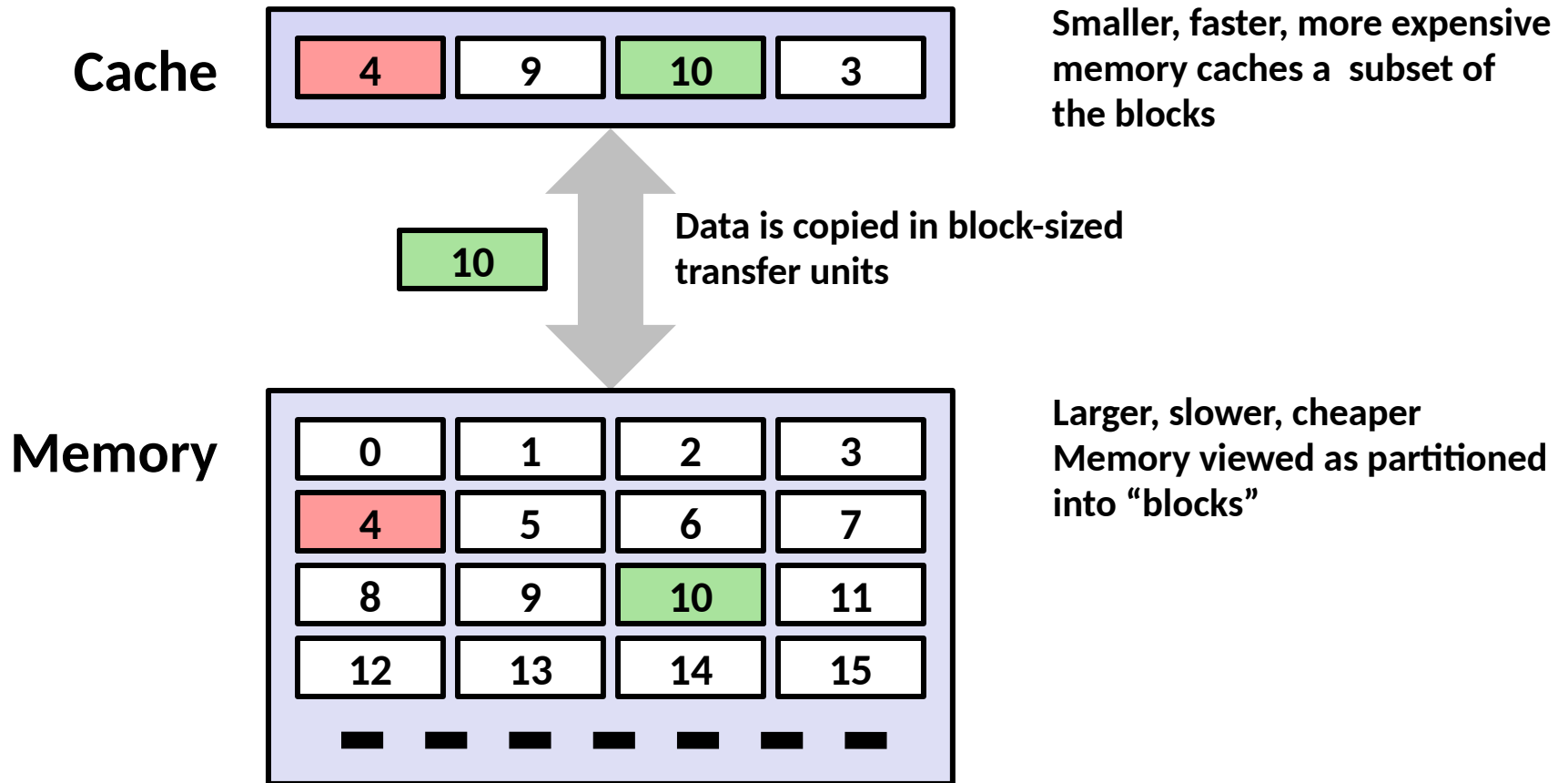  - Occurs when the set of active cache blocks (working set) is larger than the cache.

