

The Memory Hierarchy

Instructor:

R. Shathanaa

Today

- Storage technologies and trends
- Locality of reference
- Caching in 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

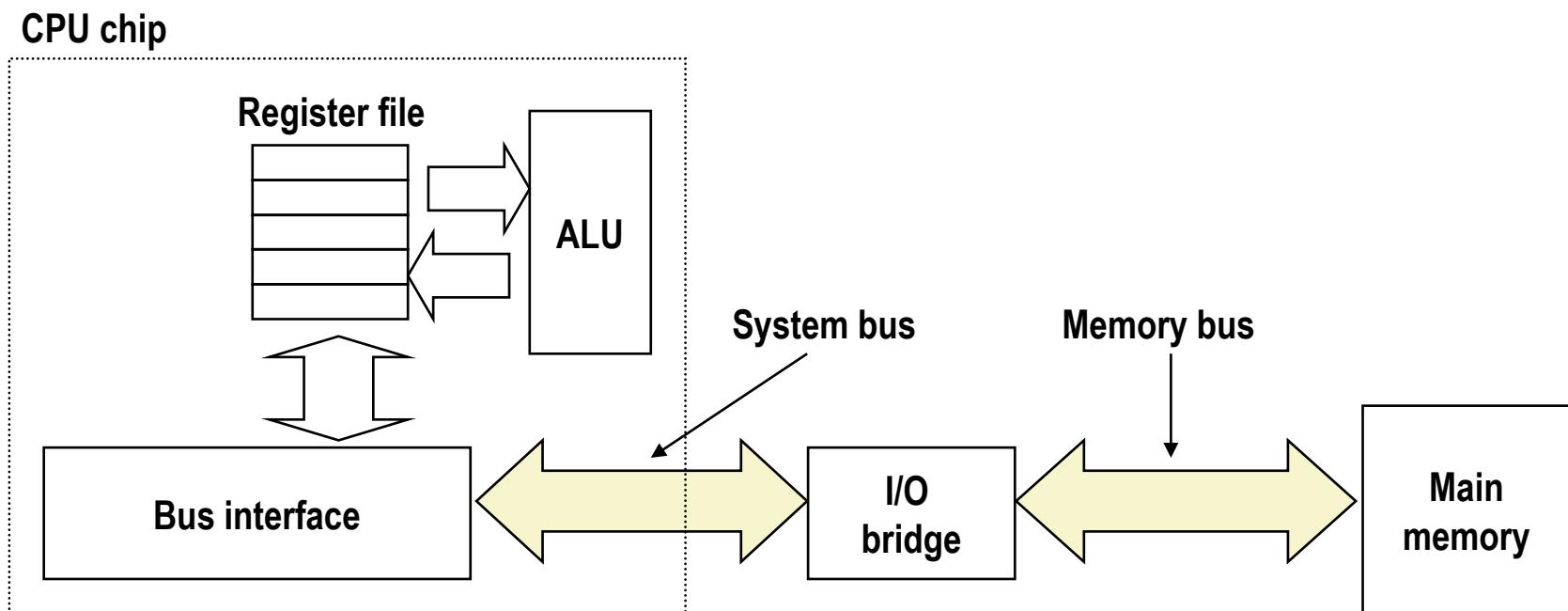
	Transistors per bit	Relative access time	Persistent?	Sensitive?	Relative cost	Applications
SRAM	6	1×	Yes	No	1,000×	Cache memory
DRAM	1	10×	No	Yes	1×	Main memory, frame buffers

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 – HDD in thumb drives, smart phones, mp3 players, tablets, laptops,...)

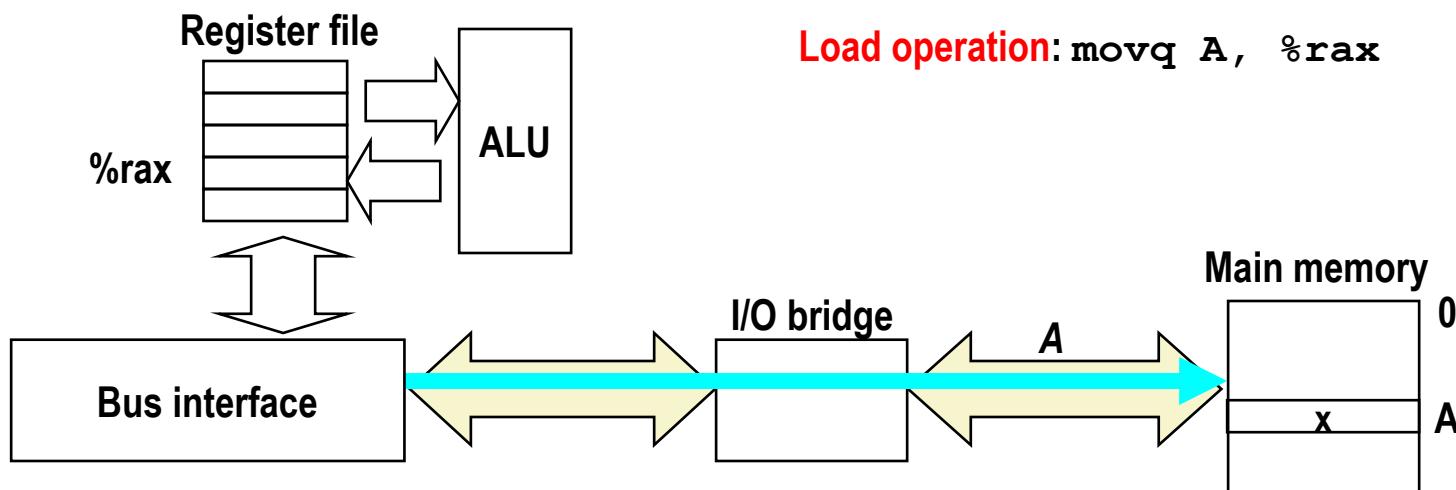
Traditional Bus Structure Connecting CPU and Memory

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



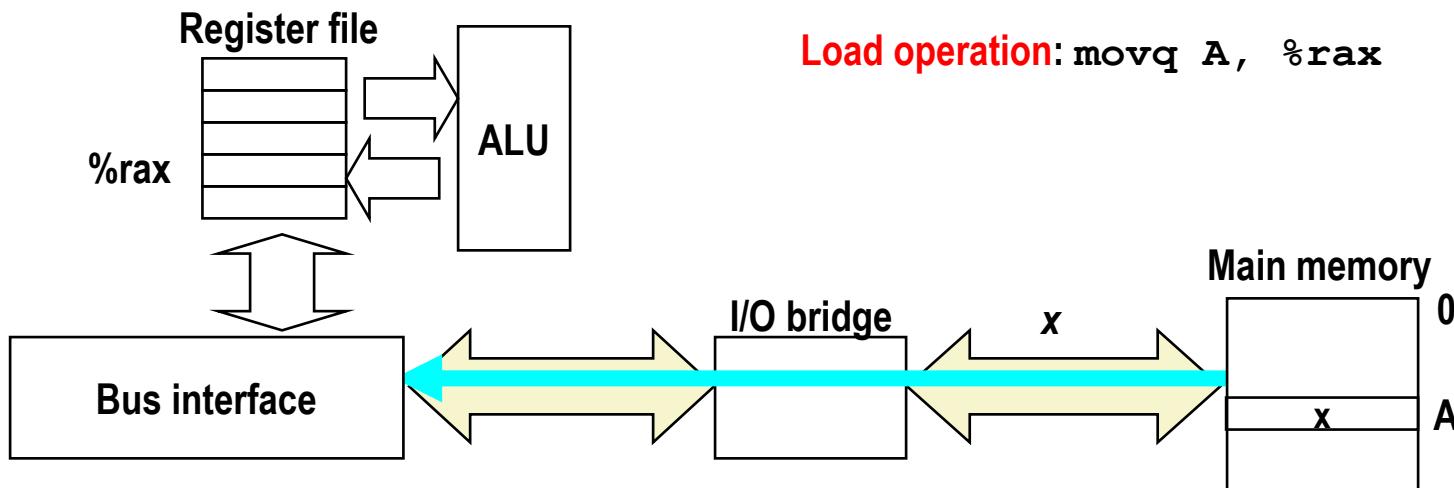
Memory Read Transaction (1)

- CPU places address A on the memory bus.



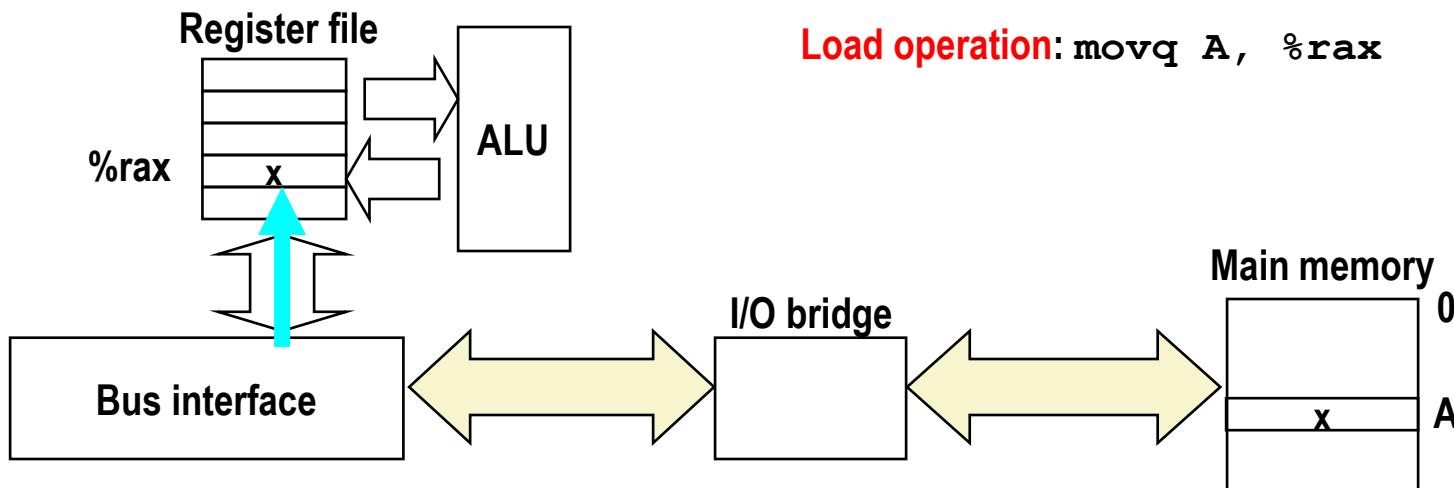
Memory Read Transaction (2)

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



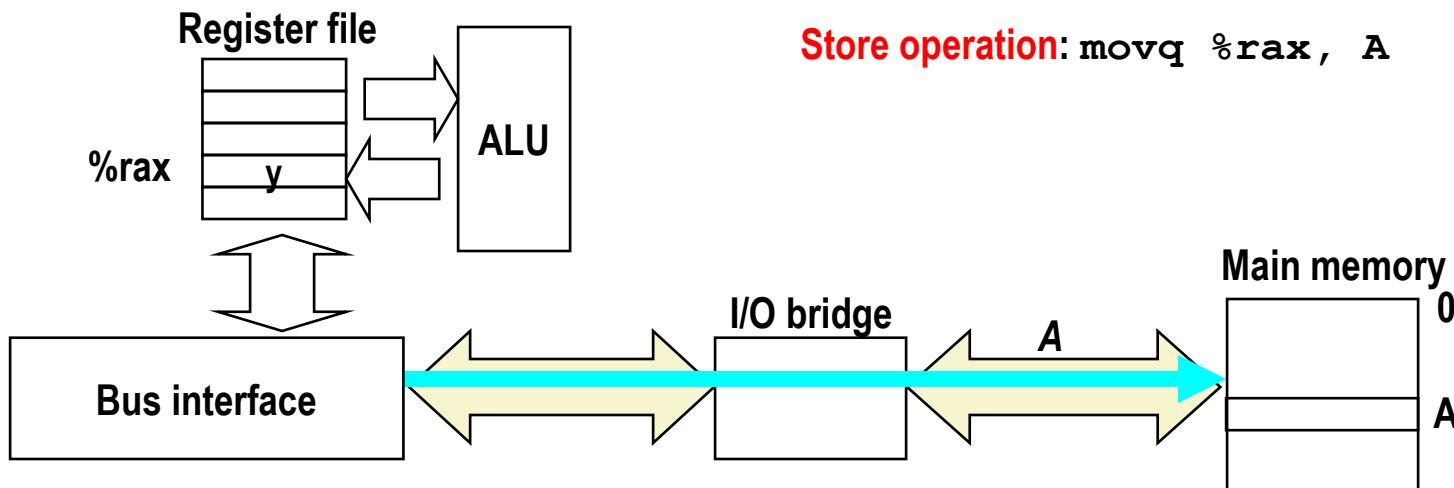
Memory Read Transaction (3)

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



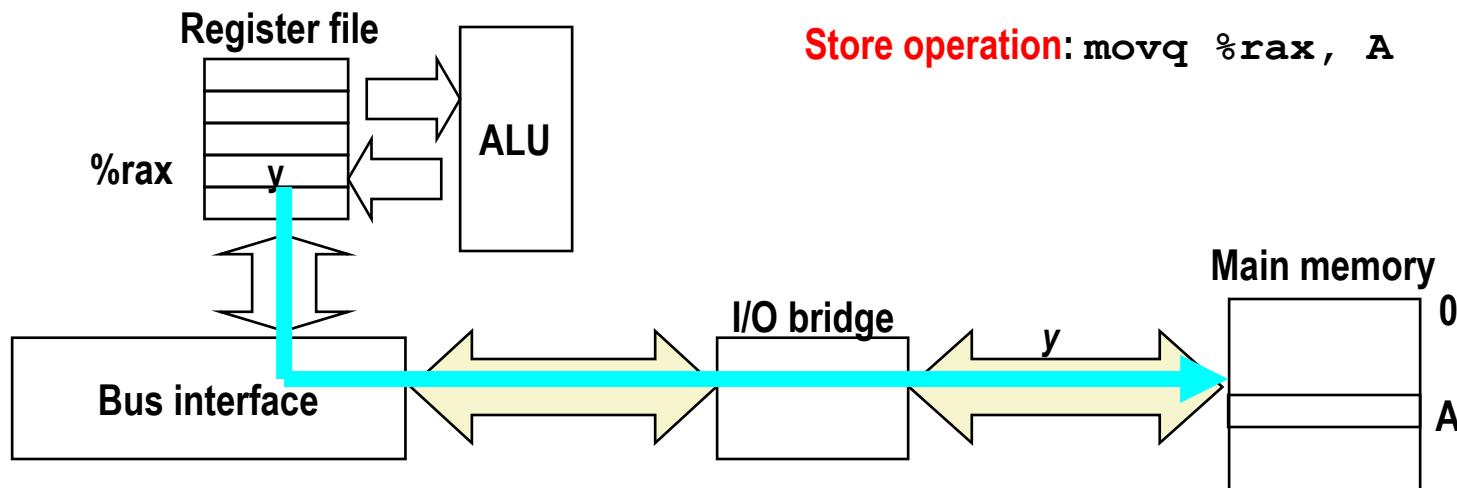
Memory Write Transaction (1)

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



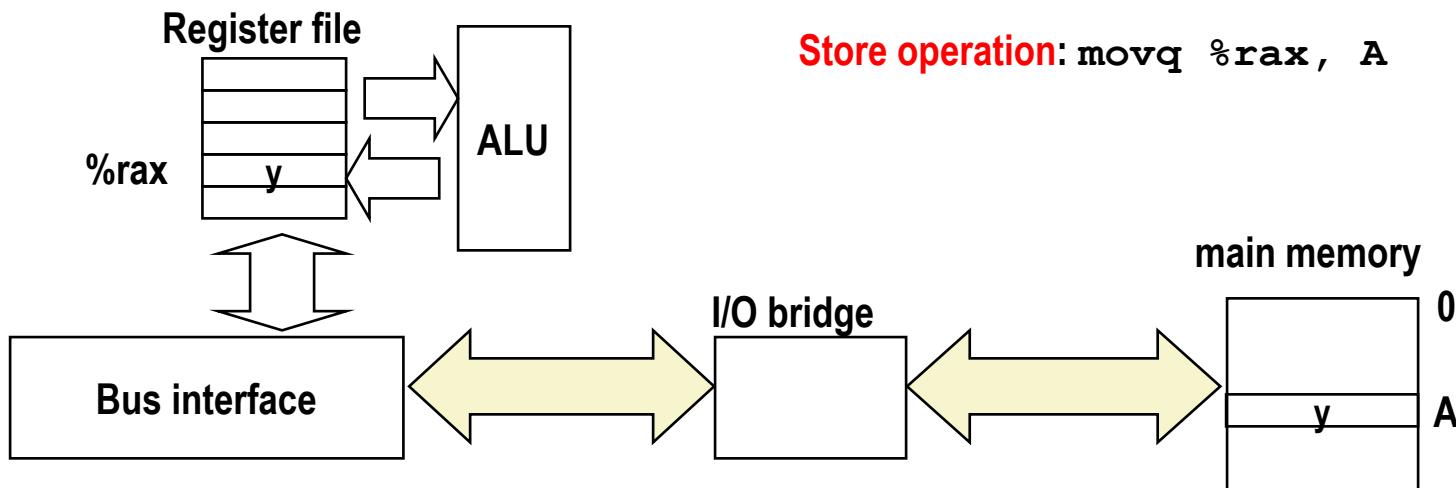
Memory Write Transaction (2)

- CPU places data word y on the bus.



Memory Write Transaction (3)

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



What's Inside A Disk Drive?

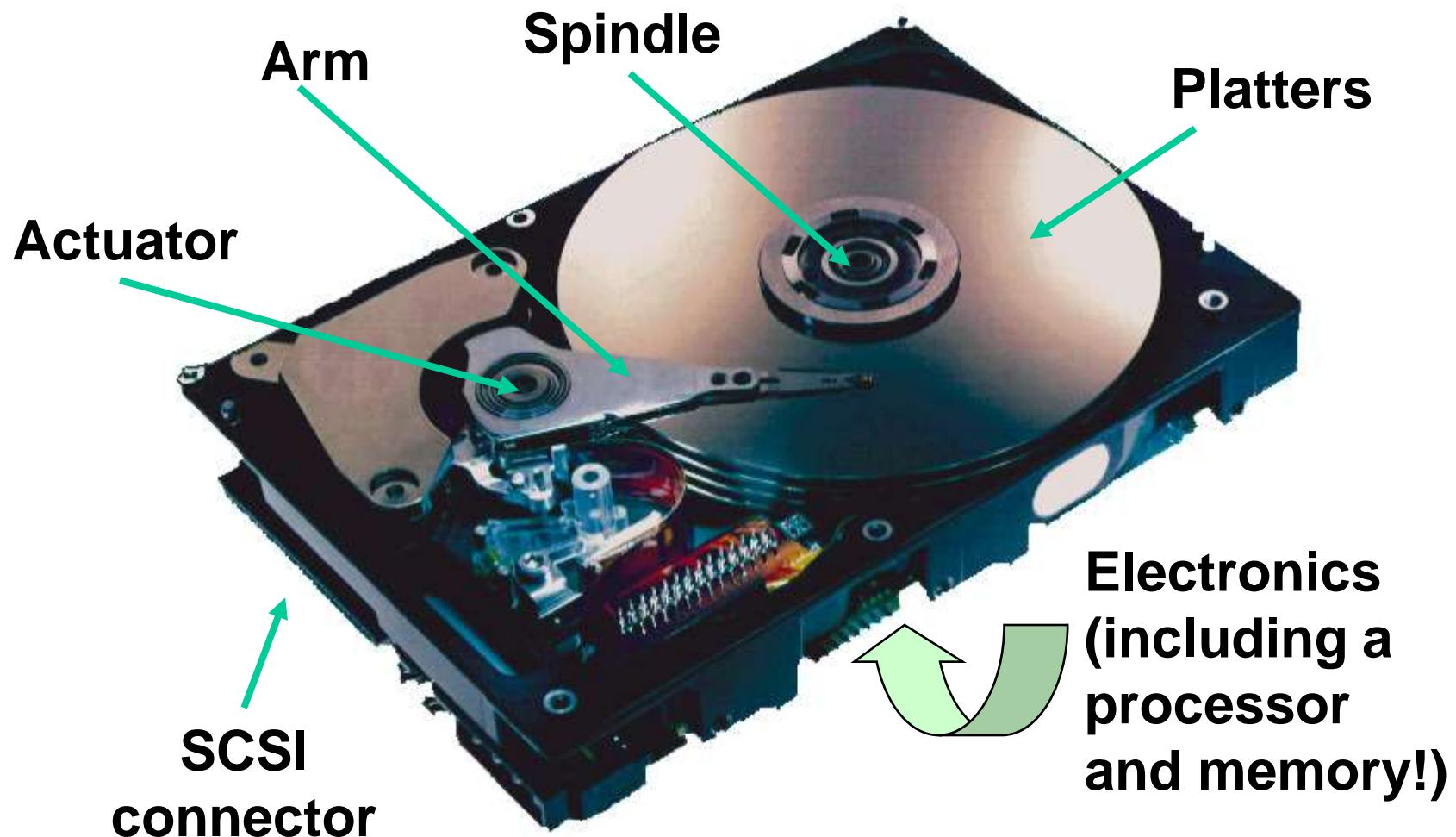
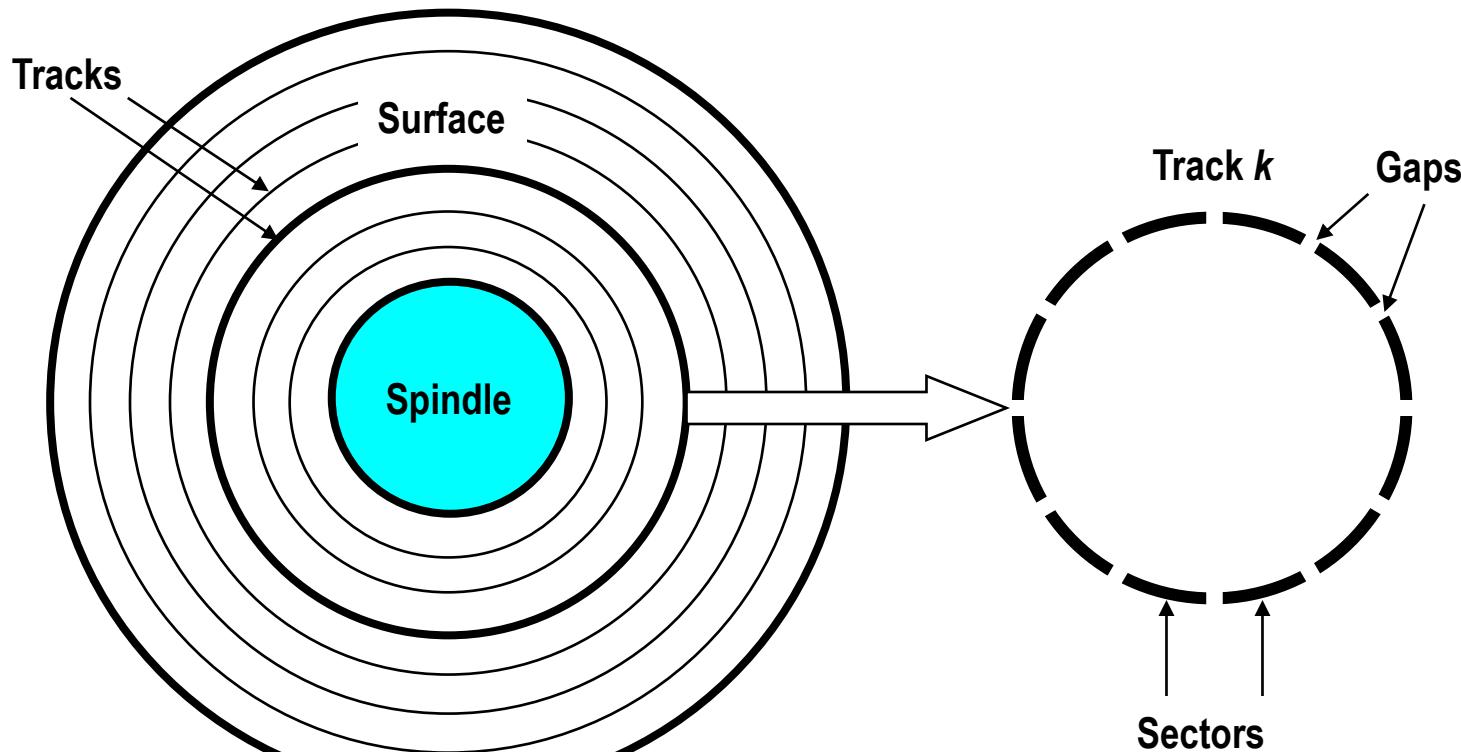


Image courtesy of Seagate Technology

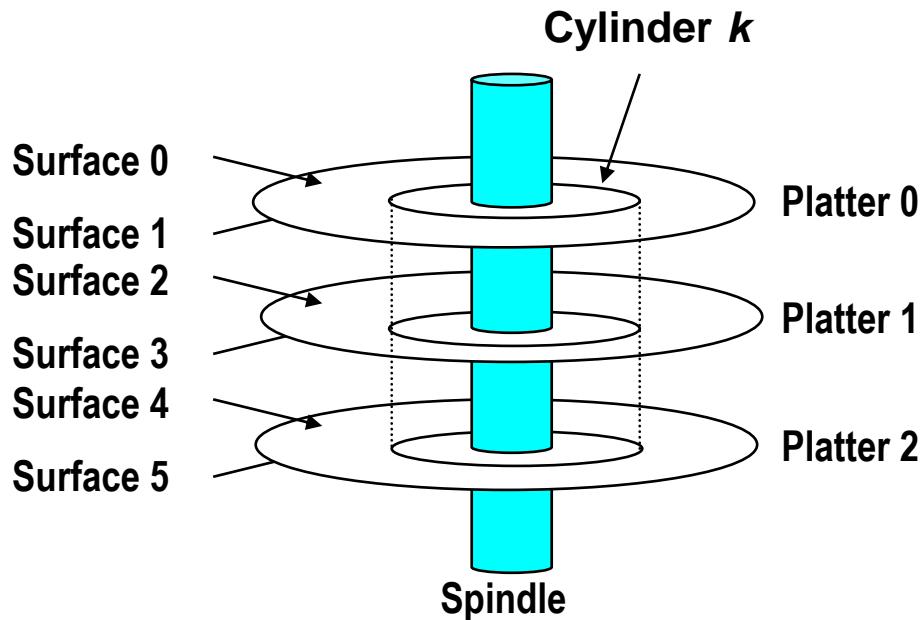
Disk Geometry

- 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 (Multiple-Platter View)

- Aligned tracks form a cylinder.

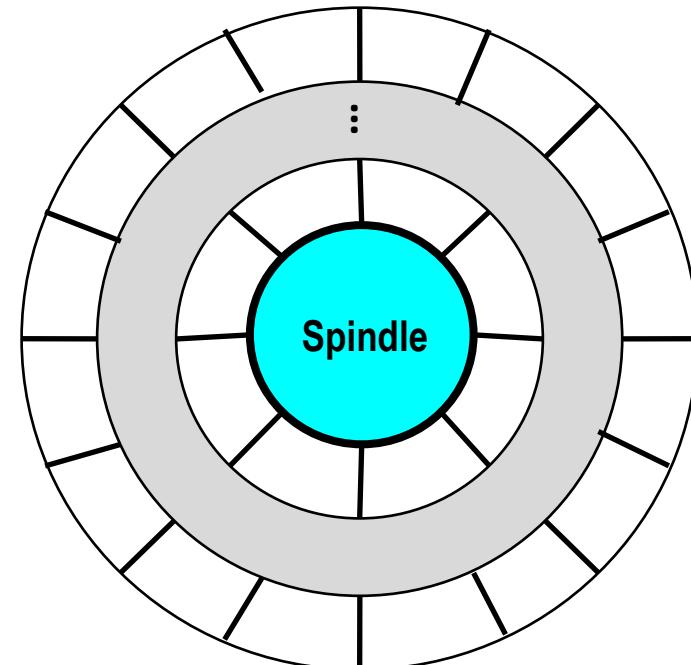


Disk Capacity

- **Capacity:** maximum number of bits that can be stored.
 - Vendors express capacity in units of gigabytes (GB), where
 $1 \text{ GB} = 10^9 \text{ Bytes}$.
- **Capacity is determined by these technology factors:**
 - **Recording density** (bits/in): number of bits that can be squeezed into a 1 inch segment of a track.
 - **Track density** (tracks/in): number of tracks that can be squeezed into a 1 inch radial segment.
 - **Areal density** (bits/in²): product of recording and track density.

Recording zones

- Modern disks partition tracks into disjoint subsets called **recording zones**
 - Each track in a zone has the same number of sectors, determined by the circumference of innermost track.
 - Each zone has a different number of sectors/track, outer zones have more sectors/track than inner zones.
 - So we use **average** number of sectors/track when computing capacity.



Computing Disk Capacity

**Capacity = (# bytes/sector) x (avg. # sectors/track) x
(# tracks/surface) x (# surfaces/platter) x
(# platters/disk)**

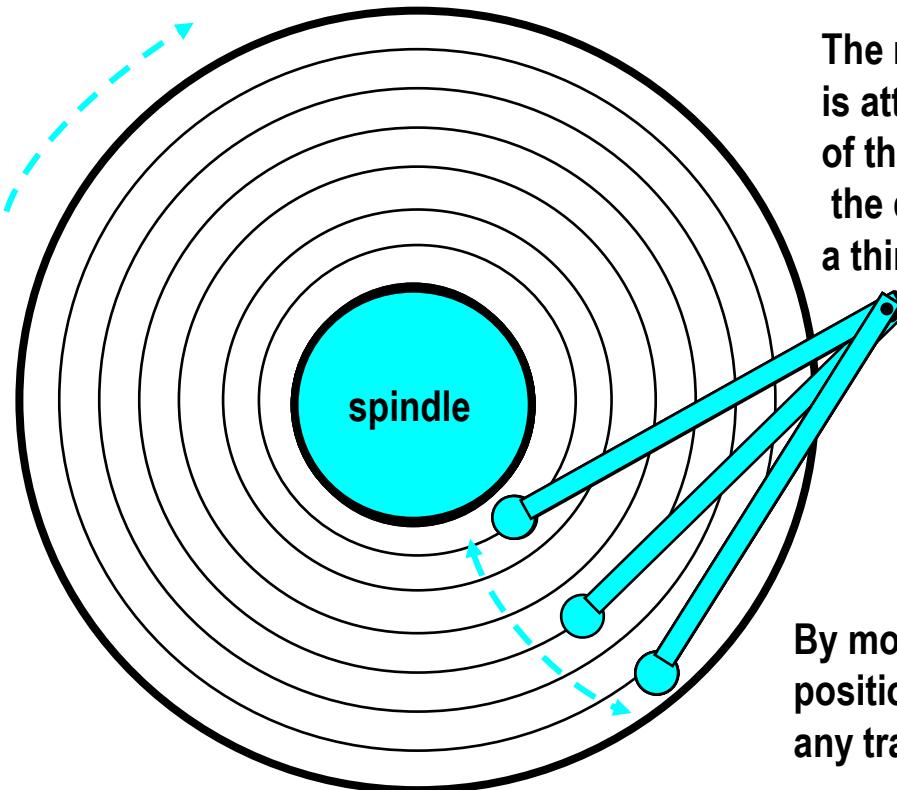
Example:

- 512 bytes/sector
- 300 sectors/track (on average)
- 20,000 tracks/surface
- 2 surfaces/platter
- 5 platters/disk

$$\begin{aligned}\text{Capacity} &= 512 \times 300 \times 20000 \times 2 \times 5 \\ &= 30,720,000,000 \\ &= 30.72 \text{ GB}\end{aligned}$$

Disk Operation (Single-Platter View)

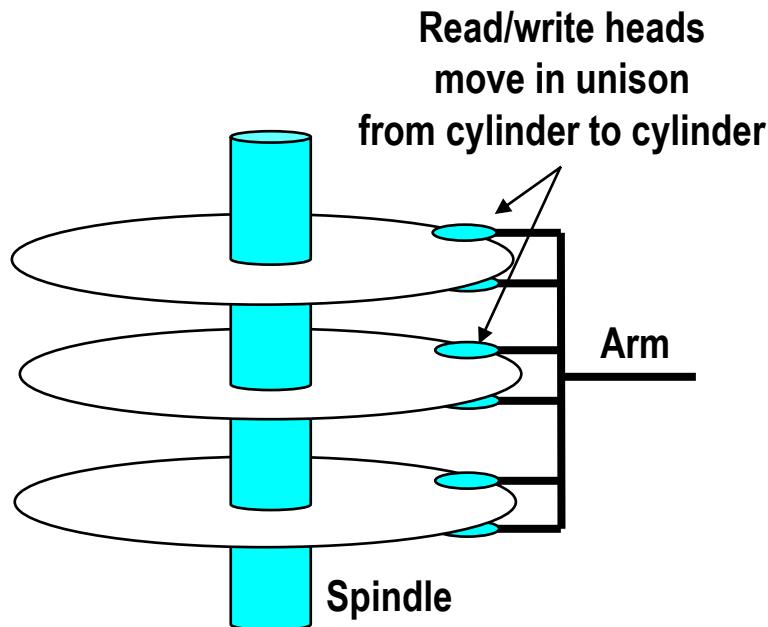
The disk surface spins at a fixed rotational rate



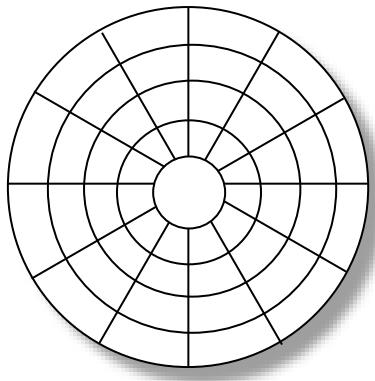
The read/write *head* is attached to the end of the *arm* and flies over the disk surface on a thin cushion of air.

