# The Memory Hierarchy and Cache Memories

Slides adapted from the CMU version of the course (thanks to Randal E. Bryant and David R. O'Hallaron)

# Today

- **Storage technologies and trends**
- **Locality of reference**
- **Concept of memory hierarchy**
- **Cache memories**

# Random-Access Memory (RAM)

- **Key features**

    - RAM is usually 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? | Cost | Applications |
|---|---|---|---|---|---|
| SRAM | 4 or 6 | 1X | No | 100x | Cache memories |
| DRAM | 1 | 10X | 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**
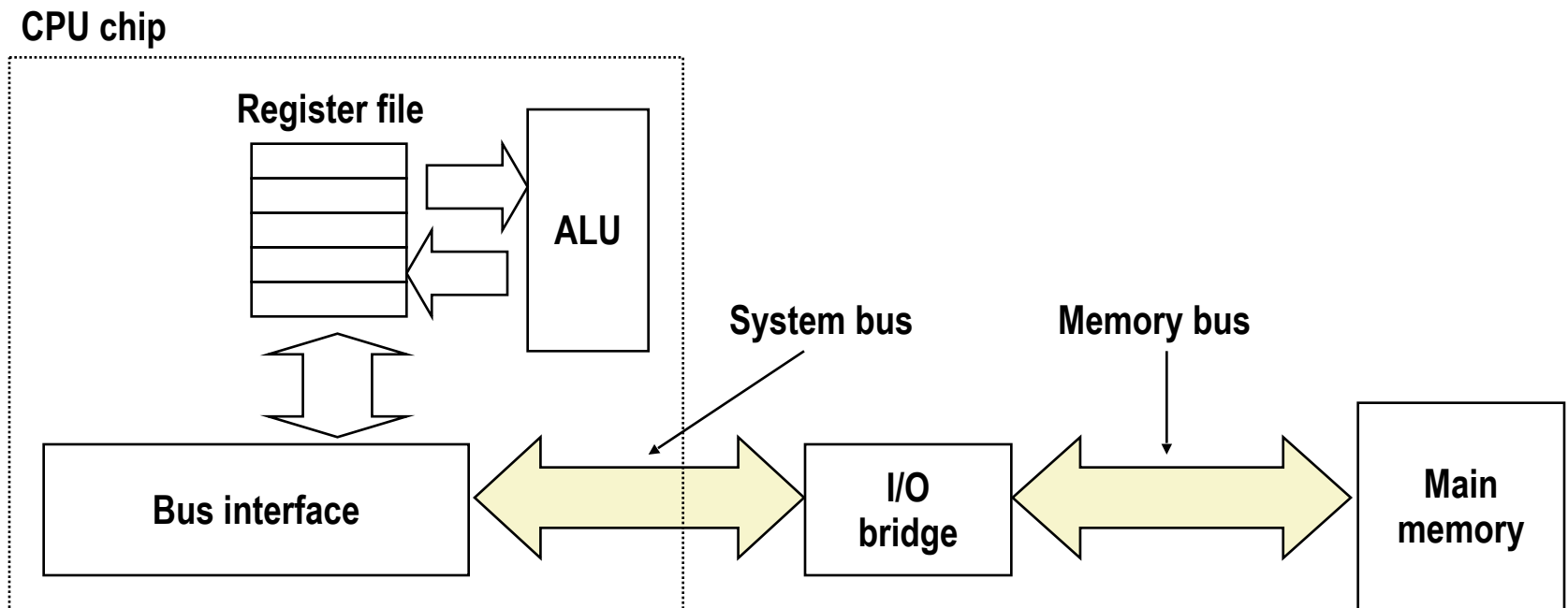  - Read-only memory (ROM): programmed during production
  - Flash memory
    - Wears out after about 100,000 erasing
- **Uses for Nonvolatile Memories**
  - Firmware programs stored in a ROM (BIOS, controllers for disks, network cards, graphics accelerators, security subsystems,…)
  - Solid state disks (replace rotating disks in thumb drives, smart phones, mp3 players, tablets, laptops,…)
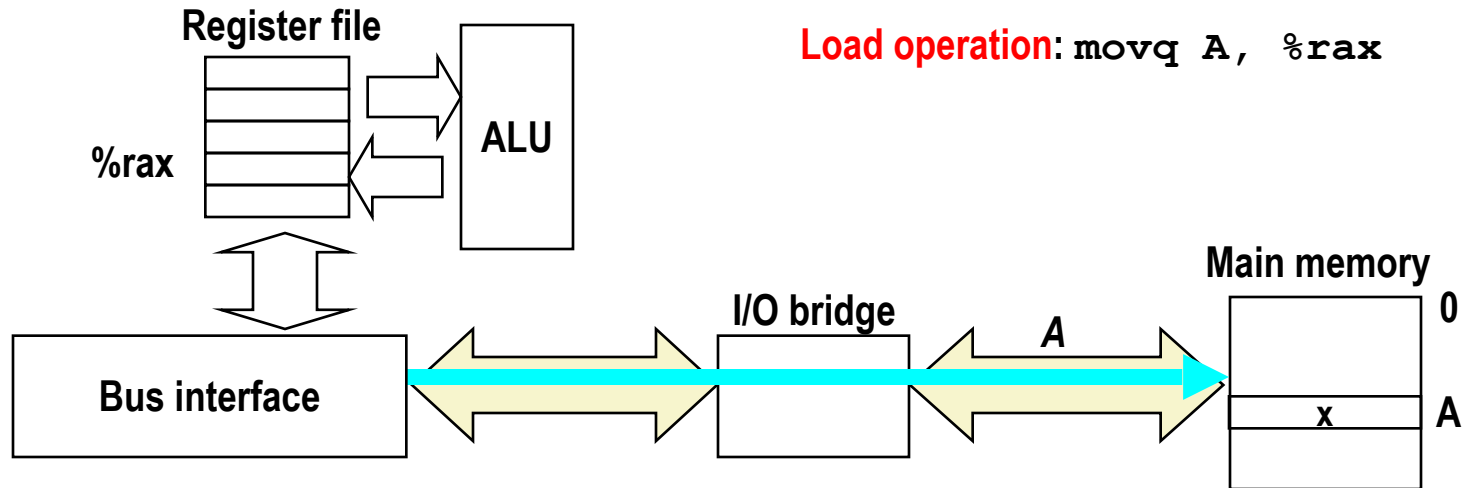  - Disk caches

# Traditional Bus Structure Connecting CPU and Memory

- A **bus** is a collection of parallel wires that carry address, data, and control signals.

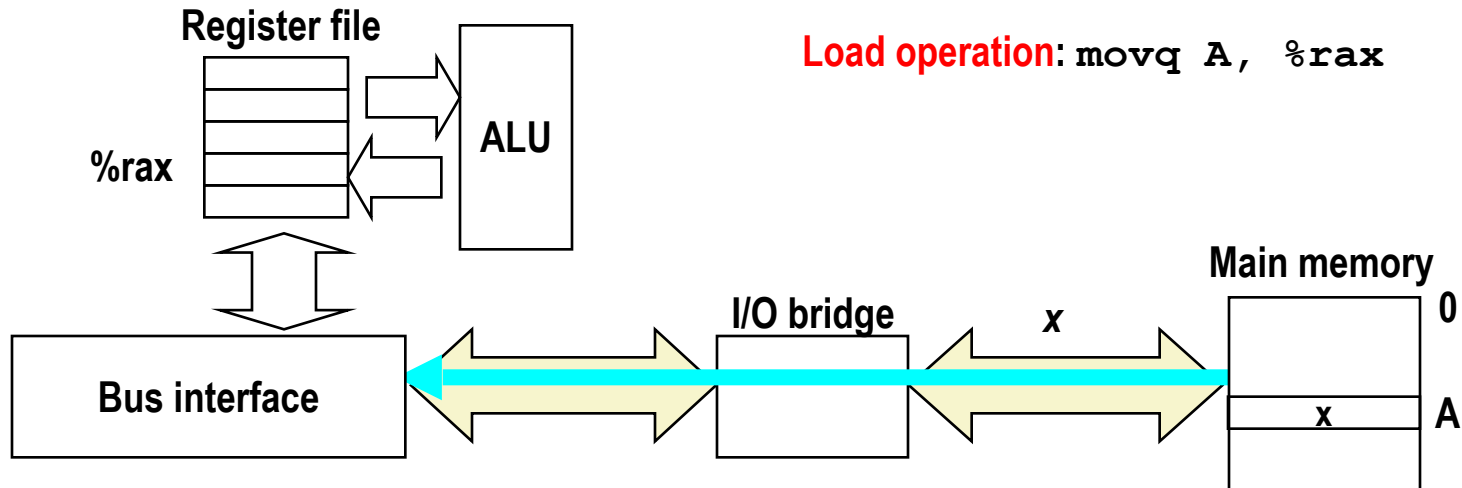- Buses are typically shared by multiple devices.