# Example Memory Hierarchy

Smaller, faster, and costlier (per byte) storage devices
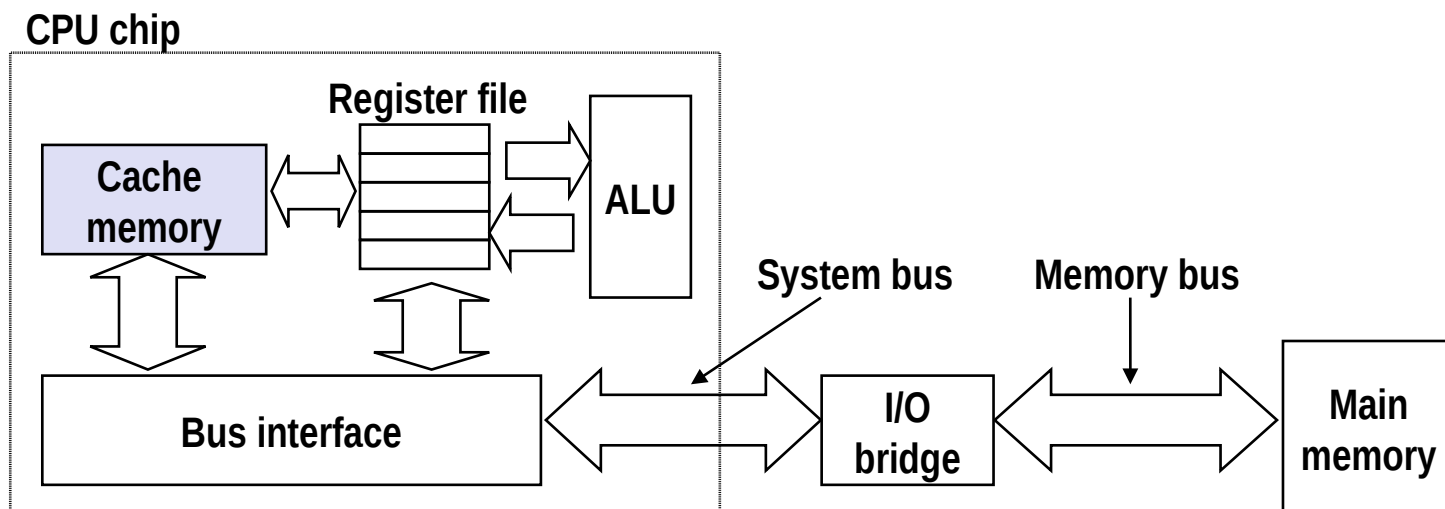
Larger, slower, and cheaper (per byte) storage devices