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

Disk Operation (Multi-Platter View)



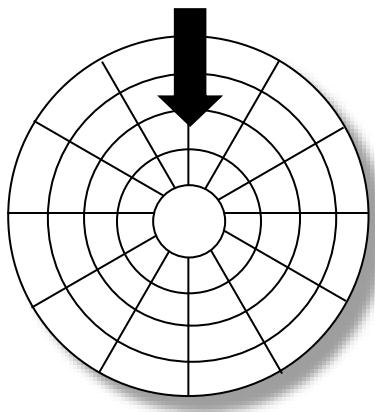
Disk Structure - top view of single platter



Surface organized into tracks

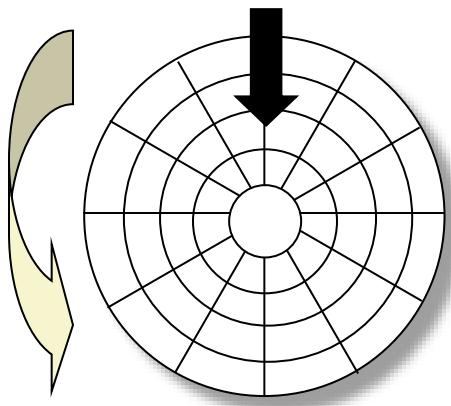
Tracks divided into sectors

Disk Access



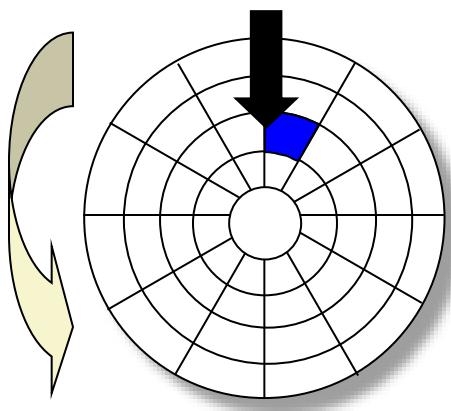
Head in position above a track

Disk Access



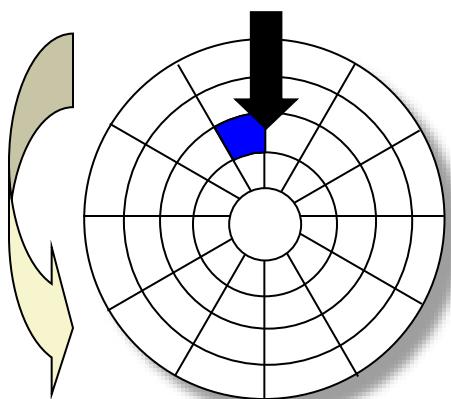
Rotation is counter-clockwise

Disk Access – Read



About to read blue sector

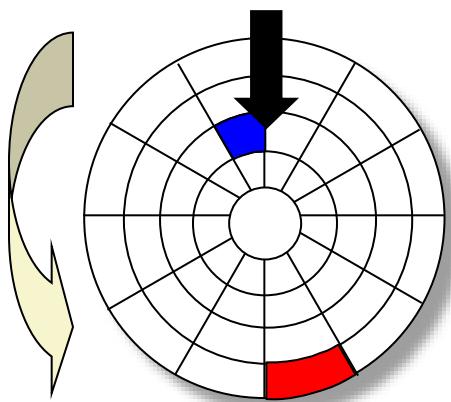
Disk Access – Read



After **BLUE** read

After reading blue sector

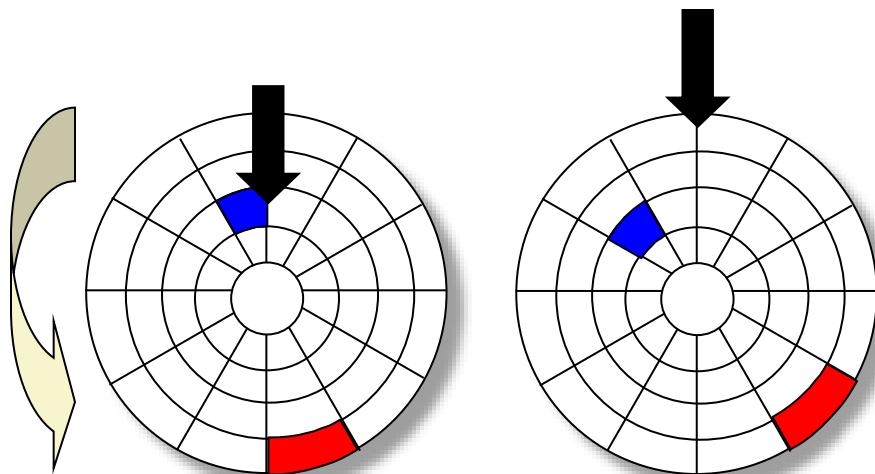
Disk Access – Read



After **BLUE** read

Red request scheduled next

Disk Access – Seek

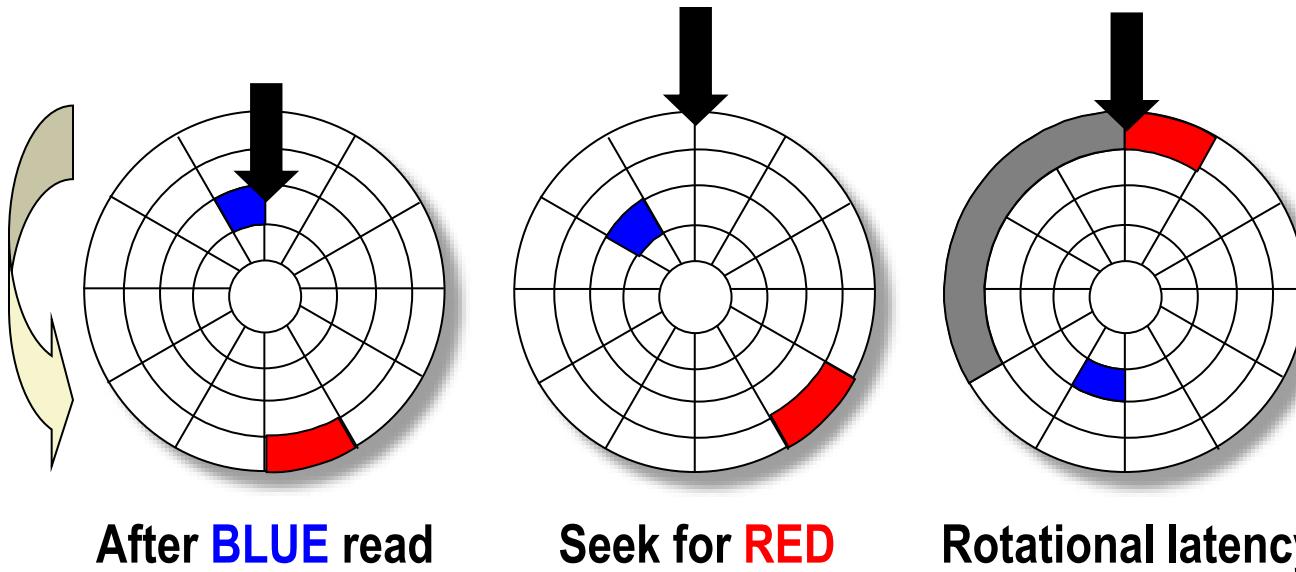


After **BLUE** read

Seek for **RED**

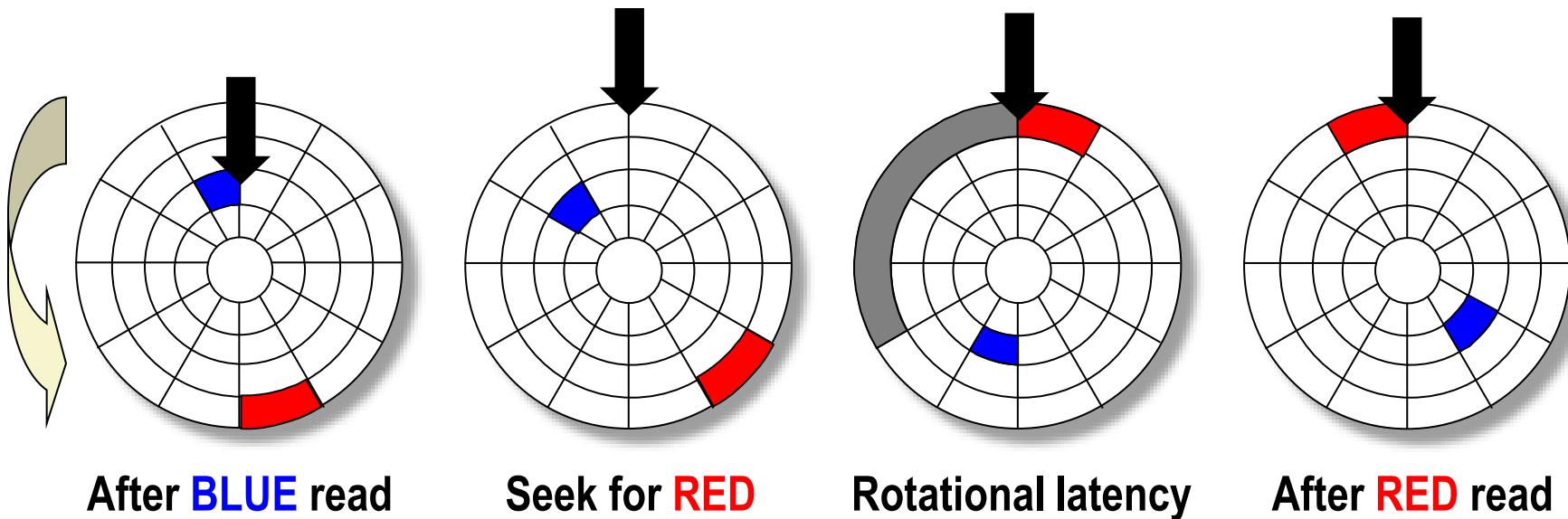
Seek to red's track

Disk Access – Rotational Latency



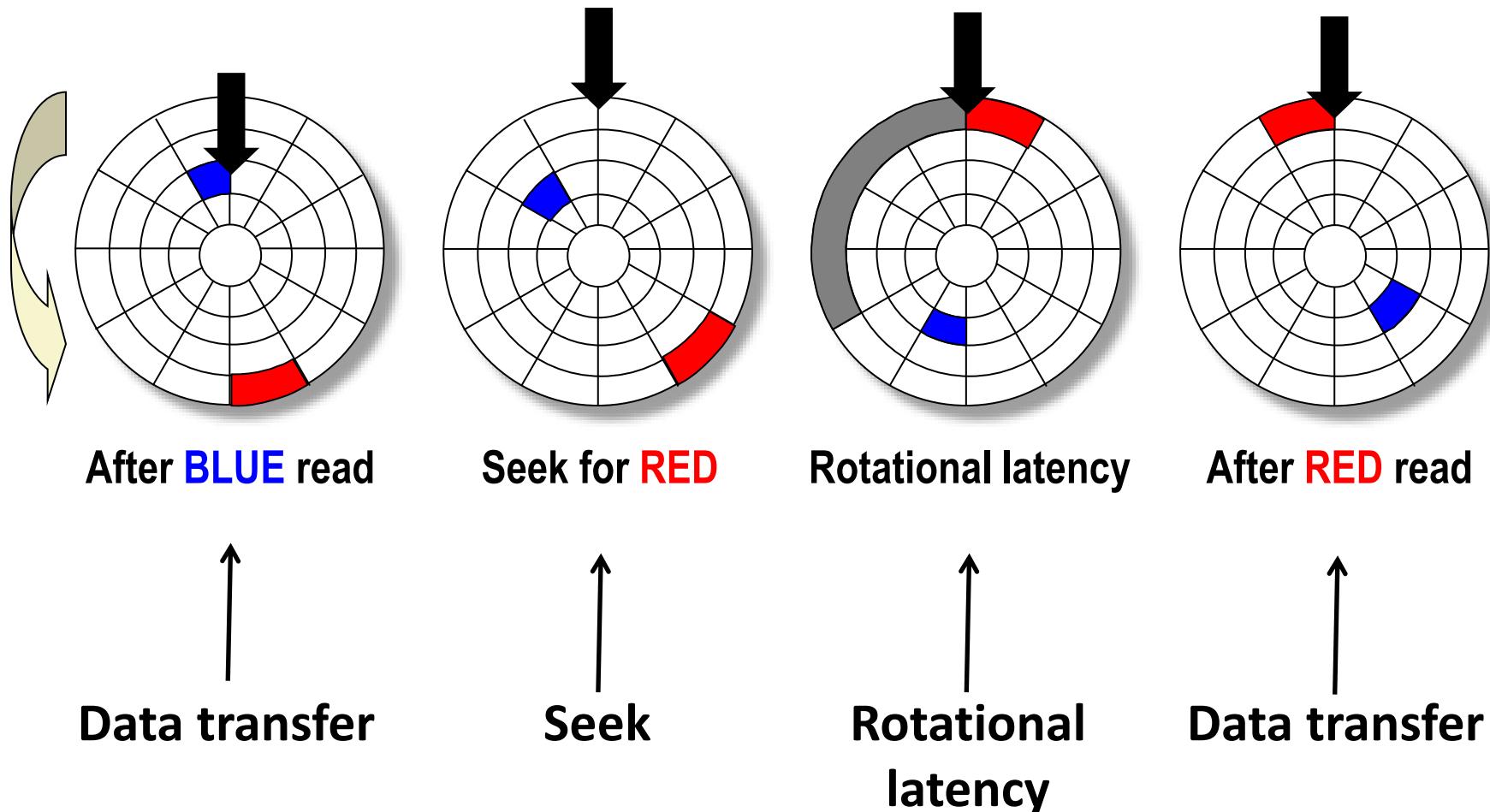
Wait for red sector to rotate around

Disk Access – Read



Complete read of red

Disk Access – Service Time Components



Disk Access Time

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

Disk Access Time Example

■ Given:

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

■ Derived:

- $T_{avg\ rotation} = 1/2 \times (60\ secs/7200\ RPM) \times 1000\ ms/sec = 4\ ms.$
- $T_{avg\ transfer} = 60/7200\ RPM \times 1/400\ secs/track \times 1000\ ms/sec = 0.02\ ms$
- $T_{access} = 9\ ms + 4\ ms + 0.02\ ms$

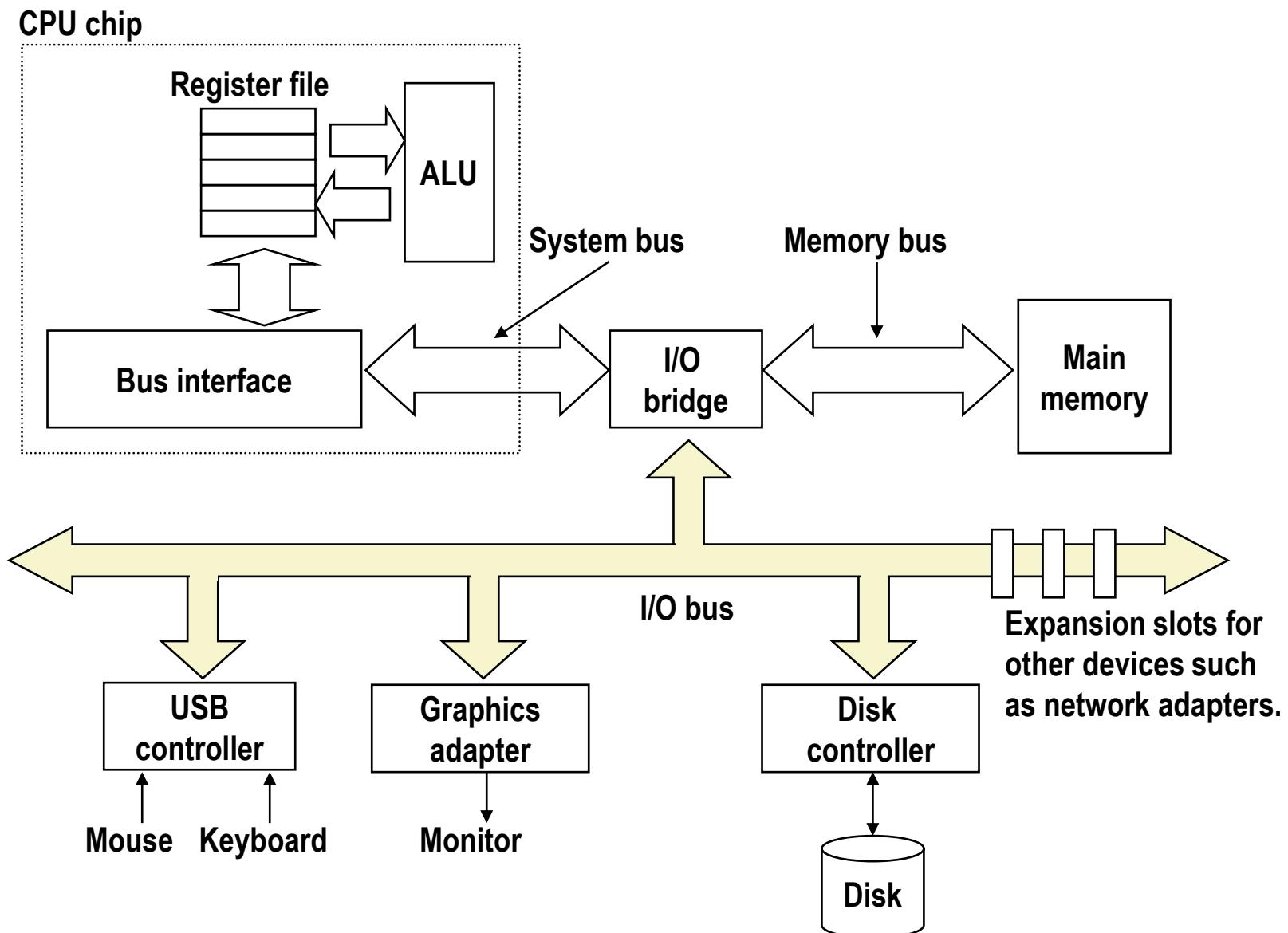
■ Important points:

- 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 logical blocks (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.

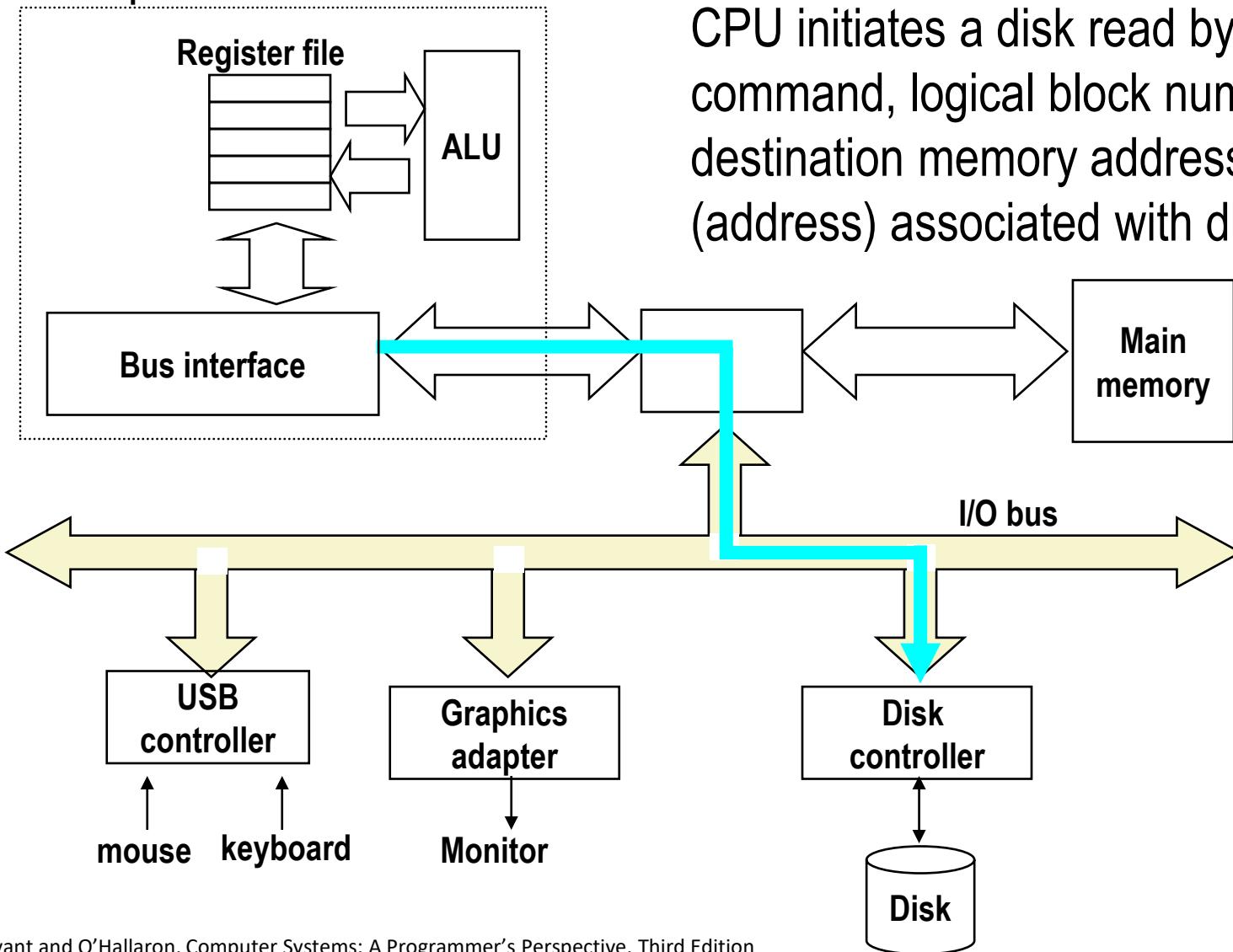
I/O Bus



Reading a Disk Sector (1)

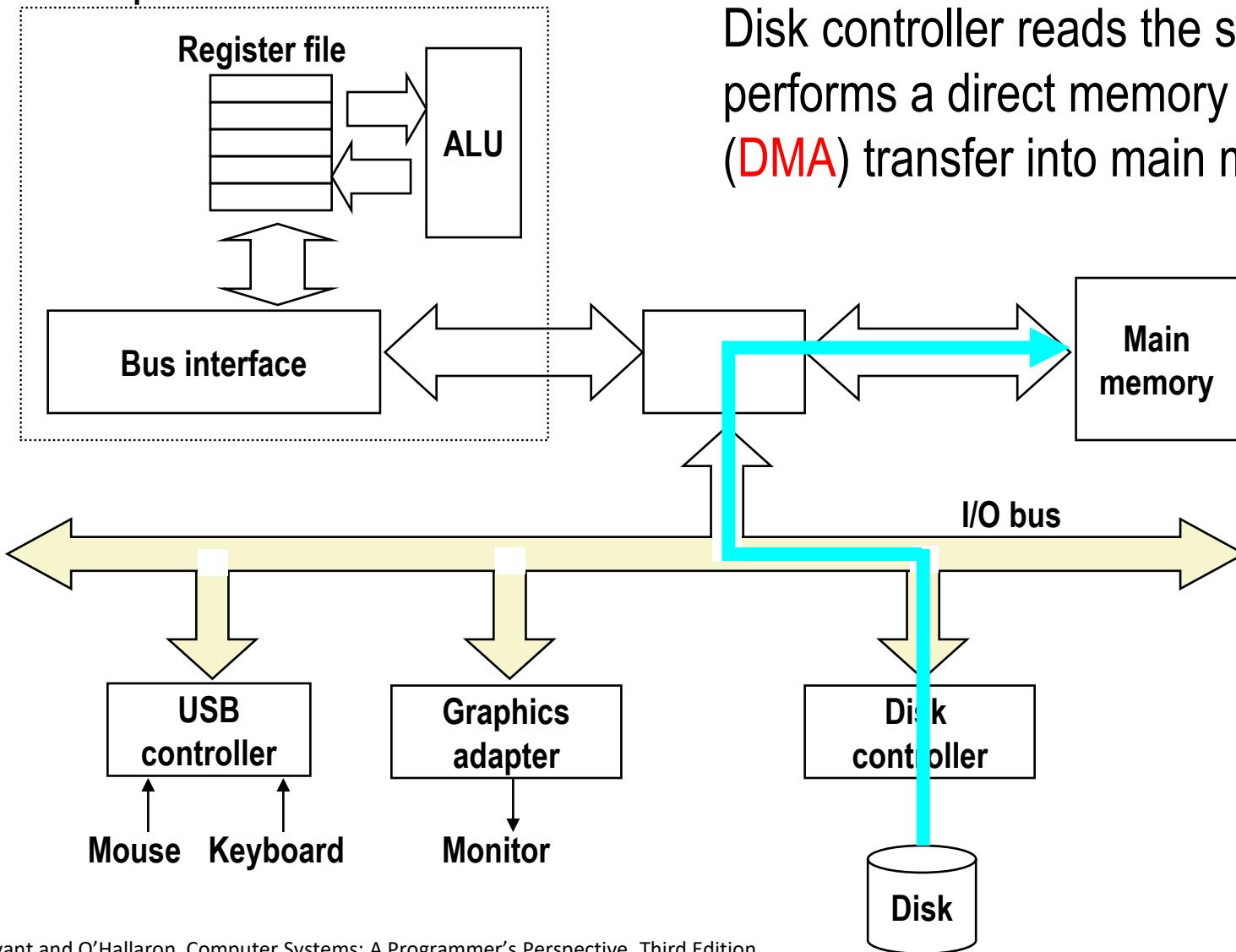
Memorymapped I/O

CPU chip



Reading a Disk Sector (2)

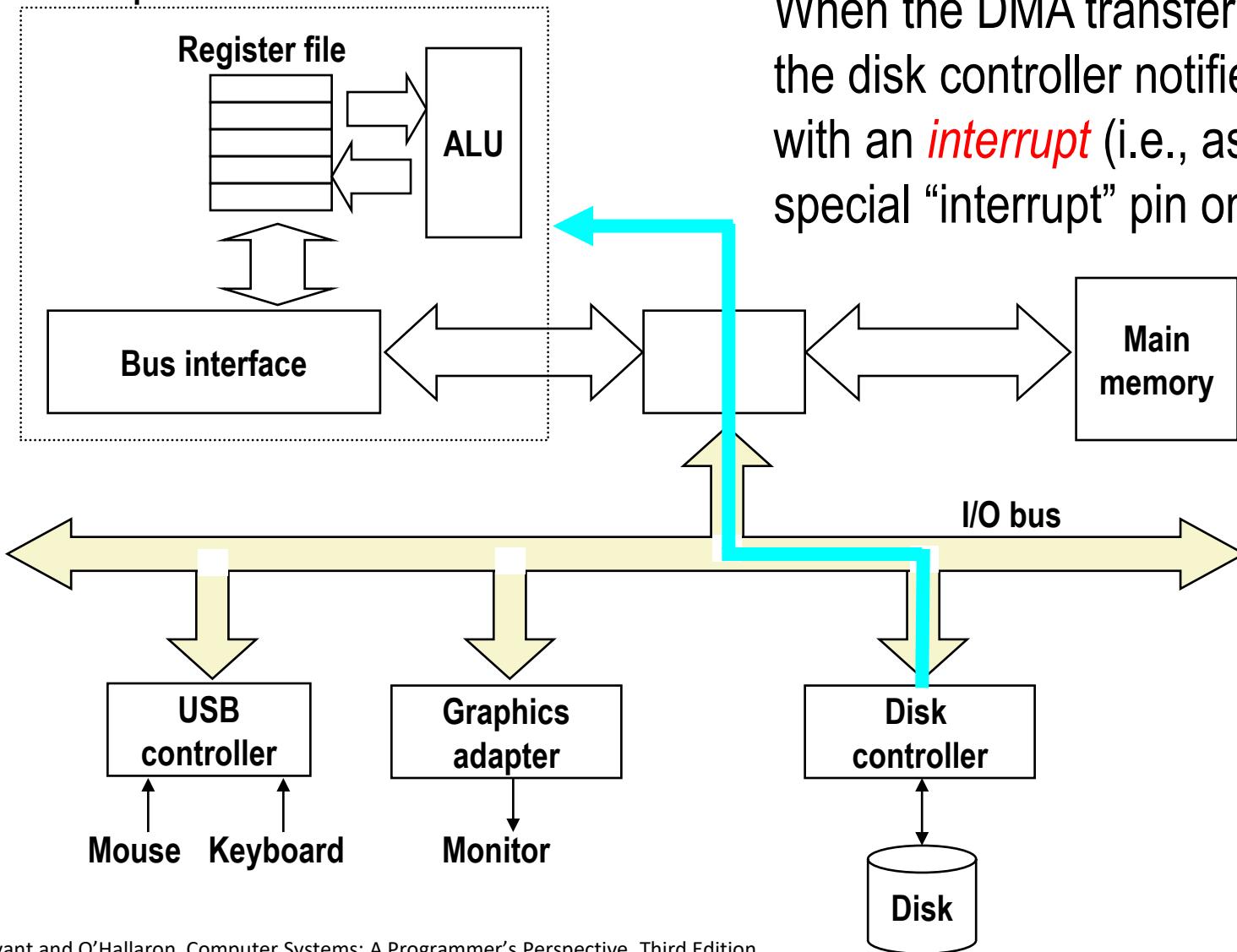
CPU chip



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

Reading a Disk Sector (3)

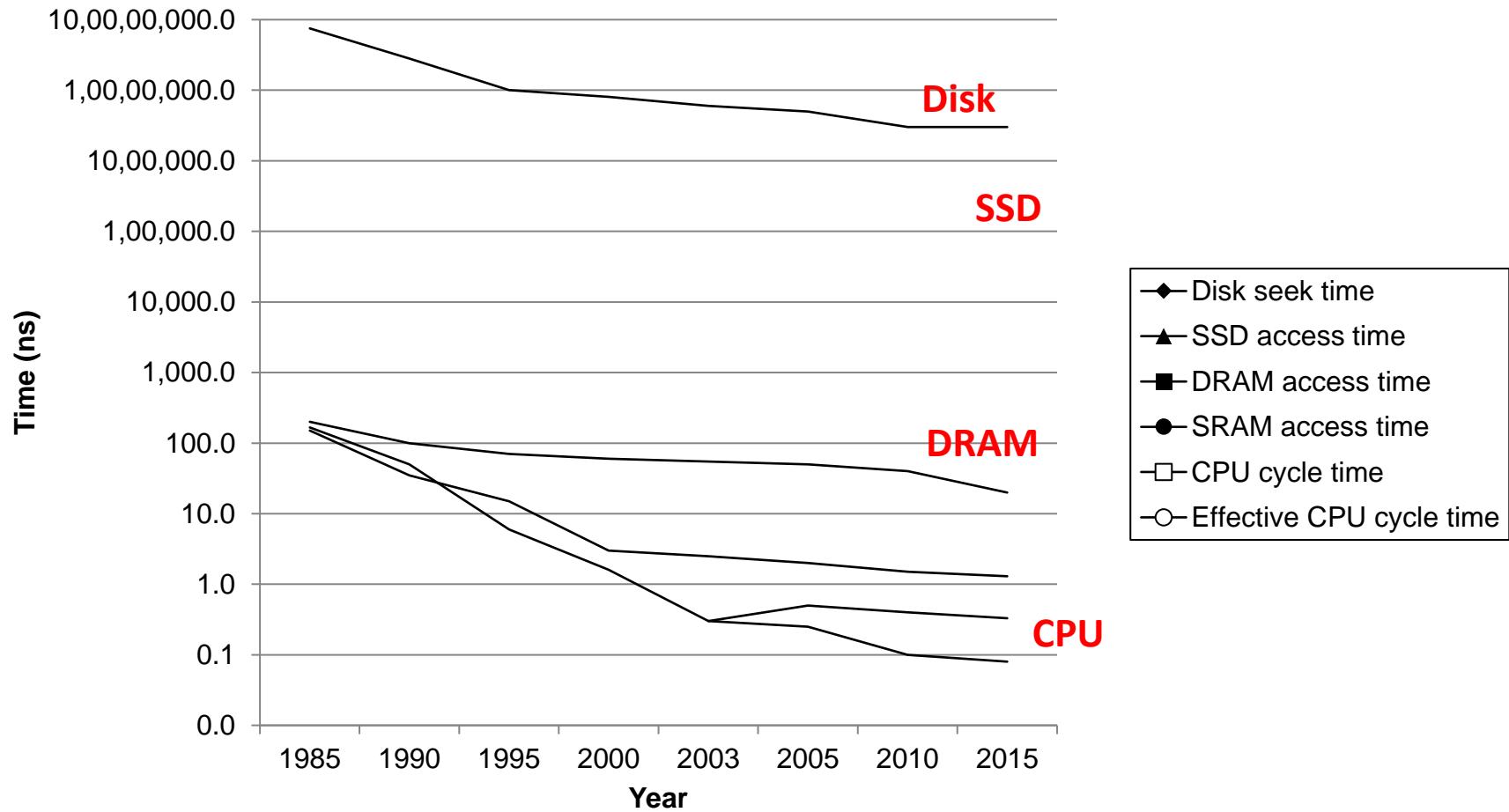
CPU chip



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

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 **locality**

Today

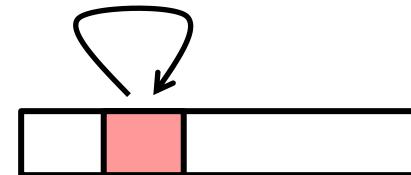
- Storage technologies and trends
- Locality of reference
- Caching in the memory hierarchy

Locality

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

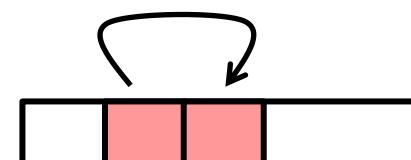
- **Temporal locality:**

- 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.

Spatial locality

Temporal locality

■ Instruction references

- Reference instructions in sequence.
- Cycle through loop repeatedly.

Spatial locality

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)?

```
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;
}
```

Memory Hierarchies

- **Some fundamental and enduring properties of hardware and software:**
 - Fast storage technologies cost more per byte, have less capacity, and require more power (heat!).
 - The gap between CPU and main memory speed is widening.
 - Well-written programs tend to exhibit good locality.
- **These fundamental properties complement each other beautifully.**
- **They suggest an approach for organizing memory and storage systems known as a **memory hierarchy**.**

Today

- Storage technologies and trends
- Locality of reference
- Caching in the memory hierarchy

Example Memory Hierarchy

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

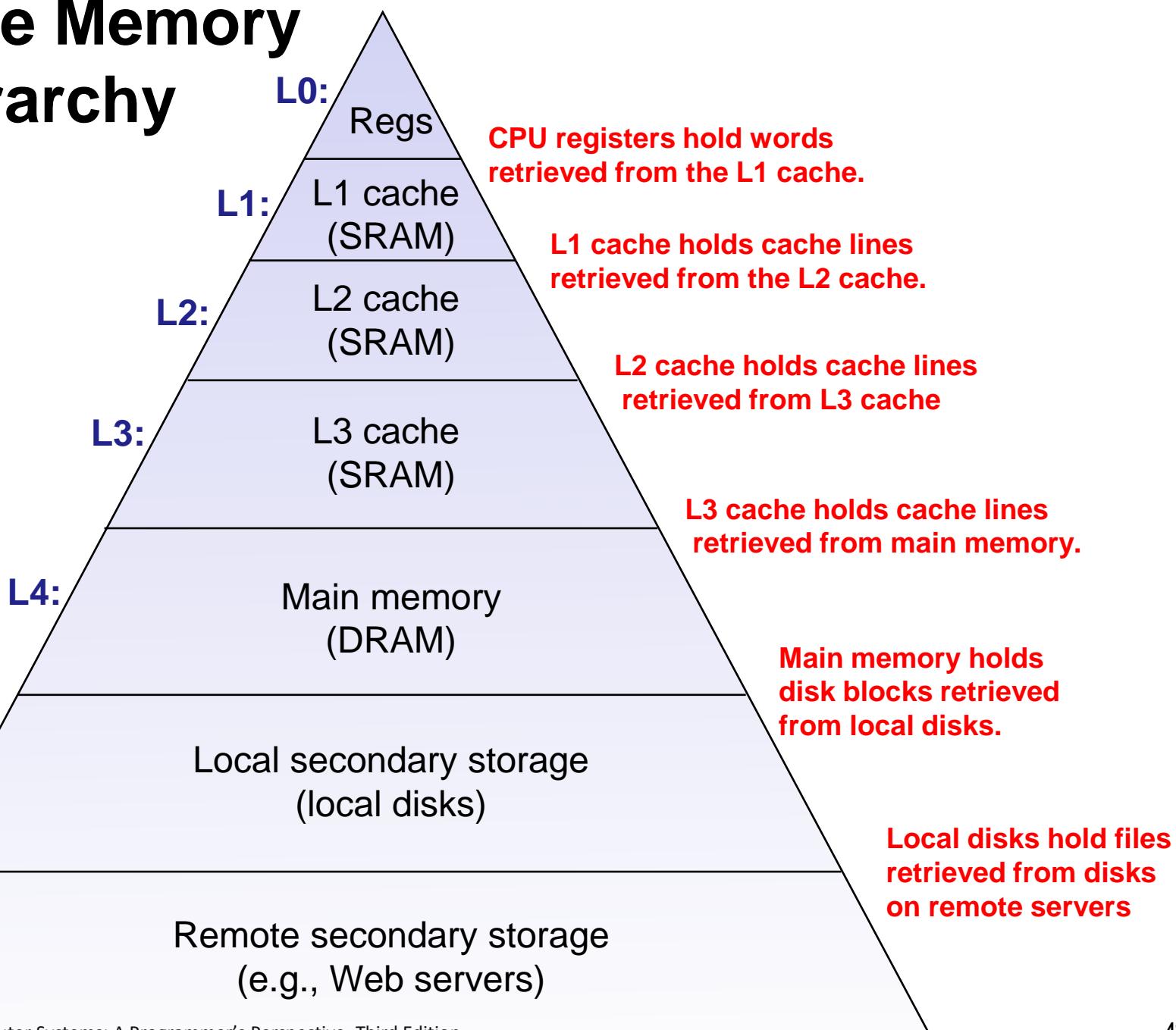
Larger,
slower,
and
cheaper
(per byte)
storage
devices

L6:

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

L5:

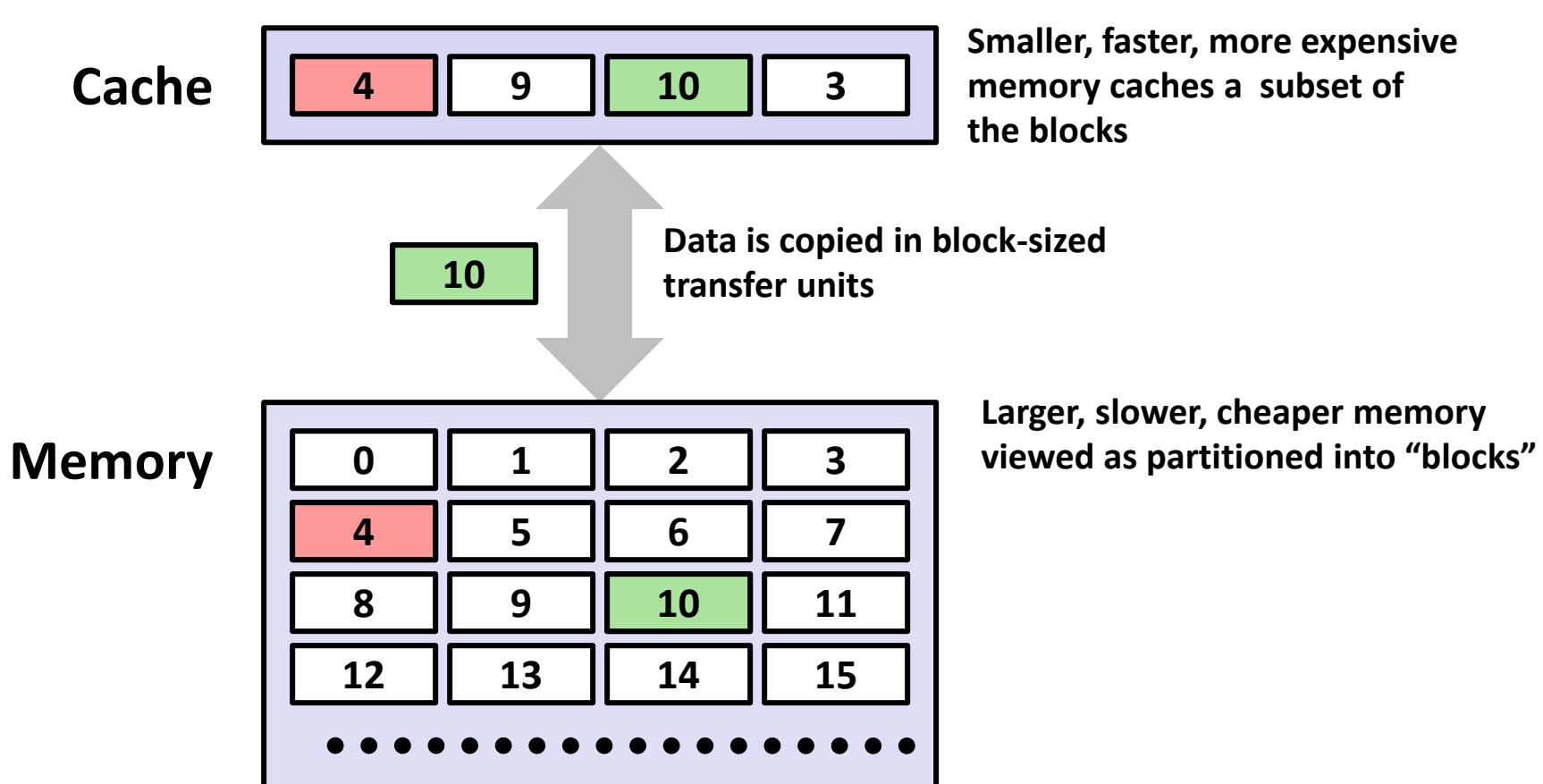
Local secondary storage
(local disks)



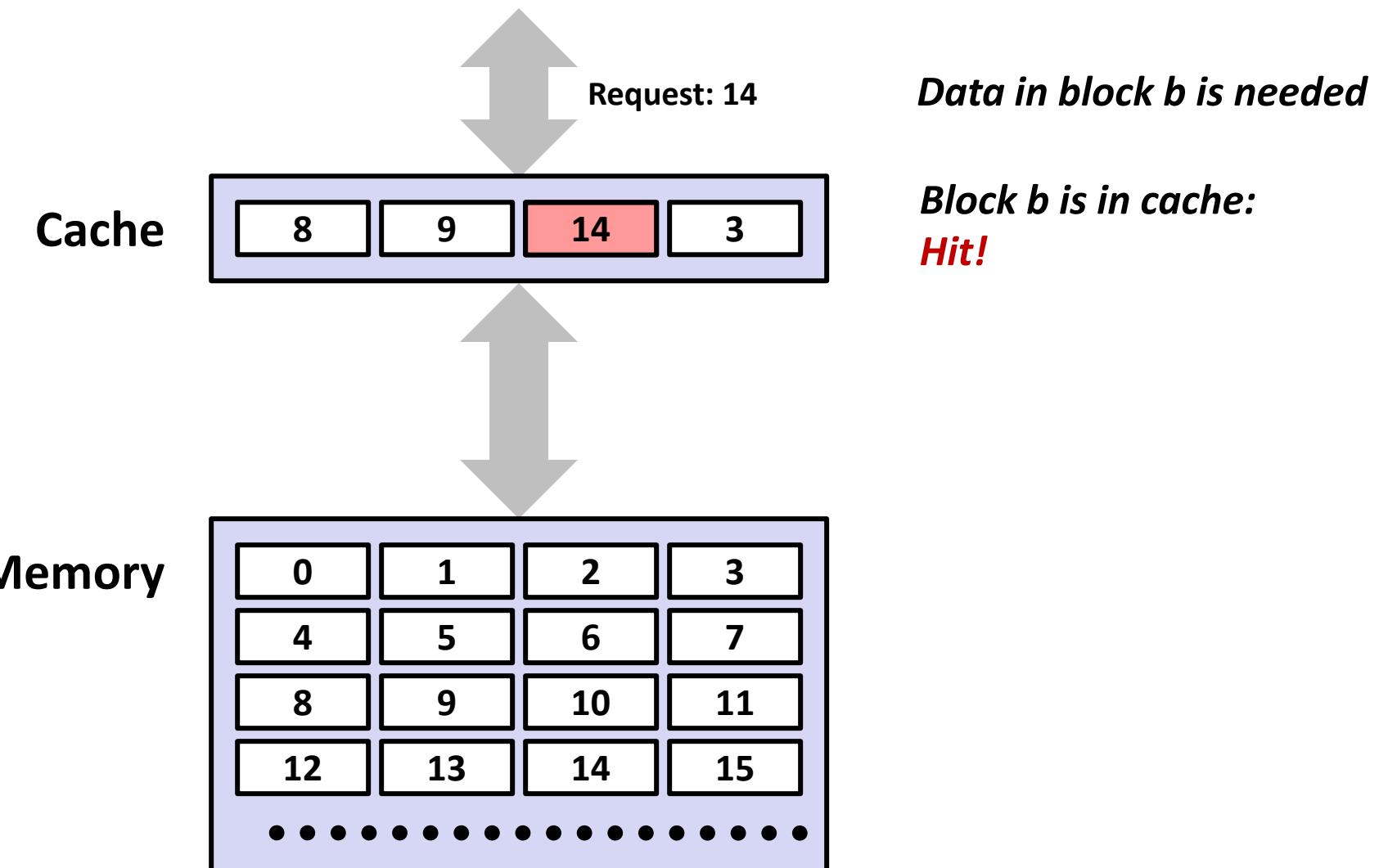
Caches

- ***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.