**CPU chip**

**Register file**

**ALU**

**System bus**

**Memory bus**

**Bus interface**

**I/O bridge**

**Main memory**

# Memory Read Transaction (1)

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

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

**Register file**
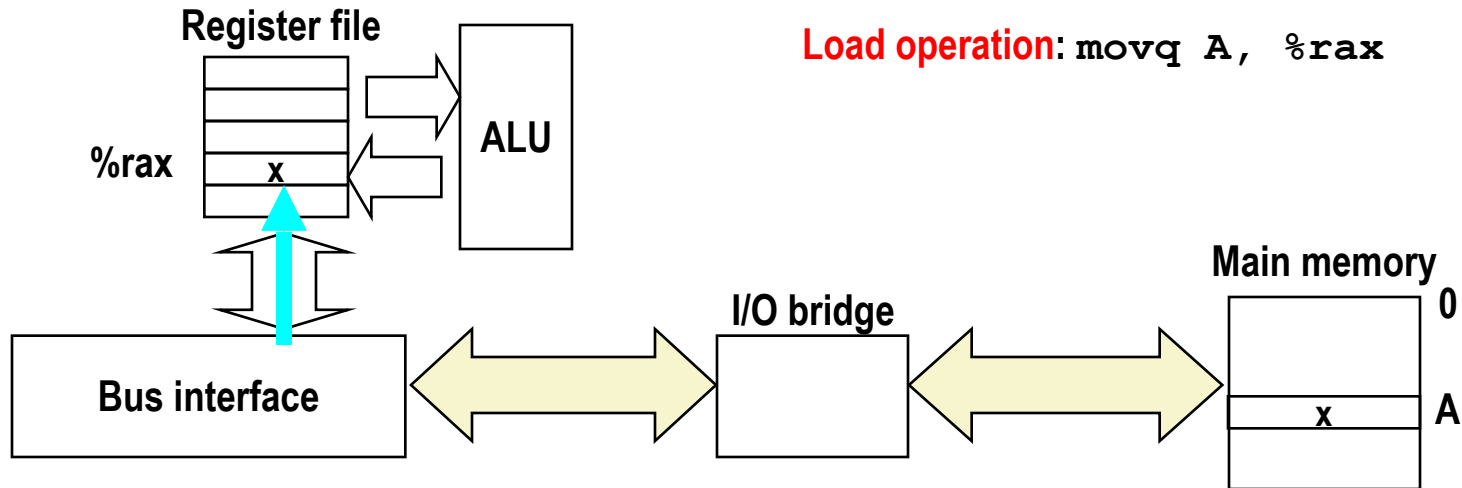
ALU

%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.**
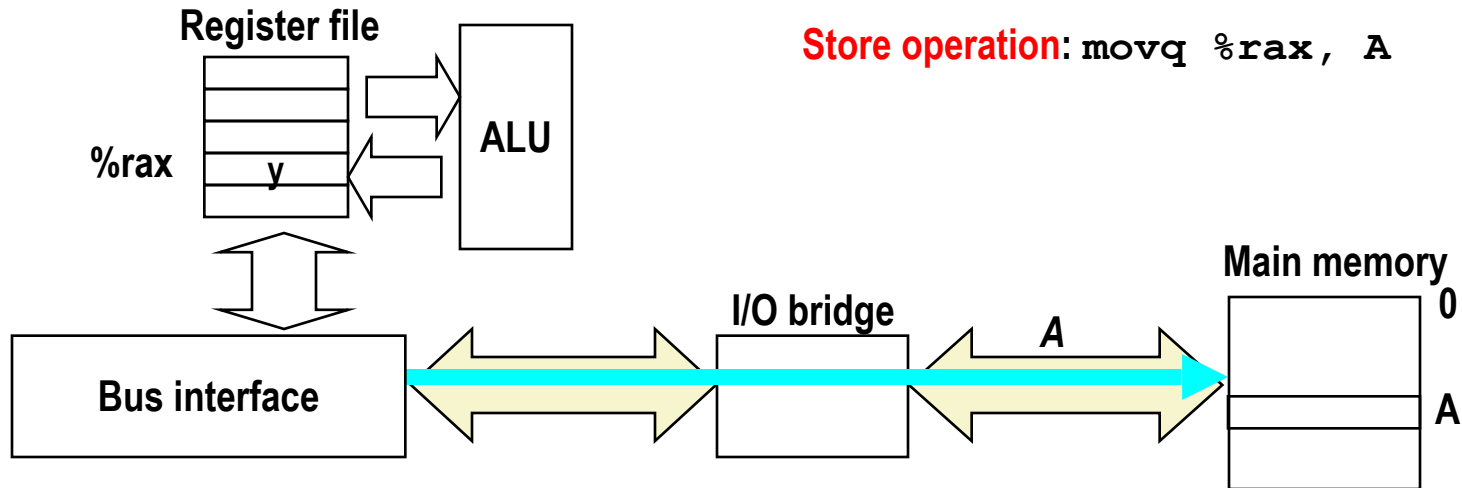


**Register file**

**%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`

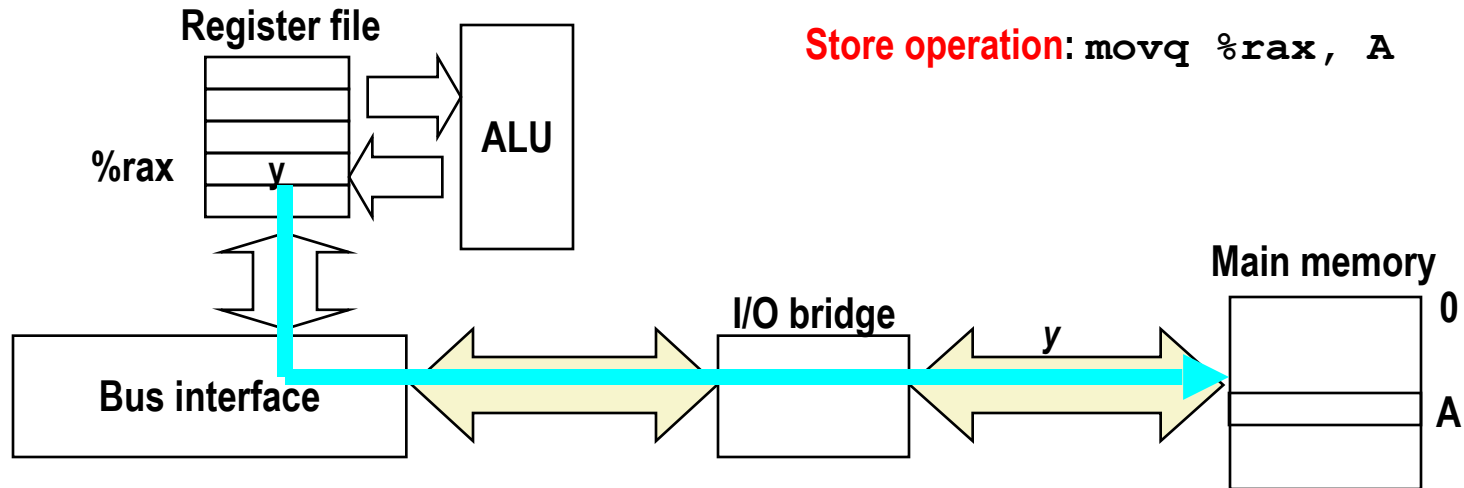%rax | x

ALU

Main memory

I/O bridge

Bus interface

0

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.