**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"
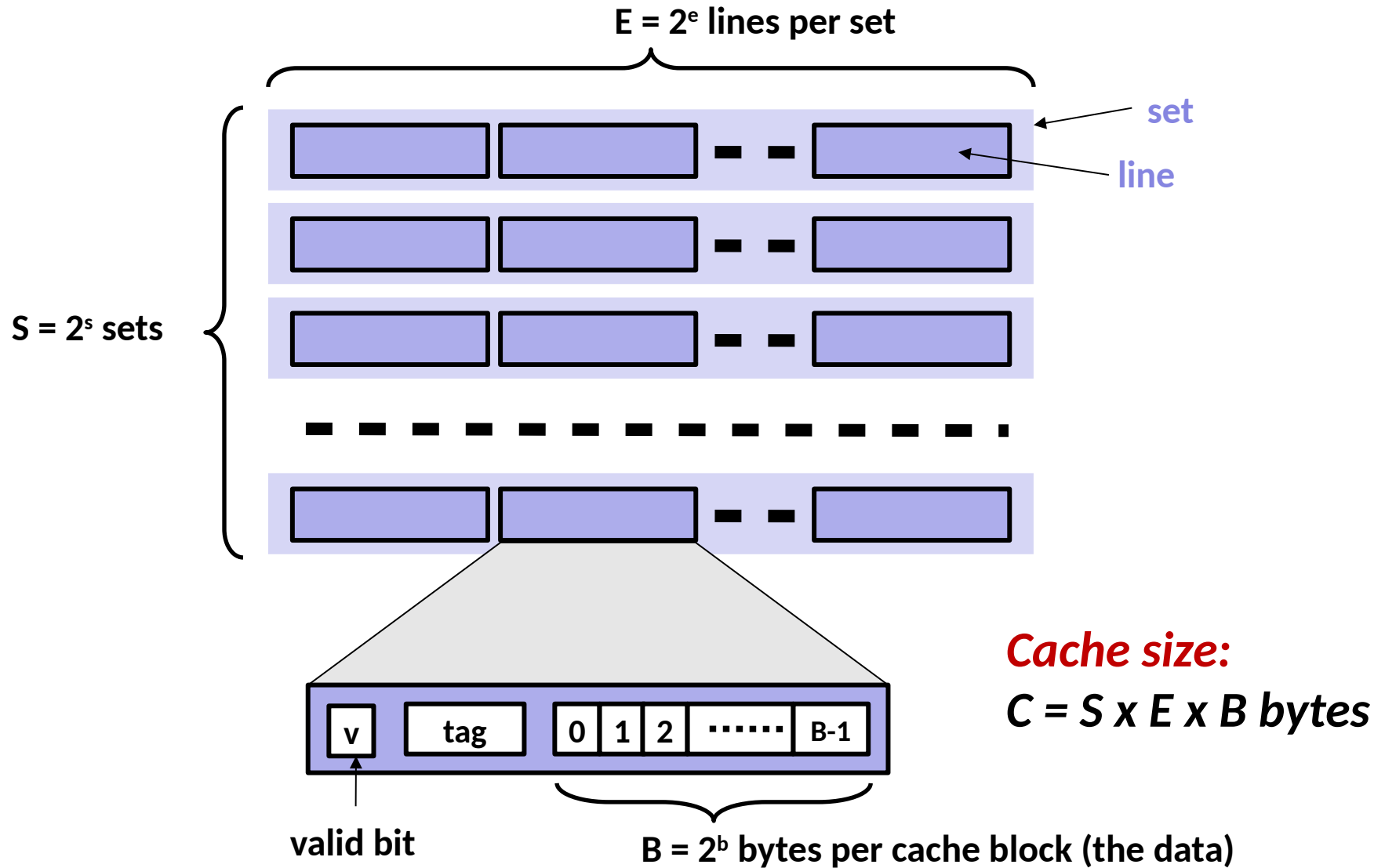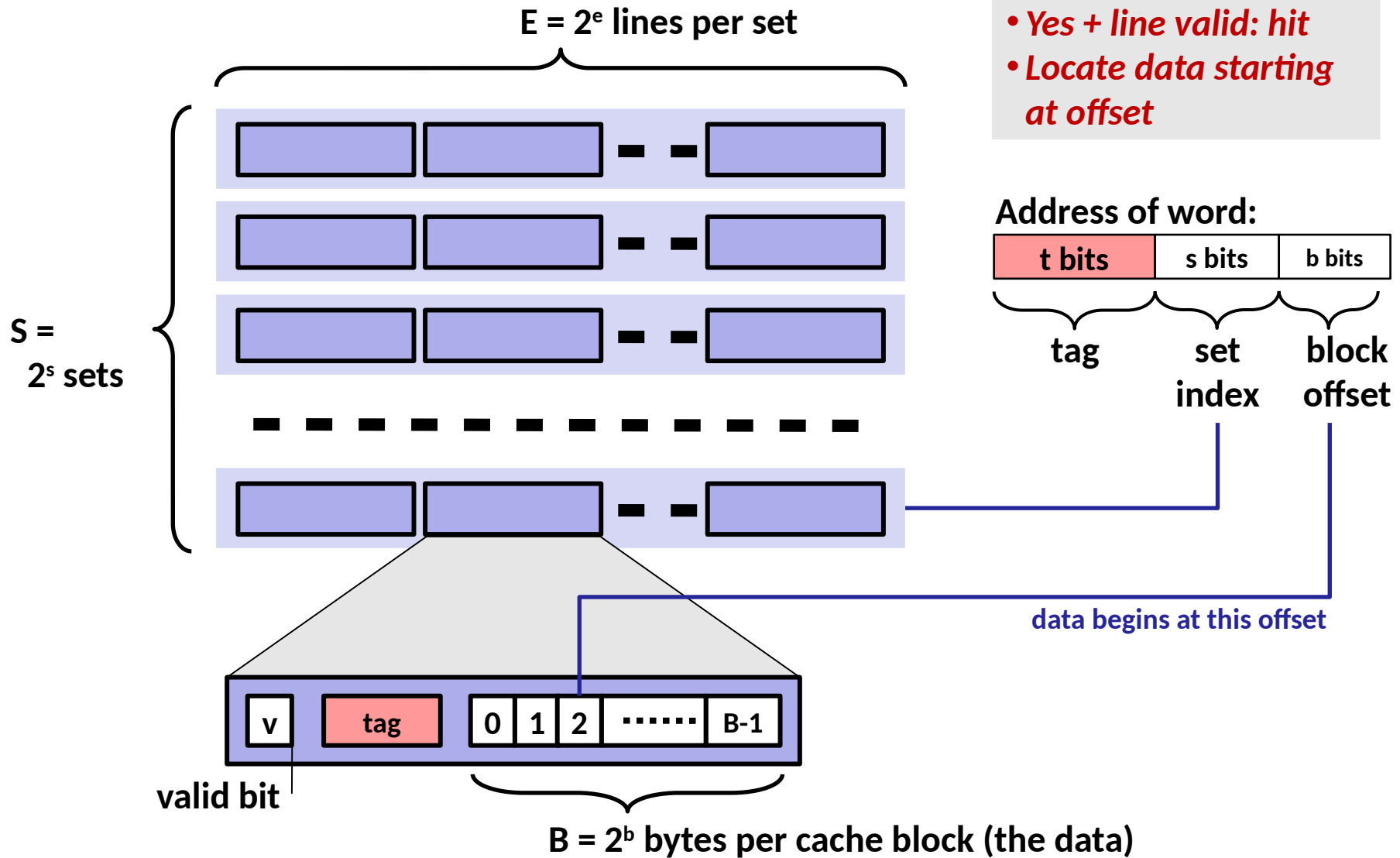
# Cache Memories

- **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:**

**CPU chip**

**Register file**

**Cache memory**

**ALU**

**System bus**   **Memory bus**

**Bus interface**

**I/O bridge**

**Main memory**

# General Cache Organization (S, E, B)

$E = 2^e$ lines per set



set

line

$S = 2^s$ sets

*Cache size:*

*C = S x E x B bytes*

| 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)