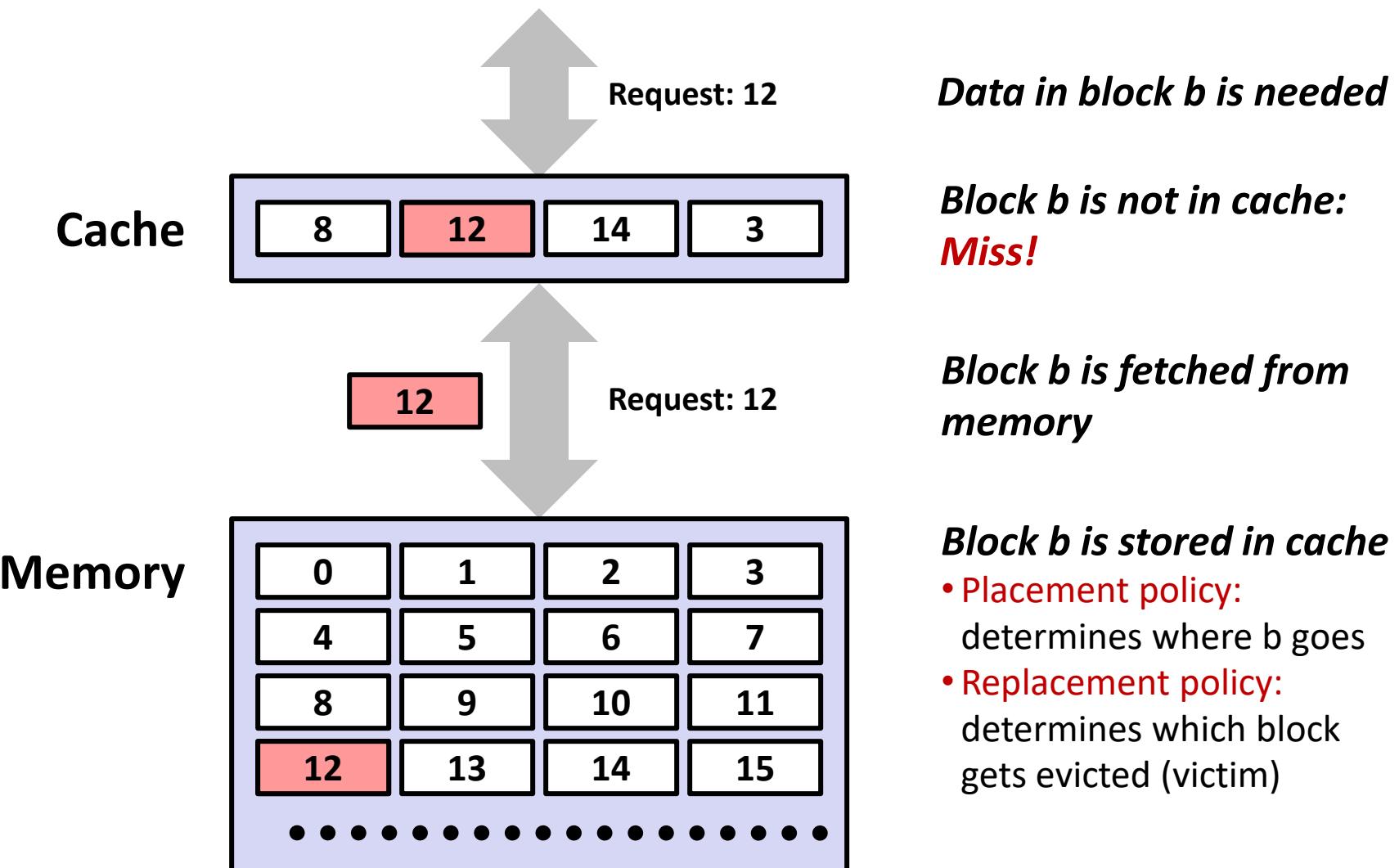
General Cache Concepts



General Cache Concepts: Hit



General Cache Concepts: Miss



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 $(i \bmod 4)$ 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.

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
TLB	Address translations	On-Chip TLB	0	Hardware MMU
L1 cache	64-byte blocks	On-Chip L1	4	Hardware
L2 cache	64-byte blocks	On-Chip L2	10	Hardware
Virtual Memory	4-KB pages	Main memory	100	Hardware + OS
Buffer cache	Parts of files	Main memory	100	OS
Disk cache	Disk sectors	Disk controller	100,000	Disk firmware
Network buffer cache	Parts of files	Local disk	10,000,000	NFS client
Browser cache	Web pages	Local disk	10,000,000	Web browser
Web cache	Web pages	Remote server disks	1,000,000,000	Web proxy server

Summary

- The speed gap between CPU, memory and mass storage continues to widen.
- Well-written programs exhibit a property called *locality*.
- Memory hierarchies based on *caching* close the gap by exploiting locality.

Cache Memories

Instructors:

R. Shathanaa

Today

- Cache memory organization and operation
- Performance impact of caches
 - Rearranging loops to improve spatial locality
 - Using blocking to improve temporal locality

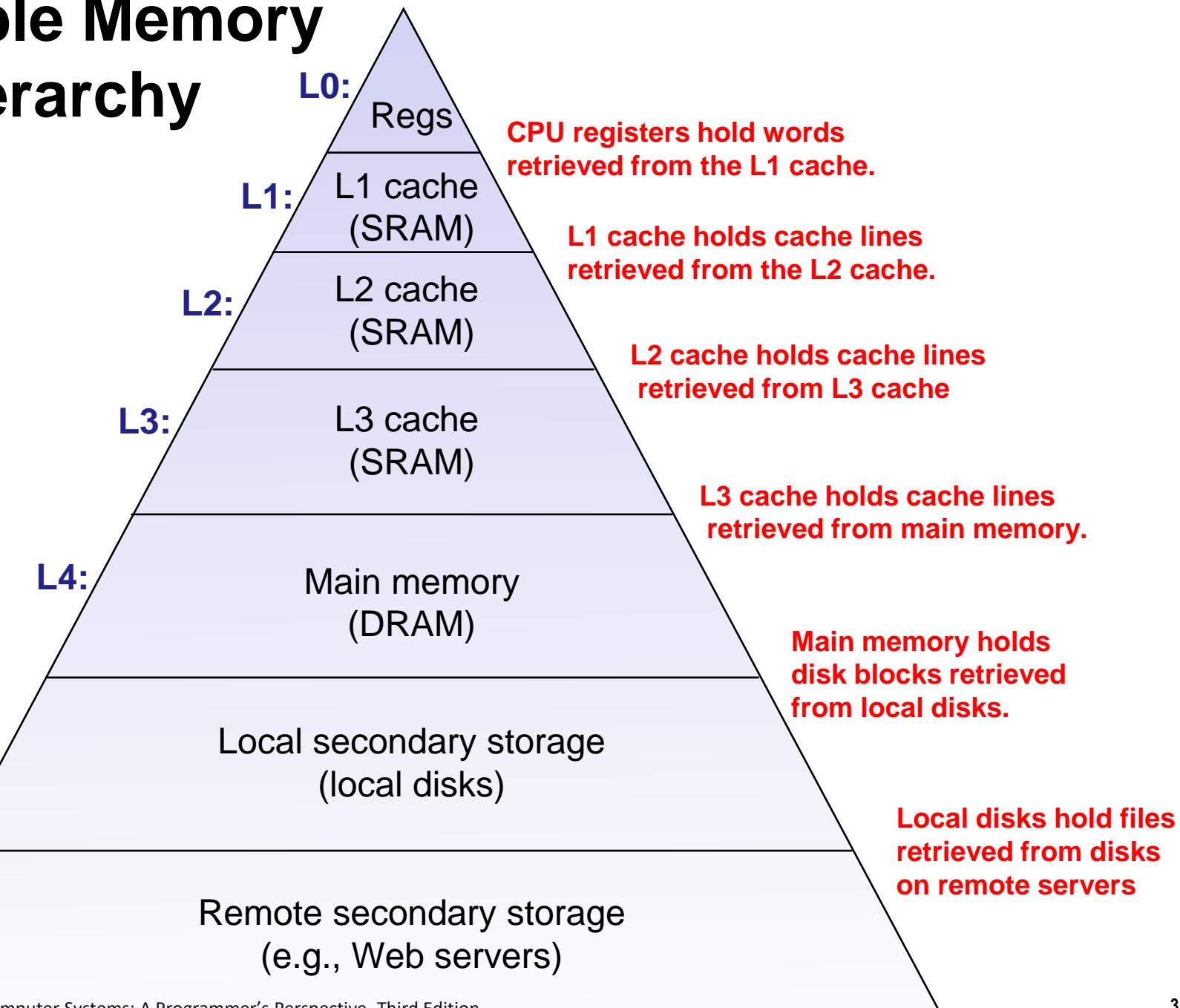
Example Memory Hierarchy

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

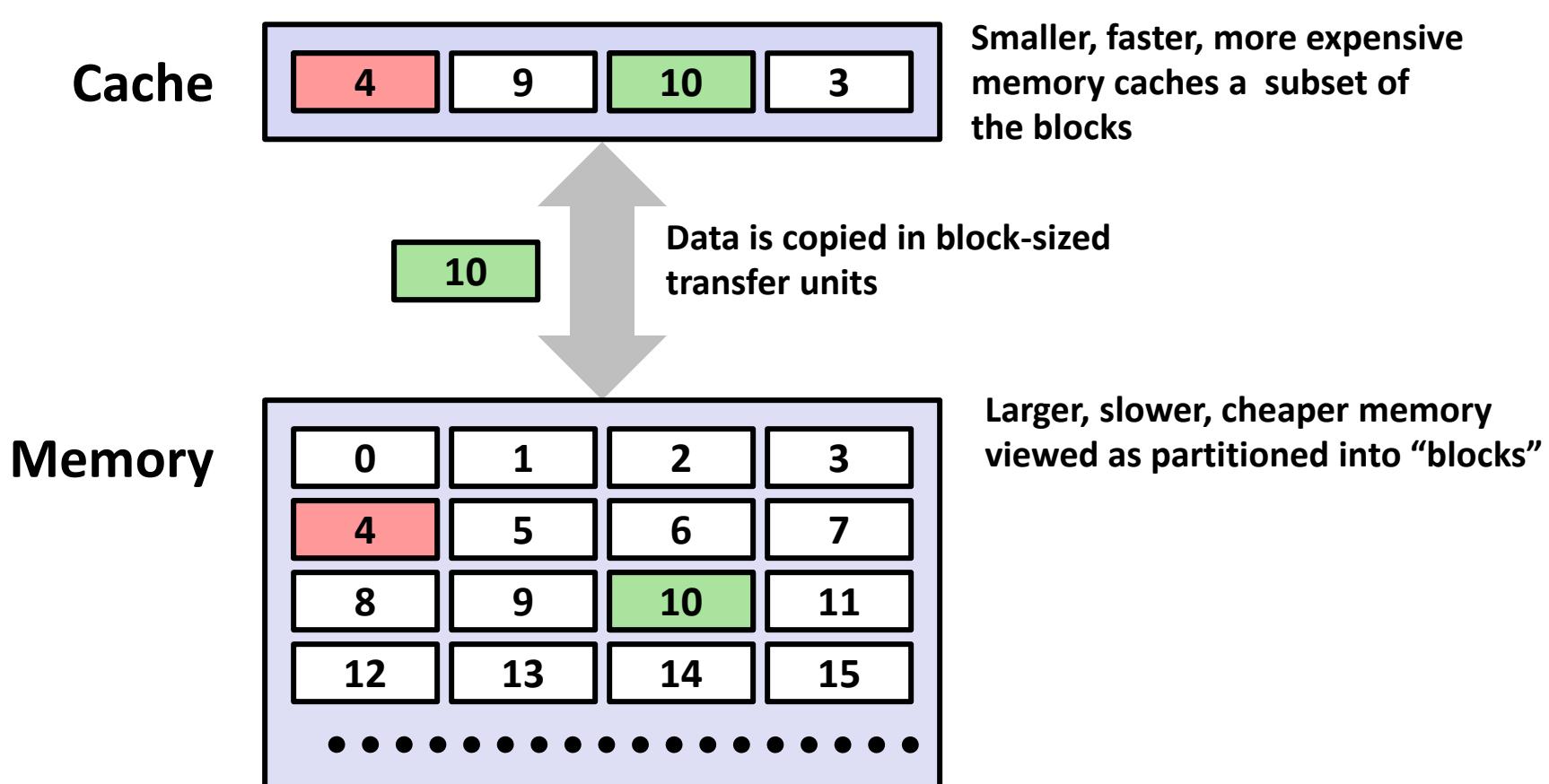
Larger,
slower,
and
cheaper
(per byte)
storage
devices

L6:

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

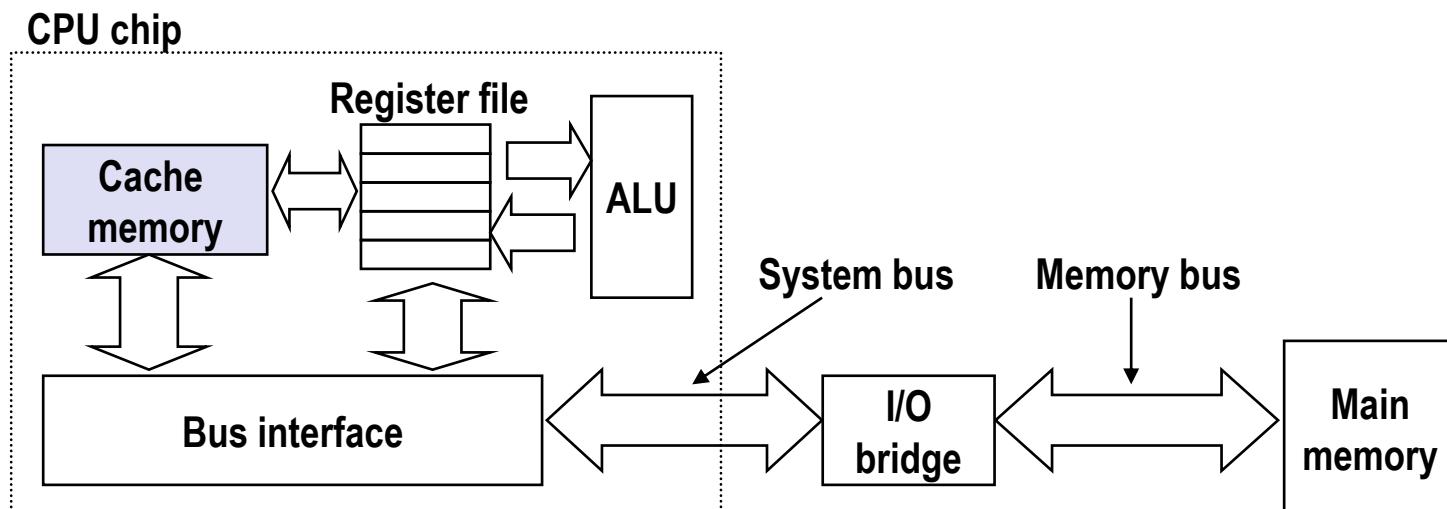


General Cache Concept



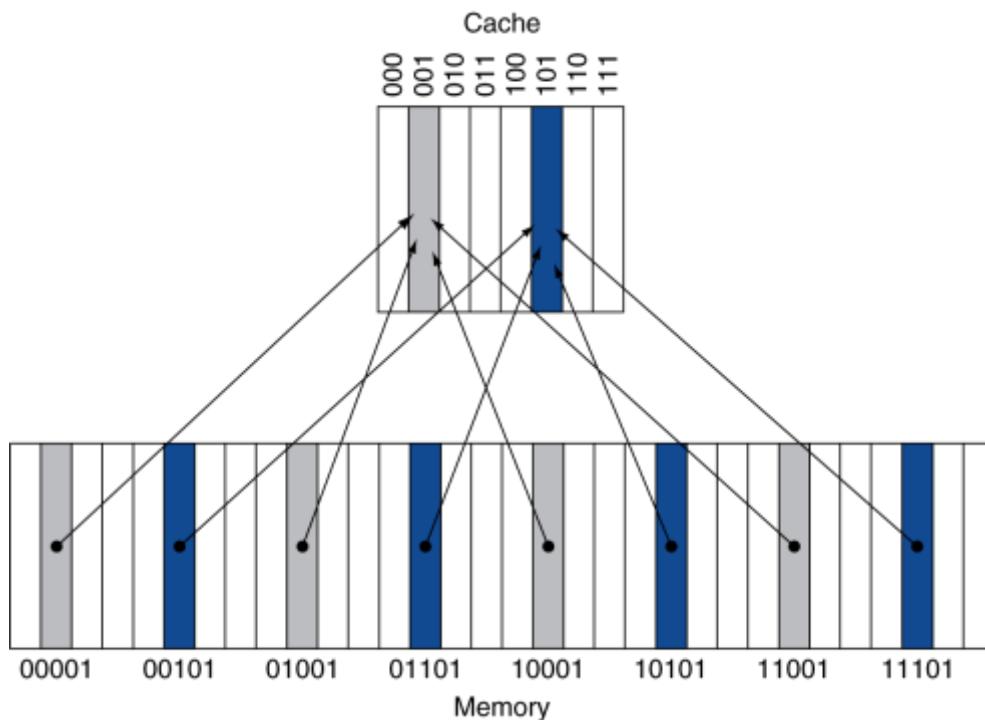
Cache Memories

- 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:



Direct Mapped Cache

- Location determined by address
- Direct mapped: only one choice
 - (Block address) modulo (#Blocks in cache)



- **#Blocks is a power of 2**
- **Use low-order address bits**

Tags and Valid Bits

- How do we know which particular block is stored in a cache location?
 - Store block address as well as the data
 - Actually, only need the high-order bits
 - Called the tag
- What if there is no data in a location?
 - Valid bit: 1 = present, 0 = not present
 - Initially 0

Cache Example

- 8-blocks, 1 word/block, direct mapped
- Initial state

Index	V	Tag	Data
000	N		
001	N		
010	N		
011	N		
100	N		
101	N		
110	N		
111	N		

Cache Example

Word addr	Binary addr	Hit/miss	Cache block
22	10 110	Miss	110

Index	V	Tag	Data
000	N		
001	N		
010	N		
011	N		
100	N		
101	N		
110	Y	10	Mem[10110]
111	N		

Cache Example

Word addr	Binary addr	Hit/miss	Cache block
26	11 010	Miss	010

Index	V	Tag	Data
000	N		
001	N		
010	Y	11	Mem[11010]
011	N		
100	N		
101	N		
110	Y	10	Mem[10110]
111	N		

Cache Example

Word addr	Binary addr	Hit/miss	Cache block
22	10 110	Hit	110
26	11 010	Hit	010

Index	V	Tag	Data
000	N		
001	N		
010	Y	11	Mem[11010]
011	N		
100	N		
101	N		
110	Y	10	Mem[10110]
111	N		

Cache Example

Word addr	Binary addr	Hit/miss	Cache block
16	10 000	Miss	000
3	00 011	Miss	011
16	10 000	Hit	000

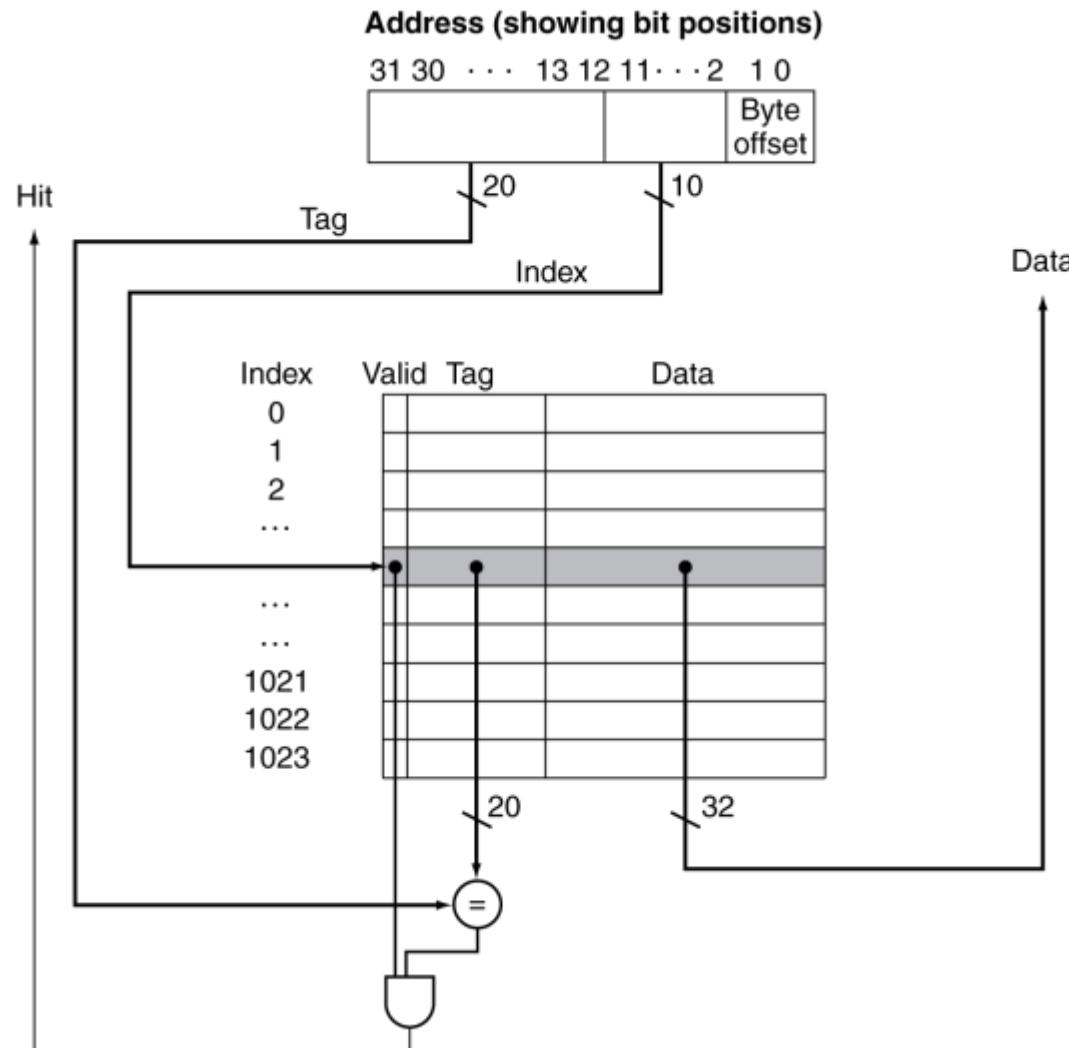
Index	V	Tag	Data
000	Y	10	Mem[10000]
001	N		
010	Y	11	Mem[11010]
011	Y	00	Mem[00011]
100	N		
101	N		
110	Y	10	Mem[10110]
111	N		

Cache Example

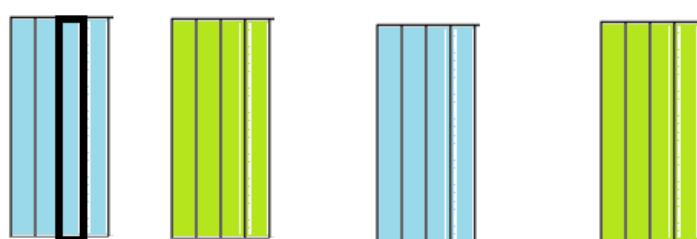
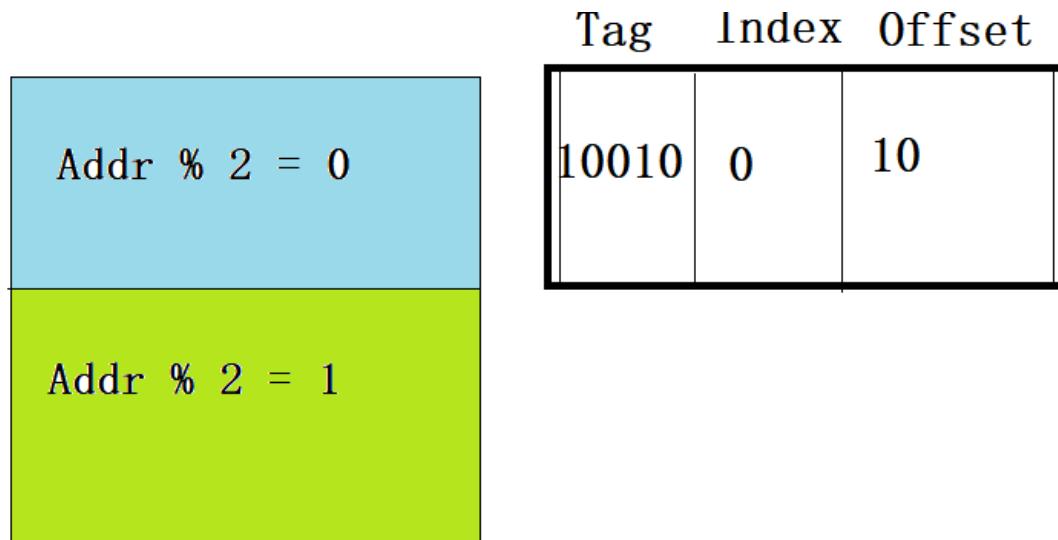
Word addr	Binary addr	Hit/miss	Cache block
18	10 010	Miss	010

Index	V	Tag	Data
000	Y	10	Mem[10000]
001	N		
010	Y	10	Mem[10010]
011	Y	00	Mem[00011]
100	N		
101	N		
110	Y	10	Mem[10110]
111	N		

Address Subdivision



Direct-Mapped Identification



1001 0000	1001 0100	1001 1000	1001 1100
1001 0001	1001 0101	1001 1001	1001 1101
1001 0010	1001 0110	1001 1010	1001 1110
1001 0011	1001 0111	1001 1011	1001 1111
100100	100101	100110	100111

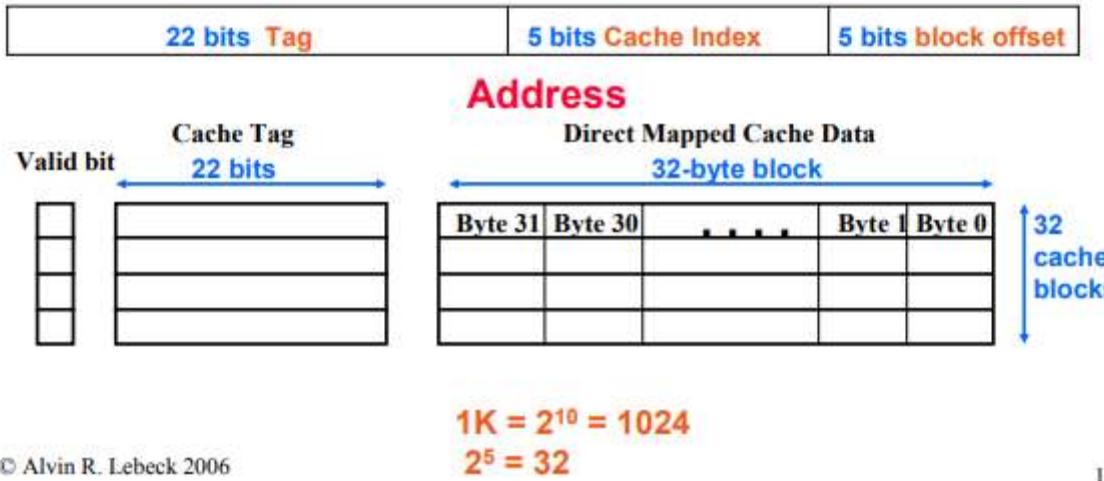
Cache Addressing

Example: 1-KB Cache with 32B blocks:

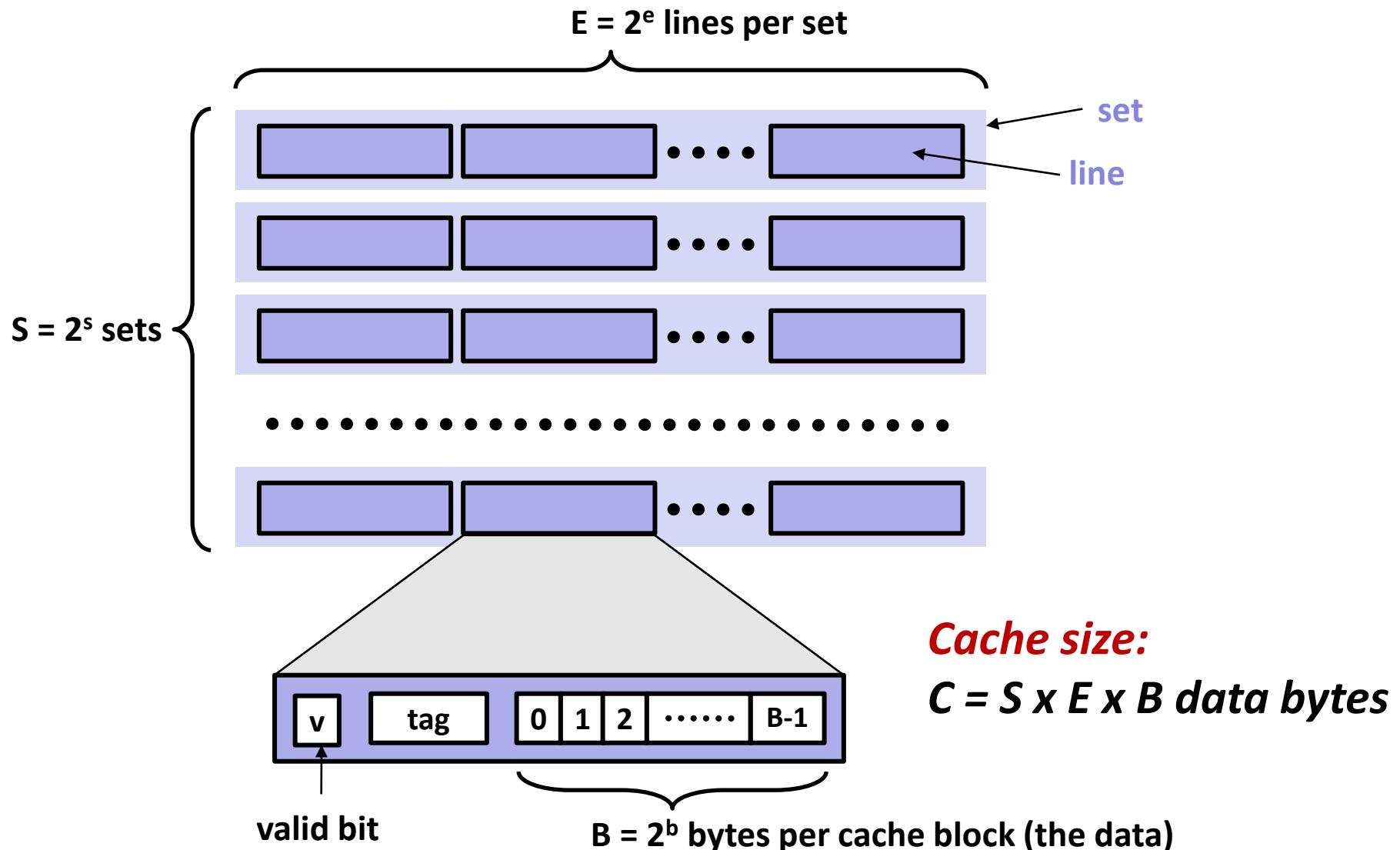
Cache Index = (**<Address>** Mod (1024))/ 32

Block-Offset = **<Address>** Mod (32)

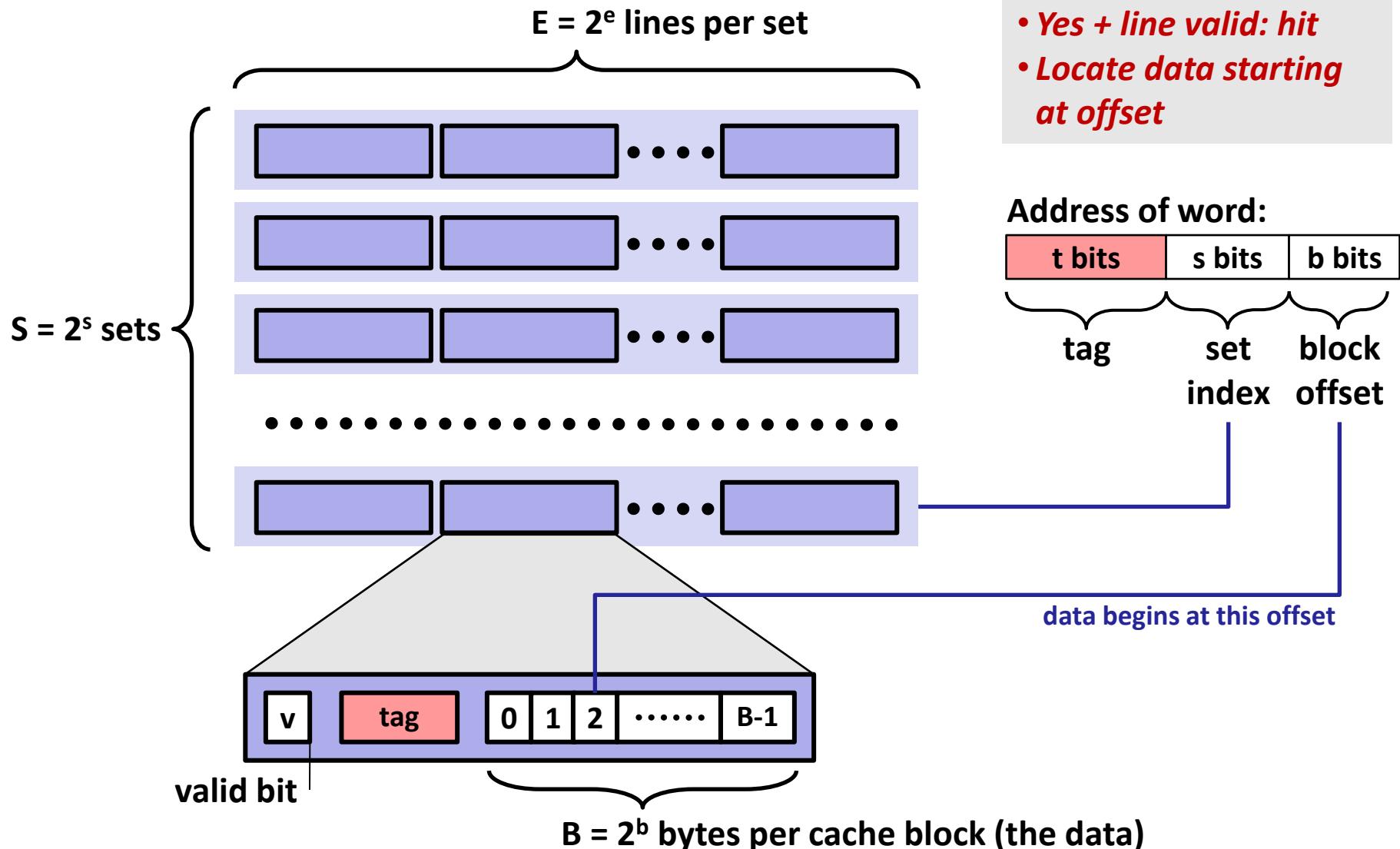
Tag = **<Address>** / (1024)



General Cache Organization (S, E, B)