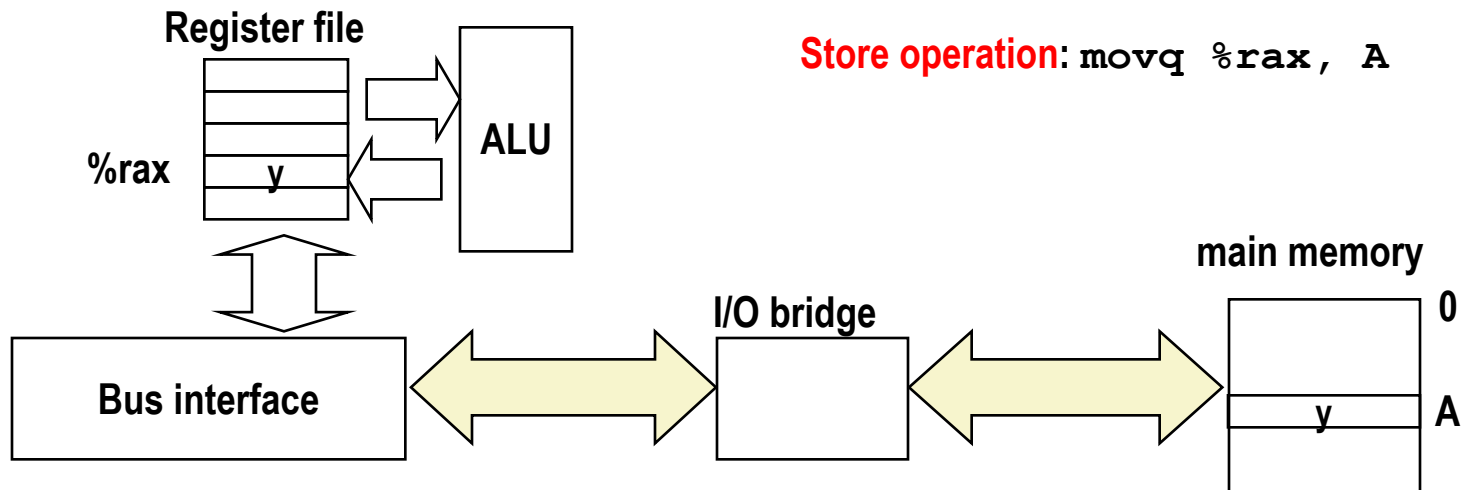**Register file**

**ALU**

**%rax** | y

**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**

**%rax**   y

**ALU**

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

**Bus interface**

**I/O bridge**   y

**Main memory**
0

A

# Memory Write Transaction (3)

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

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

Register file

%rax | y

ALU

Bus interface

I/O bridge

main memory

0

y | A

# What's Inside A Disk Drive?

Spindle

Arm

Platters

Actuator

Electronics
(including a
processor
and memory!)
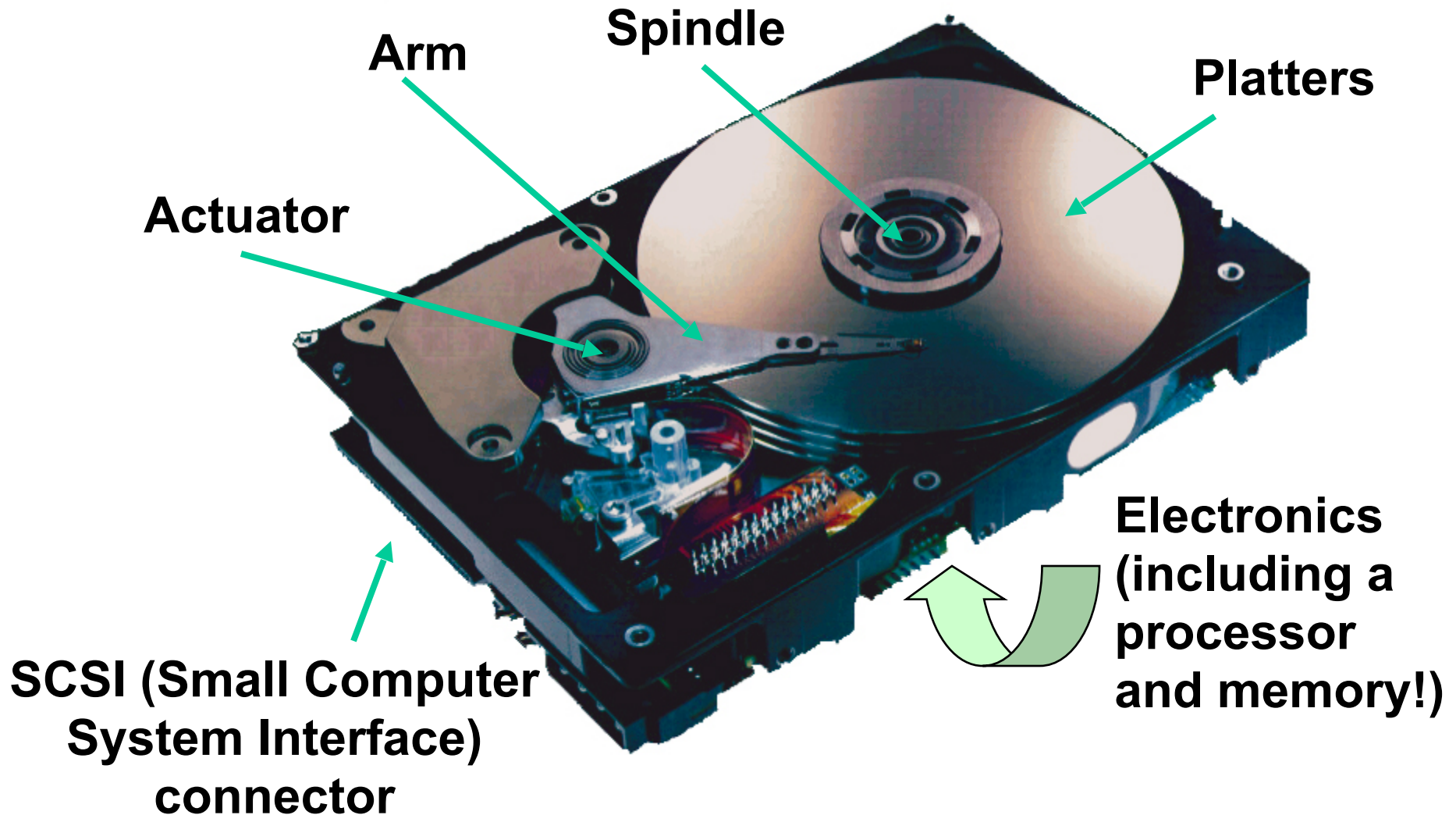
SCSI (Small Computer
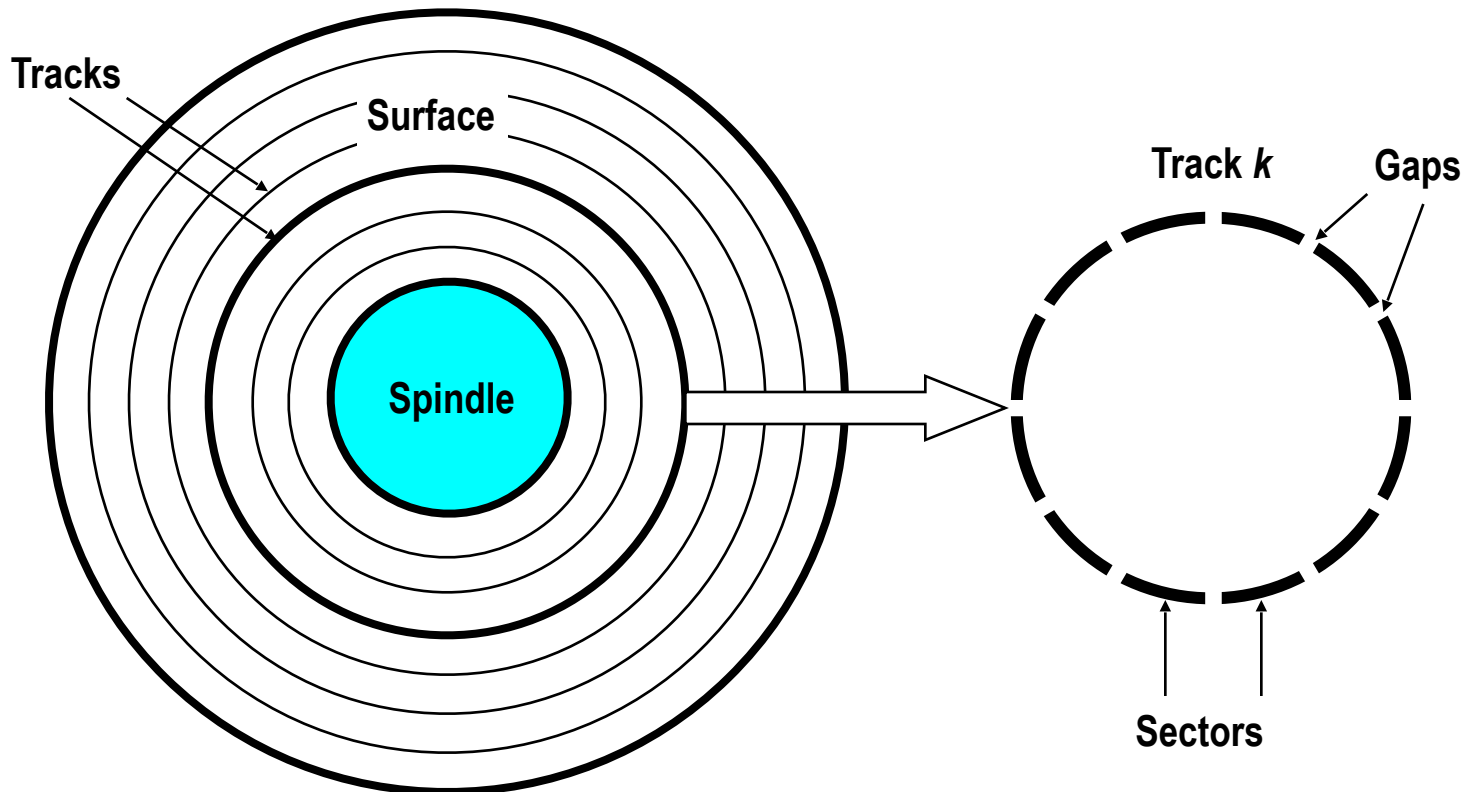System Interface)
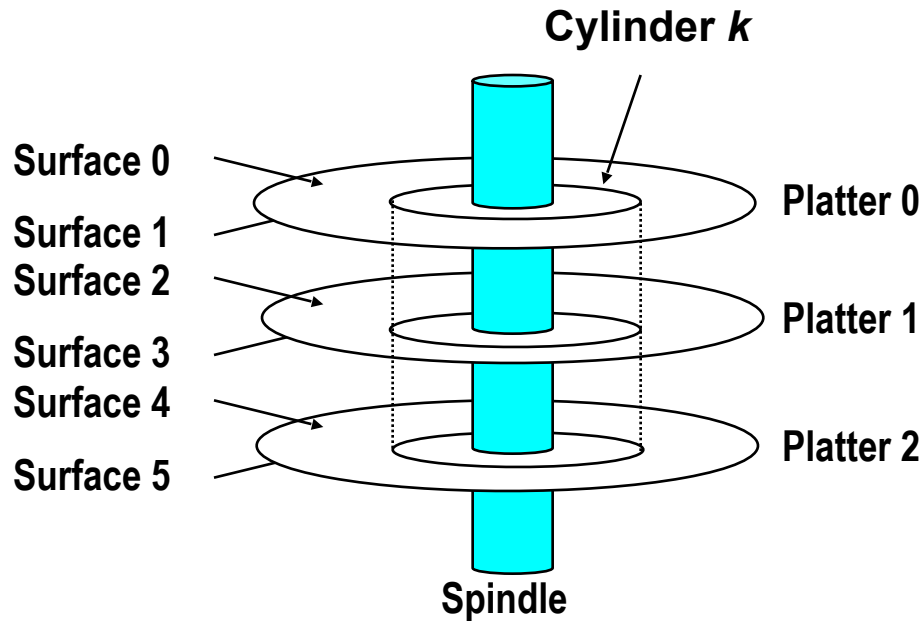connector