# Ex.: Direct Mapped Cache (E = 1)

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



S =
2$^s$ sets

**Address of int:**

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

**find set**

# Ex.: 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**

# Ex.: 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 | [$0\underline{00}0_2$], | |
| 1 | [$0\underline{00}1_2$], | miss |
| 7 | [$0\underline{11}1_2$], | hit |
| 8 | [$1\underline{00}0_2$], | miss |
| 0 | [$0\underline{00}0_2$] | miss |
| | | 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**
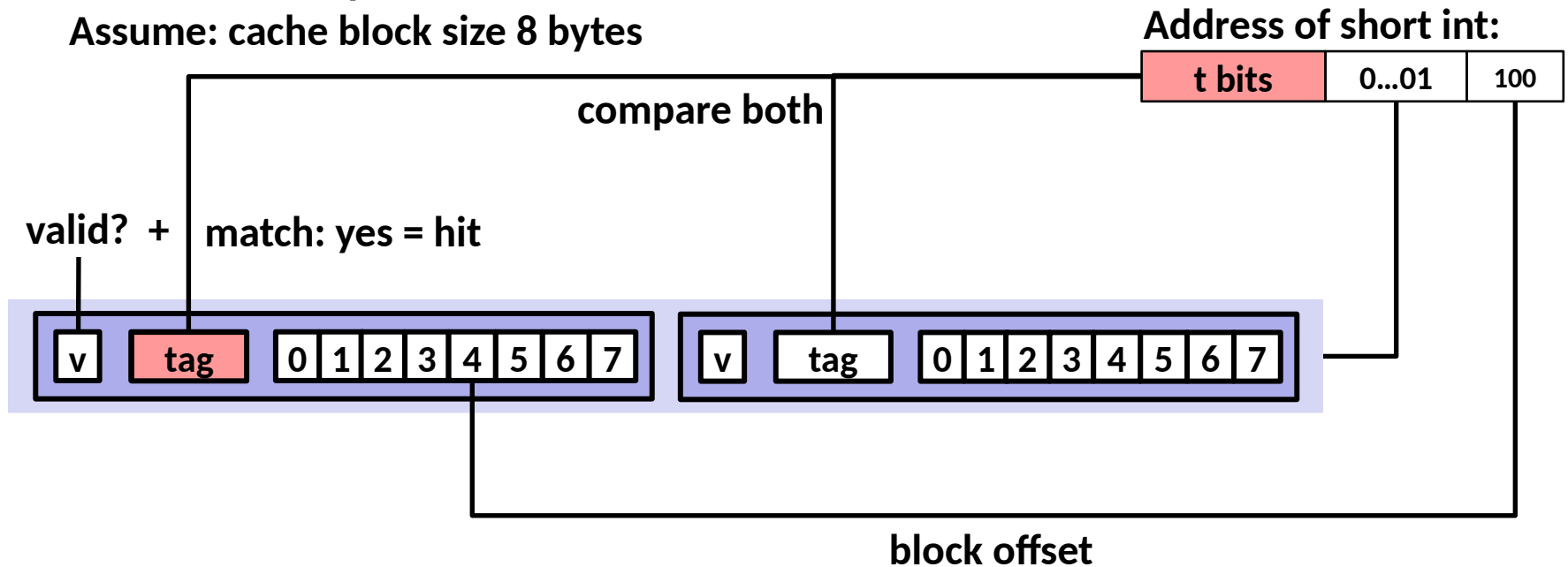**Assume: cache block size 8 bytes**

**Address of short int:**

| t bits | 0…01 | 100 |
|--------|------|-----|

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

| v | tag | 0 1 2 3 4 5 6 7 | | v | tag | 0 1 2 3 4 5 6 7 |  **find set**

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

- - - - - - - - - - - - - - - - - - - -

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

**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 |
|---|-----|---|---|---|---|---|---|---|---|

**block offset**

**short int (2 Bytes) is here**

**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 | [00$\underline{0}$0$_2$], | miss |
| 1 | [00$\underline{0}$1$_2$], | hit |
| 7 | [01$\underline{1}$1$_2$], | miss |
| 8 | [10$\underline{0}$0$_2$], | miss |
| 0 | [00$\underline{0}$0$_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 |    |        |

# Ex.: Fully Associative Cache (S = 1)

**One set containing all lines**
**Assume: cache block size 8 bytes**

One set:

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

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

| 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 | 100 |

# Cache associativity

- **Direct mapped**
  - $E = 1$ (1 line per set)

- **n-way associative**
  - n lines per set

- **Fully associative**

  - $S = 1$ (1 set containing all lines)

# Cache associativity

- **Requested address gives all information needed to locate data in the cache:**

  - Set index: which bookshelf?

  - Tag: book title

  - Block index: what page of the book?

# Cache Size

- **C = S * E * B**
  - S = #sets
  - E = #lines / set (associativity)
  - B = #bytes / block (in a line)
- **Suppose we're given:**

  - 32 KB cache

  - 64 B line

  - 4-way set associative

  - 32-bit architecture

- **How many tag, set, block index bits?**

# Cache Size

- **Suppose we're given:**

  - 32 KB cache

  - 64 B line

  - 4-way set associative

  - 32-bit architecture

- **How many tag, set, block index bits?**

  - 64 B per block = ? bits

  - How many lines in a 32 KB cache?