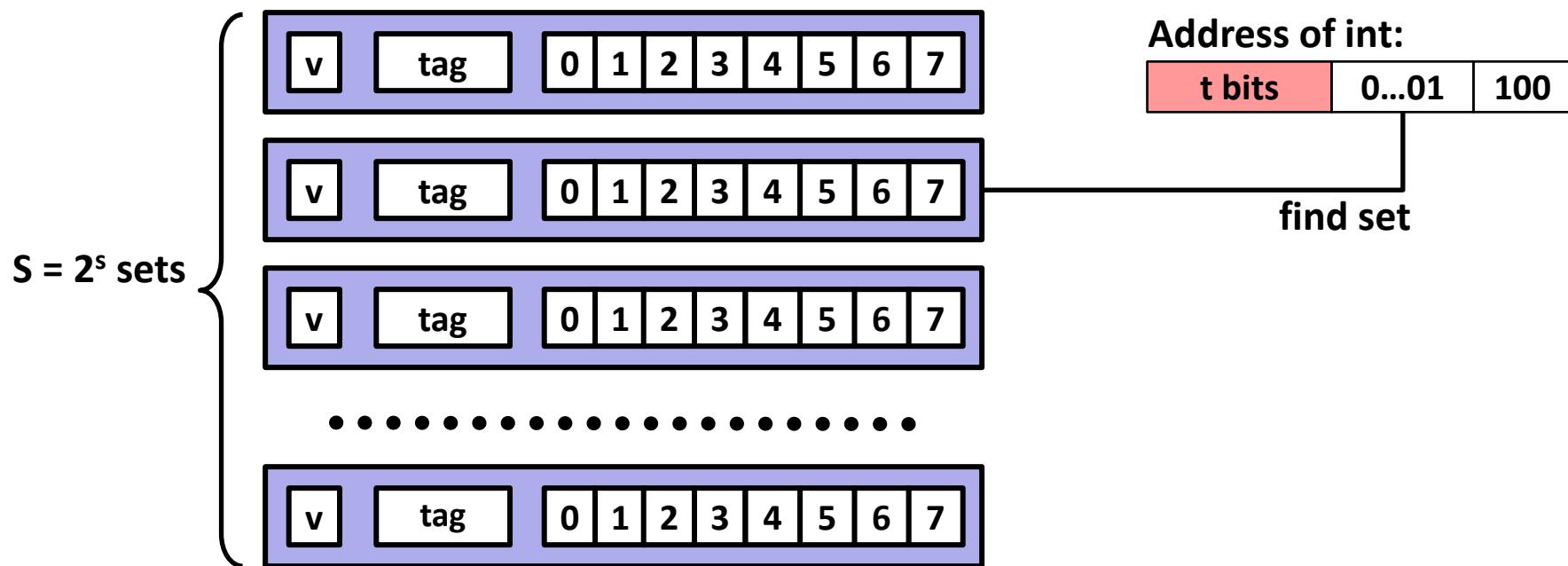
Cache Read



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

Direct mapped: One line per set

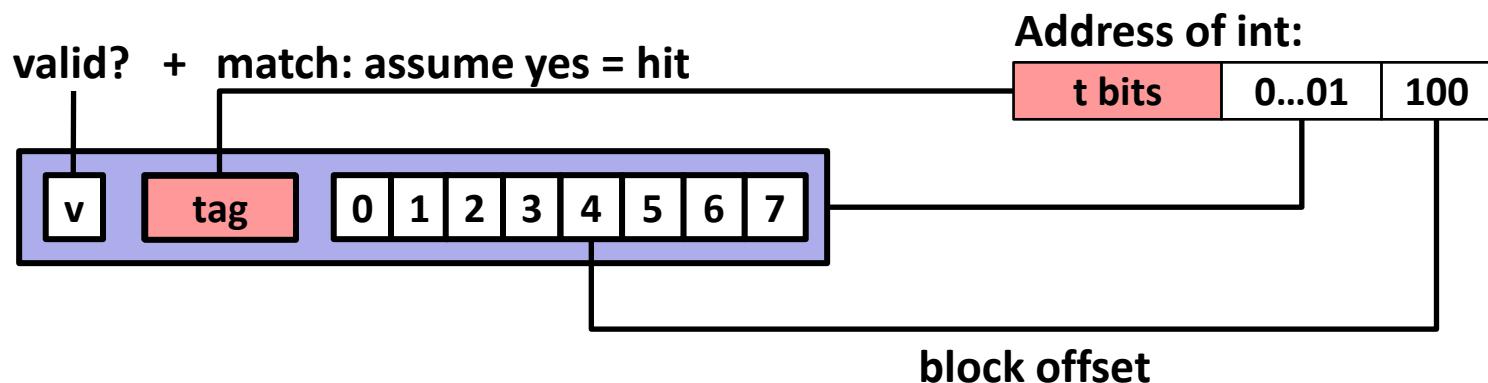
Assume: cache block size 8 bytes



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

Direct mapped: One line per set

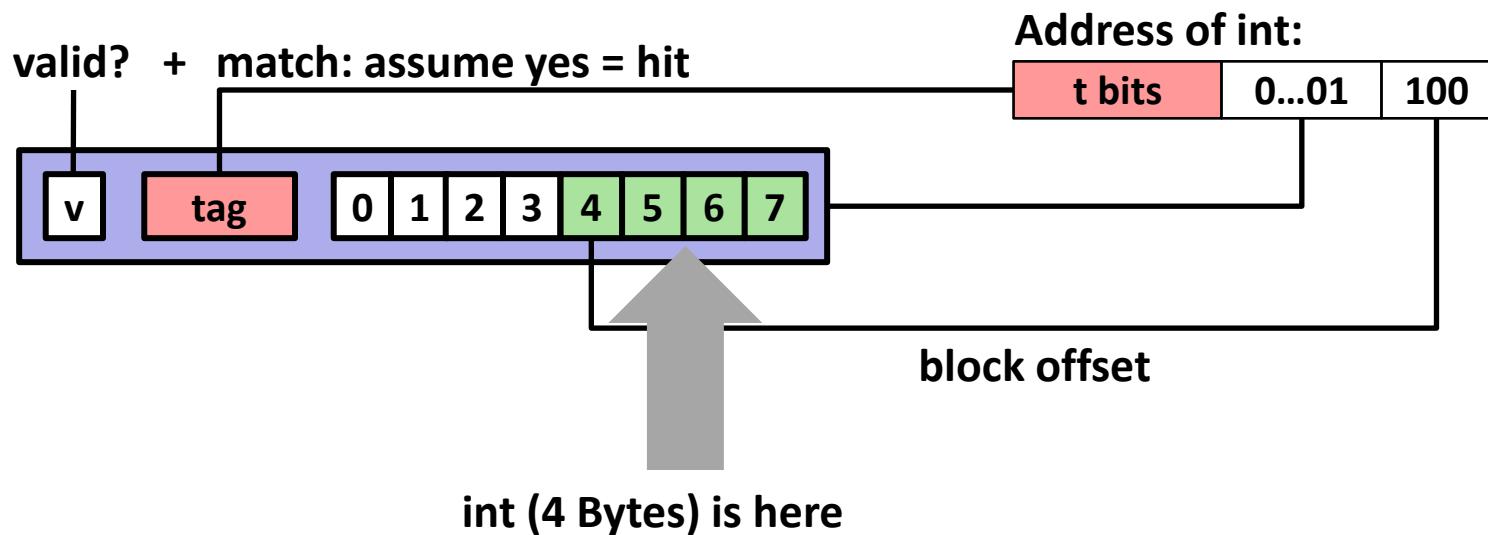
Assume: cache block size 8 bytes



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

Direct mapped: One line per set

Assume: cache block size 8 bytes



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

Direct-Mapped Cache Simulation

$t=1 \quad s=2 \quad b=1$

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

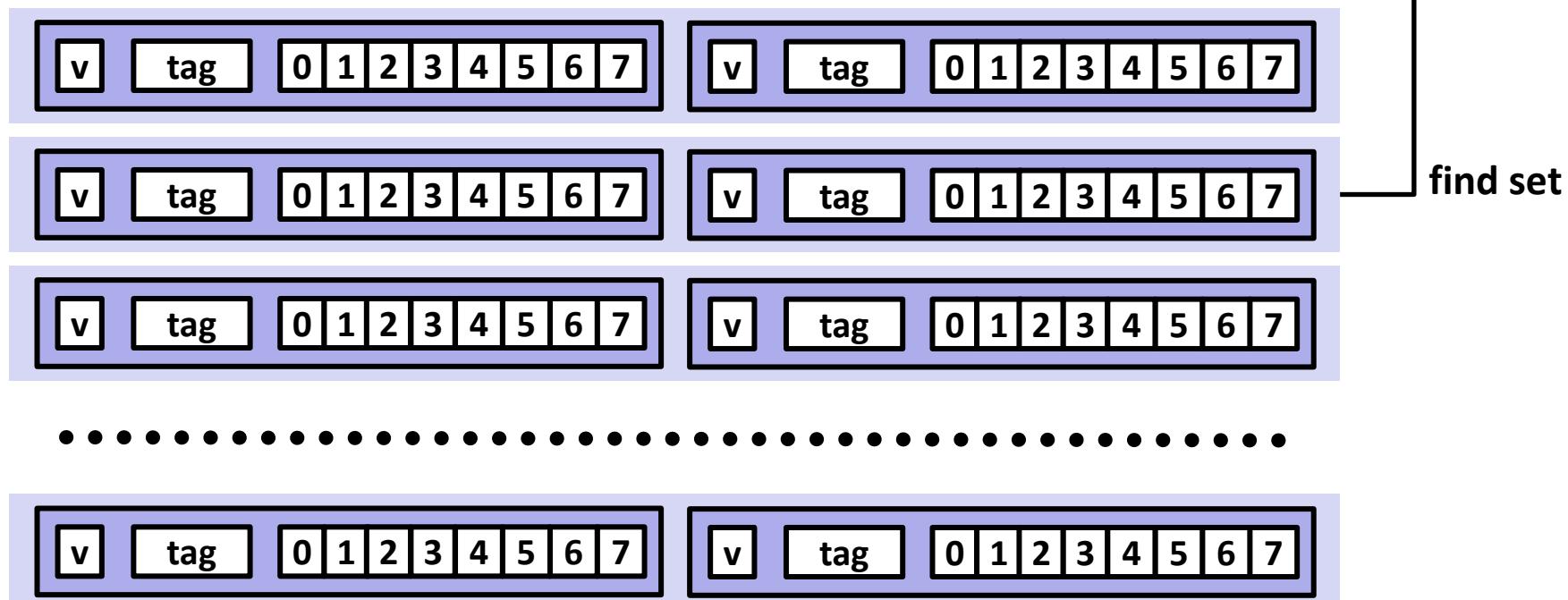
	v	Tag	Block
Set 0	1	0	$M[0-1]$
Set 1			
Set 2			
Set 3	1	0	$M[6-7]$

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

$E = 2$: Two lines per set

Assume: cache block size 8 bytes

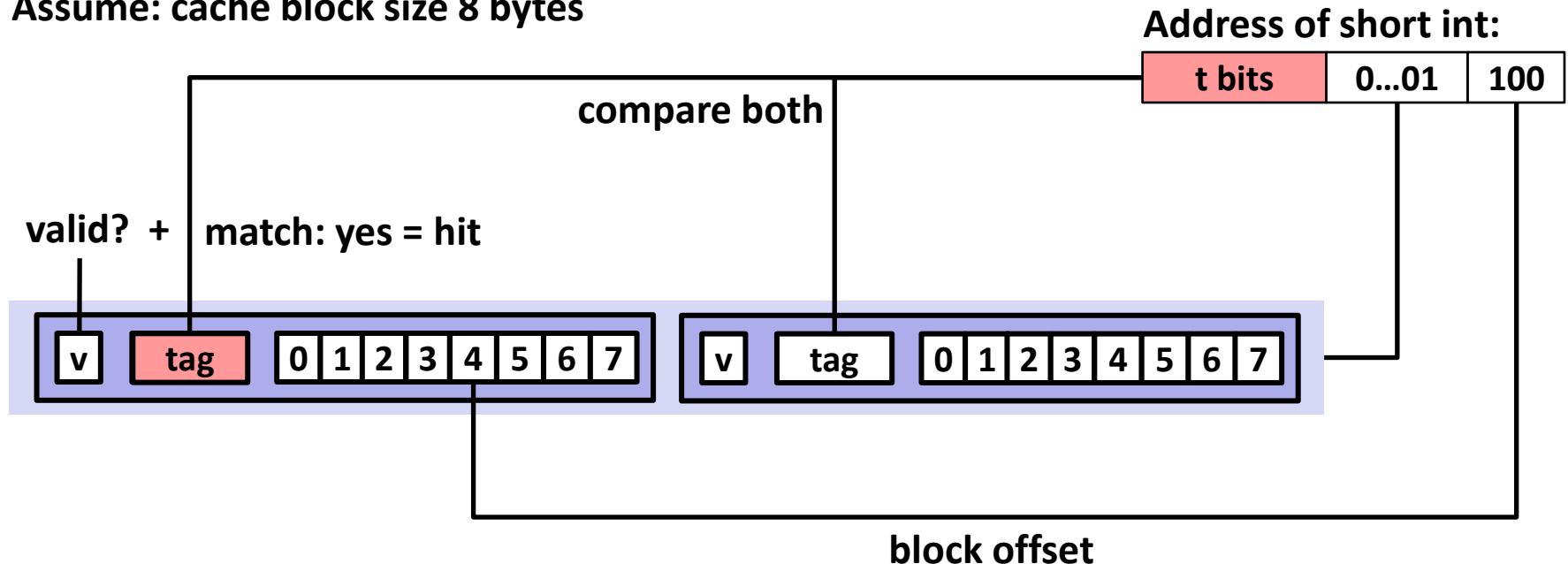
Address of short int:



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

E = 2: Two lines per set

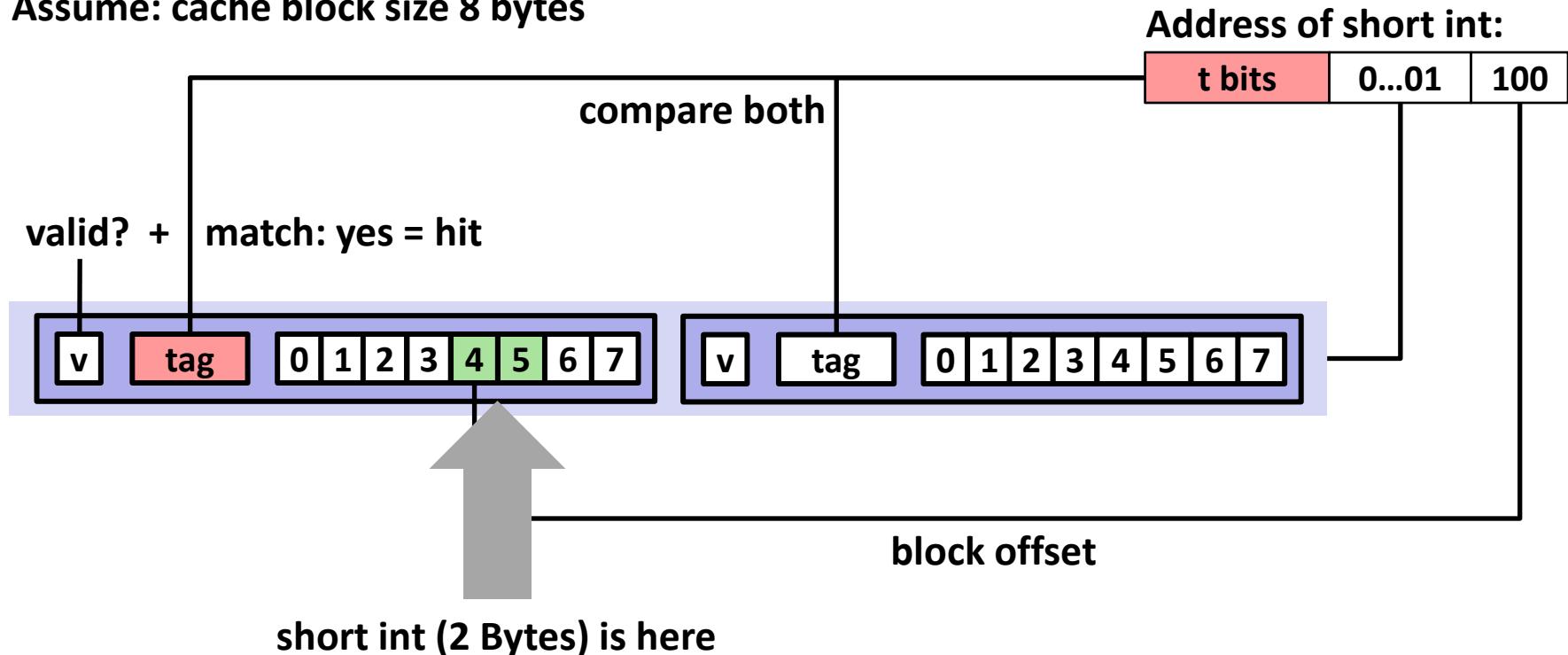
Assume: cache block size 8 bytes



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), ...

2-Way Set Associative Cache Simulation

$t=2 \quad s=1 \quad b=1$

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

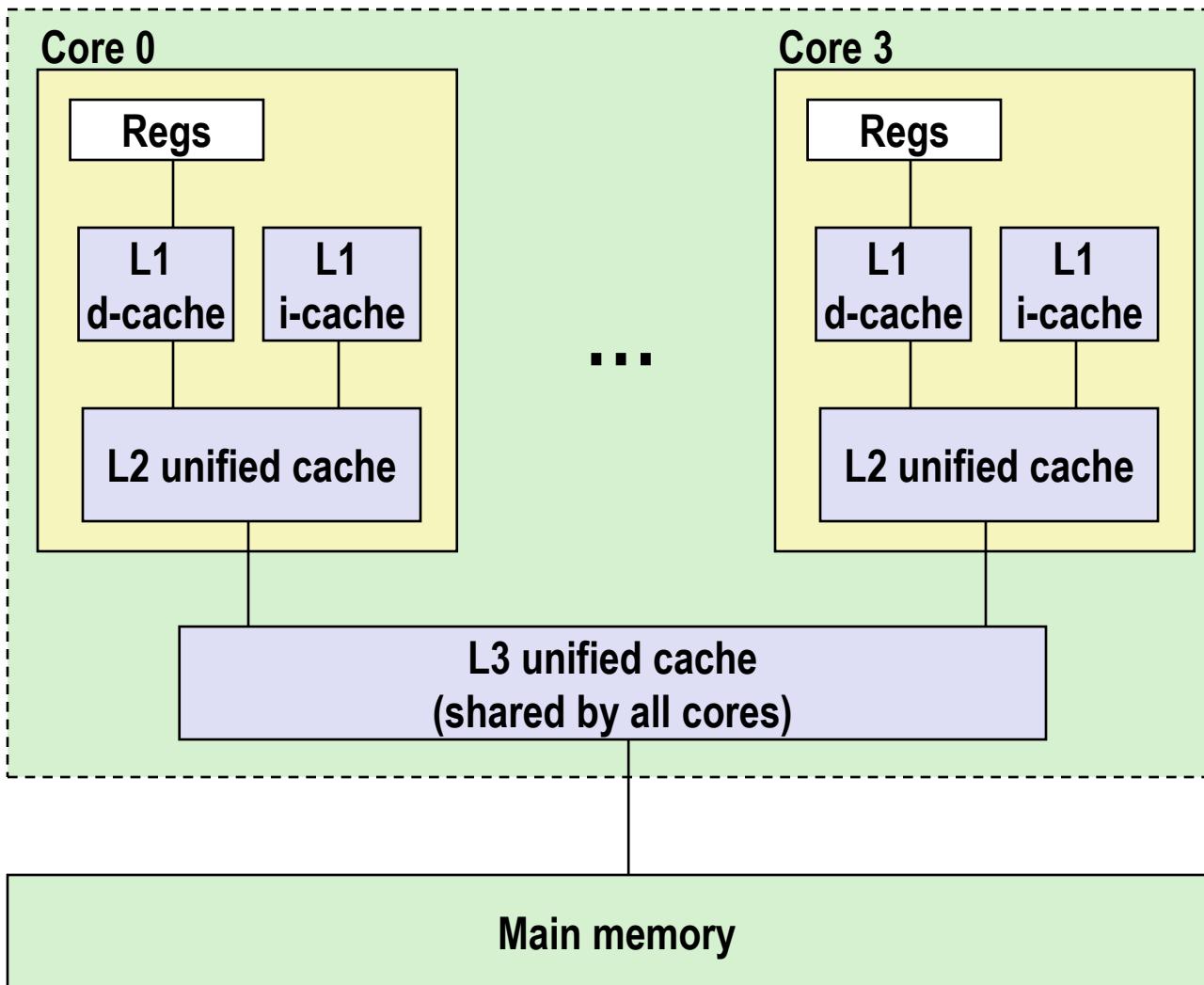
	v	Tag	Block
Set 0	1	00	M[0-1]
	1	10	M[8-9]
Set 1	1	01	M[6-7]
	0		

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)
= $1 - \text{hit rate}$
- 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

Impact of Cache Size

- On the one hand, a larger cache will tend to increase the hit rate.
- On the other hand, it is always harder to make large memories run faster. As a result, larger caches tend to increase the hit time.
- This explains why an L1 cache is smaller than an L2 cache, and an L2 cache is smaller than an L3 cache.

Impact of Block Size

- Increase the hit rate by exploiting any spatial locality that might exist in a program
- However, for a given cache size, larger blocks imply a smaller number of cache lines, which can hurt the hit rate in programs with more temporal locality than spatial locality.

Impact of Associativity

- The advantage of higher associativity (i.e., larger values of E) is that it decreases the vulnerability of the cache to **thrashing** due to conflict misses.
- However, higher associativity comes at a significant cost. Higher associativity is expensive to implement and hard to make fast. It requires more tag bits per line, additional LRU state bits per line, and additional control logic.
- Higher associativity can *increase hit time*, because of the increased complexity, and it can also *increase the miss penalty* because of the increased complexity of choosing a victim line.

Impact of Write Strategy

- Write-through caches are simpler to implement and can use a *write buffer* that works independently of the cache to update memory.
- Furthermore, read misses are less expensive because they do not trigger a memory write.
- On the other hand, write-back caches result in fewer transfers.
- Further, reducing the number of transfers becomes increasingly important as we move down the hierarchy and the transfer times increase.
- In general, caches further down the hierarchy are more likely to use write-back than write-through.

Today

- Cache organization and operation
- Performance impact of caches
 - Rearranging loops to improve spatial locality
 - Using blocking to improve temporal locality

Matrix Multiplication Example

■ Description:

- Multiply $N \times N$ 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

```
/* 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

Variable sum held in register

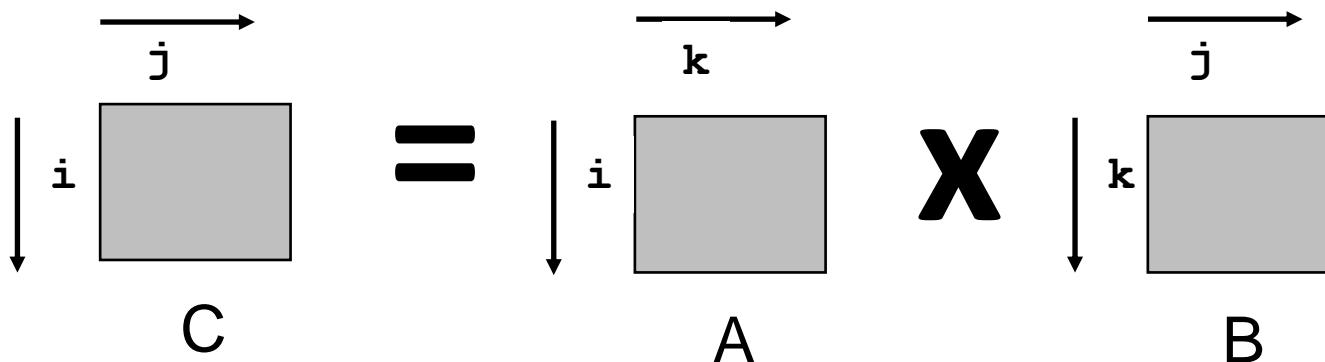
Miss Rate Analysis for Matrix Multiply

■ 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 ( $B$ ) >  $\text{sizeof}(a_{ij})$  bytes, exploit spatial locality
  - miss rate =  $\text{sizeof}(a_{ij}) / B$

## ■ Stepping through rows in one column:

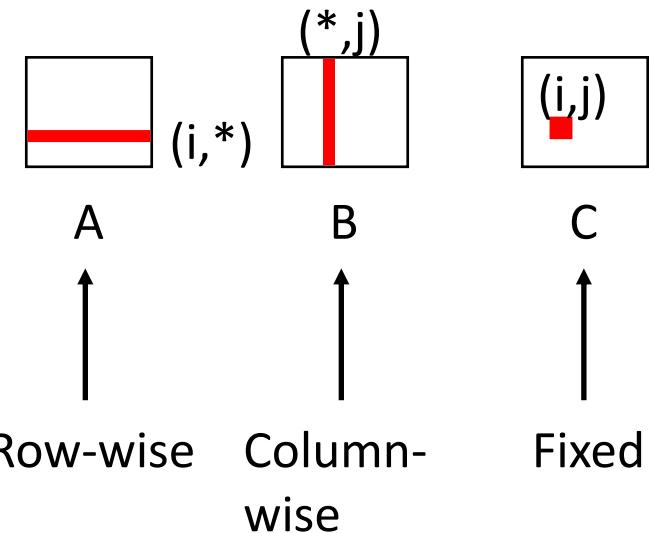
- ```
for (i = 0; i < n; i++)
    sum += a[i][0];
```
- accesses distant elements
- no spatial locality!
 - miss rate = 1 (i.e. 100%)

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:



Misses per inner loop iteration:

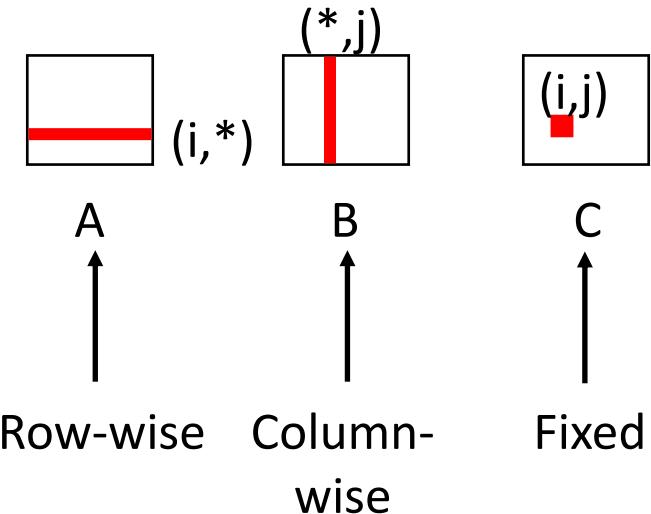
<u>A</u>	<u>B</u>	<u>C</u>
0.25	1.0	0.0

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:



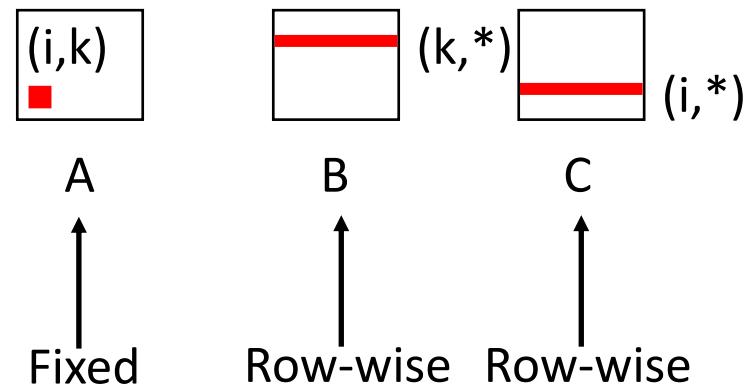
Misses per inner loop iteration:

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:



Misses per inner loop iteration:

A
0.0

B
0.25

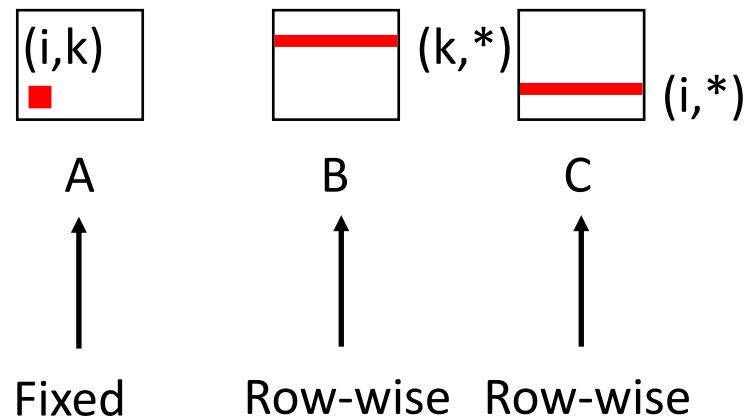
C
0.25

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:



Misses per inner loop iteration:

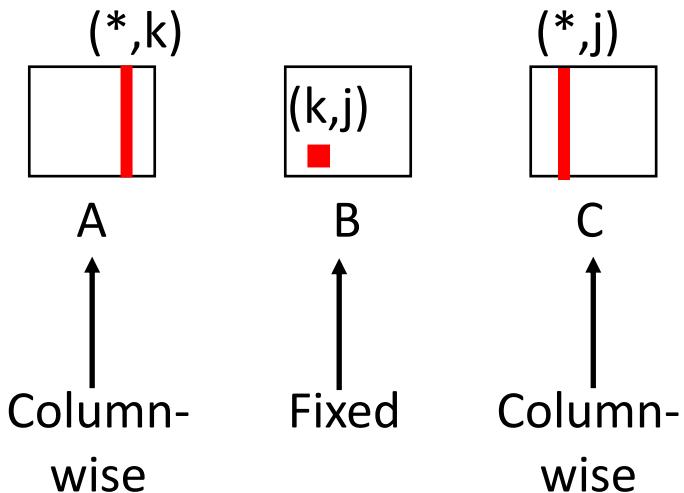
A	B	C
0.0	0.25	0.25

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:



Misses per inner loop iteration:

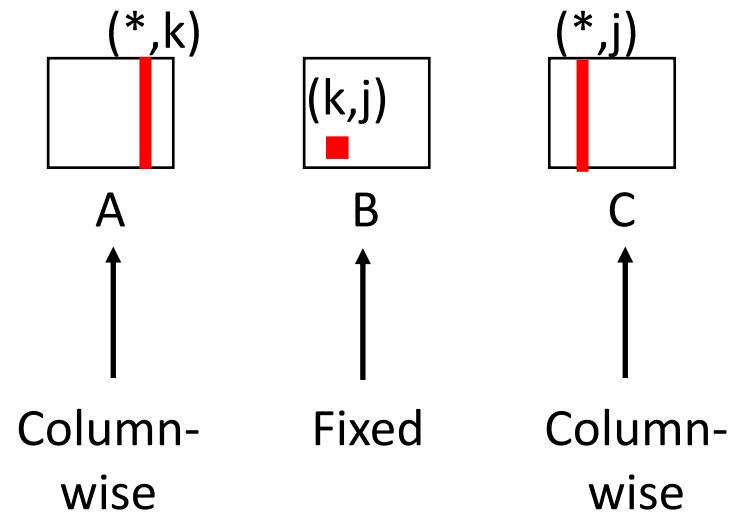
<u>A</u>	<u>B</u>	<u>C</u>
1.0	0.0	1.0

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:



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;
    }
}

```

```

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];
    }
}

```

```

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;
    }
}

```

ijk (& jik):

- 2 loads, 0 stores
- misses/iter = **1.25**

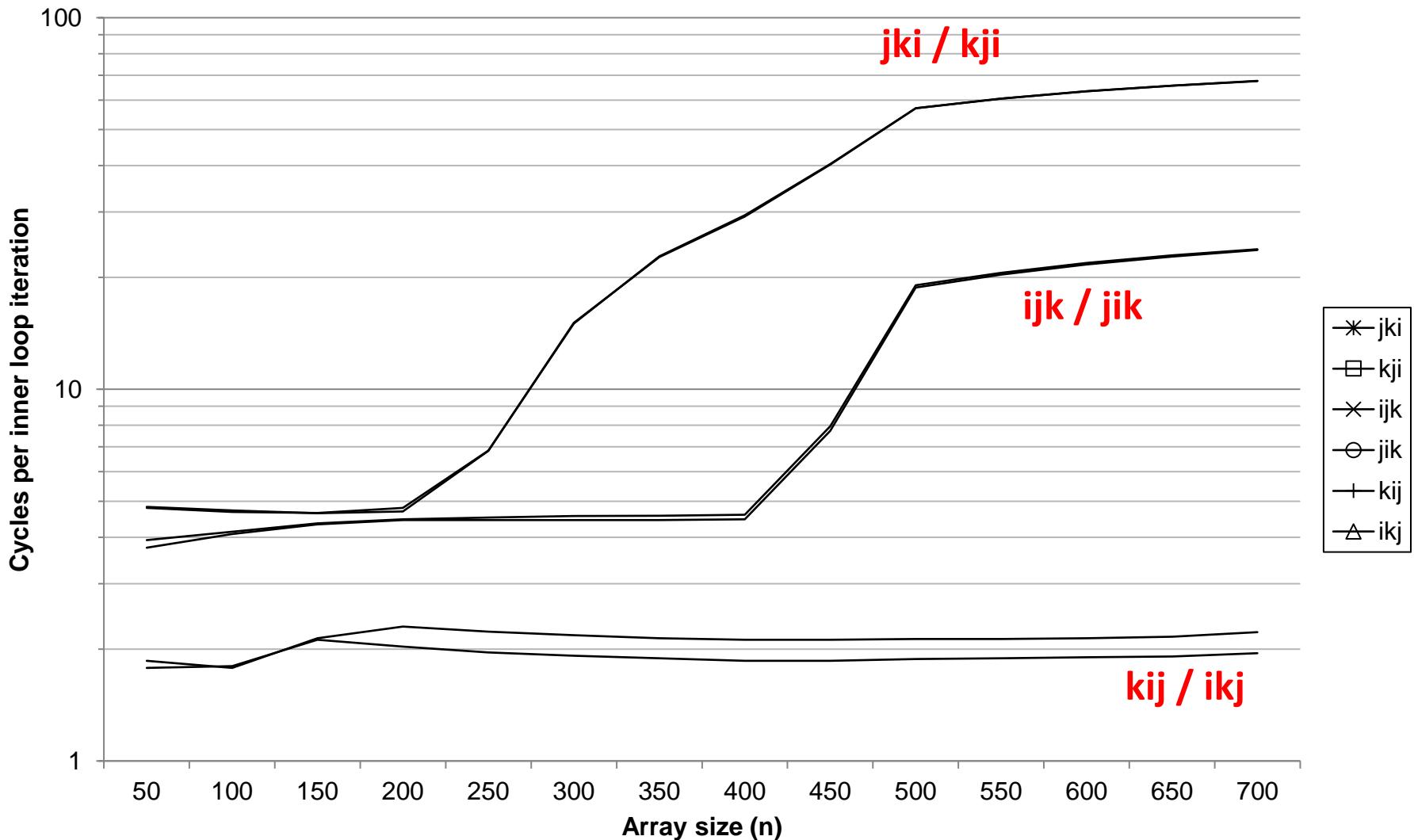
kij (& ikj):

- 2 loads, 1 store
- misses/iter = **0.5**

jki (& kji):

- 2 loads, 1 store
- misses/iter = **2.0**

Core i7 Matrix Multiply Performance

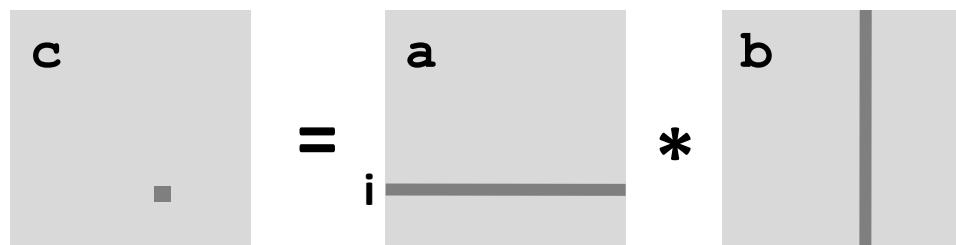


Today

- Cache organization and operation
- Performance impact of caches
 - Rearranging loops to improve spatial locality
 - Using blocking to improve temporal locality

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];  
}
```



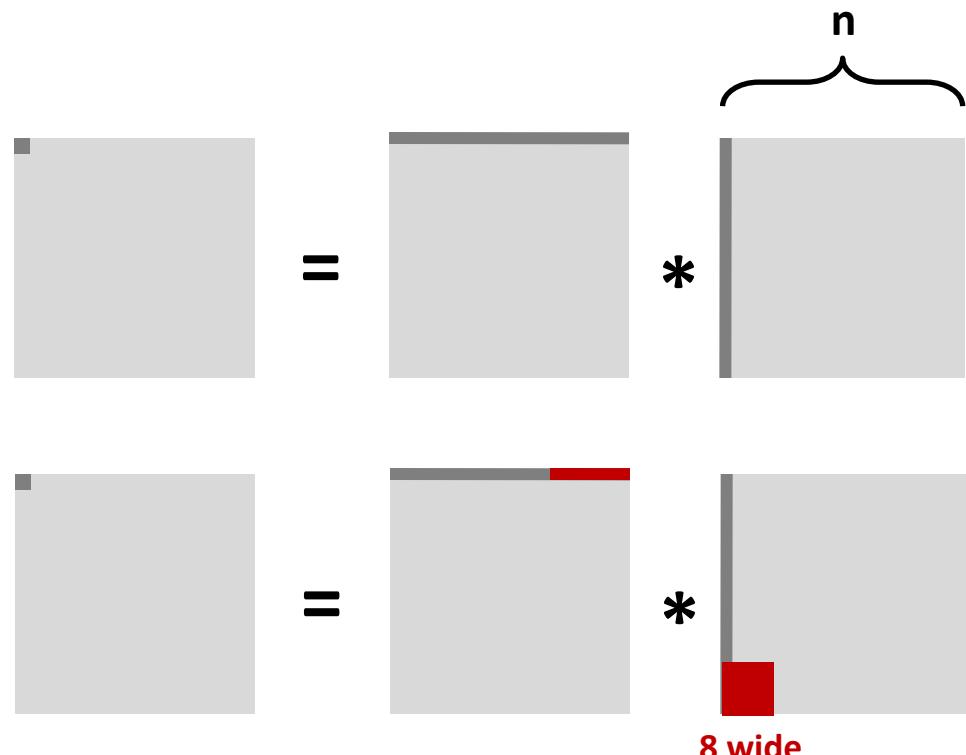
Cache Miss Analysis

■ Assume:

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

■ First iteration:

- $n/8 + n = 9n/8$ misses



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

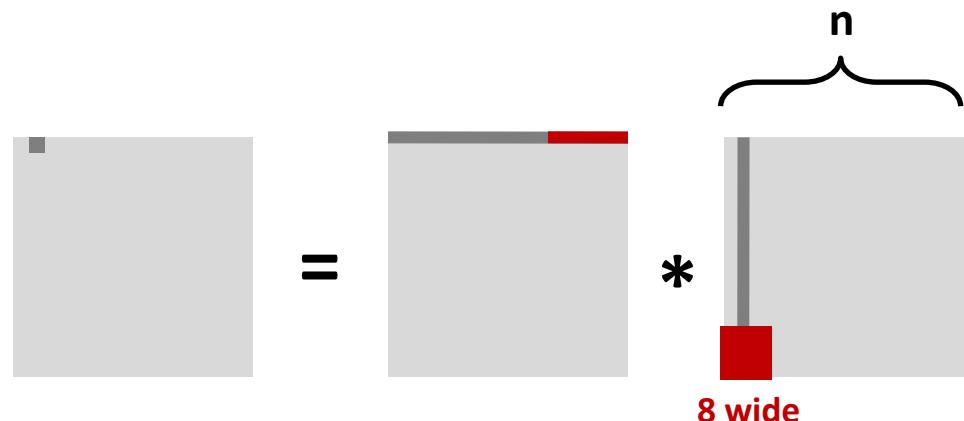
Cache Miss Analysis

■ Assume:

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

■ Second iteration:

- Again:
 $n/8 + n = 9n/8$ misses



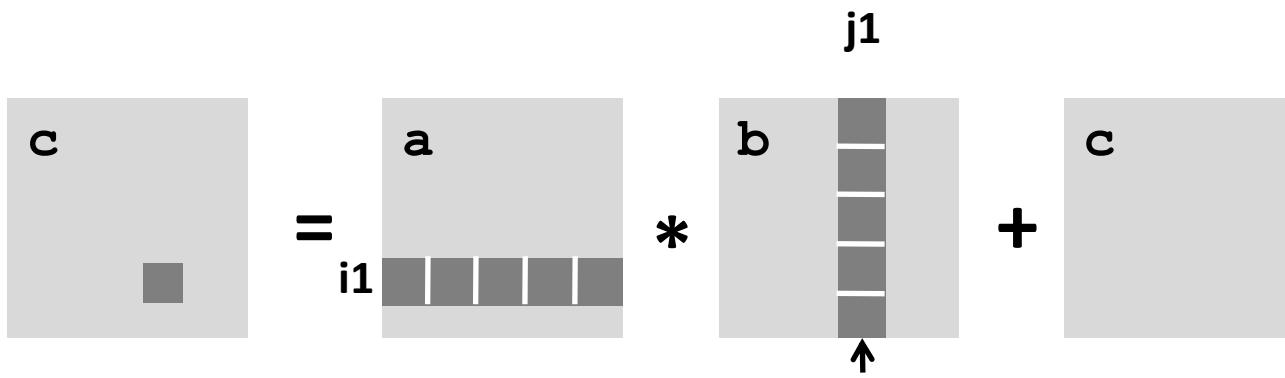
■ Total misses:

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

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; i++)
                    for (j1 = j; j1 < j+B; j++)
                        for (k1 = k; k1 < k+B; k++)
                            c[i1*n+j1] += a[i1*n + k1]*b[k1*n + j1];
}
                                            matmult/bmm.c
```



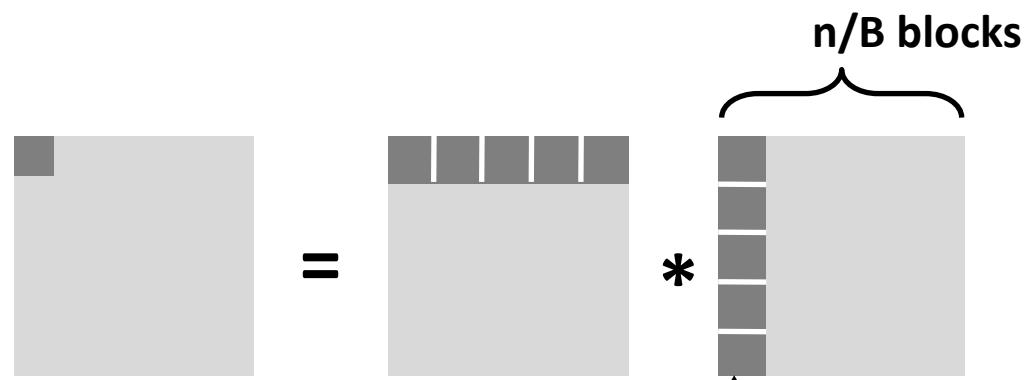
Cache Miss Analysis

■ Assume:

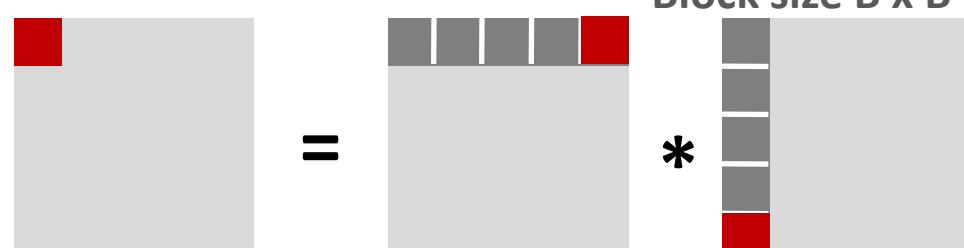
- Cache block = 8 doubles
- Cache size $C \ll n$ (much smaller than n)
- Three blocks ■ fit into cache: $3B^2 < C$

■ First (block) iteration:

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



- Afterwards in cache
(schematic)



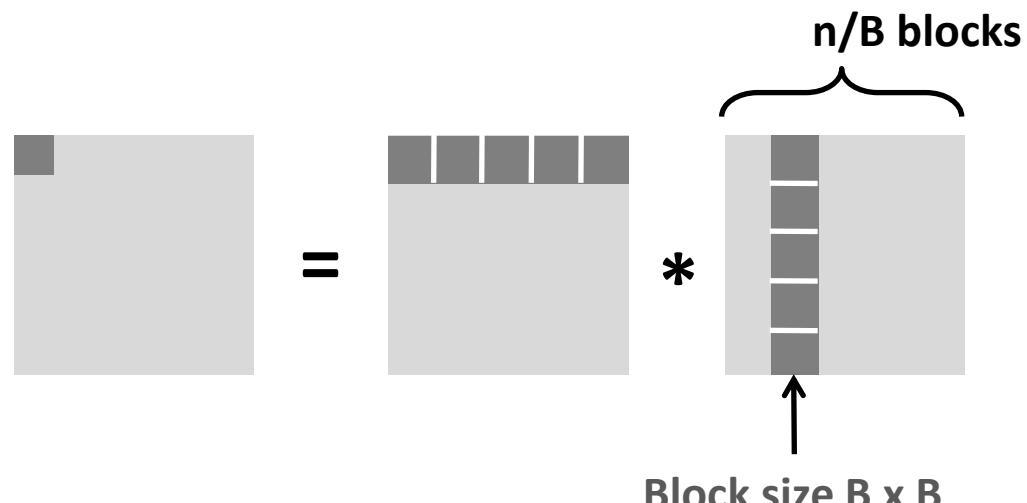
Cache Miss Analysis

■ Assume:

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

■ Second (block) iteration:

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



■ Total misses:

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

Blocking Summary

- No blocking: $(9/8) * n^3$
- Blocking: $1/(4B) * n^3$
- Suggest 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 elements used $O(n)$ times!
 - But program has to be written properly

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 sequentially with stride 1.
 - Try to maximize temporal locality by using a data object as often as possible once it's read from memory.