*Image courtesy of Seagate Technology*

# Disk Geometry

- **Disks consist of platters, each with two surfaces.**
- **Each surface consists of concentric rings called tracks.**
- **Each track consists of sectors separated by gaps.**

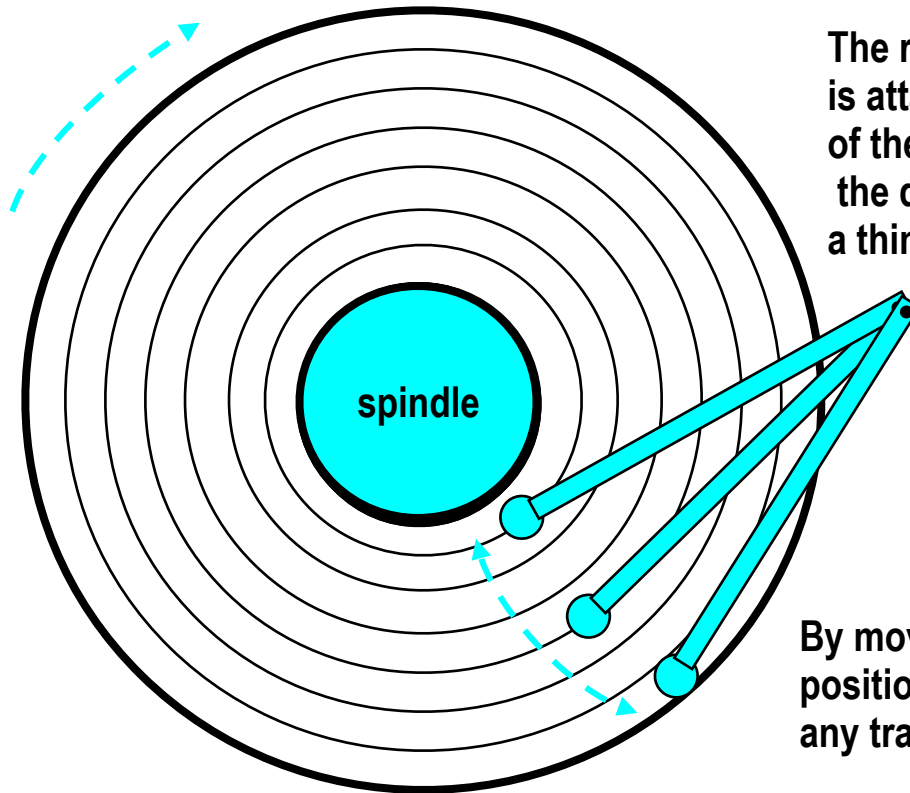# Disk Geometry (Muliple-Platter View)

- **Aligned tracks form a cylinder.**

**Cylinder *k***

**Surface 0**
**Surface 1**
**Surface 2**
**Surface 3**
**Surface 4**
**Surface 5**

**Platter 0**
**Platter 1**
**Platter 2**

**Spindle**

# Disk Operation (Single-Platter View)
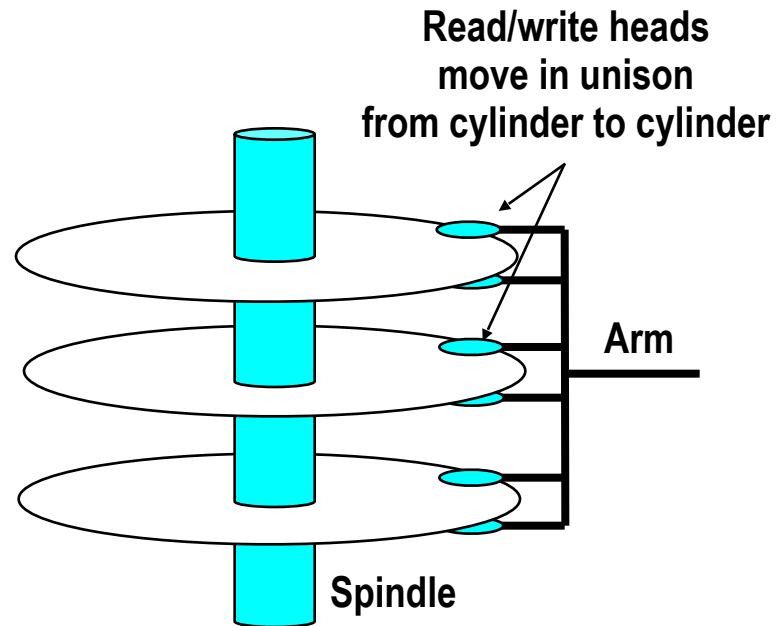
**The disk surface spins at a fixed rotational rate**