  - How many sets to hold those lines?

# Cache Size

- **Suppose we're given:**

  - 32 KB cache

  - 64 B line

  - 4-way set associative

  - 32-bit architecture

- **How many tag, set, block index bits?**

  - 64 B per block = 6 bits

  - 32 KB cache / 64 B per line = 2^9 lines

  - 2^9 lines / 4 lines per set = 2^7 sets

  - So 6 block index bits, 7 set index bits, ??? tag bits

# Cache Size

- **Suppose we're given:**

  - 32 KB cache

  - 64 B line

  - 4-way set associative

  - 32-bit architecture

- **How many tag, set, block index bits?**

  - 64 B per block = 6 bits

  - 32 KB cache / 64 B per line = 2^9 lines

  - 2^9 lines / 4 lines per set = 2^7 sets

  - So 6 block index bits, 7 set index bits

  - 32 – 6 – 7 = 19 tag bits

# Cache Size

- **Suppose we're given:**

  - Read 0x2b74ce2d

- **What are the tag, set, block index?**

  - (19 tag, 7 set, 6 block)

# Cache Size

- **Suppose we're given:**

  - Read 0x2b74ce2d

- **What are the tag, set, block index?**

  - (19 tag, 7 set, 6 block)

  - 0x2b74ce2d =

    0010 1011 0111 0100 1100 1110 0010 1101

# Cache Size

- **Suppose we're given:**

  - Read 0x2b74ce2d

- **What are the tag, set, block index?**

  - (19 tag, 7 set, 6 block)

  - 0x2b74ce2d =

    0010 1011 0111 0100 1100 1110 0010 1101

  - Tag = 0010101101110100110

  - Set index = 0111000 (56)

  - Block index = 101101 (45)

# 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**

Core 0

Regs

L1 d-cache | L1 i-cache

L2 unified cache

. . .

Core 3

Regs

L1 d-cache | L1 i-cache

L2 unified cache

L3 unified cache
(shared by all cores)

Main memory

**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!)

# Cache Performance Metrics

- **Average memory access time (AMAT)**
  - Each request may hit or miss.
  - What is the average time over many requests?

- **AMAT = H + MR * AMP**
  - H = hit latency
  - MR = miss rate
  - AMP = average miss penalty

# Cache Performance Metrics

- **AMAT = H + MR * AMP**

  - H = hit latency

  - MR = miss rate

  - AMP = average miss penalty

- **Suppose we're given:**

  - Hit latency of L1 = 1 cycle

  - Hit rate of L1 = 30%

  - Latency of main memory = 200 cycles

- **What is AMAT?**

# Cache Performance Metrics

- **AMAT = H + MR * AMP**

  - H = hit latency

  - MR = miss rate

  - AMP = average miss penalty

- **Suppose we're given:**

  - Hit latency of L1 = 1 cycle

  - Hit rate of L1 = 30%

  - Latency of main memory = 200 cycles

- **What is AMAT?**

  - 1 + 0.70 * 200 = 141 cycles

# Cache Performance Metrics

- **What about multi-level caches?**
  - AMAT(L1) = H(L1) + MR(L1) * AMAT(L2)
  - AMAT(L2) = H(L2) + MR(L2) * AMP

- **Suppose we're given:**
  - L1: 1 cycle, 30% hit rate
  - L2: 10 cycles, 55% hit rate
  - RAM: 200 cycles

- **What is AMAT?**

# Cache Performance Metrics

- **What about multi-level caches?**
    - AMAT(L1) = H(L1) + MR(L1) * AMAT(L2)
    - AMAT(L2) = H(L2) + MR(L2) * AMP

- **Suppose we're given:**

    - L1: 1 cycle, 30% hit rate

    - L2: 10 cycles, 55% hit rate

    - RAM: 200 cycles

- **What is AMAT?**

    - 1 + 0.70 * AMAT(L2)

# Cache Performance Metrics

- **What about multi-level caches?**
  - AMAT(L1) = H(L1) + MR(L1) * AMAT(L2)
  - AMAT(L2) = H(L2) + MR(L2) * AMP

- **Suppose we're given:**

  - L1: 1 cycle, 30% hit rate

  - L2: 10 cycles, 55% hit rate

  - RAM: 200 cycles

- **What is AMAT?**

  - 1 + 0.70 * (10 + 0.45 * 200) = 71 cycles

# 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 (<span style="color:red">temporal locality</span>)
  - Stride-1 reference patterns are good (<span style="color:red">spatial locality</span>)

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

# Cache Performance

- **What is the hit rate of this code?**
  - 64 B block
  - int is 4 B

```
int a[m][n];
for (int i = 0; i < m; i++)
    for (int j = 0; j < n; j++)
        a[i][j] += 1;
```

# Cache Performance

- **What is the hit rate of this code?**
  - 64 B block
  - Int is 4 B

  ```
  int a[m][n];
  for (int i = 0; i < m; i++)
      for (int j = 0; j < n; j++)
          a[i][j] += 1;
  ```