The Processor

Adapted and Supplemented by,
Dr. R. Shathanaa

MIPS ISA

- Microprocessor without Interlocked Pipelined Stages

Arithmetic and Logical Instructions

Instruction	Operation
add \$d, \$s, \$t	\$d = \$s + \$t
addu \$d, \$s, \$t	\$d = \$s + \$t
addi \$t, \$s, i	\$t = \$s + SE(i)
addiu \$t, \$s, i	\$t = \$s + SE(i)
and \$d, \$s, \$t	\$d = \$s & \$t
andi \$t, \$s, i	\$t = \$s & ZE(i)
div \$s, \$t	lo = \$s / \$t; hi = \$s % \$t
divu \$s, \$t	lo = \$s / \$t; hi = \$s % \$t
mult \$s, \$t	hi:lo = \$s * \$t
multu \$s, \$t	hi:lo = \$s * \$t
nor \$d, \$s, \$t	\$d = ~(\$s \$t)
or \$d, \$s, \$t	\$d = \$s \$t
ori \$t, \$s, i	\$t = \$s ZE(i)
sll \$d, \$t, a	\$d = \$t << a
sllv \$d, \$t, \$s	\$d = \$t << \$s
sra \$d, \$t, a	\$d = \$t >> a
sraw \$d, \$t, \$s	\$d = \$t >> \$s
srl \$d, \$t, a	\$d = \$t >>> a
srlv \$d, \$t, \$s	\$d = \$t >>> \$s

Comparison Instructions

Instruction	Operation
slt \$d, \$s, \$t	\$d = (\$s < \$t)
sltu \$d, \$s, \$t	\$d = (\$s < \$t)
slti \$t, \$s, i	\$t = (\$s < SE(i))
sltiu \$t, \$s, i	\$t = (\$s < SE(i))

Instruction	Operation
beq \$s, \$t, label	if (\$s == \$t) pc += i << 2
bgtz \$s, label	if (\$s > 0) pc += i << 2
blez \$s, label	if (\$s <= 0) pc += i << 2
bne \$s, \$t, label	if (\$s != \$t) pc += i << 2

Jump Instructions

Instruction	Operation
j label	pc += i << 2
jal label	\$31 = pc; pc += i << 2
jalr \$s	\$31 = pc; pc = \$s
jr \$s	pc = \$s

Load Instructions

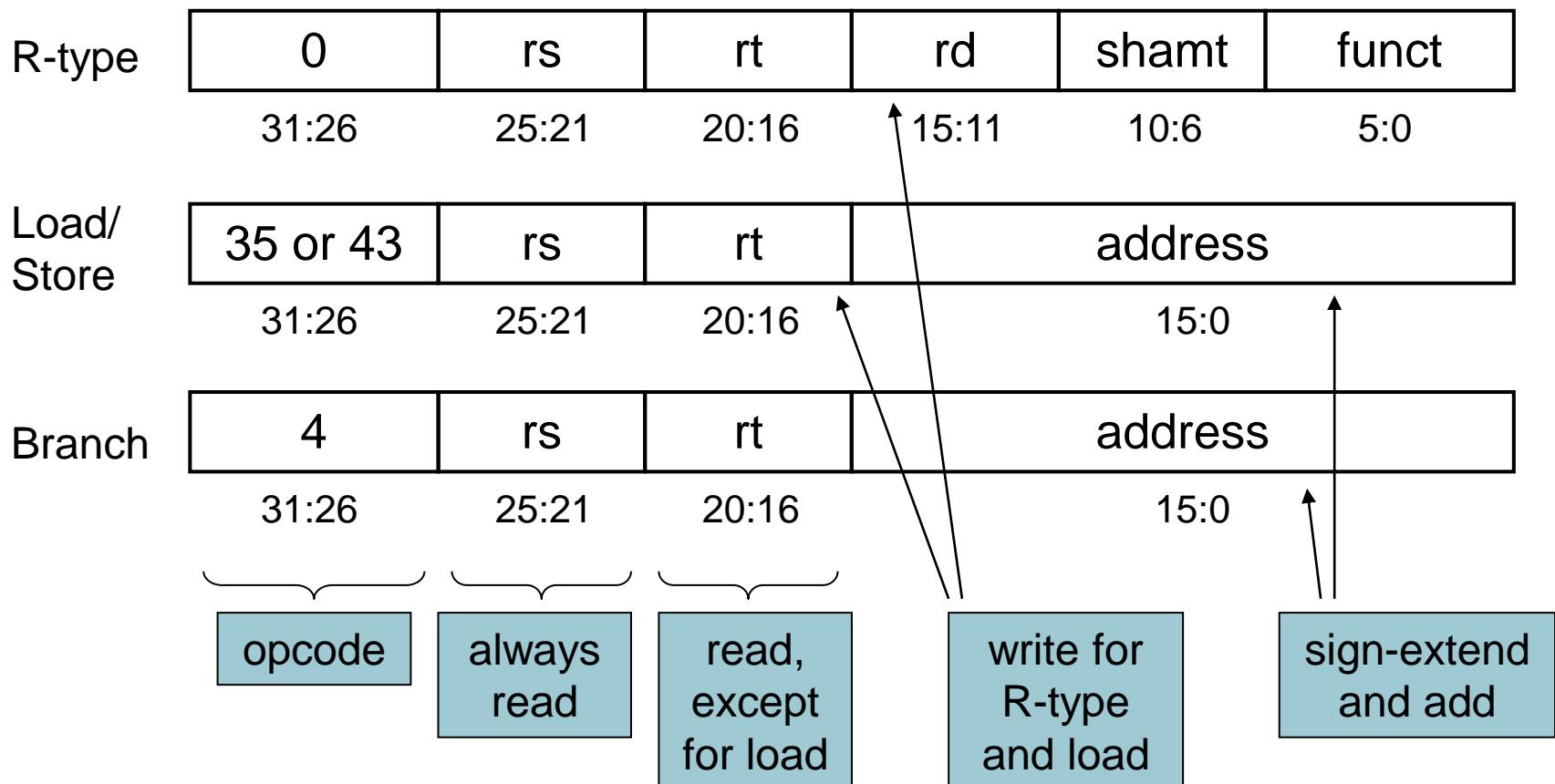
Instruction	Operation
lb \$t, i(\$s)	\$t = SE (MEM [\$s + i]:1)
lbu \$t, i(\$s)	\$t = ZE (MEM [\$s + i]:1)
lh \$t, i(\$s)	\$t = SE (MEM [\$s + i]:2)
lhu \$t, i(\$s)	\$t = ZE (MEM [\$s + i]:2)
lw \$t, i(\$s)	\$t = MEM [\$s + i]:4

Store Instructions

Instruction	Operation
sb \$t, i(\$s)	MEM [\$s + i]:1 = LB (\$t)
sh \$t, i(\$s)	MEM [\$s + i]:2 = LH (\$t)
sw \$t, i(\$s)	MEM [\$s + i]:4 = \$t

The Main Control Unit

Control signals derived from instruction



Introduction

- CPU performance factors
 - Instruction count
 - Determined by ISA and compiler
 - CPI and Cycle time
 - Determined by CPU hardware
- We will examine two MIPS implementations
 - A simplified version
 - A more realistic pipelined version
- Simple subset, shows most aspects
 - Memory reference: lw, sw
 - Arithmetic/logical: add, sub, and, or, slt
 - Control transfer: beq, j

Instruction Execution

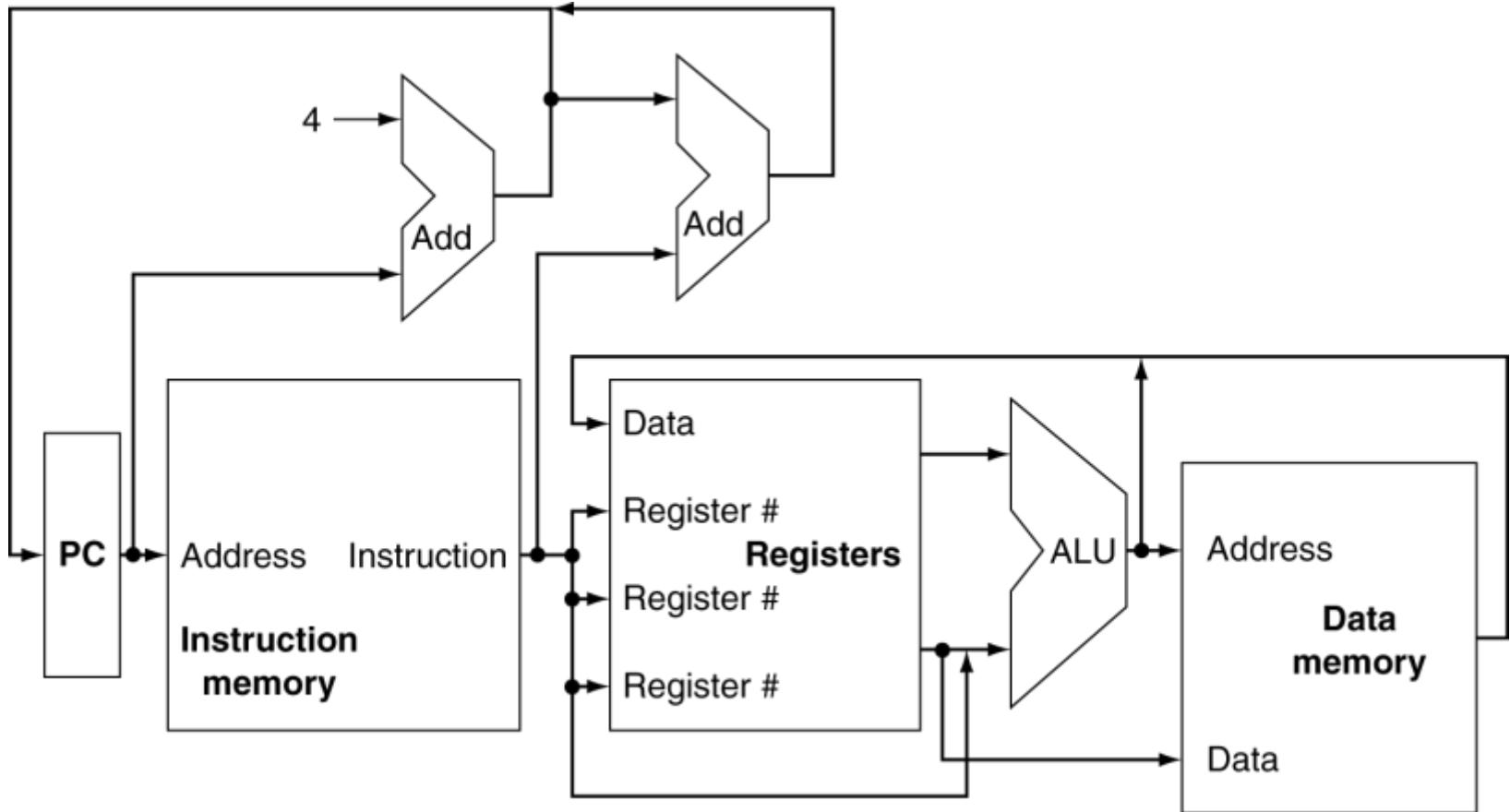
- PC → instruction memory, fetch instruction
- Register numbers → register file, read registers
- Depending on instruction class
 - Use ALU to calculate
 - Arithmetic result
 - Memory address for load/store
 - Branch target address
 - Access data memory for load/store
 - PC ← target address or PC + 4

Instruction Execution

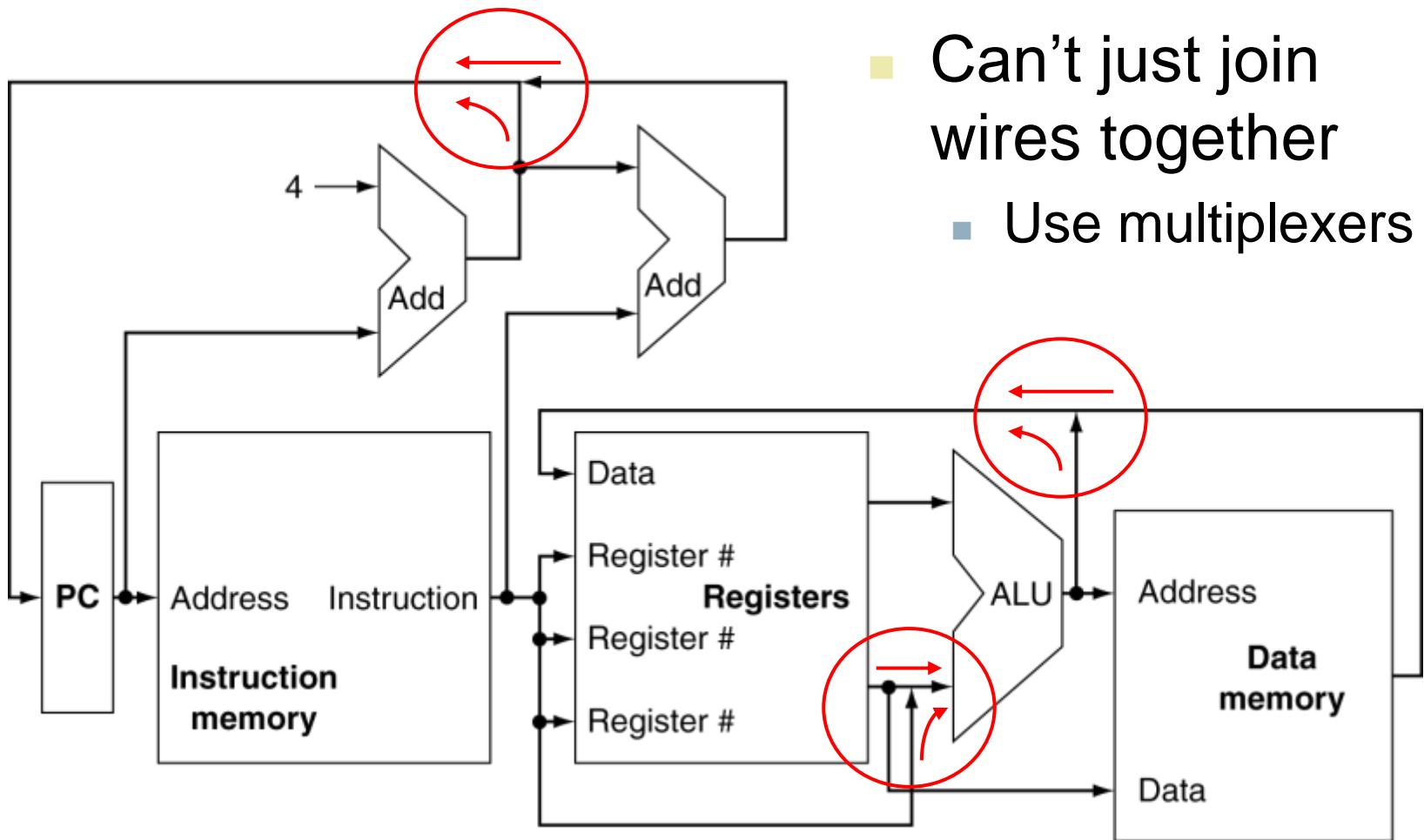
1. Read instruction from instruction memory
 2. Decode instruction and read operands
- Arithmetic (add \$s1, \$s2, \$s3)
 3. Perform arithmetic operation (add, sub, etc.)
 4. Write the result to destination register
 - Load / Store (lw \$s1, 0(\$t1))
 3. Calculate Effective Address
 4. Read from memory (Write to memory incase of store)
 5. Write the read value to the destination register (For load alone)
 - Branch (bne \$t1,\$t2, loop)
 3. Calculate branch outcome and branch target address

Final step: Update PC with PC+4 or branch target address.

CPU Overview

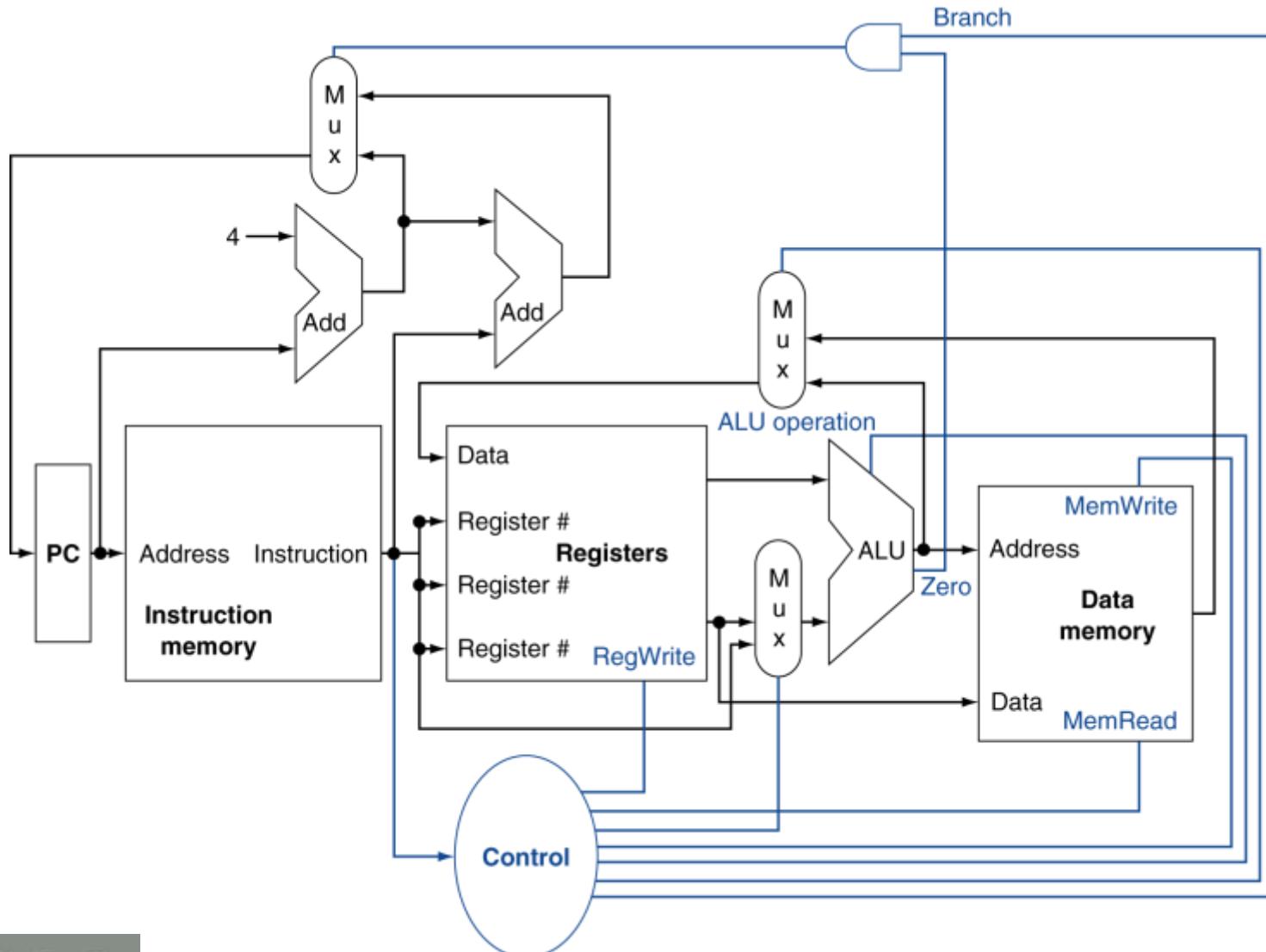


Multiplexers



- Can't just join wires together
 - Use multiplexers

Control

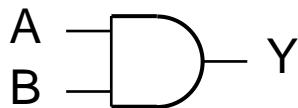


Logic Design Basics

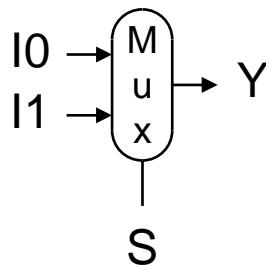
- Information encoded in binary
 - Low voltage = 0, High voltage = 1
 - One wire per bit
 - Multi-bit data encoded on multi-wire buses
- Combinational element
 - Operate on data
 - Output is a function of input
- State (sequential) elements
 - Store information

Combinational Elements

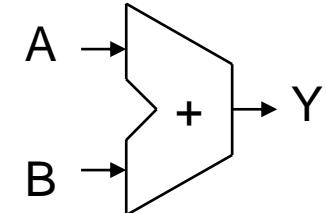
- AND-gate
 - $Y = A \& B$



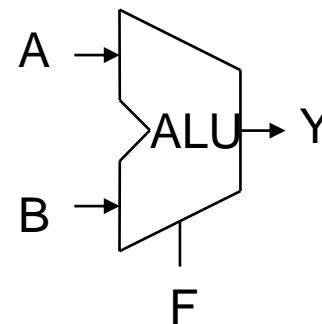
- Multiplexer
 - $Y = S ? I_1 : I_0$



- Adder
 - $Y = A + B$



- Arithmetic/Logic Unit
 - $Y = F(A, B)$



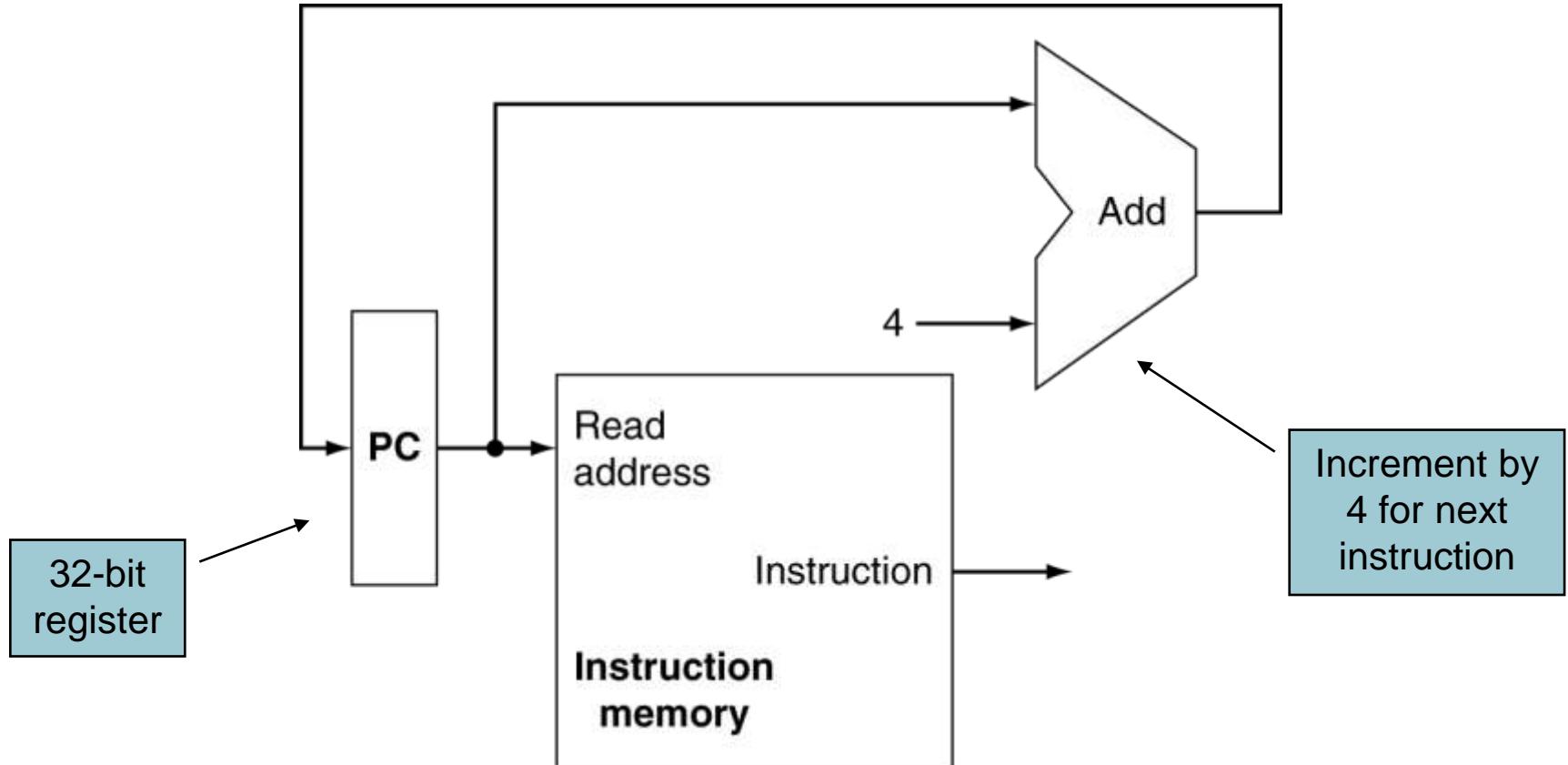
Sequential Elements

- Register: stores data in a circuit
 - Uses a clock signal to determine when to update the stored value
 - Edge-triggered: update when Clk changes from 0 to 1
- Register with write control
 - Only updates on clock edge when write control input is 1
 - Used when stored value is required later

Building a Datapath

- Datapath
 - Elements that process data and addresses in the CPU
 - Registers, ALUs, mux's, memories, ...
- We will build a MIPS datapath incrementally
 - Refining the overview design

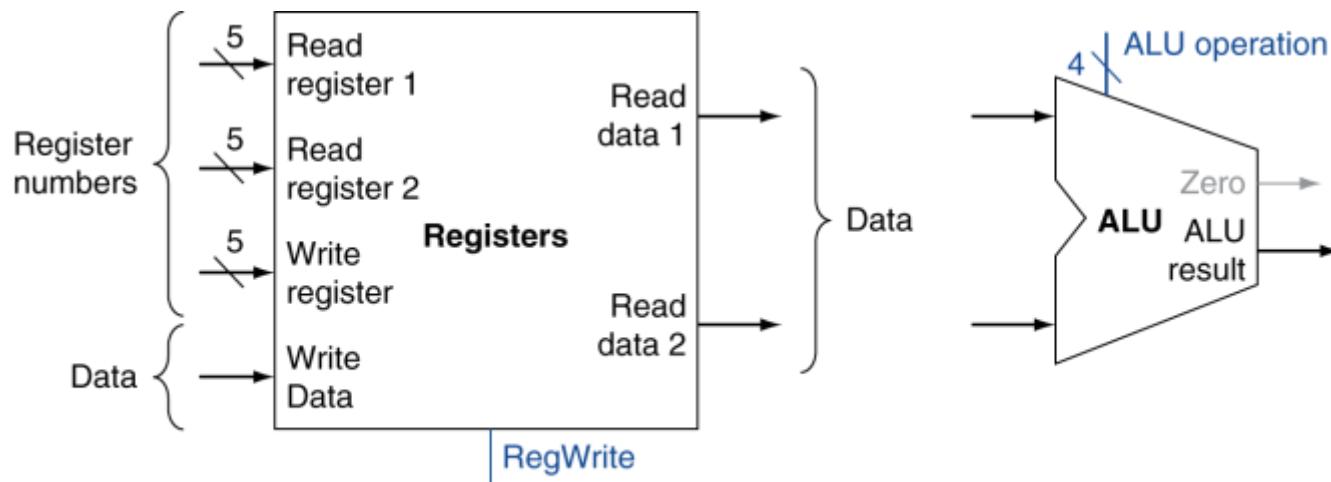
Instruction Fetch



R-Format Instructions

- Read two register operands ADD Perform arithmetic/logical operation
- Write register result

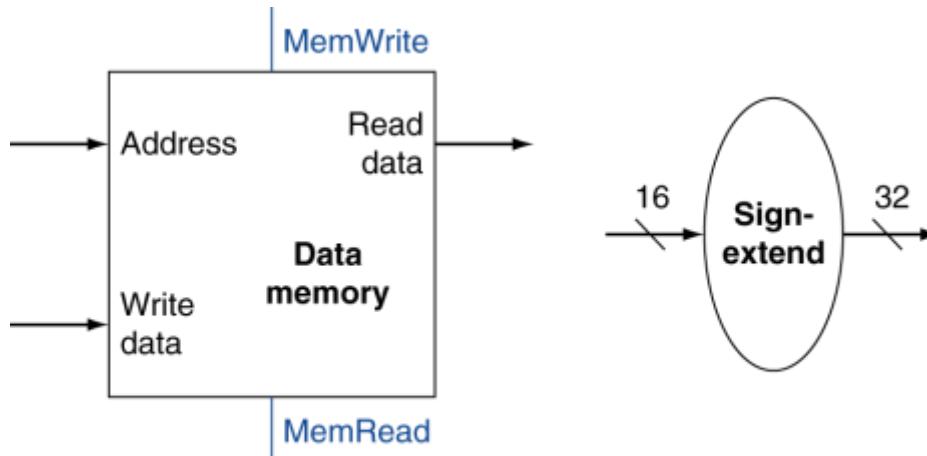
ADD \$S1,\$S2,\$S3



Load/Store Instructions

- Read register operands
- Calculate address using 16-bit offset
 - Use ALU, but sign-extend offset
- Load: Read memory and update register
- Store: Write register value to memory

LD \$T1, 8(\$S1)



a. Data memory unit

b. Sign extension unit

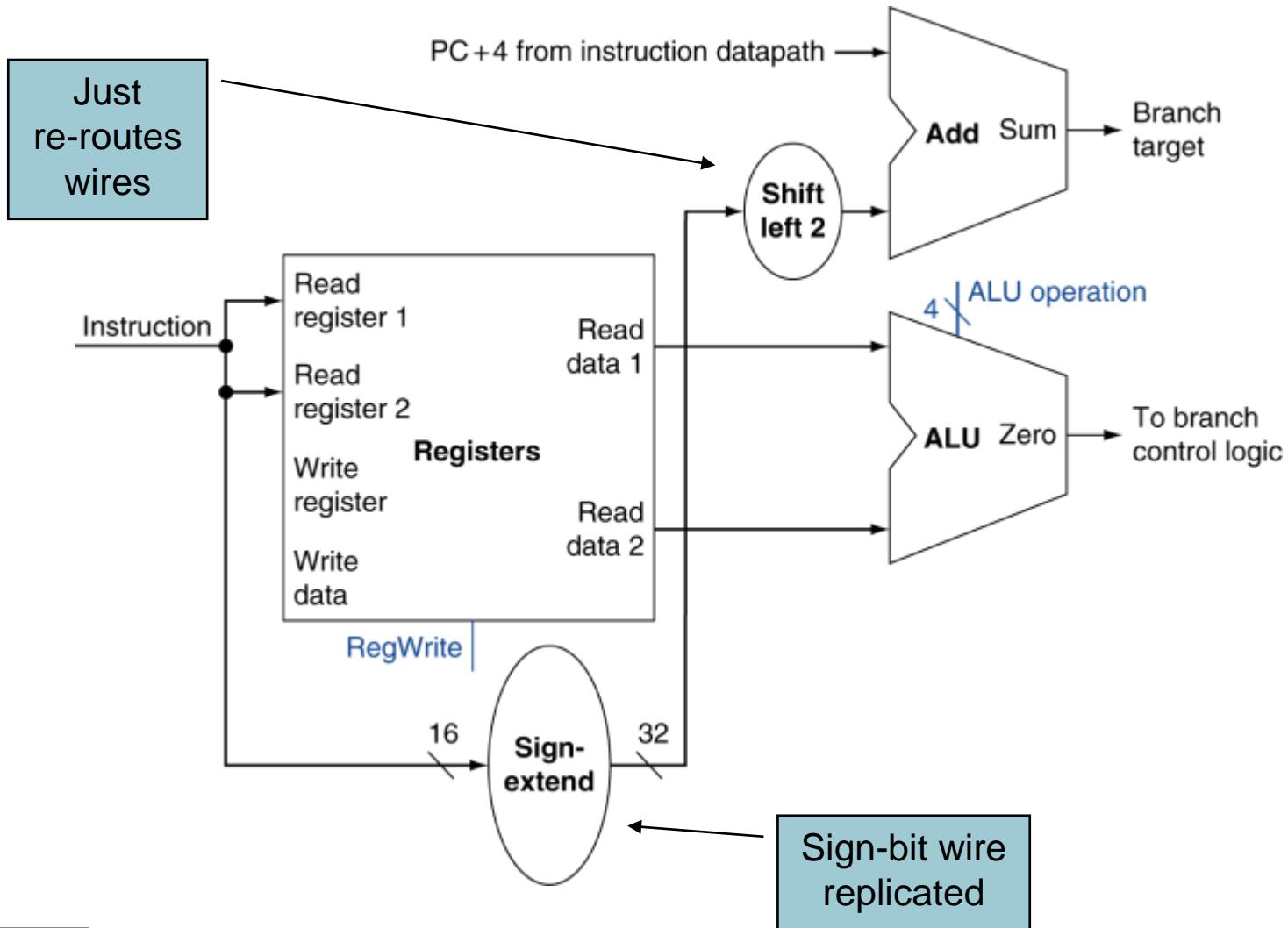
Branch Instructions

- Read register operands
- Compare operands
 - Use ALU, subtract and check Zero output
- Calculate target address
 - Sign-extend displacement
 - Shift left 2 places (word displacement)
 - Add to PC + 4
 - Already calculated by instruction fetch