**The read/write *head* is attached to the end of the *arm* and flies over the disk surface on a thin cushion of air.**

**spindle**

**By moving radially, the arm can position the read/write head over any track.**

# Disk Operation (Multi-Platter View)

Read/write heads
move in unison
from cylinder to cylinder

Arm

Spindle

# Disk Structure - top view of single platter

**Surface organized into tracks**

**Tracks divided into sectors**

# Disk Access



**Head in position above a track**

# Disk Access



**Rotation is counter-clockwise**

# Disk Access – Read

**About to read blue sector**

# Disk Access – Read



After BLUE read

## After reading blue sector

# Disk Access – Read



After **BLUE** read

**Red request scheduled next**

# Disk Access – Seek



**After BLUE read**          **Seek for RED**

**Seek to red's track**

# Disk Access – Rotational Latency



After **BLUE** read          Seek for **RED**          Rotational latency

## Wait for red sector to rotate around

# Disk Access – Read



After **BLUE** read          Seek for **RED**          **Rotational latency**          After **RED** read

## Complete read of red

# Disk Access – Service Time Components



After **BLUE** read          Seek for **RED**          **Rotational latency**          After **RED** read

**Data transfer**          **Seek**          **Rotational latency**          **Data transfer**

# Disk Access Time

- **Average time to access some target sector approximated by :**
  - Taccess = Tavg seek + Tavg rotation + Tavg transfer
- **Seek time (Tavg seek)**
  - Time to position heads over cylinder containing target sector.
  - Typical Tavg seek is 3—9 ms
- **Rotational latency (Tavg rotation)**
  - Time waiting for first bit of target sector to pass under r/w head.
  - Tavg rotation = 1/2 x 1/RPMs x 60 sec/1 min
  - Typical Tavg rotation = 7200 RPMs
- **Transfer time (Tavg transfer)**
  - Time to read the bits in the target sector.
  - Tavg transfer = 1/RPM x 1/(avg # sectors/track) x 60 secs/1 min.

# Disk Access Time Example

■ **Given:**

- Rotational rate = 7,200 RPM
- Average seek time = 9 ms.
- Avg # sectors/track = 400.

■ **Derived:**

- Tavg rotation = 1/2 x (60 secs/7200 RPM) x 1000 ms/sec = 4 ms.
- Tavg transfer = 60/7200 RPM x 1/400 secs/track x 1000 ms/sec = 0.02 ms
- Taccess  = 9 ms + 4 ms + 0.02 ms

■ **Important points:**

- Access time dominated by seek time and rotational latency.
- First bit in a sector is the most expensive, the rest are free.
- SRAM access time is about  4 ns/doubleword, DRAM about  60 ns
  - Disk is about 40,000 times slower than SRAM,
  - 2,500 times slower than DRAM.

# Logical Disk Blocks

- **Modern disks present a simpler abstract view of the complex sector geometry:**
  - The set of available sectors is modeled as a sequence of b-sized <span style="color:red">logical blocks</span> (0, 1, 2, …)
- **Mapping between logical blocks and actual (physical) sectors**
  - Maintained by hardware/firmware device called disk controller.
  - Converts requests for logical blocks into (surface,track,sector) triples.
- **Allows controller to set aside spare cylinders for each zone.**
  - Accounts for the difference in "formatted capacity" and "maximum capacity".

# I/O Bus

**CPU chip**

**Register file**

**ALU**

**System bus**

**Memory bus**

**Bus interface**

**I/O bridge**

**Main memory**

**I/O bus**

**USB controller**

**Graphics adapter**

**Disk controller**

**Expansion slots for other devices such as network adapters.**

**Mouse**  **Keyboard**

**Monitor**

**Disk**

# Reading a Disk Sector (1)

**CPU chip**

**Register file**

**ALU**

**Bus interface**

CPU initiates a disk read by writing a command, logical block number, and destination memory address to a port (address) associated with disk controller.

**Main memory**

**I/O bus**

**USB controller**

mouse    keyboard

**Graphics adapter**

Monitor

**Disk controller**

**Disk**

# Reading a Disk Sector (2)



Disk controller reads the sector and performs a direct memory access (DMA) transfer into main memory.

**CPU chip**

Register file

ALU

Bus interface

Main memory

I/O bus

USB controller

Graphics adapter

Disk controller

Mouse    Keyboard

Monitor

Disk

# Reading a Disk Sector (3)

**CPU chip**

**Register file**

**ALU**

**Bus interface**

**Main memory**

When the DMA transfer completes, the disk controller notifies the CPU with an *interrupt* (i.e., asserts a special "interrupt" pin on the CPU)

**I/O bus**

**USB controller**

**Graphics adapter**

**Disk controller**

**Mouse**   **Keyboard**

**Monitor**

**Disk**

# Solid State Disks (SSDs)

I/O bus

Requests to read and
write logical disk blocks

Solid State Disk (SSD)

Flash
translation layer

Flash memory

Block 0

| Page 0 | Page 1 | · · · | Page P-1 |

· · ·

Block  B-1

| Page 0 | Page 1 | · · · | Page P-1 |

- **Pages: 512KB to 4KB, Blocks: 32 to 128 pages**
- **Data read/written in units of pages.**
- **Page can be written only after its block has been erased**
- **A block wears out after about 100,000 repeated writes.**

# SSD Performance Characteristics

| | |
|---|---|
| **Avg seq read time** | **50 micro seconds** |
| **Avg seq write time** | **60 micro seconds** |

## Writes are somewhat slower

- Erasing a block takes a long time (~1 ms)
- Modifying a block page requires all other pages to be copied to new block
- In earlier SSDs, the read/write gap was much larger.

**Source: Intel SSD 730 product specification.**

# SSD Tradeoffs vs Rotating Disks

- **Advantages of SSDs**
  - No moving parts → faster, less power

- **Disadvantages**
  - Have the potential to wear out
    - Mitigated by "wear leveling logic" in flash translation layer
    - E.g. Intel SSD 730 guarantees 128 petabyte ($128 \times 10^{15}$ bytes) of writes before they wear out
  - About 30 times more expensive per byte

- **Applications**
  - MP3 players, smart phones
  - Laptops and servers

# The CPU-Memory Gap

**The gap widens between DRAM, disk, and CPU speeds.**

# Locality to the Rescue!

**The key to bridging this CPU-Memory gap is a fundamental property of computer programs known as <span style="color:red">locality</span>**

# Today

- **Storage technologies and trends**
- **Locality of reference**
- **Concept of memory hierarchy**
- **Cache memories**

# 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).     **Spatial locality**

- Reference variable `sum` each iteration.     **Temporal locality**

■ **Instruction references**

- Reference instructions in sequence.     **Spatial locality**

- Cycle through loop repeatedly.     **Temporal locality**

# Qualitative Estimates of Locality

- **Claim:** Being able to look at code and get a qualitative sense of its locality is a key skill for a professional programmer.