  - For each cache line (64 B), it holds 16 ints
  - First will miss, next 15 will hit
  - Then another miss, then another 15 hits, …
  - Hit rate = 15/16

# Matrix Multiplication Example

- **Description:**
  - Multiply N x N matrices
  - Matrix elements are doubles (8 bytes)
  - $O(N^3)$ total operations
  - N reads per source element
  - N values summed per destination
    - but may be able to hold in register

*Variable **sum** held in register*

```
/* ijk */
for (i=0; i<n; i++)  {
  for (j=0; j<n; j++) {
    sum = 0.0;
    for (k=0; k<n; k++)
      sum += a[i][k] * b[k][j];
    c[i][j] = sum;
  }
}
```

*matmult/mm.c*

# Miss Rate Analysis for Matrix Multiply

- **Assume:**
  - Block size = 32B (big enough for four doubles)
  - Matrix dimension (N) is very large
    - Approximate 1/N as 0.0
  - Cache is not even big enough to hold multiple rows
- **Analysis Method:**
  - Look at access pattern of inner loop

# Layout of C Arrays in Memory (review)

- **C arrays allocated in row-major order**
  - each row in contiguous memory locations
- **Stepping through columns in one row:**
  - ```
    for (i = 0; i < N; i++)
      sum += a[0][i];
    ```
  - accesses successive elements
  - if block size (B) > sizeof($a_{ij}$) bytes, exploit spatial locality
    - miss rate = sizeof($a_{ij}$) / B
- **Stepping through rows in one column:**
  - ```
    for (i = 0; i < n; i++)
      sum += a[i][0];
    ```
  - accesses distant elements
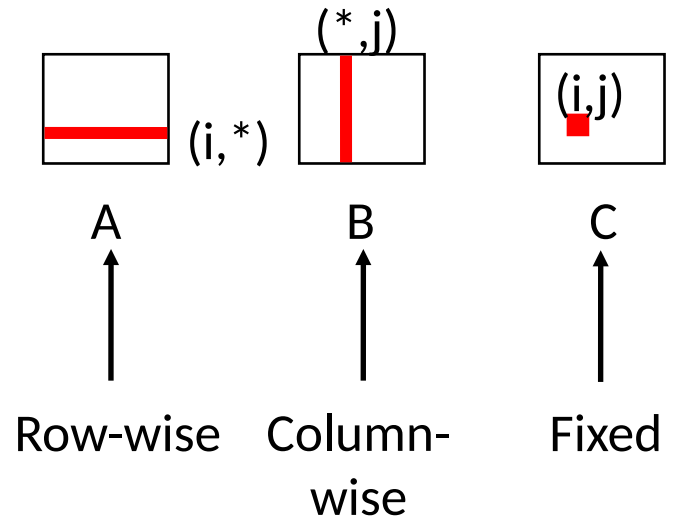  - no spatial locality!
    - miss rate = 1 (i.e. 100%)

# Matrix Multiplication (ijk)

```
/* ijk */
for (i=0; i<n; i++)  {
  for (j=0; j<n; j++) {
    sum = 0.0;
    for (k=0; k<n; k++)
      sum += a[i][k] * b[k][j];
    c[i][j] = sum;
  }
}
```
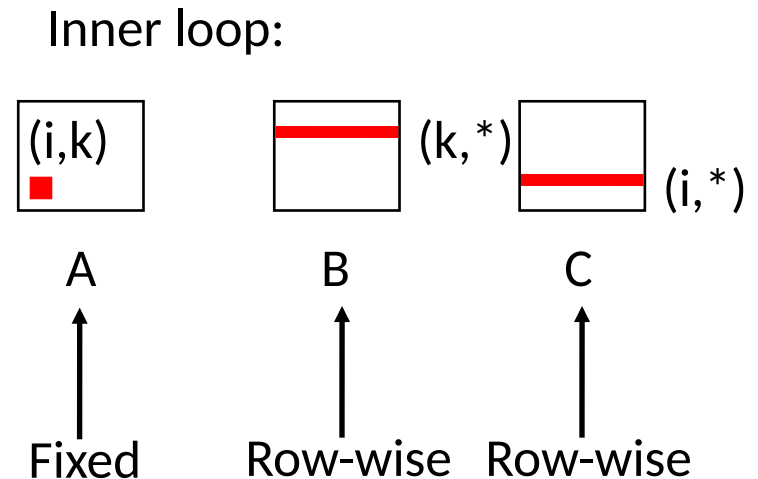*matmult/mm.c*

Inner loop:



| A | B | C |
|---|---|---|
| Row-wise | Column-wise | Fixed |

Misses per inner loop iteration:

| A | B | C |
|------|-----|-----|
| 0.25 | 1.0 | 0.0 |

# Matrix Multiplication (jik)

```
/* jik */
for (j=0; j<n; j++) {
  for (i=0; i<n; i++) {
    sum = 0.0;
    for (k=0; k<n; k++)
      sum += a[i][k] * b[k][j];
    c[i][j] = sum
  }
}
```
*matmult/mm.c*