BEQ \$T1, \$T2, LOOP

Branch Instructions

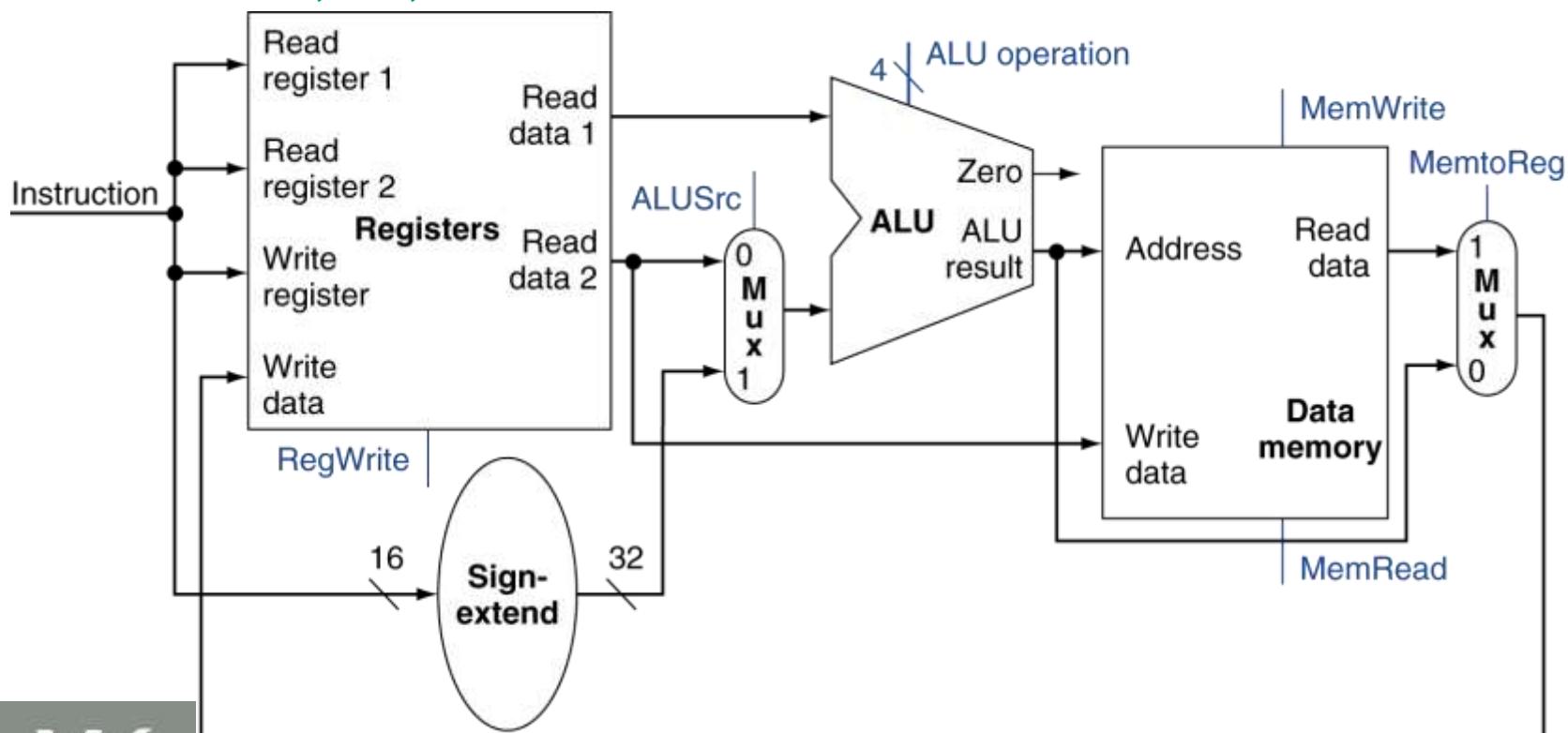
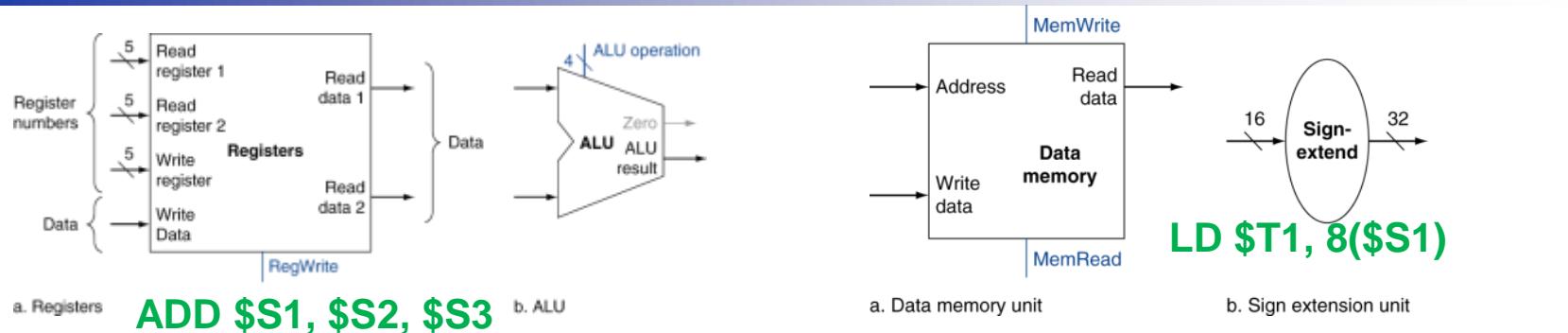
BEQ \$T1, \$T2, LOOP



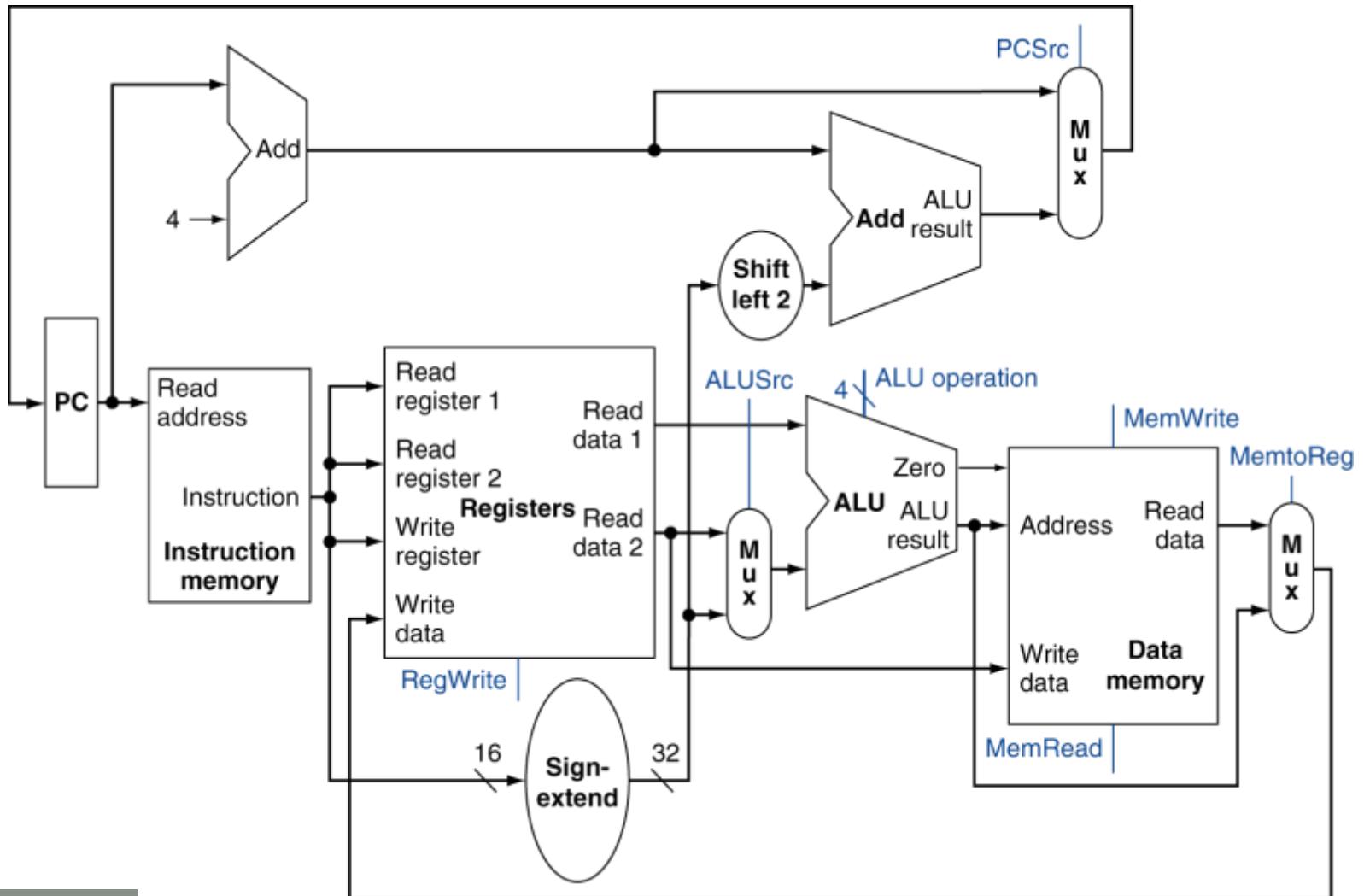
Composing the Elements

- First-cut data path does an instruction in one clock cycle
 - Each datapath element can only do one function at a time
 - Hence, we need separate instruction and data memories
- Reuse units and use multiplexers where alternate data sources are used for different instructions

R-Type/Load/Store Datapath



Full Datapath



ALU Control

Memory reference: lw, sw
Arithmetic/logical: add, sub,
and, or, slt
Control transfer: beq, j

■ ALU used for

- Load/Store: F = add
- Branch: F = subtract
- R-type: F depends on funct field

ALU control	Function
0000	AND
0001	OR
0010	add
0110	subtract
0111	set-on-less-than
1100	NOR

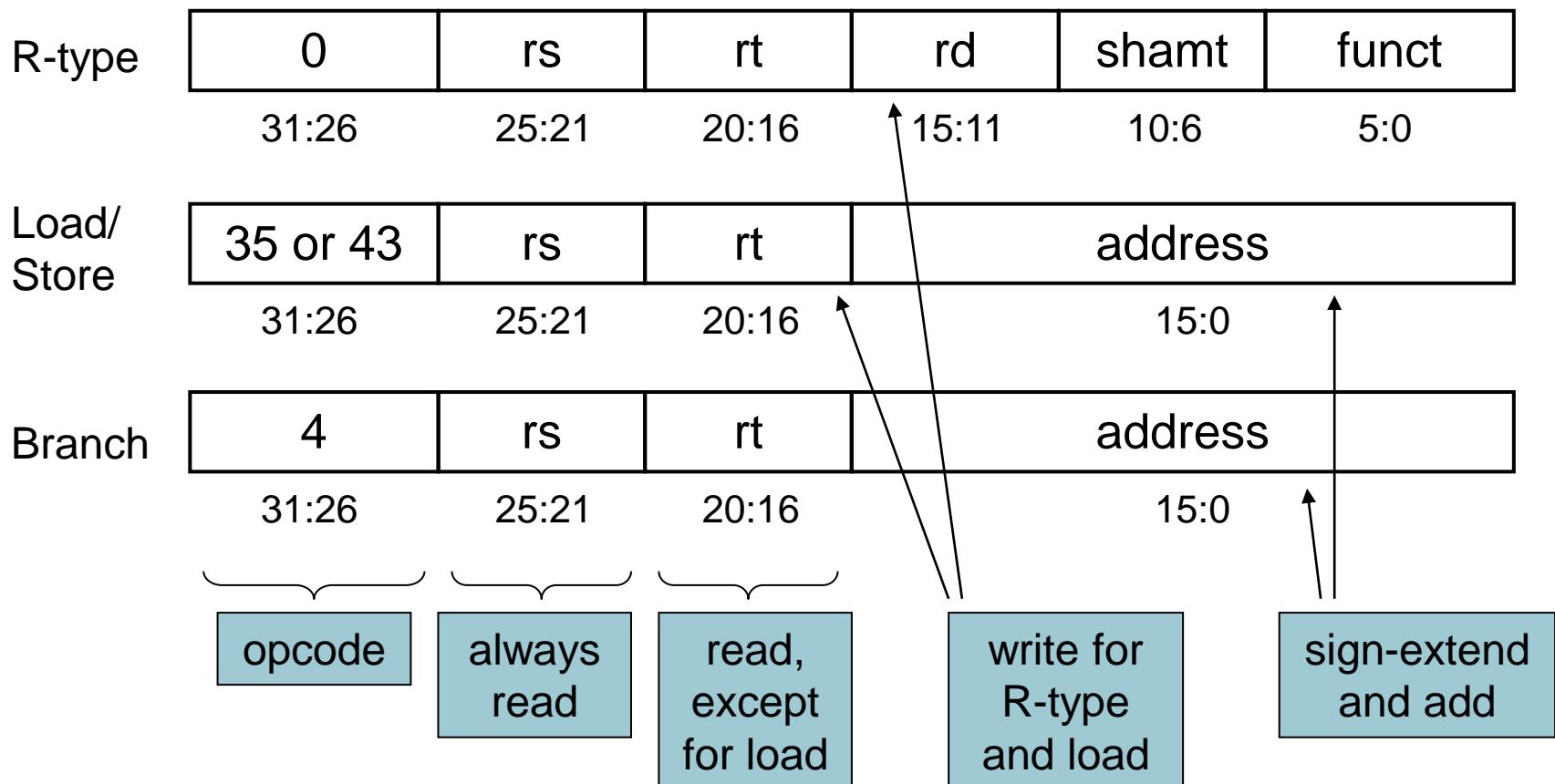
ALU Control

- Assume 2-bit ALUOp derived from opcode
 - Combinational logic derives ALU control
 - Multiple levels of control – reduces size of control

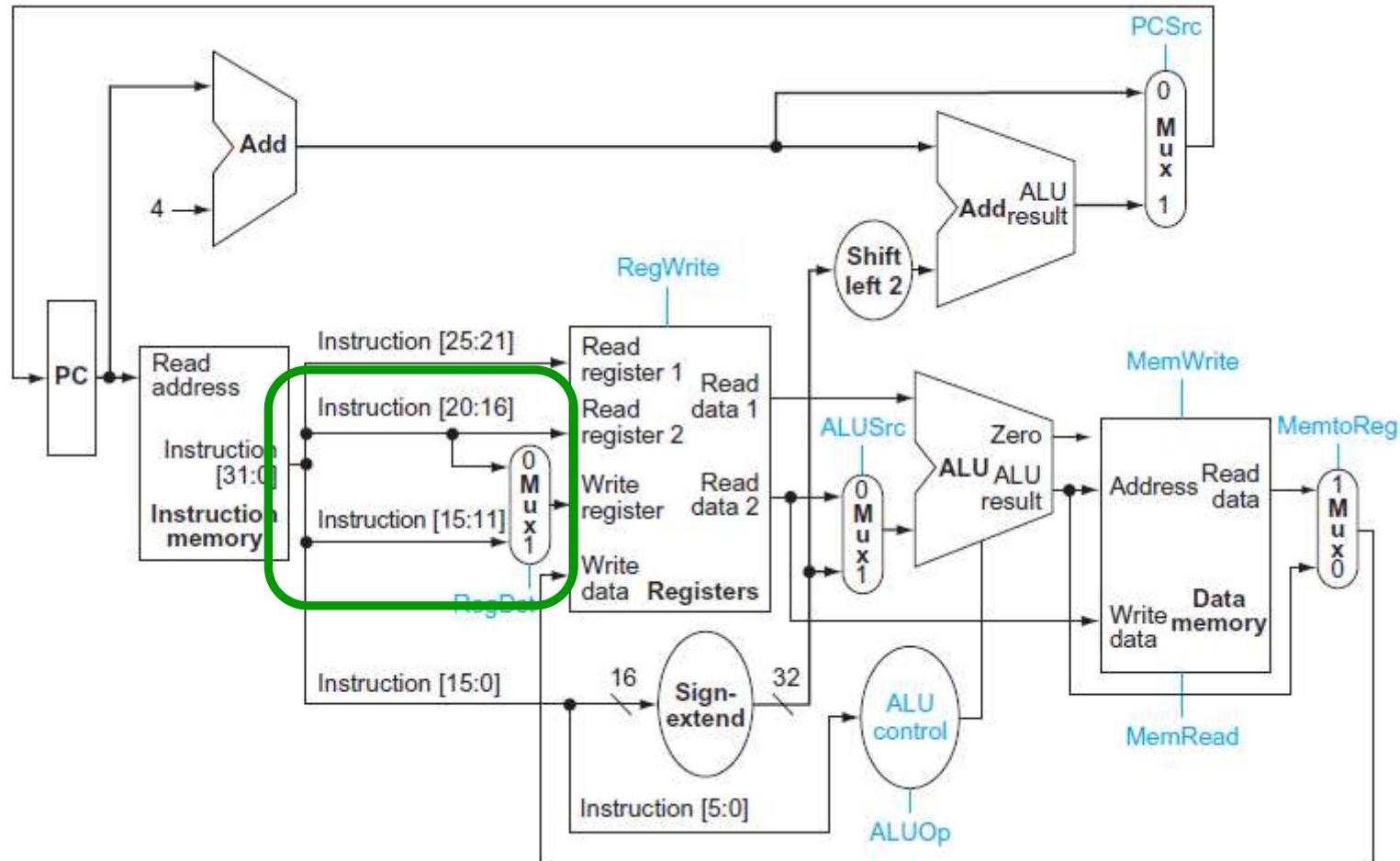
opcode	ALUOp	Operation	funct	ALU function	ALU control
lw	00	load word	xxxxxx	add	0010
sw	00	store word	xxxxxx	add	0010
beq	01	branch equal	xxxxxx	subtract	0110
R-type	10	add	100000	add	0010
		subtract	100010	subtract	0110
		AND	100100	AND	0000
		OR	100101	OR	0001
		set-on-less-than	101010	set-on-less-than	0111

The Main Control Unit

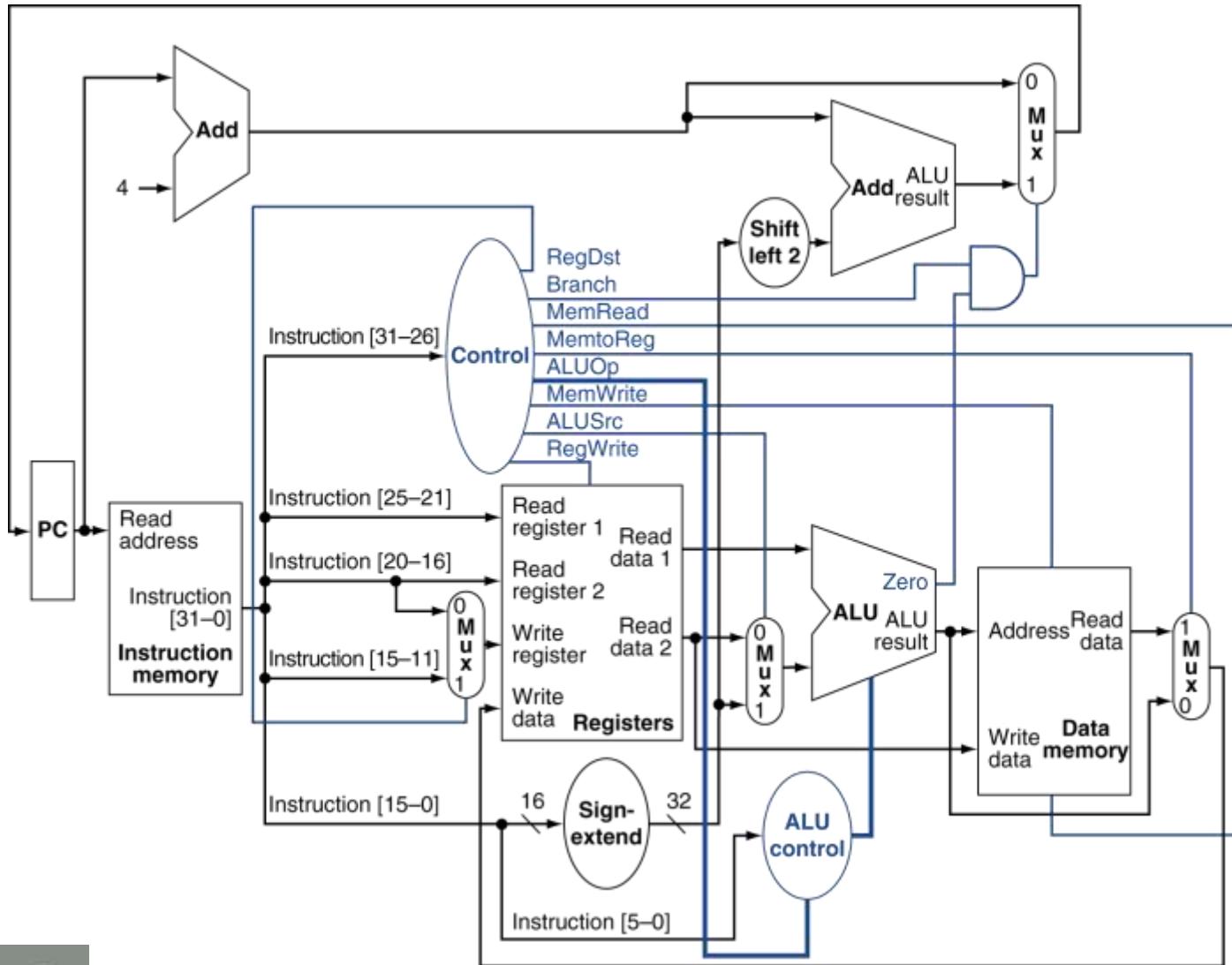
Control signals derived from instruction



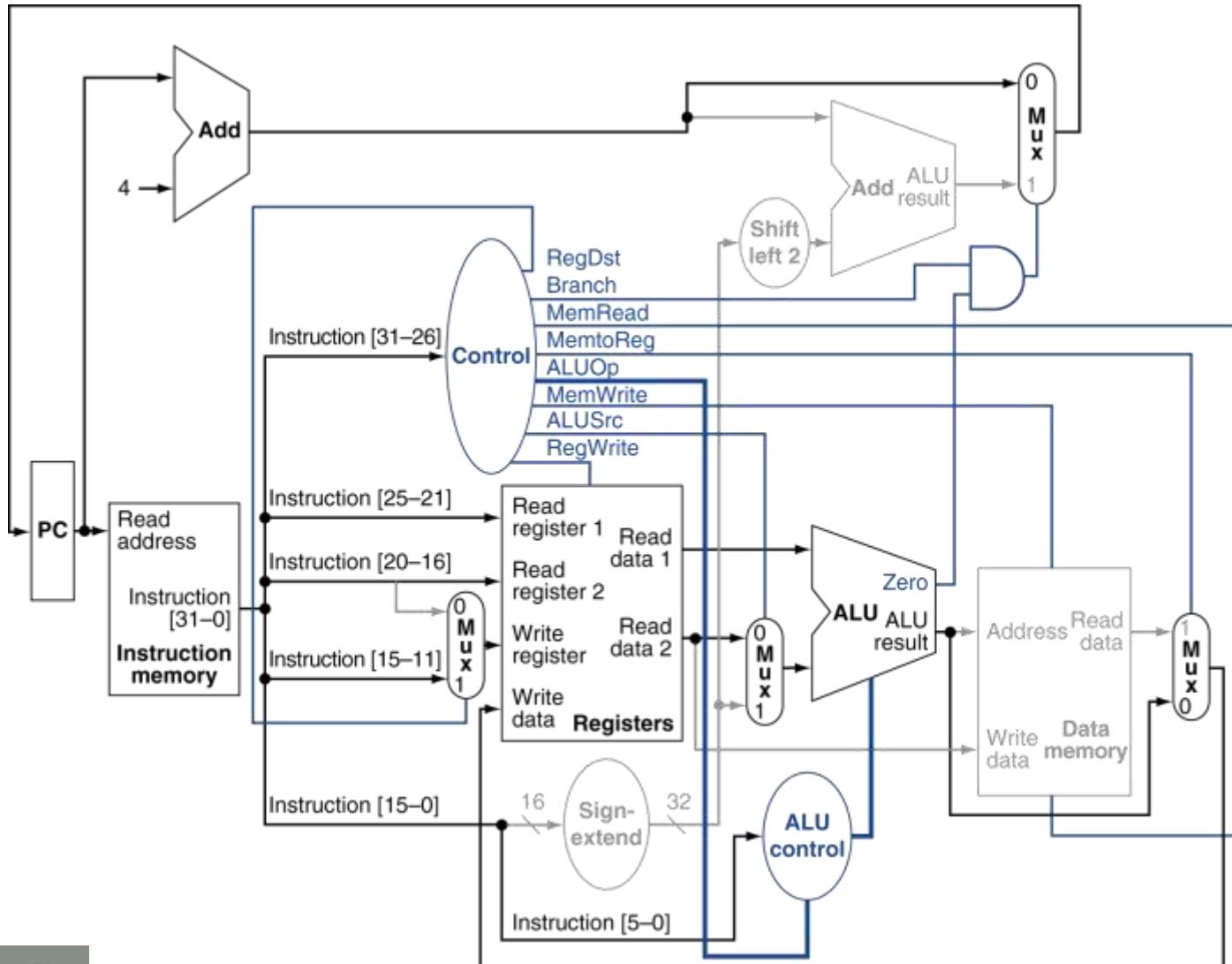
Datapath with Instruction Fields



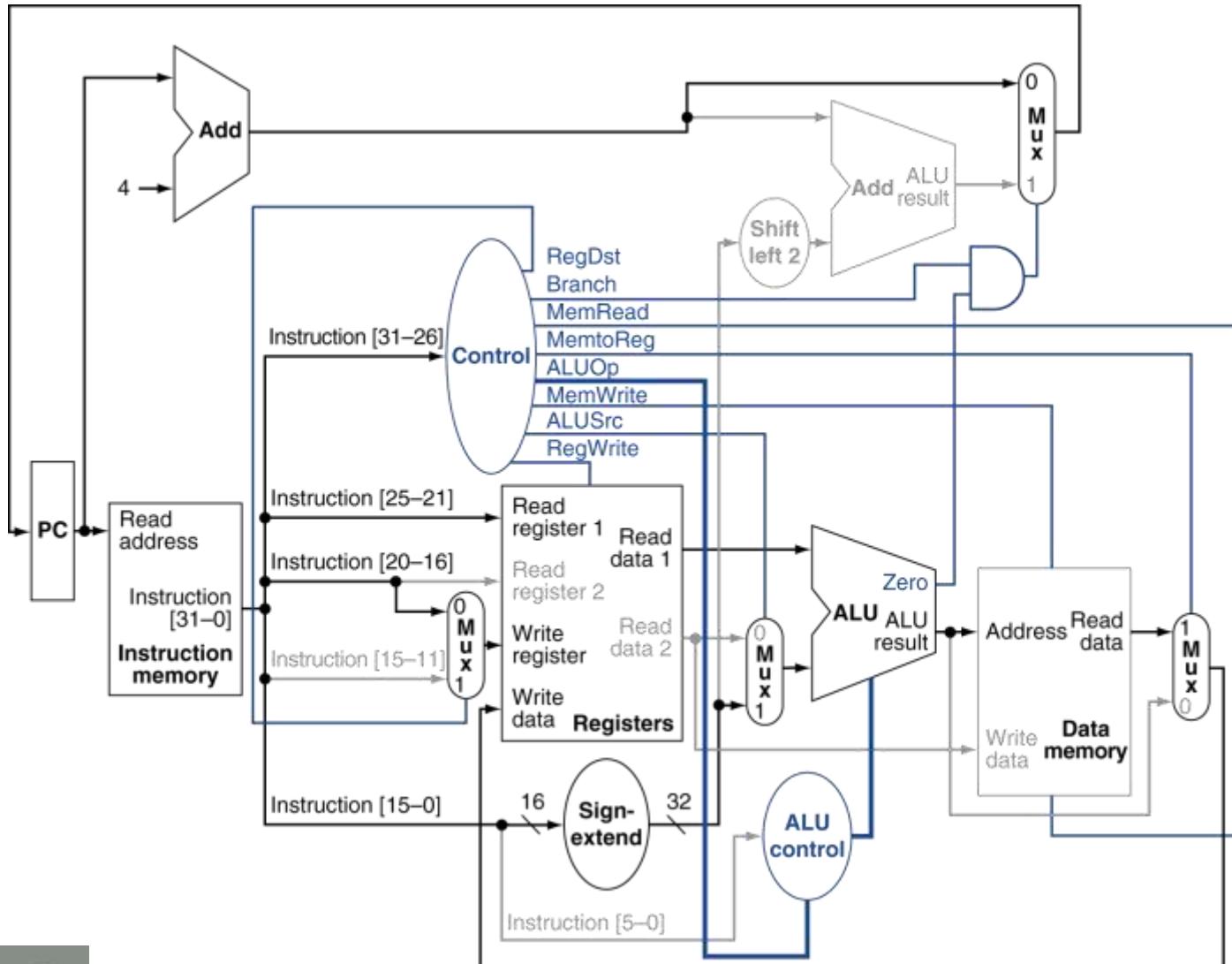
Datapath With Control



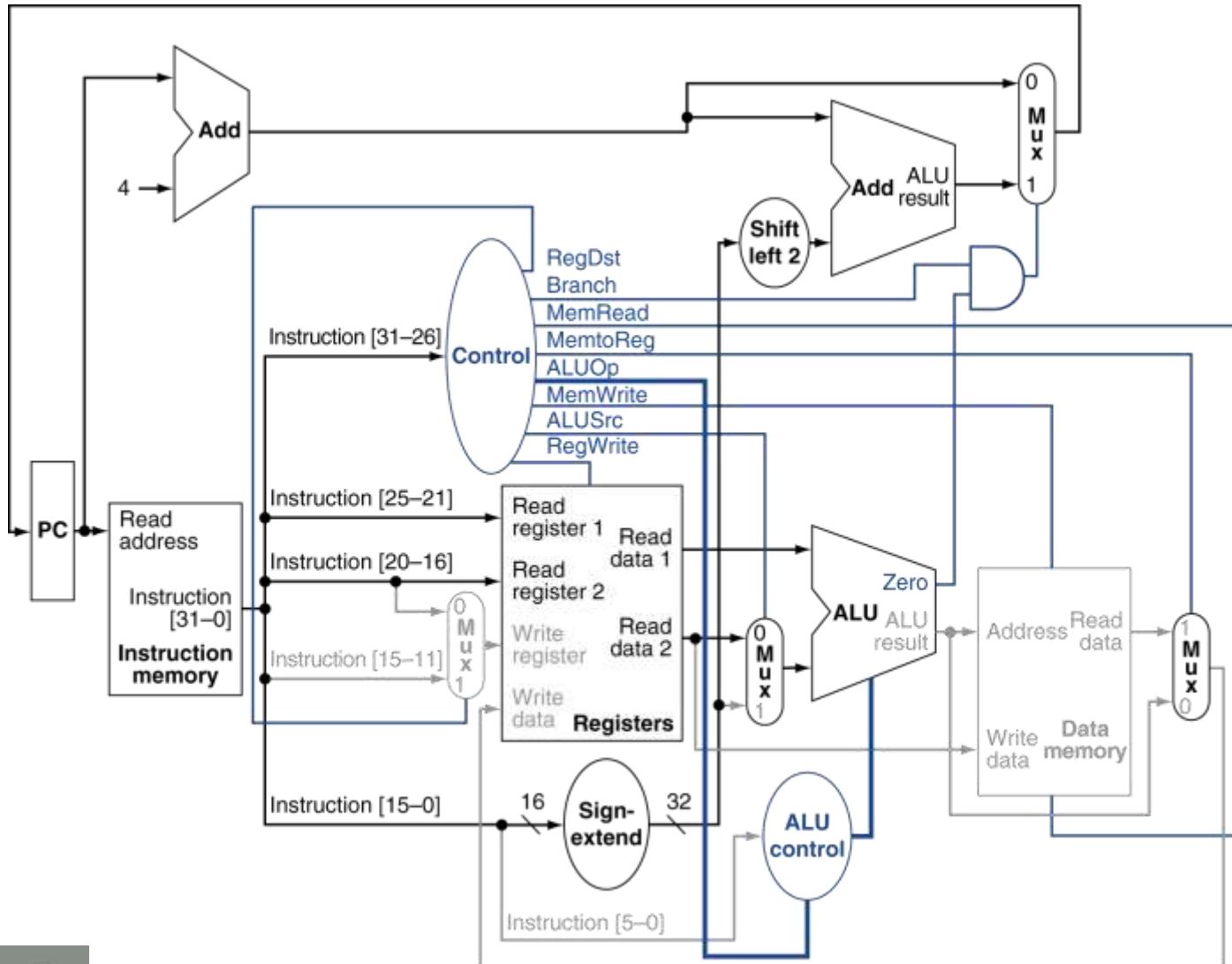
R-Type Instruction



Load Instruction



Branch-on-Equal Instruction



Performance Issues

- Longest delay determines clock period
 - Critical path: load instruction
 - Instruction memory → register file → ALU → data memory → register file
- Not feasible to vary period for different instructions
- Violates design principle
 - Making the common case fast
- We will improve performance by pipelining

Next Topic ...

- Pipelining



COMPUTER ORGANIZATION AND DESIGN

The Hardware/Software Interface

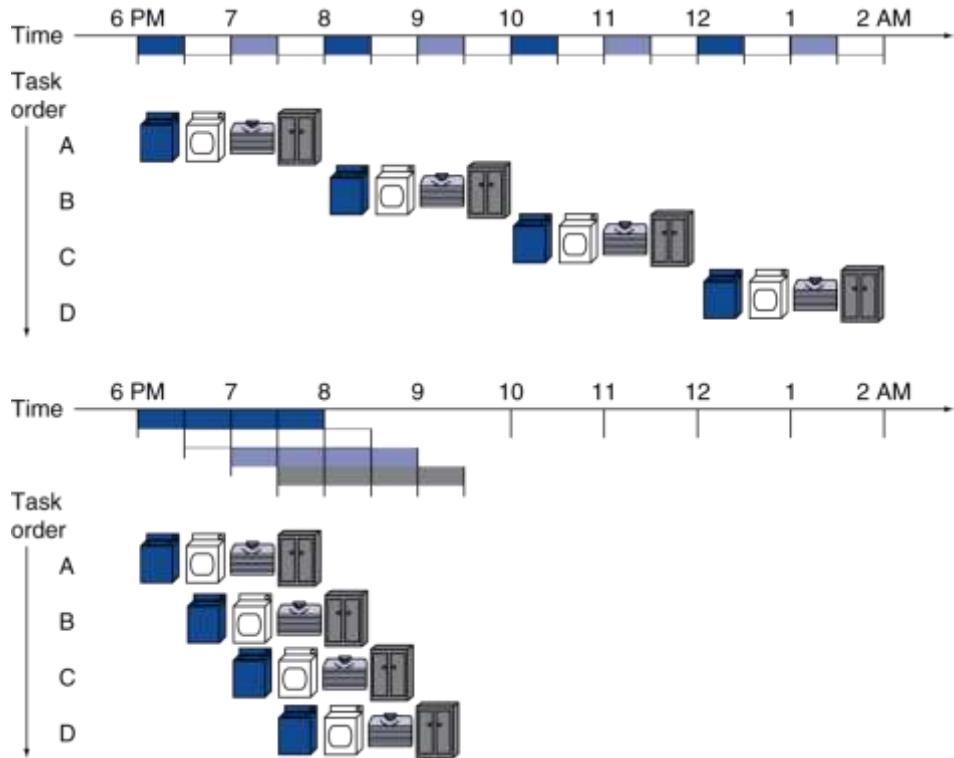


Processor Pipelining

Adapted and Supplemented by,
Dr. R. Shathanaa

Pipelining Analogy

- Pipelined laundry: overlapping execution
 - Parallelism improves performance



- Four loads:
 - Speedup
 $= 8/3.5 = 2.3$
- Non-stop:
 - Speedup
 $= 2n/0.5n + 1.5 \approx 4$
= number of stages

MIPS Pipeline

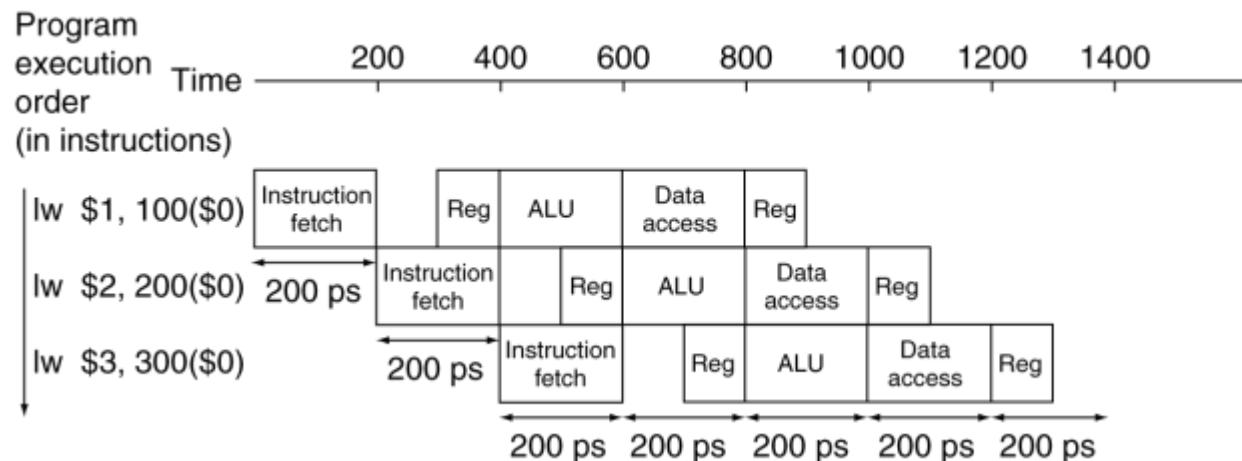
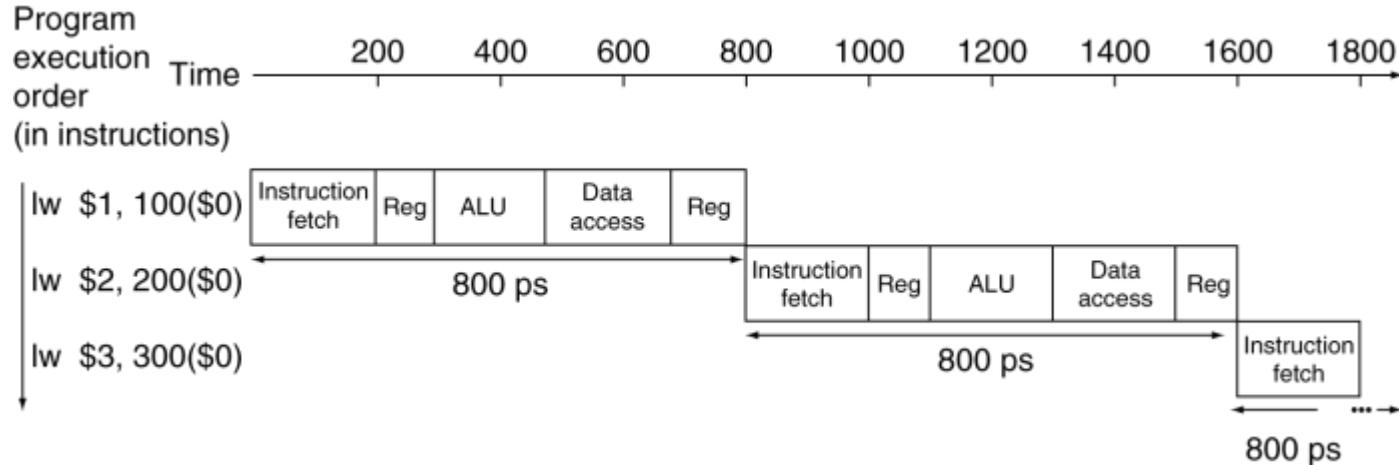
- Five stages, one step per stage
 1. IF: Instruction fetch from memory
 2. ID: Instruction decode & register read
 3. EX: Execute operation or calculate address
 4. MEM: Access memory operand
 5. WB: Write result back to register

Pipeline Performance

- Assume time for stages is
 - 100ps for register read or write
 - 200ps for other stages
- Compare pipelined datapath with single-cycle datapath

Instr	Instr fetch	Register read	ALU op	Memory access	Register write	Total time
lw	200ps	100 ps	200ps	200ps	100 ps	800ps
sw	200ps	100 ps	200ps	200ps		700ps
R-format	200ps	100 ps	200ps		100 ps	600ps
beq	200ps	100 ps	200ps			500ps

Pipeline Performance



Pipeline Speedup

- If all stages are balanced
 - i.e., all take the same time
 - Time between instructions_{pipelined}
= Time between instructions_{nonpipelined}

Number of stages
- If not balanced, speedup is less
- Speedup due to increased throughput
 - Latency (time for each instruction) does not decrease