- **Question:** Does this function have good locality with respect to array $a$?

```
int sum_array_rows(int a[M][N])
{
    int i, j, sum = 0;

    for (i = 0; i < M; i++)
        for (j = 0; j < N; j++)
            sum += a[i][j];
    return sum;
}
```

# Locality Example

- **Question:** Does this function have good locality with respect to array $a$?

```
int sum_array_cols(int a[M][N])
{
    int i, j, sum = 0;

    for (j = 0; j < N; j++)
        for (i = 0; i < M; i++)
            sum += a[i][j];
    return sum;
}
```

# Locality Example

■ **Question: Can you permute the loops so that the function scans the 3-d array** $a$ **with a stride-1 reference pattern (and thus has good spatial locality)?**

```c
int sum_array_3d(int a[M][N][N])
{
    int i, j, k, sum = 0;

    for (i = 0; i < M; i++)
        for (j = 0; j < N; j++)
            for (k = 0; k < N; k++)
                sum += a[k][i][j];
    return sum;
}
```

# Today

- **Storage technologies and trends**
- **Locality of reference**
- **Concept of memory hierarchy**
- **Cache memories**

# 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

47

# Today

- **Storage technologies and trends**
- **Locality of reference**
- **Concept of memory hierarchy**
- **Cache memories**

# 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.**

# 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 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**

**Cache**

| 8 | 9 | 14 | 3 |
|---|---|----|---|

**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 in cache:*
*Hit!*

# General Cache Concepts: Miss

**Request: 12**

**Cache**

| 8 | 12 | 14 | 3 |

*Data in block b is needed*

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

**12**    **Request: 12**

*Block b is fetched from memory*

**Memory**

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

*Block b is stored in cache*

- Replacement policy: determines which block gets evicted (victim)

# Types of Cache Misses

- **Cold (compulsory) miss**
  - Happens when cache is empty.
- **Capacity miss**
  - Happens when the working set is larger than the cache.
- **Conflict miss**
  - Happens when multiple data objects map to the same cache block.
    - E.g. Block i at level k+1 must be placed in block (i mod 4) at level k
    - So referencing blocks 0, 8, 0, 8, 0, 8, ... would miss every time

# Examples of Caching in the Mem. Hierarchy

| Cache Type | What is Cached? | Where is it Cached? | Latency (cycles) | Managed By |
|---|---|---|---|---|
| Registers | 4-8 bytes words | CPU core | 0 | Compiler |
| L1 cache | 64-byte blocks | On-Chip L1 | 4 | Hardware |
| L2 cache | 64-byte blocks | On-Chip L2 | 10 | Hardware |
| Main Memory | 4-KB pages | Main memory | 100 | Hardware + OS |
| Buffer cache | Parts of files | Main memory | 100 | OS |
| Browser cache | Web pages | Local disk | 10,000,000 | Web browser |

# General Cache Organization (S, E, B)

E    lines per set

set

line

S    sets

*Cache size:*
*C = S x E x B data bytes*

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

valid bit    B    bytes per cache block (the data)

# Another View



1 valid bit per line   $t$ tag bits per line   $B$ bytes per cache block

Set 0:
| Valid | Tag | 0 | 1 | $\cdots$ | $B{-}1$ |
| Valid | Tag | 0 | 1 | $\cdots$ | $B{-}1$ |

$E$ lines per set

Set 1:
| Valid | Tag | 0 | 1 | $\cdots$ | $B{-}1$ |
| Valid | Tag | 0 | 1 | $\cdots$ | $B{-}1$ |

Set $S{-}1$:
| Valid | Tag | 0 | 1 | $\cdots$ | $B{-}1$ |
| Valid | Tag | 0 | 1 | $\cdots$ | $B{-}1$ |

$S$ sets

Cache size:   $C = B \times E \times S$ data bytes

(a)

Address:   $t$ bits   $s$ bits   $b$ bits

$m{-}1$ ........................ 0

Tag   Set index   Block offset

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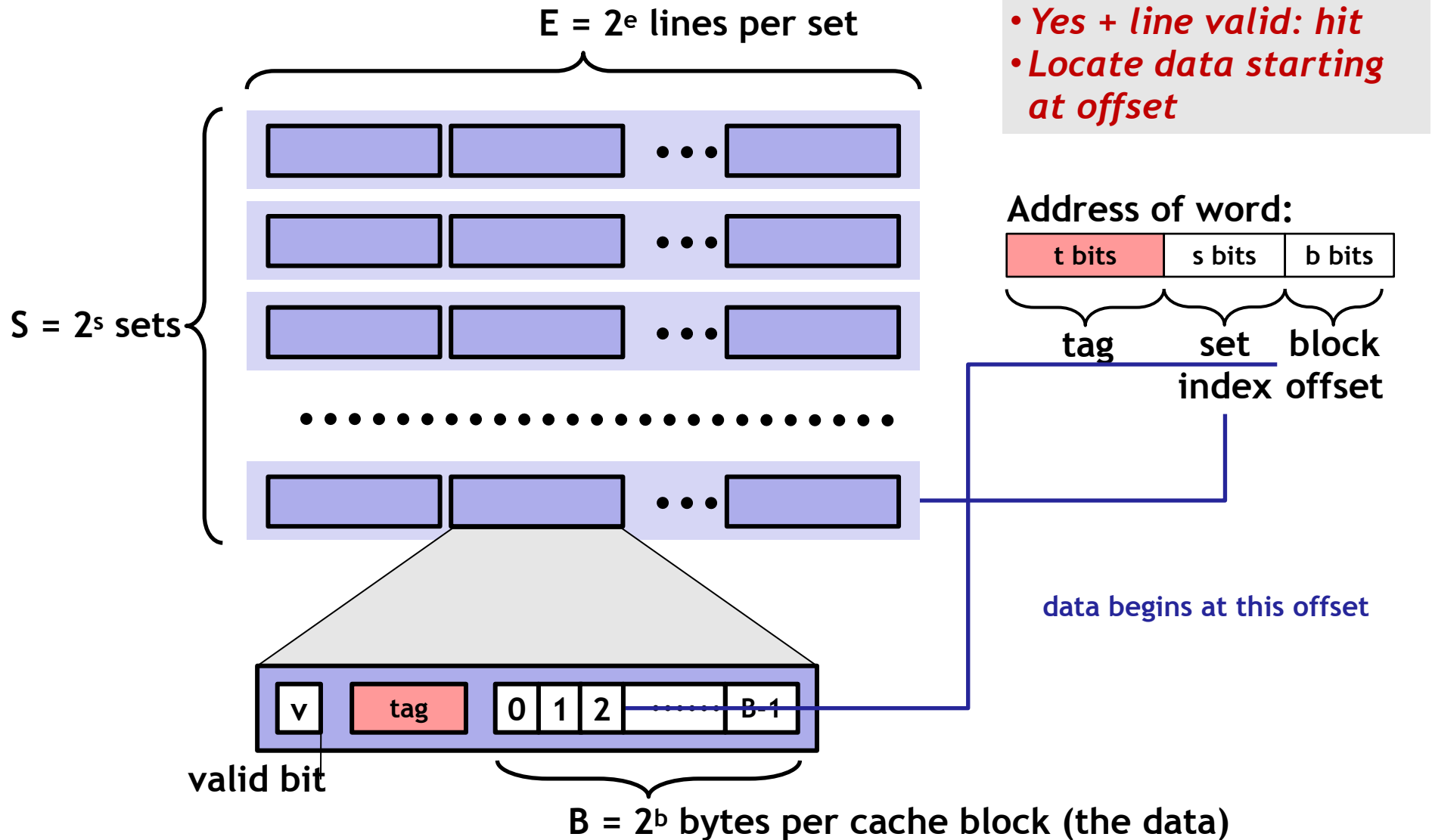
# Example: Direct Mapped Cache (E = 1)