Inner loop:

(*,j)

(i,*)

(i,j)

A         B         C

Row-wise   Column-wise   Fixed

Misses per inner loop iteration:

| A | B | C |
|---|---|---|
| 0.25 | 1.0 | 0.0 |

# Matrix Multiplication (kij)

```
/* kij */
for (k=0; k<n; k++) {
  for (i=0; i<n; i++) {
    r = a[i][k];
    for (j=0; j<n; j++)
      c[i][j] += r * b[k][j];

  }
}                          matmult/mm.c
```

Inner loop:



|  |  |  |
|------|-----------|-----------|
| A | B | C |
| Fixed | Row-wise | Row-wise |

Misses per inner loop iteration:

| A | B | C |
|-----|------|------|
| 0.0 | 0.25 | 0.25 |

# Matrix Multiplication (ikj)

```
/* ikj */
for (i=0; i<n; i++) {
  for (k=0; k<n; k++) {
    r = a[i][k];
    for (j=0; j<n; j++)
      c[i][j] += r * b[k][j];
  }
}
```
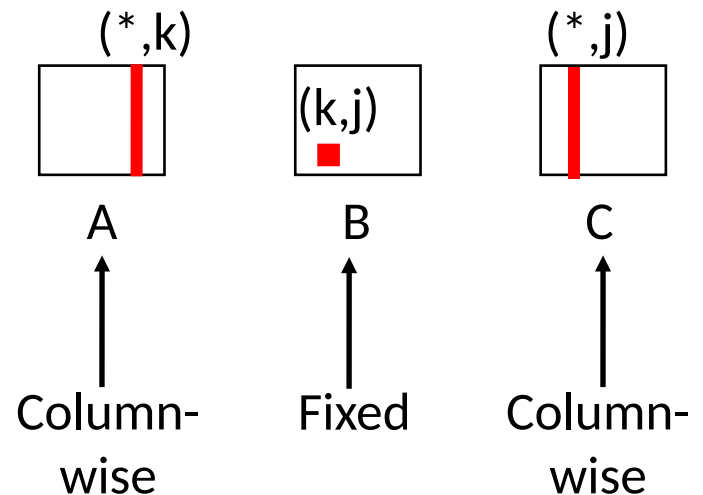*matmult/mm.c*

Inner loop:



| A | B | C |
|---|---|---|
| Fixed | Row-wise | Row-wise |

Misses per inner loop iteration:

| A | B | C |
|---|---|---|
| 0.0 | 0.25 | 0.25 |

# Matrix Multiplication (jki)

```
/* jki */
for (j=0; j<n; j++) {
  for (k=0; k<n; k++) {
    r = b[k][j];
    for (i=0; i<n; i++)
      c[i][j] += a[i][k] * r;
  }
}
```
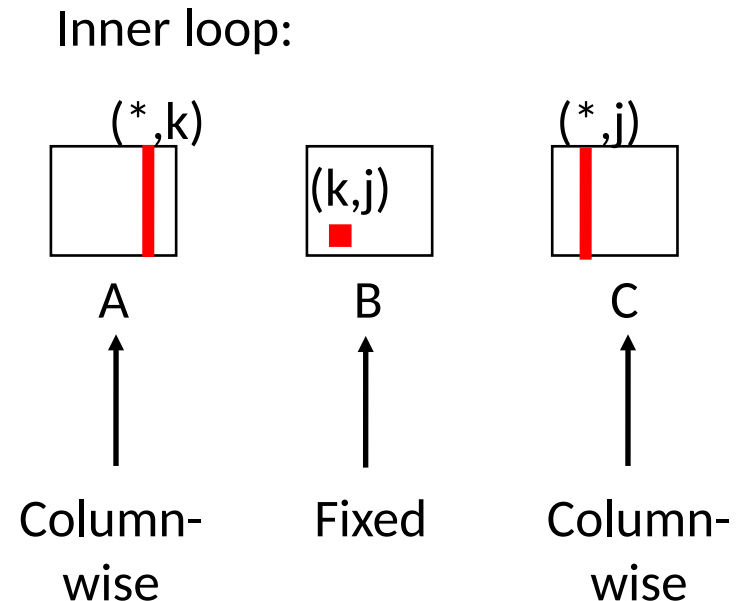*matmult/mm.c*

Inner loop:



(*,k)          (*,j)
A      B      C

Column-    Fixed    Column-
wise                wise

Misses per inner loop iteration:

A      B      C

1.0 0.0 1.0

# Matrix Multiplication (kji)

```
/* kji */
for (k=0; k<n; k++) {
  for (j=0; j<n; j++) {
    r = b[k][j];
    for (i=0; i<n; i++)
      c[i][j] += a[i][k] * r;
  }
}
```

matmult/mm.c

Inner loop:



|  | (*,k) |  | (k,j) |  | (*,j) |
| A | | | B | | | C |

Column-wise    Fixed    Column-wise

Misses per inner loop iteration:

| A | B | C |
| --- | --- | --- |
| 1.0 | 0.0 | 1.0 |

# Summary of Matrix Multiplication

```
for (i=0; i<n; i++) {
  for (j=0; j<n; j++) {
    sum = 0.0;
    for (k=0; k<n; k++)
      sum += a[i][k] * b[k][j];
    c[i][j] = sum;
  }
}
```