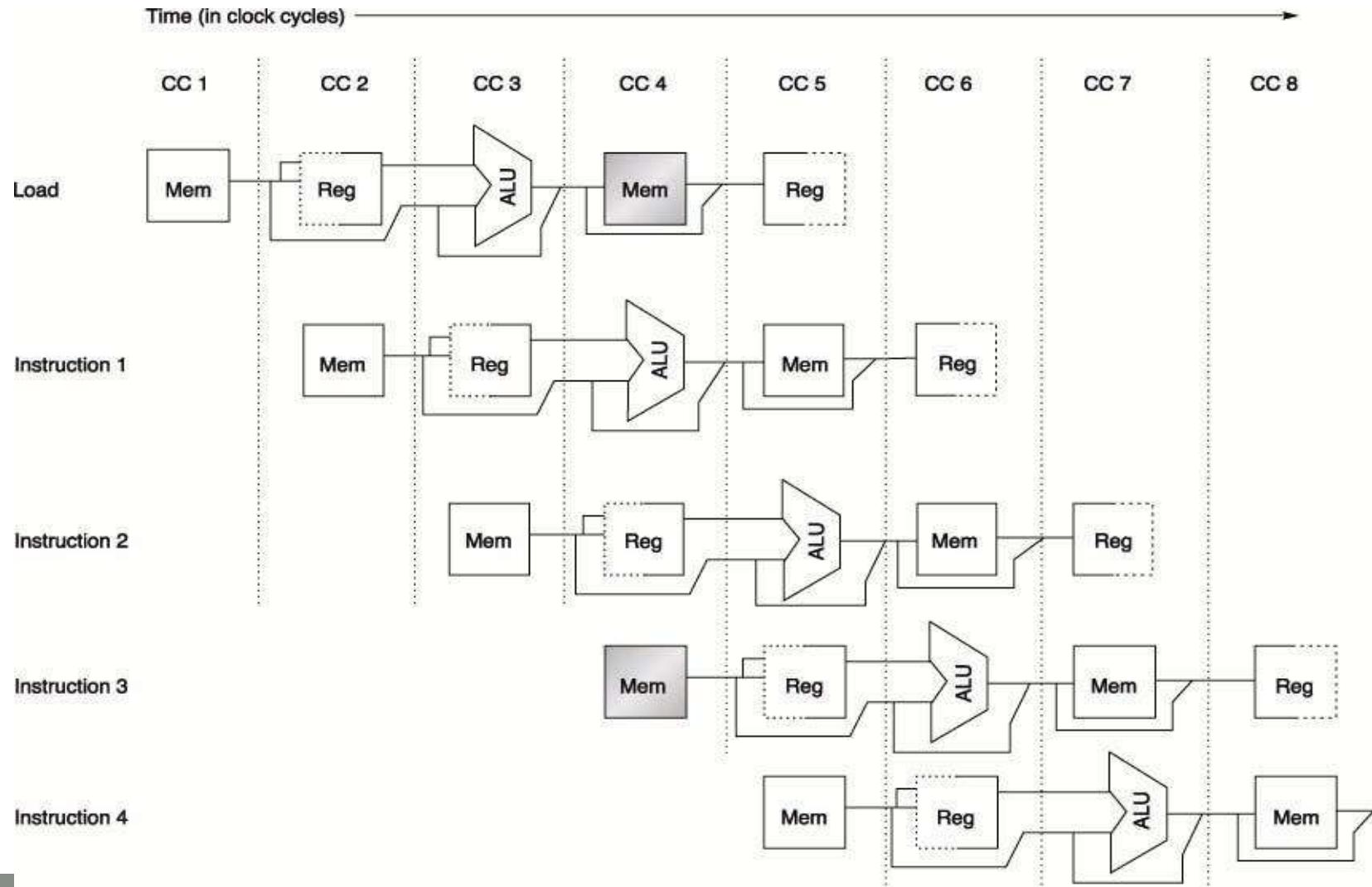
Hazards

- Situations that prevent starting the next instruction in the next cycle
- Structural hazard
 - A required resource is busy
- Data hazard
 - Need to wait for previous instruction to complete its data read/write
- Control hazard
 - Deciding on control action depends on previous instruction

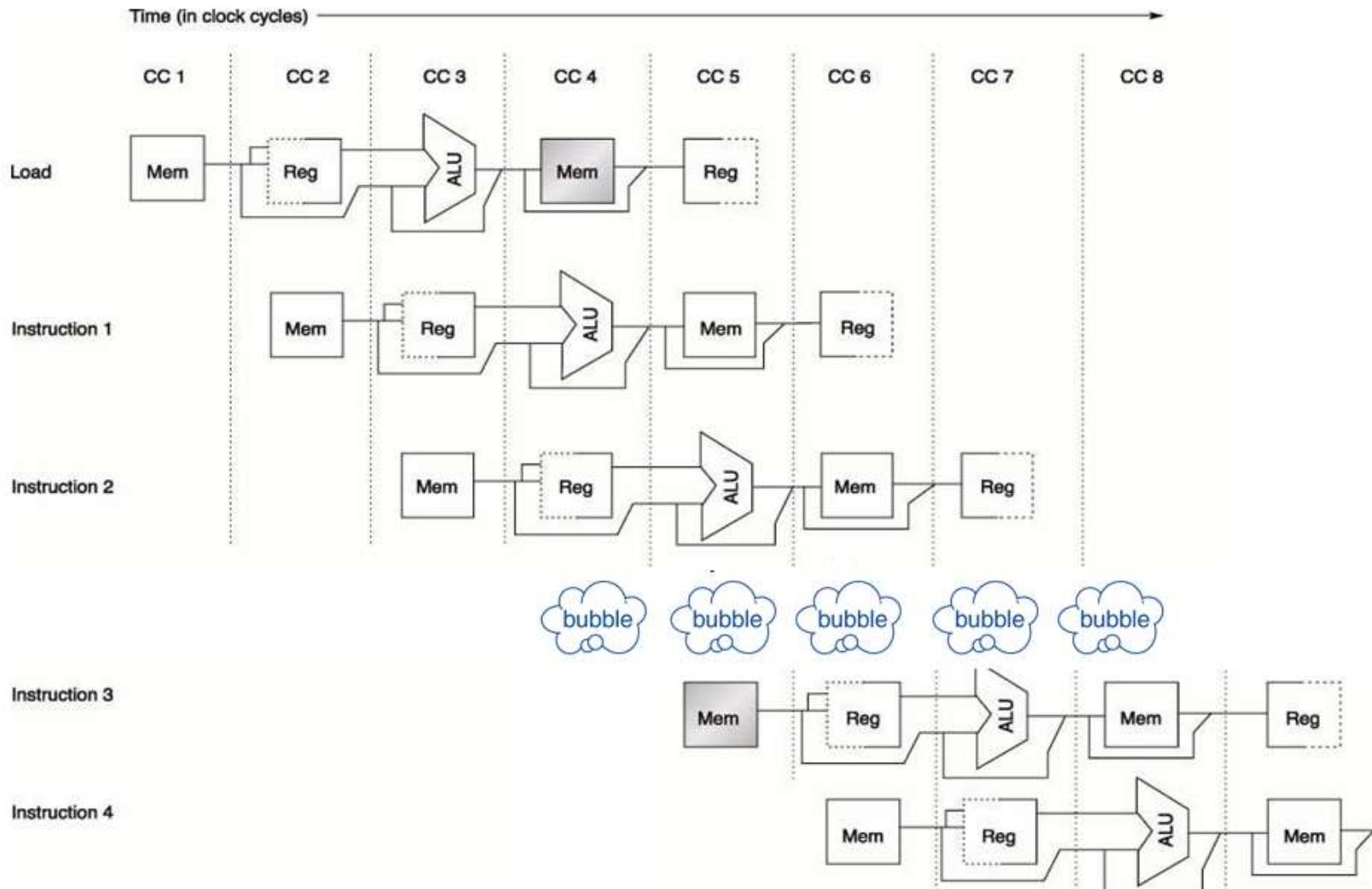
Structural Hazard

- Conflict for use of a resource
- In MIPS pipeline with a single memory
 - Load/store requires data access
 - Instruction fetch would have to *stall* for that cycle
 - Would cause a pipeline “bubble”
- Hence, pipelined datapaths require separate instruction/data memories
 - Or separate instruction/data caches
 - In general, resources need to be replicated to avoid structural hazard and stall

Structural Hazard

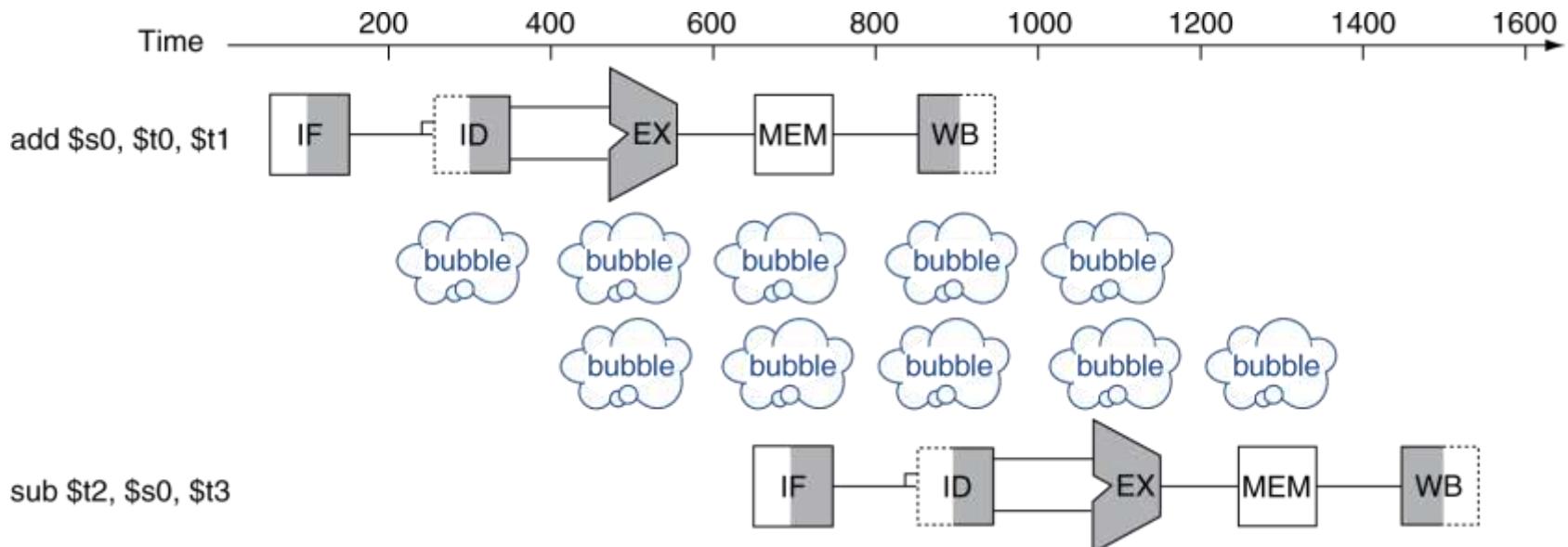


Structural Hazard



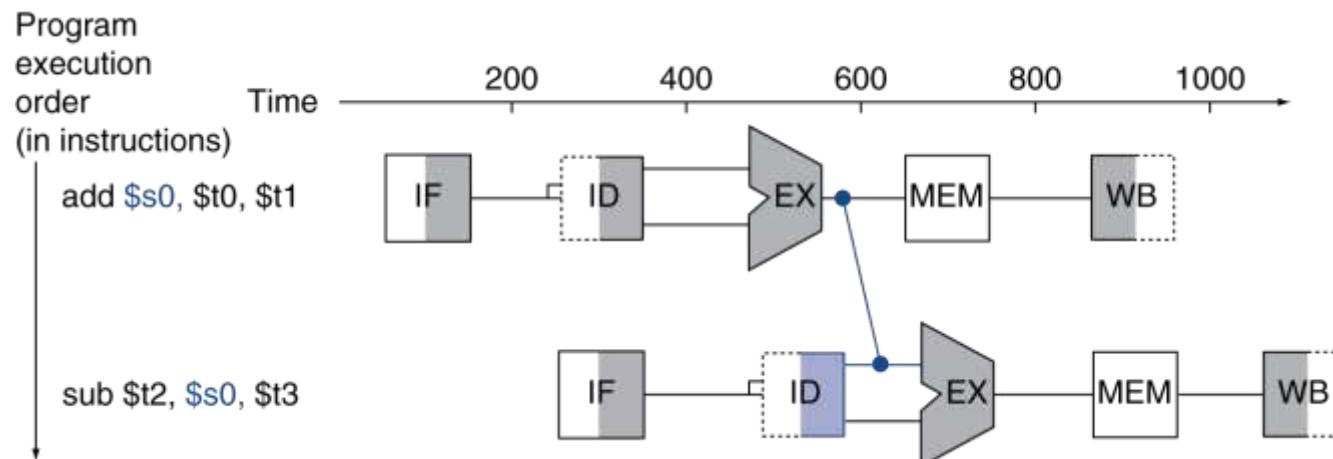
Data Hazards

- An instruction depends on completion of data access by a previous instruction
 - add \$s0, \$t0, \$t1
 - sub \$t2, \$s0, \$t3



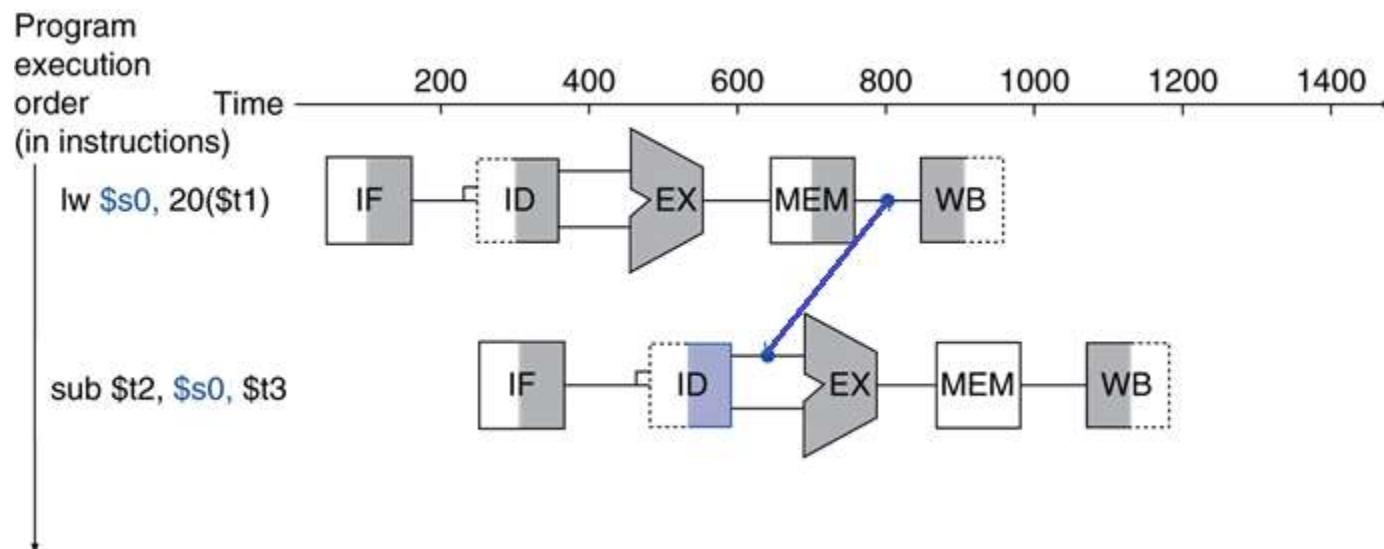
Forwarding (aka Bypassing)

- Use result when it is computed
 - Don't wait for it to be stored in a register
 - Requires extra connections in the datapath



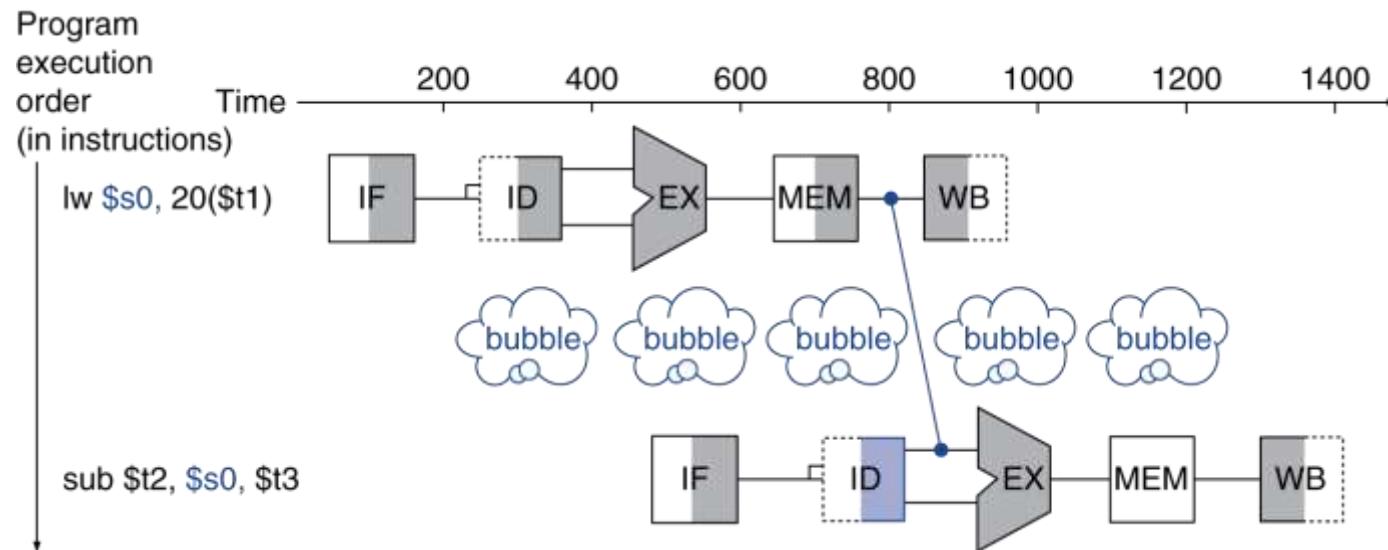
Load-Use Data Hazard

- Can't always avoid stalls by forwarding
 - If value not computed when needed
 - Can't forward backward in time!



Load-Use Data Hazard

- Can't always avoid stalls by forwarding
 - If value not computed when needed
 - Can't forward backward in time!



Code Scheduling to Avoid Stalls

- Reorder code to avoid use of load result in the next instruction
- C code for $A = B + E; C = B + F;$

The diagram illustrates the process of code scheduling to eliminate stalls. It shows two versions of the same assembly-like code, separated by a blue arrow pointing from left to right.

Initial Code (Left):

```
lw    $t1, 0($t0)
lw    $t2, 4($t0)
stall
add $t3, $t1, $t2
sw    $t3, 12($t0)
lw    $t4, 8($t0)
stall
add $t5, $t1, $t4
sw    $t5, 16($t0)
```

Optimized Code (Right):

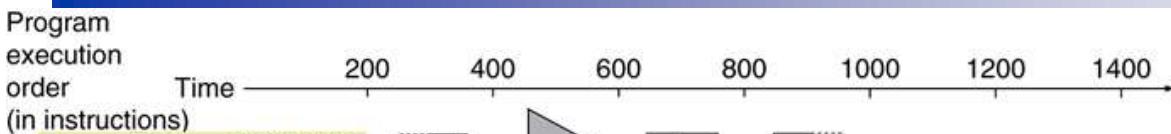
```
lw    $t1, 0($t0)
lw    $t2, 4($t0)
lw    $t4, 8($t0)
add $t3, $t1, $t2
sw    $t3, 12($t0)
add $t5, $t1, $t4
sw    $t5, 16($t0)
```

Performance Metrics:

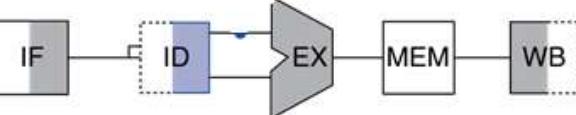
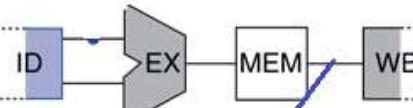
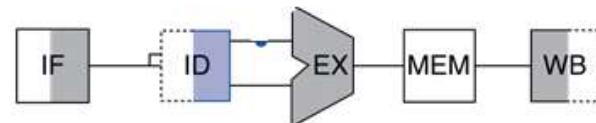
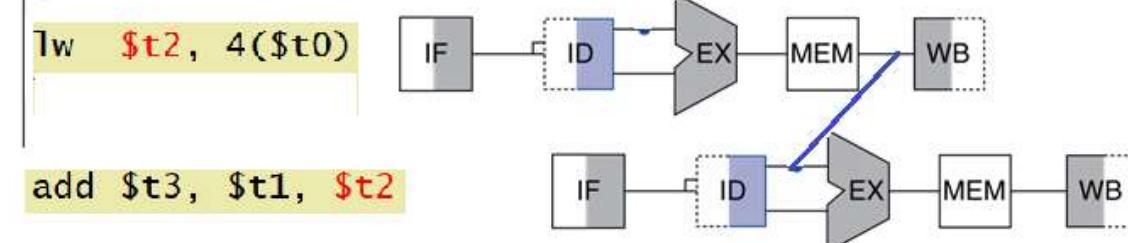
- Initial Code (Left):** 13 cycles
- Optimized Code (Right):** 11 cycles

11 cycles

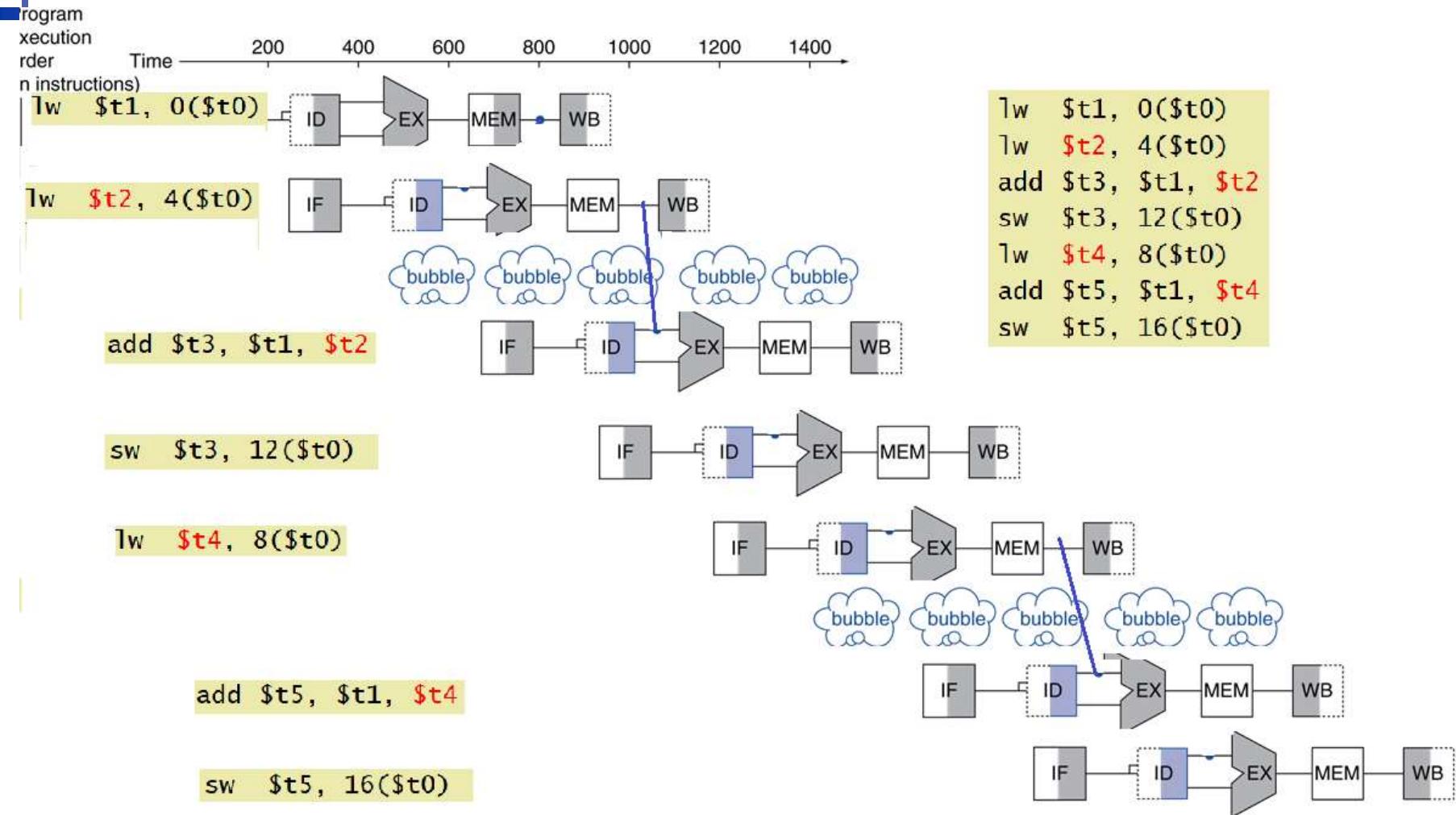
Code Scheduling to Avoid Stalls



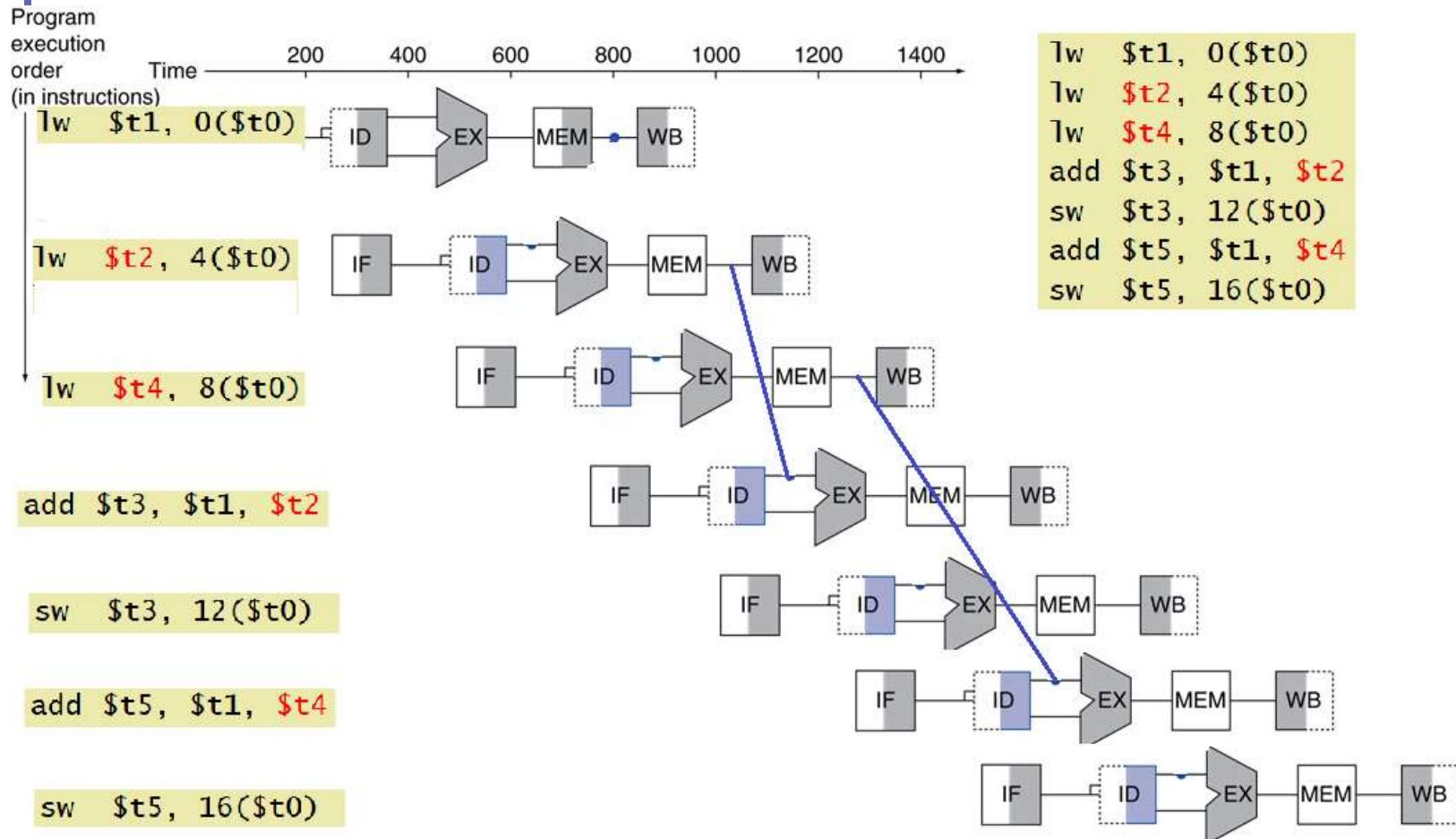
```
lw $t1, 0($t0)
lw $t2, 4($t0)
add $t3, $t1, $t2
sw $t3, 12($t0)
lw $t4, 8($t0)
add $t5, $t1, $t4
sw $t5, 16($t0)
```



Code Scheduling to Avoid Stalls



Code Scheduling to Avoid Stalls



Next Topic ...

- Control Hazards



COMPUTER ORGANIZATION AND DESIGN

The Hardware/Software Interface



Chapter 4

The Processor

Adapted and Supplemented by,
Dr. R. Shathanaa

Control Hazards

- Branch determines flow of control
 - Fetching next instruction depends on branch outcome
 - Pipeline can't always fetch correct instruction
 - Still working on ID stage of branch
- In MIPS pipeline
 - Need to compare registers and compute target early in the pipeline
 - Add hardware to do it in ID stage

Stall on Branch

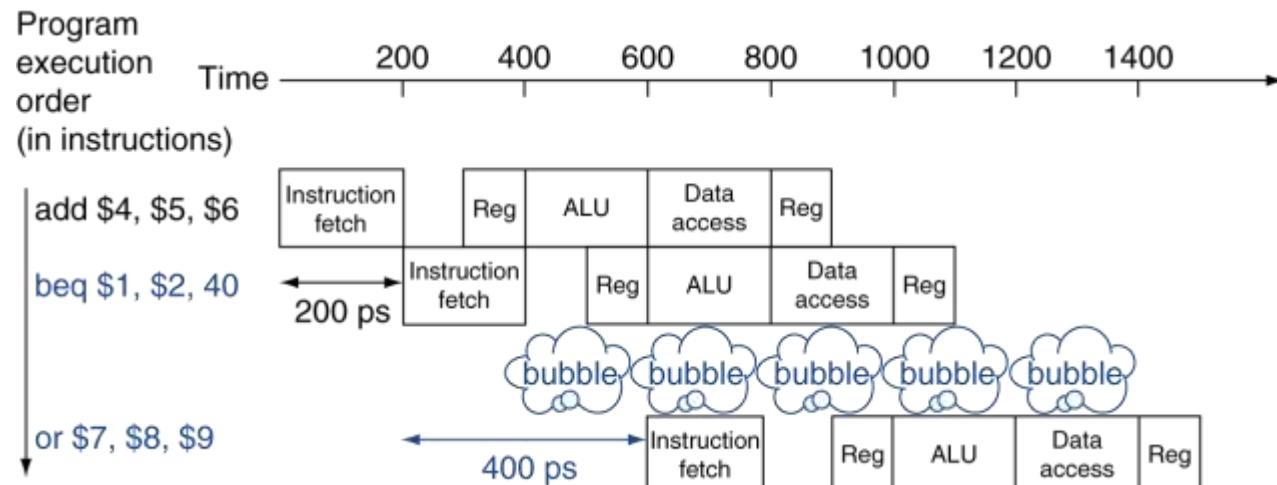
- Wait until branch outcome determined (during Decode stage) before fetching next instruction

add \$4,\$5,\$6

beq \$1,\$2,40

lw \$3,300(\$0)

40: or \$7,\$8,\$9



Performance of “Stall on Branch”

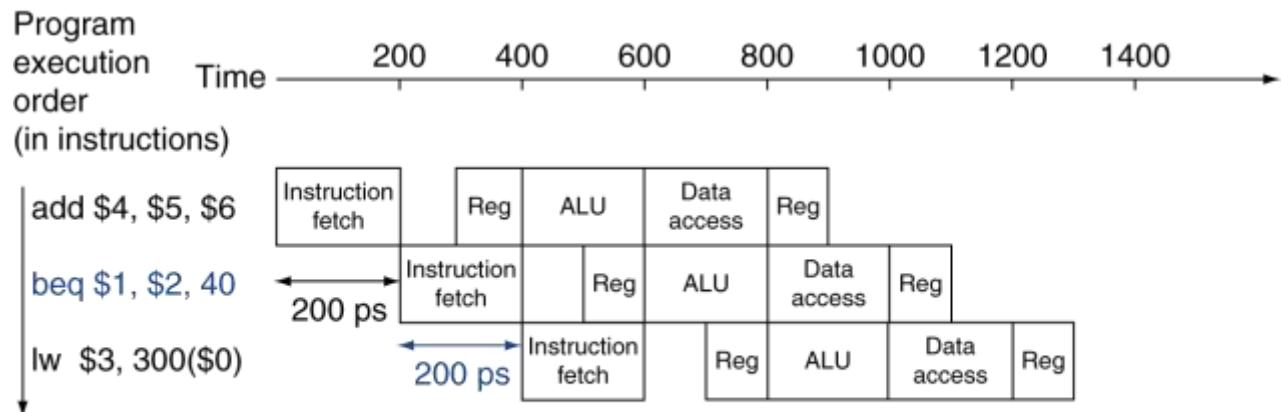
- Estimate the impact on the *clock cycles per instruction* (CPI) of stalling on branches. Assume all other instructions have a CPI of 1 and branches occur 17% of the time.
- Solution
 - Since the other instructions run have a CPI of 1 and branches occur 17% of time, they will incur latency only for that 17% of time
 - new CPI = old CPI + Branch penalty
$$\begin{aligned} &= 1 + (0.17 \times 1 \text{ cycle}) \\ &= 1.17 \end{aligned}$$

Branch Prediction

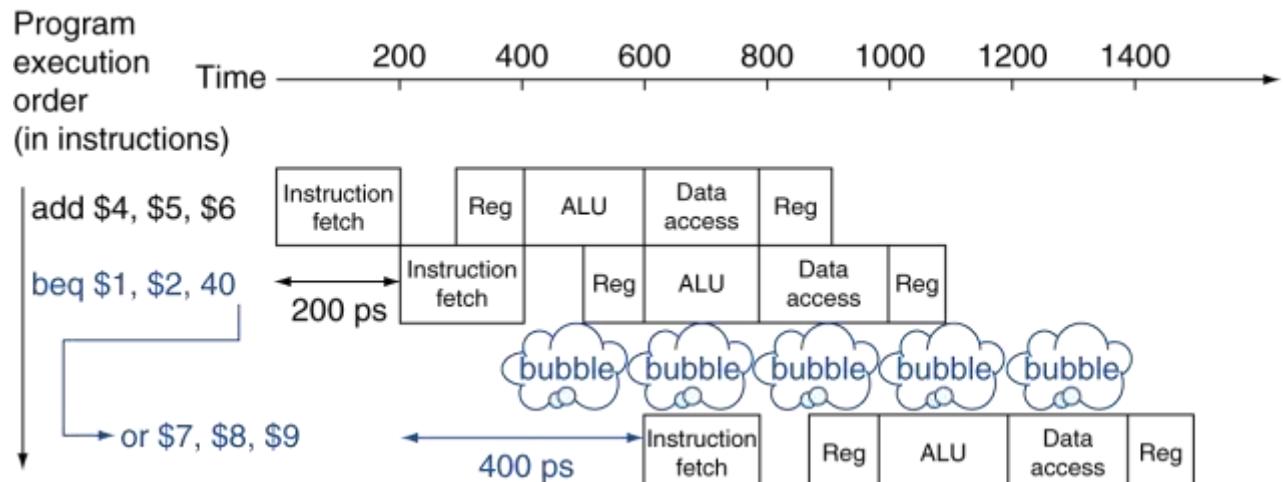
- Longer pipelines can't readily determine branch outcome early
 - Stall penalty becomes unacceptable
- Predict outcome of branch
 - Only stall if prediction is wrong
- In MIPS pipeline
 - Can predict branches not taken
 - Fetch instruction after branch, with no delay

MIPS with Predict Not Taken

Prediction correct



Prediction incorrect



More-Realistic Branch Prediction

- Static branch prediction
 - Based on typical branch behavior
 - Example: loop and if-statement branches
 - Predict backward branches taken
 - Predict forward branches not taken
- Dynamic branch prediction
 - Hardware measures actual branch behavior
 - e.g., **record recent history** of each branch
 - Assume future behavior will continue the trend
 - When wrong, stall while re-fetching, and update history

Pipeline Summary

The BIG Picture

- Pipelining improves performance by increasing instruction throughput
 - Executes multiple instructions in parallel
 - Each instruction has the same latency
- Subject to hazards
 - Structure, data, control
- Instruction set design affects complexity of pipeline implementation

Program Optimization

Instructors:

R. Shathanaa

Today

- **Overview**
- **Generally Useful Optimizations**
 - Code motion/precomputation
 - Strength reduction
 - Sharing of common subexpressions
 - Removing unnecessary procedure calls
- **Optimization Blockers**
 - Procedure calls
 - Memory aliasing
- **Exploiting Instruction-Level Parallelism**
- **Dealing with Conditionals**

Performance Realities

- *There's more to performance than asymptotic complexity*
- Constant factors matter too!
 - Easily see 10:1 performance range depending on how code is written
 - Must optimize at multiple levels:
 - algorithm, data representations, procedures, and loops
- Must understand system to optimize performance
 - How programs are compiled and executed
 - How modern processors + memory systems operate
 - How to measure program performance and identify bottlenecks
 - How to improve performance without destroying code modularity and generality

Optimizing Compilers

- **Provide efficient mapping of program to machine**
 - register allocation
 - code selection and ordering (scheduling)
 - dead code elimination
 - eliminating minor inefficiencies
- **Don't (usually) improve asymptotic efficiency**
 - up to programmer to select best overall algorithm
 - big-O savings are (often) more important than constant factors
 - but constant factors also matter
- **Have difficulty overcoming “optimization blockers”**
 - potential memory aliasing
 - potential procedure side-effects

Generally Useful Optimizations

- Optimizations that you or the compiler should do regardless of processor / compiler

- Code Motion

- Reduce frequency with which computation performed
 - If it will always produce same result
 - Especially moving code out of loop

```
void set_row(double *a, double *b,
    long i, long n)
{
    long j;
    for (j = 0; j < n; j++)
        a[n*i+j] = b[j];
}
```



```
long j;
int ni = n*i;
for (j = 0; j < n; j++)
    a[ni+j] = b[j];
```

Compiler-Generated Code Motion (-O1)

```
void set_row(double *a, double *b,
    long i, long n)
{
    long j;
    for (j = 0; j < n; j++)
        a[n*i+j] = b[j];
}
```

```
long j;
long ni = n*i;
double *rowp = a+ni;
for (j = 0; j < n; j++)
    *rowp++ = b[j];
```

```
set_row:
    testq    %rcx, %rcx          # Test n
    jle     .L1                  # If 0, goto done
    imulq    %rcx, %rdx          # ni = n*i
    leaq     (%rdi,%rdx,8), %rdx # rowp = A + ni*8
    movl    $0, %eax            # j = 0
.L3:
    movsd    (%rsi,%rax,8), %xmm0 # loop:
    movsd    %xmm0, (%rdx,%rax,8) # t = b[j]
    addq    $1, %rax             # M[A+ni*8 + j*8] = t
    cmpq    %rcx, %rax           # j++
    jne     .L3                  # j:n
    .L1:
        rep ; ret               # if !=, goto loop
                                # done:
```

Reduction in Strength

- Replace costly operation with simpler one
- Shift, add instead of multiply or divide