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



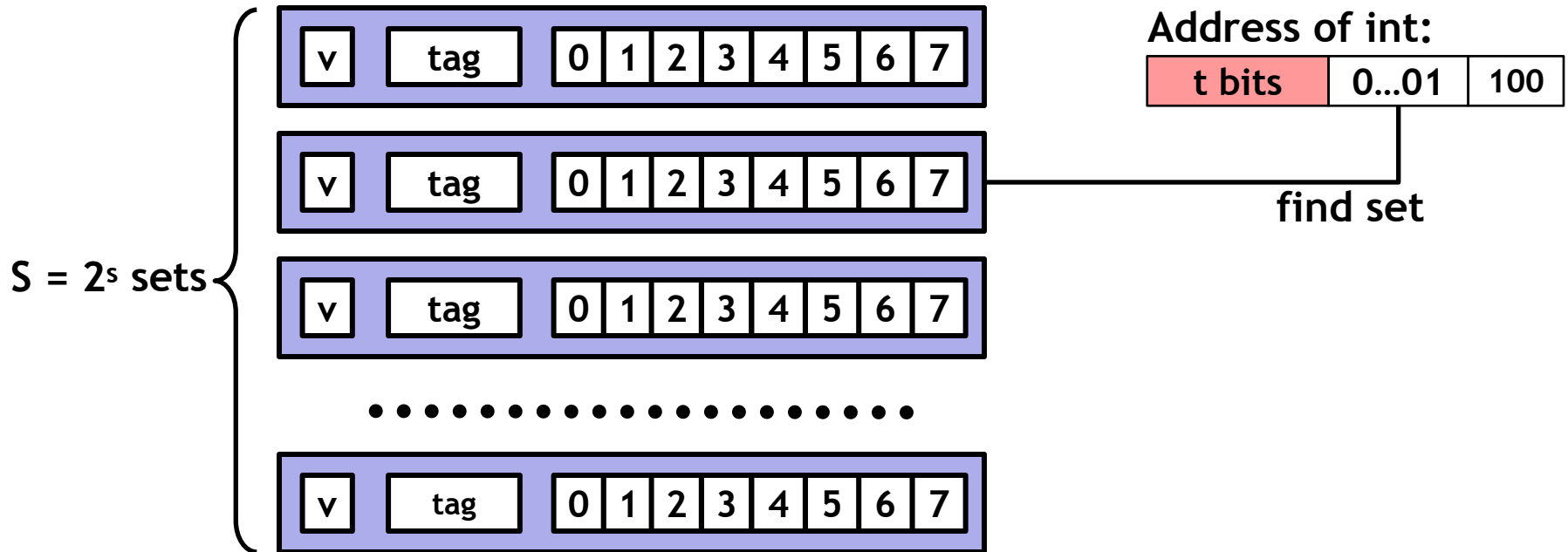S = $2^s$ sets

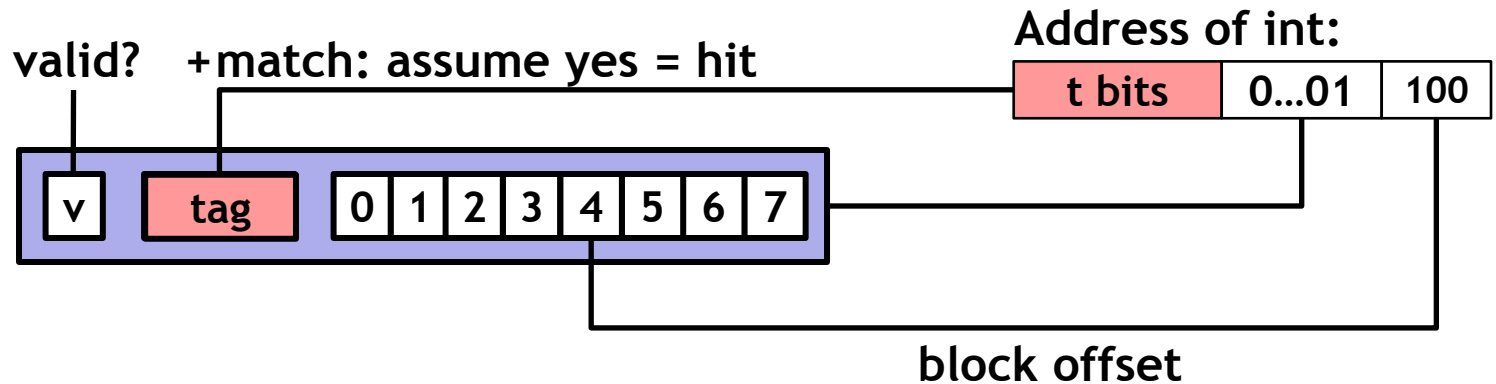Address of int:

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

find set

# Example: Direct Mapped Cache (E = 1)
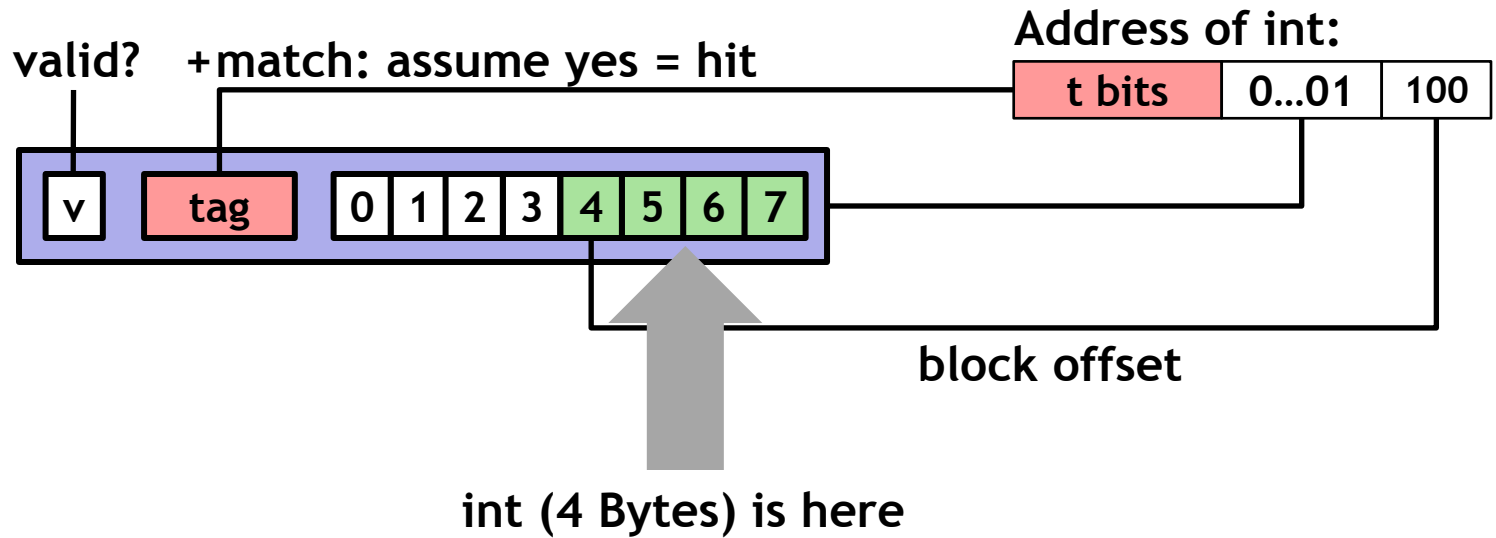
Direct mapped: One line per set
Assume: cache block size 8 bytes

# 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 line/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 | [$0000_2$] | 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:**

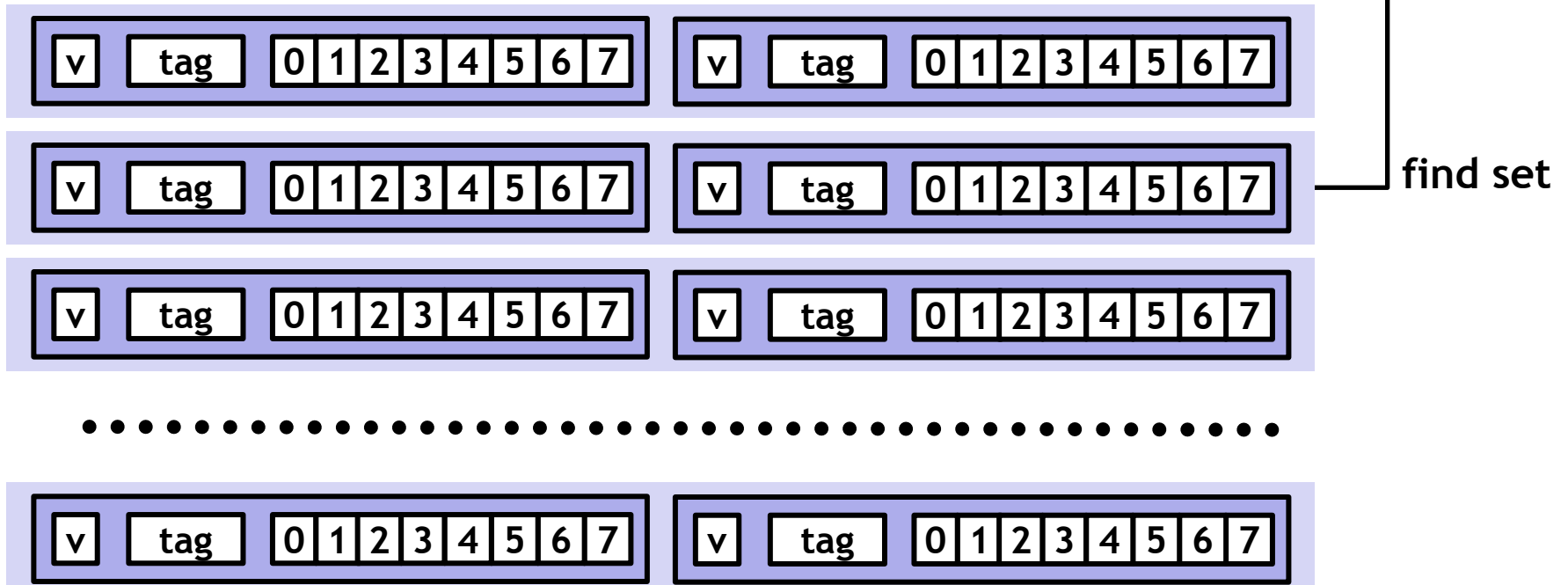| t bits | 0...01 | 100 |
|--------|--------|-----|