**ijk (& jik):**
- 2 loads, 0 stores
- misses/iter = **1.25**

```
for (k=0; k<n; k++) {
 for (i=0; i<n; i++) {
  r = a[i][k];
  for (j=0; j<n; j++)
   c[i][j] += r * b[k][j];
 }
}
```
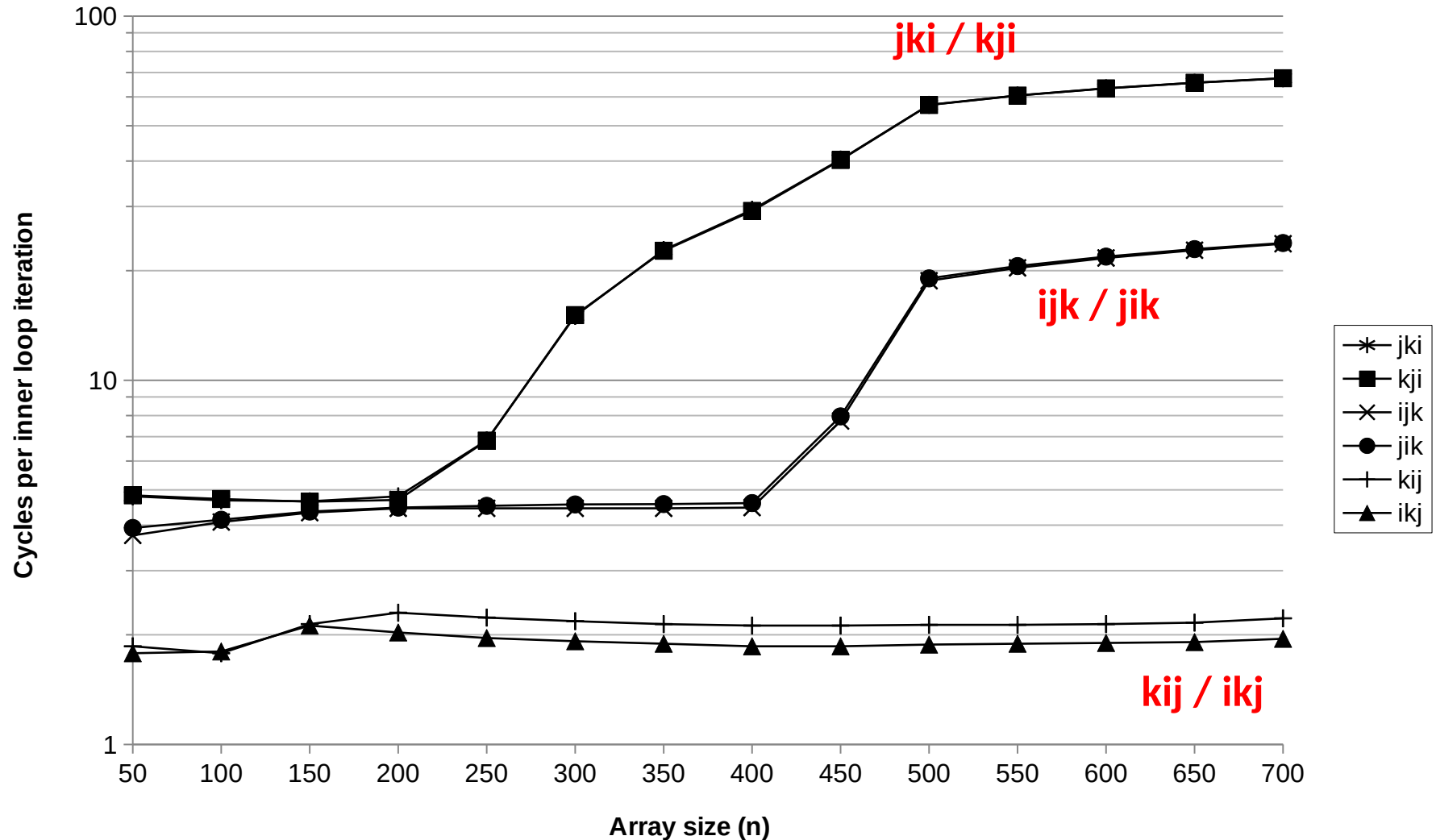
**kij (& ikj):**
- 2 loads, 1 store
- misses/iter = **0.5**

```
for (j=0; j<n; j++) {
 for (k=0; k<n; k++) {
   r = b[k][j];
   for (i=0; i<n; i++)
    c[i][j] += a[i][k] * r;
 }
}
```

**jki (& kji):**
- 2 loads, 1 store
- misses/iter = **2.0**

# Core i7 Matrix Multiply Performance

# Cache Summary

- **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.