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:

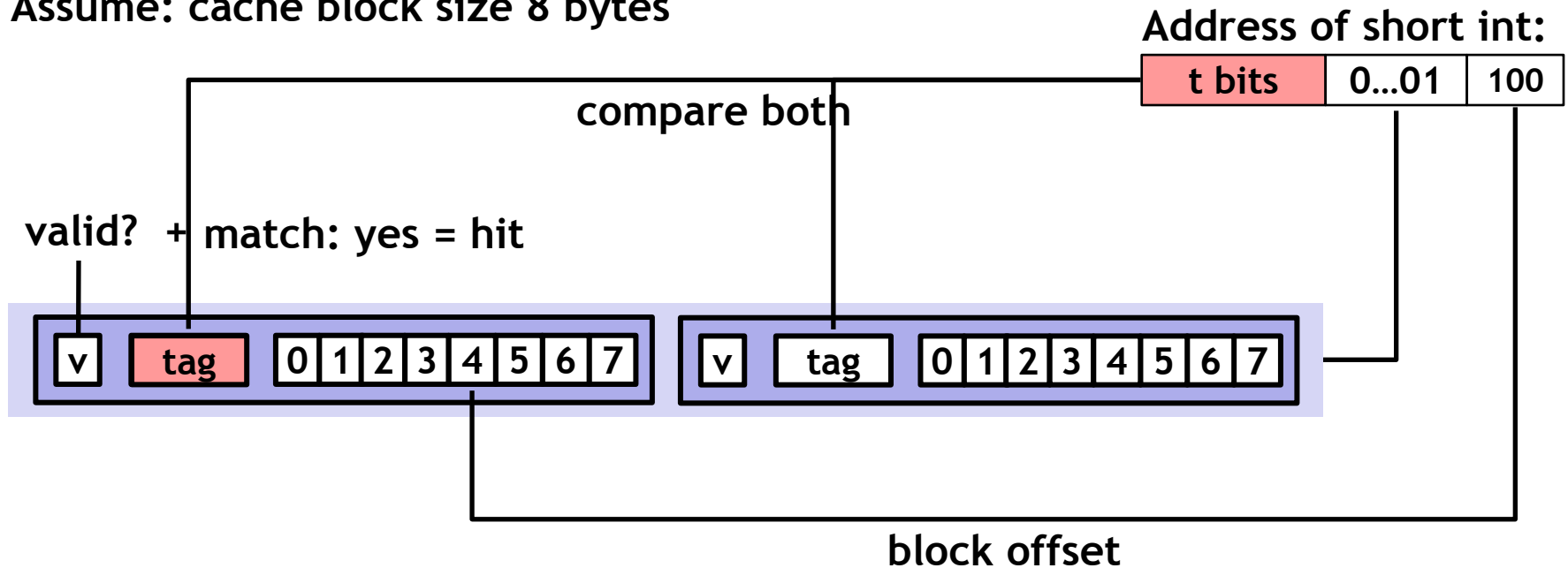| t bits | 0...01 | 100 |

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 lines/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 |     |        |

# 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 lines/set

Address trace (reads, one byte per read):

| 0 | [01$\underline{0}$0$_2$], |
| 1 | [00$\underline{0}$1$_2$], |
| 7 | [00$\underline{0}$1$_2$], |
| 8 | [11$\underline{0}$0$_2$], |
| 0 | [00$\underline{0}$0$_2$] |

**v      Tag        Block**

Set 0

Set 1

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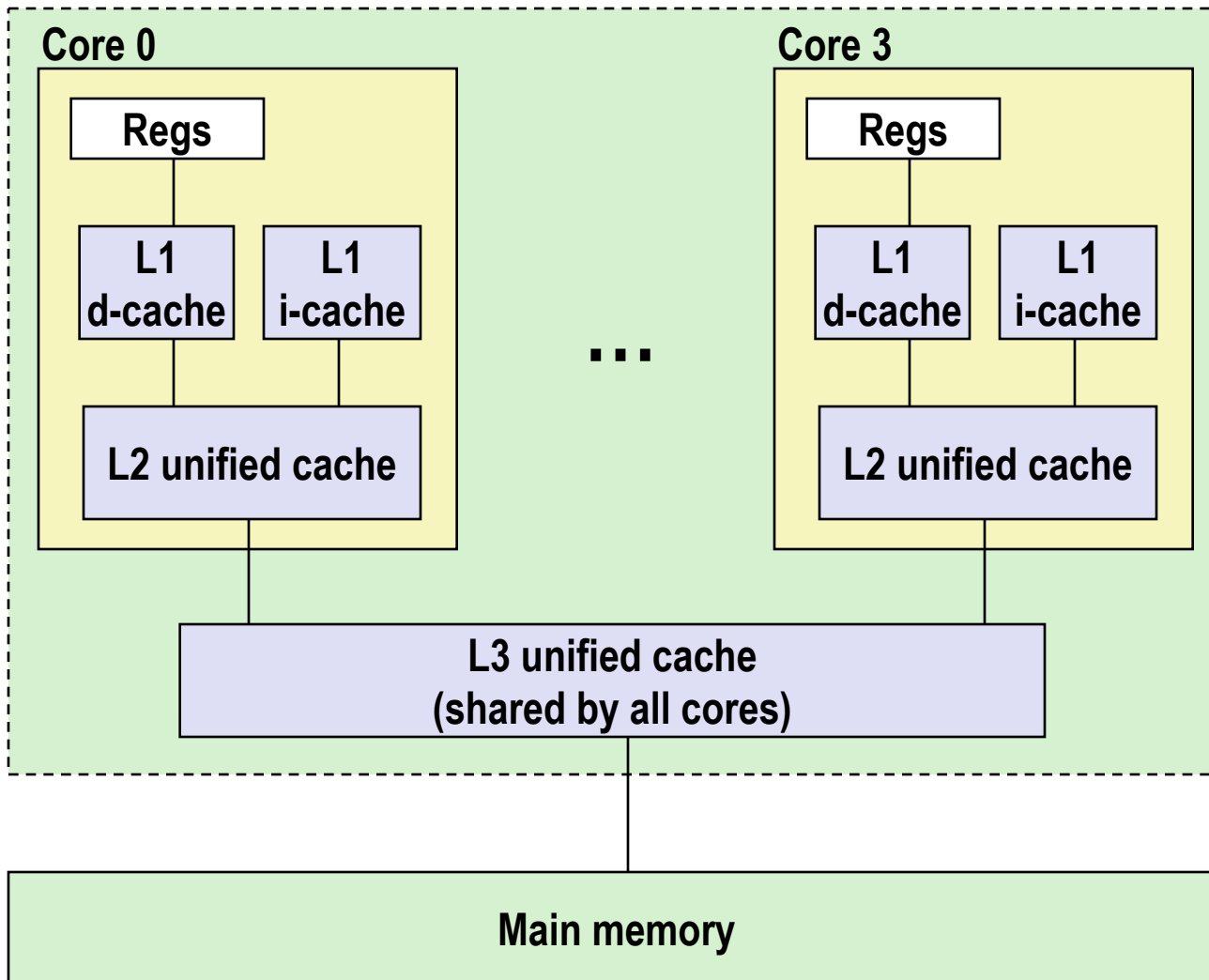- **What to do on a write-miss?**
  - Write-allocate (load into cache, update line in cache)
    - Good if more writes to the location follow
  - No-write-allocate (writes straight to memory, does not load into cache)
- **Typical**
  - Write-through + No-write-allocate
  - **Write-back + Write-allocate**

# 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.

# 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)