$16 * x \rightarrow x \ll 4$

- Utility machine dependent
- Depends on cost of multiply or divide instruction
 - On Intel Nehalem, integer multiply requires 3 CPU cycles
- Recognize sequence of products

```
for (i = 0; i < n; i++) {  
    int ni = n*i;  
    for (j = 0; j < n; j++)  
        a[ni + j] = b[j];  
}
```



```
int ni = 0;  
for (i = 0; i < n; i++) {  
    for (j = 0; j < n; j++)  
        a[ni + j] = b[j];  
    ni += n;  
}
```

Share Common Subexpressions

- Reuse portions of expressions
- GCC will do this with -O1

```
/* Sum neighbors of i,j */
up = val[(i-1)*n + j];
down = val[(i+1)*n + j];
left = val[i*n      + j-1];
right = val[i*n     + j+1];
sum = up + down + left + right;
```

```
long inj = i*n + j;
up = val[inj - n];
down = val[inj + n];
left = val[inj - 1];
right = val[inj + 1];
sum = up + down + left + right;
```

3 multiplications: $i*n$, $(i-1)*n$, $(i+1)*n$

```
leaq  1(%rsi), %rax # i+1
leaq -1(%rsi), %r8  # i-1
imulq %rcx, %rsi   # i*n
imulq %rcx, %rax   # (i+1)*n
imulq %rcx, %r8    # (i-1)*n
addq  %rdx, %rsi   # i*n+j
addq  %rdx, %rax   # (i+1)*n+j
addq  %rdx, %r8    # (i-1)*n+j
```

1 multiplication: $i*n$

```
imulq  %rcx, %rsi # i*n
addq  %rdx, %rsi # i*n+j
movq  %rsi, %rax # i*n+j
subq  %rcx, %rax # i*n+j-n
leaq  (%rsi,%rcx), %rcx # i*n+j+n
```

Limitations of Optimizing Compilers

- **Operate under fundamental constraint**
 - Must not cause any change in program behavior
- **Behavior that may be obvious to the programmer can be obfuscated by languages and coding styles**
 - e.g., Data ranges may be more limited than variable types suggest
- **Most analysis is performed only within procedures**
 - Whole-program analysis is too expensive in most cases
 - Newer versions of GCC do interprocedural analysis within individual files
 - But, not between code in different files
- **Most analysis is based only on *static* information**
 - Compiler has difficulty anticipating run-time inputs
- **When in doubt, the compiler must be conservative**

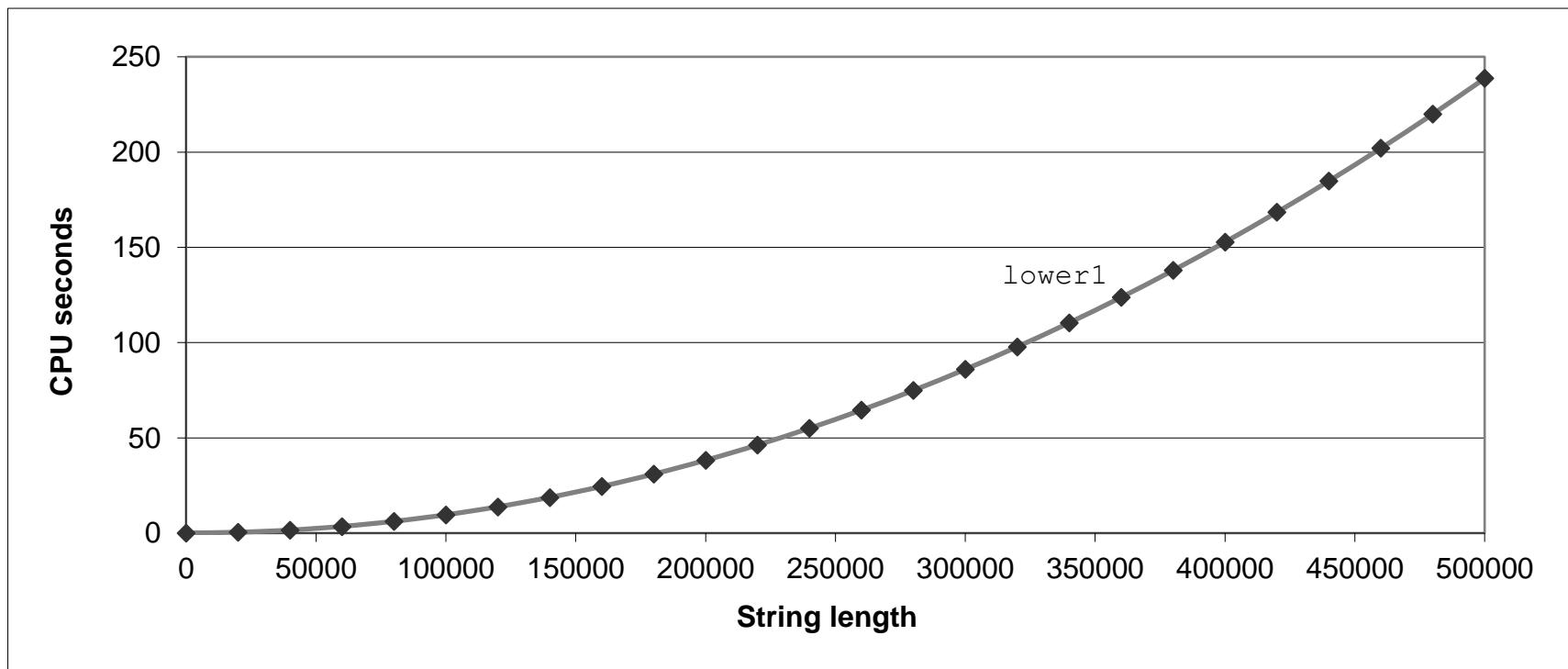
Optimization Blocker #1:

■ Procedure to Convert String to Lower Case

```
void lower(char *s)
{
    size_t i;
    for (i = 0; i < strlen(s); i++)
        if (s[i] >= 'A' && s[i] <= 'Z')
            s[i] -= ('A' - 'a');
}
```

Lower Case Conversion Performance

- Time quadruples when double string length
- Quadratic performance



Convert Loop To Goto Form

```
void lower(char *s)
{
    size_t i = 0;
    if (i >= strlen(s))
        goto done;
loop:
    if (s[i] >= 'A' && s[i] <= 'Z')
        s[i] -= ('A' - 'a');
    i++;
    if (i < strlen(s))
        goto loop;
done:
}
```

- `strlen` executed every iteration

Calling Strlen

```
/* My version of strlen */
size_t strlen(const char *s)
{
    size_t length = 0;
    while (*s != '\0') {
        s++;
        length++;
    }
    return length;
}
```

■ Strlen performance

- Only way to determine length of string is to scan its entire length, looking for null character.

■ Overall performance, string of length N

- N calls to strlen
- Require times N, N-1, N-2, ..., 1
- Overall $O(N^2)$ performance

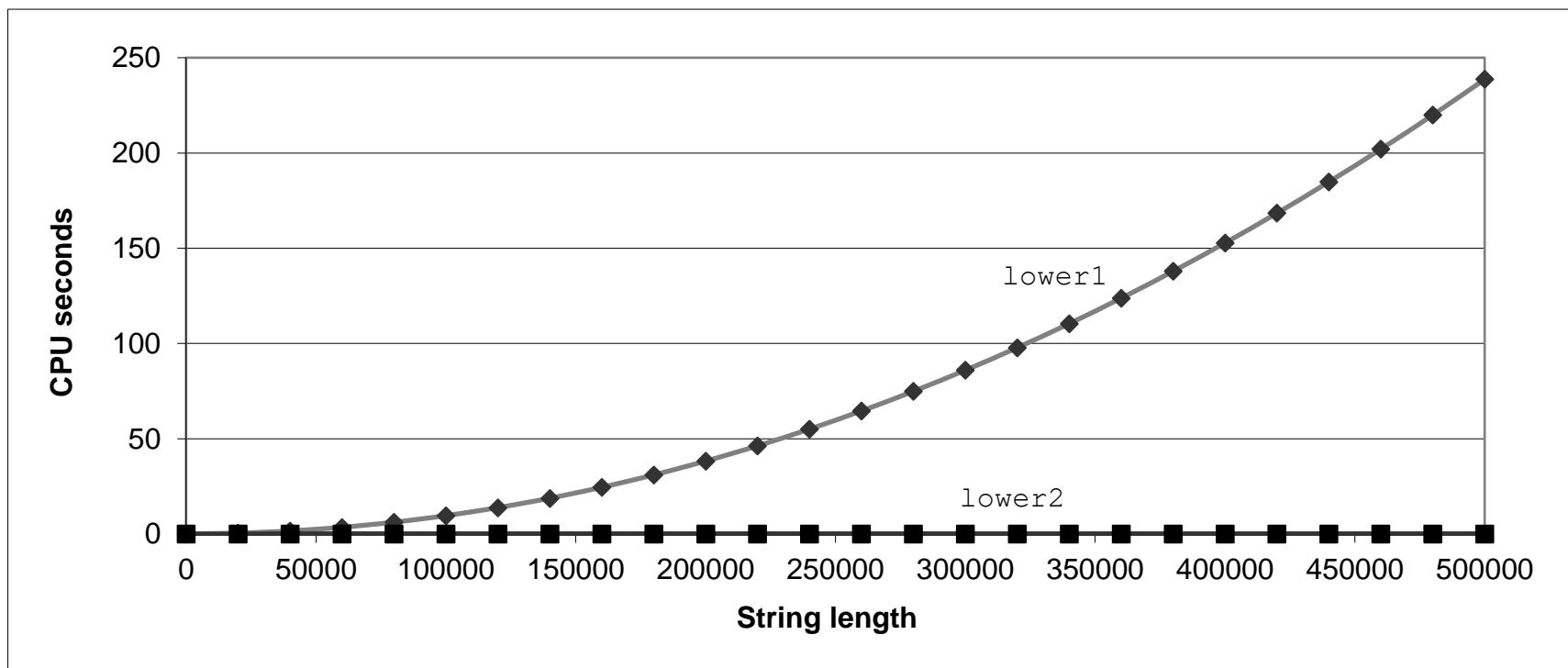
Improving Performance

```
void lower(char *s)
{
    size_t i;
    size_t len = strlen(s);
    for (i = 0; i < len; i++)
        if (s[i] >= 'A' && s[i] <= 'Z')
            s[i] -= ('A' - 'a');
}
```

- Move call to `strlen` outside of loop
- Since result does not change from one iteration to another
- Form of code motion

Lower Case Conversion Performance

- Time doubles when double string length
- Linear performance of lower2



Optimization Blocker: Procedure Calls

■ *Why couldn't compiler move strlen out of inner loop?*

- Procedure may have side effects
 - Alters global state each time called
- Function may not return same value for given arguments
 - Depends on other parts of global state
 - Procedure `lower` could interact with `strlen`

■ Warning:

- Compiler treats procedure call as a black box
- Weak optimizations near them

■ Remedies:

- Use of inline functions
 - GCC does this with `-O1`
 - Within single file
- Do your own code motion

```
size_t lencnt = 0;
size_t strlen(const char *s)
{
    size_t length = 0;
    while (*s != '\0') {
        s++; length++;
    }
    lencnt += length;
    return length;
}
```

Memory Matters

```
/* Sum rows of n X n matrix a
   and store in vector b */
void sum_rows1(double *a, double *b, long n) {
    long i, j;
    for (i = 0; i < n; i++) {
        b[i] = 0;
        for (j = 0; j < n; j++)
            b[i] += a[i*n + j];
    }
}
```

```
# sum_rows1 inner loop
.L4:
    movsd    (%rsi,%rax,8), %xmm0      # FP load
    addsd    (%rdi), %xmm0              # FP add
    movsd    %xmm0, (%rsi,%rax,8)       # FP store
    addq    $8, %rdi
    cmpq    %rcx, %rdi
    jne     .L4
```

- Code updates $b[i]$ on every iteration
- Why couldn't compiler optimize this away?

```
1 void twiddle1(long *xp, long *yp)
2 {
3     *xp += *yp;
4     *xp += *yp;
5 }
6
7 void twiddle2(long *xp, long *yp)
8 {
9     *xp += 2* *yp;
10 }
```

At first glance, both procedures seem to have identical behavior. They both add twice the value stored at the location designated by pointer yp to that designated by pointer xp. On the other hand, function twiddle2 is more efficient. It requires only three memory references (read *xp, read *yp, write *xp), whereas twiddle1 requires six (two reads of *xp, two reads of *yp, and two writes of *xp). Hence, if a compiler is given procedure twiddle1 to compile, one might think it could generate more efficient code based on the computations performed by twiddle2.

Consider, however, the case in which xp and yp are equal. Then function twiddle1 will perform the following computations:

```
3     *xp += *xp; /* Double value at xp */
4     *xp += *xp; /* Double value at xp */
```

The result will be that the value at xp will be increased by a factor of 4. On the other hand, function twiddle2 will perform the following computation:

```
9     *xp += 2* *xp; /* Triple value at xp */
```

Memory Aliasing

```
/* Sum rows is of n X n matrix a
   and store in vector b */
void sum_rows1(double *a, double *b, long n) {
    long i, j;
    for (i = 0; i < n; i++) {
        b[i] = 0;
        for (j = 0; j < n; j++)
            b[i] += a[i*n + j];
    }
}
```

```
double A[9] =
{ 0, 1, 2,
  4, 8, 16,
  32, 64, 128};

double B[3] = A+3;

sum_rows1(A, B, 3);
```

Value of B:

init: [4, 8, 16]

i = 0: [3, 8, 16]

i = 1: [3, 28, 16]

i = 2: [3, 28, 224]

- Code updates `b[i]` on every iteration
- Must consider possibility that these updates will affect program behavior

Removing Aliasing

```
/* Sum rows is of n X n matrix a
   and store in vector b  */
void sum_rows2(double *a, double *b, long n) {
    long i, j;
    for (i = 0; i < n; i++) {
        double val = 0;
        for (j = 0; j < n; j++)
            val += a[i*n + j];
        b[i] = val;
    }
}
```

Optimization Blocker: Memory Aliasing

■ Aliasing

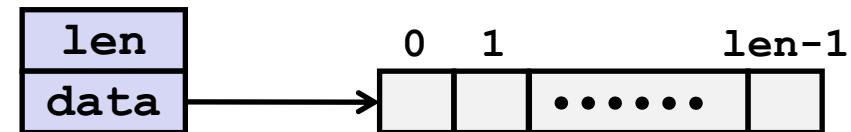
- Two different memory references specify single location
- Easy to have happen in C
 - Since allowed to do address arithmetic
 - Direct access to storage structures
- Get in habit of introducing local variables
 - Accumulating within loops
 - **Your way of telling compiler not to check for aliasing**

Exploiting Instruction-Level Parallelism

- **Need general understanding of modern processor design**
 - Hardware can execute multiple instructions in parallel
- **Performance limited by data dependencies**
- **Simple transformations can yield dramatic performance improvement**
 - Compilers often cannot make these transformations
 - Lack of associativity and distributivity in floating-point arithmetic

Benchmark Example: Data Type for Vectors

```
/* data structure for vectors */
typedef struct{
    size_t len;
    data_t *data;
} vec;
```



■ Data Types

- Use different declarations for `data_t`
- `int`
- `long`
- `float`
- `double`

```
/* retrieve vector element
   and store at val */
int get_vec_element
    (*vec v, size_t idx, data_t *val)
{
    if (idx >= v->len)
        return 0;
    *val = v->data[idx];
    return 1;
}
```

Benchmark Computation

```
void combinel(vec_ptr v, data_t *dest)
{
    long int i;
    *dest = IDENT;
    for (i = 0; i < vec_length(v); i++) {
        data_t val;
        get_vec_element(v, i, &val);
        *dest = *dest OP val;
    }
}
```

Compute sum or product of vector elements

■ Data Types

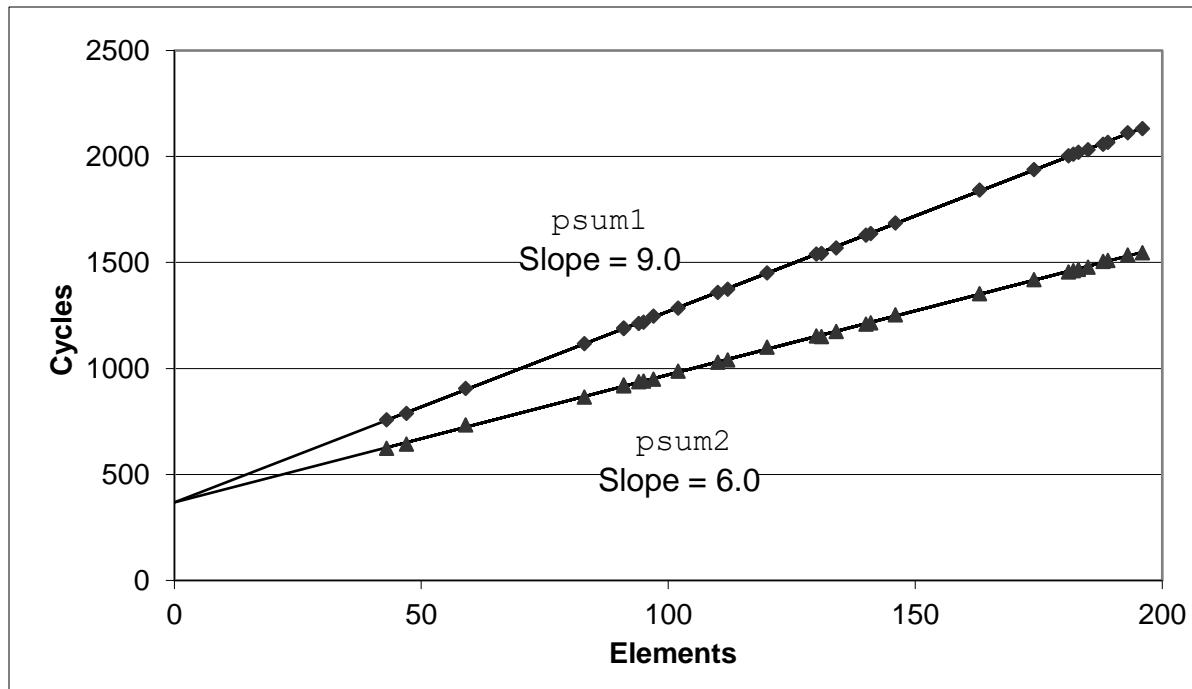
- Use different declarations for `data_t`
- `int`
- `long`
- `float`
- `double`

■ Operations

- Use different definitions of `OP` and `IDENT`
 - `+` / `0`
 - `*` / `1`

Cycles Per Element (CPE)

- Convenient way to express performance of program that operates on vectors or lists
- Length = n
- In our case: **CPE = cycles per OP**
- $T = \text{CPE} * n + \text{Overhead}$
 - CPE is slope of line



Benchmark Performance

```

void combine1(vec_ptr v, data_t *dest)
{
    long int i;
    *dest = IDENT;
    for (i = 0; i < vec_length(v); i++) {
        data_t val;
        get_vec_element(v, i, &val);
        *dest = *dest OP val;
    }
}

```

Compute sum or product of vector elements

Method	Integer		Double FP	
Operation	Add	Mult	Add	Mult
Combine1 unoptimized	22.68	20.02	19.98	20.18
Combine1 -O1	10.12	10.12	10.17	11.14

O1. Eliminating Loop Inefficiencies

Code Motion

```

1  /* Implementation with maximum use of data abstraction */
2  void combine1(vec_ptr v, data_t *dest)
3  {
4      long i;
5
6      *dest = IDENT;
7      for (i = 0; i < vec_length(v); i++) {
8          data_t val;
9          get_vec_element(v, i, &val);
10         *dest = *dest OP val;
11     }
12 }

1  /* Move call to vec_length out of loop */
2  void combine2(vec_ptr v, data_t *dest)
3  {
4      long i;
5      long length = vec_length(v);
6
7      *dest = IDENT;
8      for (i = 0; i < length; i++) {
9          data_t val;
10         get_vec_element(v, i, &val);
11         *dest = *dest OP val;
12     }

```

Function	Page	Method	Integer		Floating point	
			+	*	+	*
combine1	543	Abstract -01	10.12	10.12	10.17	11.14
combine2	545	Move vec_length	7.02	9.03	9.02	11.03

O2. Reducing Procedure Calls

code/opt/vec.c

```
1 data_t *get_vec_start(vec_ptr v)
2 {
3     return v->data;
4 }
```

code/opt/vec.c

```
1 /* Direct access to vector data */
2 void combine3(vec_ptr v, data_t *dest)
3 {
4     long i;
5     long length = vec_length(v);
6     data_t *data = get_vec_start(v);
7
8     *dest = IDENT;
9     for (i = 0; i < length; i++) {
10         *dest = *dest OP data[i];
11     }
12 }
```

Function	Page	Method	Integer		Floating point	
			+	*	+	*
combine2	545	Move <code>vec_length</code>	7.02	9.03	9.02	11.03
combine3	549	Direct data access	7.17	9.02	9.02	11.03

Surprisingly, there is no apparent performance improvement. Indeed, the performance for integer sum has gotten slightly worse.

Evidently, other operations in the inner loop are forming a bottleneck that limits the performance more than the call to `get_vec_element`. We will return to this function later and see why the repeated bounds checking by `combine2` does not incur a performance penalty.

For now, we can view this transformation as one of a series of steps that will ultimately lead to greatly improved performance.

O3. Eliminating Unneeded Memory References

```
1  /* Accumulate result in local variable */
2  void combine4(vec_ptr v, data_t *dest)
3  {
4      long i;
5      long length = vec_length(v);
6      data_t *data = get_vec_start(v);
7      data_t acc = IDENT;
8
9      for (i = 0; i < length; i++) {
10         acc = acc OP data[i];
11     }
12     *dest = acc;
13 }
```

Function	Page	Method	Integer		Floating point	
			+	*	+	*
combine3	549	Direct data access	7.17	9.02	9.02	11.03
combine4	551	Accumulate in temporary	1.27	3.01	3.01	5.01

Basic Optimizations

```
void combine4(vec_ptr v, data_t *dest)
{
    long i;
    long length = vec_length(v);
    data_t *d = get_vec_start(v);
    data_t t = IDENT;
    for (i = 0; i < length; i++)
        t = t OP d[i];
    *dest = t;
}
```

- Move `vec_length` out of loop
- Avoid bounds check on each cycle
- Accumulate in temporary

Effect of Basic Optimizations

```
void combine4(vec_ptr v, data_t *dest)
{
    long i;
    long length = vec_length(v);
    data_t *d = get_vec_start(v);
    data_t t = IDENT;
    for (i = 0; i < length; i++)
        t = t OP d[i];
    *dest = t;
}
```

Method	Integer		Double FP	
Operation	Add	Mult	Add	Mult
Combine1 -O1	10.12	10.12	10.17	11.14
Combine4	1.27	3.01	3.01	5.01

- Eliminates sources of overhead in loop

Program Optimization

Instructors:

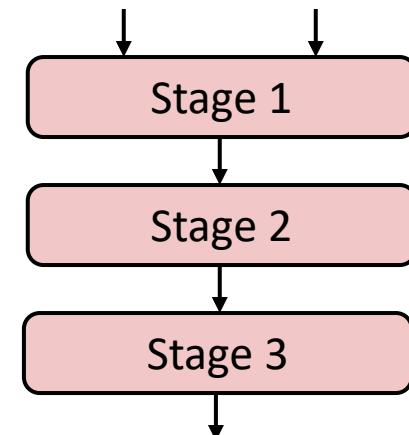
R. Shathanaa

Superscalar Processor

- **Definition:** A superscalar processor can issue and execute *multiple instructions in one cycle*. The instructions are retrieved from a sequential instruction stream and are usually scheduled dynamically.
- **Benefit:** without programming effort, superscalar processor can take advantage of the *instruction level parallelism* that most programs have
- Most modern CPUs are superscalar.
- Intel: since Pentium (1993)

Pipelined Functional Units

```
long mult_eg(long a, long b, long c) {
    long p1 = a*b;
    long p2 = a*c;
    long p3 = p1 * p2;
    return p3;
}
```



	Time							
	1	2	3	4	5	6	7	
Stage 1	a*b	a*c			p1*p2			
Stage 2		a*b	a*c			p1*p2		
Stage 3			a*b	a*c				p1*p2

- Divide computation into stages
- Pass partial computations from stage to stage
- Stage i can start on new computation once values passed to $i+1$
- E.g., complete 3 multiplications in 7 cycles, even though each requires 3 cycles

Haswell CPU

- 8 Total Functional Units
- **Multiple instructions can execute in parallel**

2 load, with address computation
1 store, with address computation
4 integer
2 FP multiply
1 FP add
1 FP divide

- **Some instructions take > 1 cycle, but can be pipelined**

<i>Instruction</i>	<i>Latency</i>	<i>Cycles/Issue</i>
Load / Store	4	1
Integer Multiply	3	1
Integer/Long Divide	3-30	3-30
Single/Double FP Multiply	5	1
Single/Double FP Add	3	1
Single/Double FP Divide	3-15	3-15

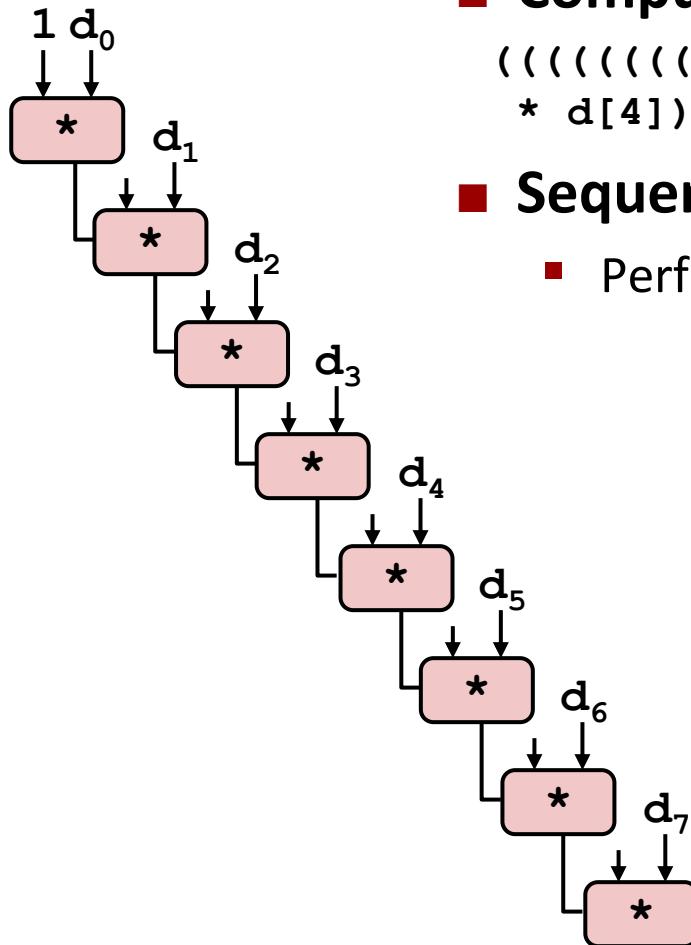
x86-64 Compilation of Combine4

■ Inner Loop (Case: Integer Multiply)

```
void combine4(vec_ptr v, data_t *dest)
{
    long i;
    long length = vec_length(v);
    data_t *d = get_vec_start(v);
    data_t t = IDENT;
    for (i = 0; i < length; i++)
        t = t OP d[i];
    *dest = t;
}
```

Method	Integer		Double FP	
Operation	Add	Mult	Add	Mult
Combine4	1.27	3.01	3.01	5.01
Latency Bound	1.00	3.00	3.00	5.00

Combine4 = Serial Computation (OP = *)



■ Computation (length=8)

```
(((((1 * d[0]) * d[1]) * d[2]) * d[3])  
 * d[4]) * d[5]) * d[6]) * d[7])
```

■ Sequential dependence

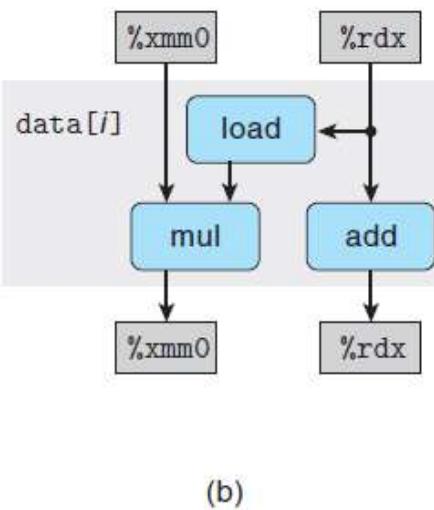
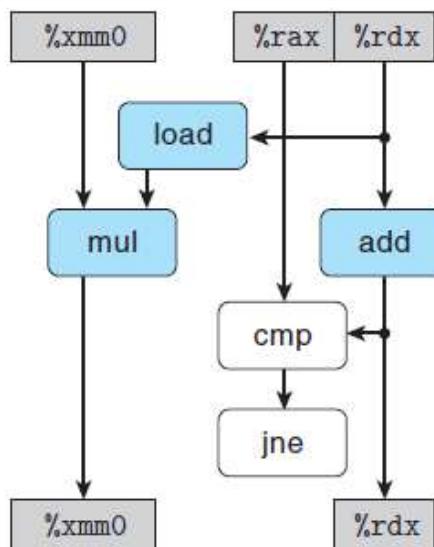
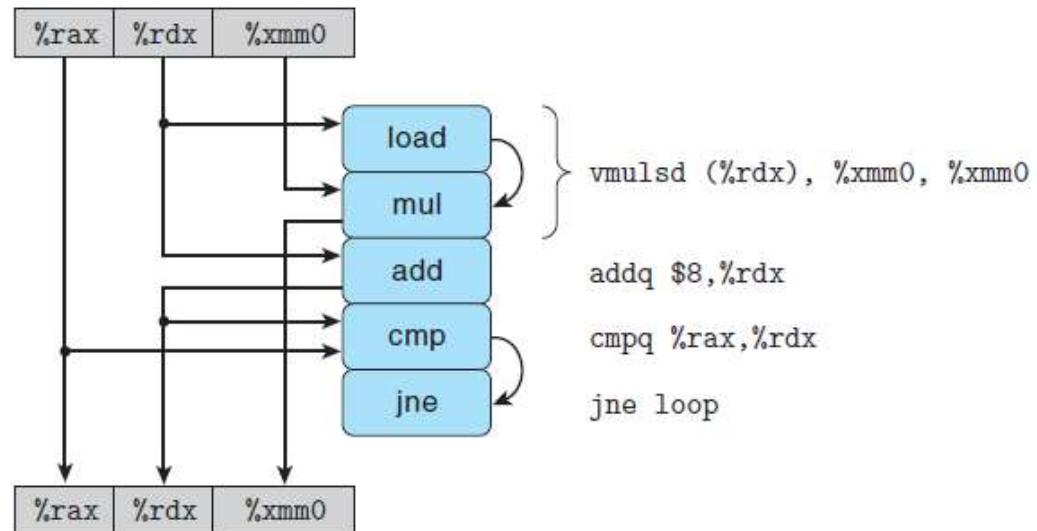
- Performance: determined by latency of OP

Data-Flow Graphs

- A way to visualize how the data dependencies in a program dictate its performance.

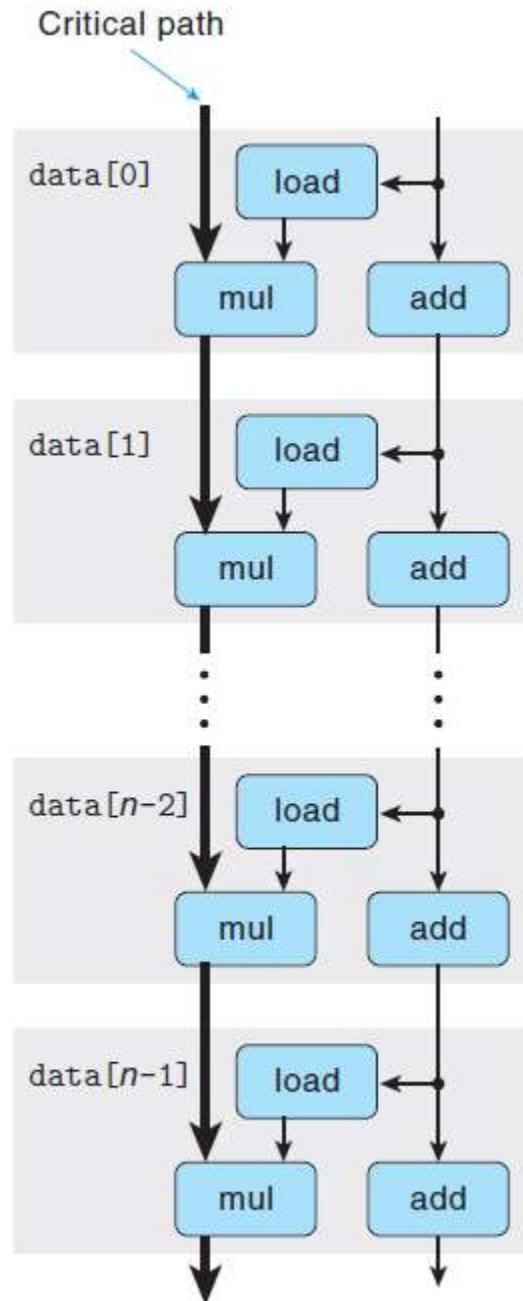
```
Inner loop of combine4.  data_t = double, OP = *
acc in %xmm0, data+i in %rdx, data+length in %rax
1 .L25:                                loop:
2   vmulsd  (%rdx), %xmm0, %xmm0      Multiply acc by data[i]
3   addq    $8, %rdx                   Increment data+i
4   cmpq    %rax, %rdx                 Compare to data+length
5   jne     .L25                      If !=, goto loop
```

Data-Flow Graph



Data-Flow Graphs

- For a code segment forming a loop, we can classify the registers that are accessed into four categories:
 - *Read-only*. used as source values, not modified within the loop. **rax**
 - *Write-only*. used as the destinations of data-movement operations. **none**
 - *Local*. These are updated and used within the loop, but there is no dependency from one iteration to another. **condition flags**
 - *Loop*. These are used both as source values and as destinations for the loop, with the value generated in one iteration being used in another. **rdx, xmm0**
 - **Assume mul – 5 clk and add - 1 clk**



Loop Unrolling (2x1)

```
void unroll2a_combine(vec_ptr v, data_t *dest)
{
    long length = vec_length(v);
    long limit = length-1;
    data_t *d = get_vec_start(v);
    data_t x = IDENT;
    long i;
    /* Combine 2 elements at a time */
    for (i = 0; i < limit; i+=2) {
        x = (x OP d[i]) OP d[i+1];
    }
    /* Finish any remaining elements */
    for (; i < length; i++) {
        x = x OP d[i];
    }
    *dest = x;
}
```

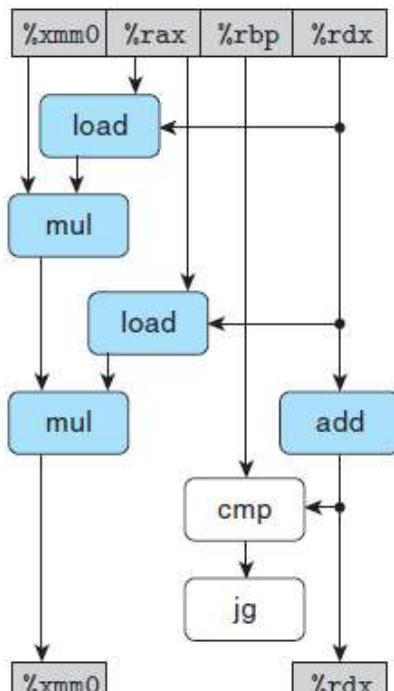
- Perform 2x more useful work per iteration

Effect of Loop Unrolling

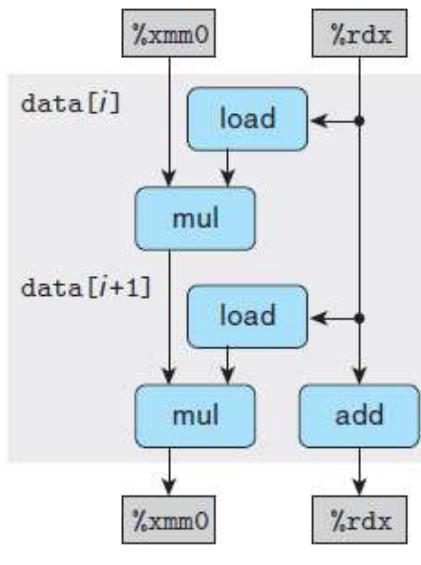
Method	Integer		Double FP	
Operation	Add	Mult	Add	Mult
Combine4	1.27	3.01	3.01	5.01
Unroll 2x1	1.01	3.01	3.01	5.01
Latency Bound	1.00	3.00	3.00	5.00

- Helps integer add
 - Achieves latency bound
- Others don't improve. *Why?*
 - Still sequential dependency

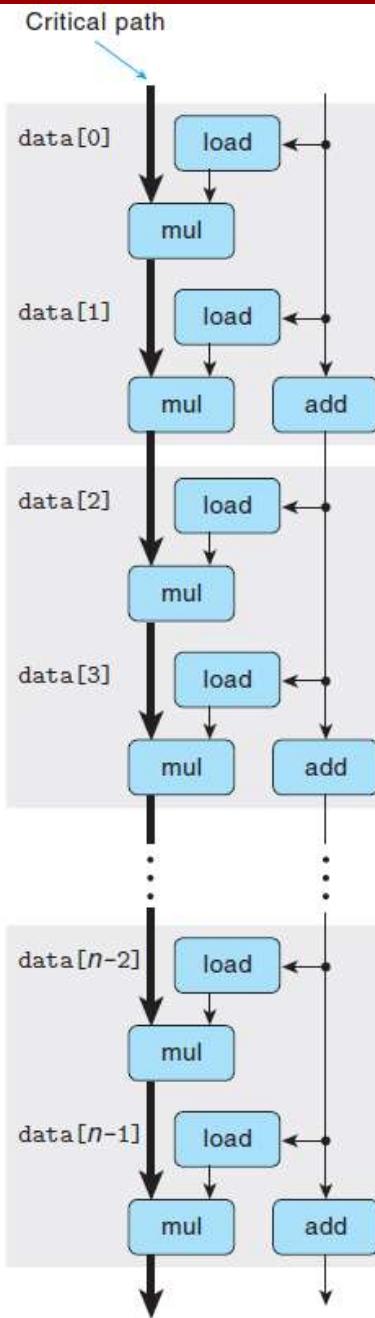
```
x = (x OP d[i]) OP d[i+1];
```



(a)



(b)



Loop Unrolling with Reassociation (2x1a)

```
void unroll2aa_combine(vec_ptr v, data_t *dest)
{
    long length = vec_length(v);
    long limit = length-1;
    data_t *d = get_vec_start(v);
    data_t x = IDENT;
    long i;
    /* Combine 2 elements at a time */
    for (i = 0; i < limit; i+=2) {
        x = x OP (d[i] OP d[i+1]);
    }
    /* Finish any remaining elements */
    for (; i < length; i++) {
        x = x OP d[i];
    }
    *dest = x;
}
```

Compare to before

$x = (x \text{ OP } d[i]) \text{ OP } d[i+1];$

- Can this change the result of the computation?
- Yes, for FP. *Why?*

Effect of Reassociation

Method	Integer		Double FP	
Operation	Add	Mult	Add	Mult
Combine4	1.27	3.01	3.01	5.01
Unroll 2x1	1.01	3.01	3.01	5.01
Unroll 2x1a	1.01	1.51	1.51	2.51
Latency Bound	1.00	3.00	3.00	5.00
Throughput Bound	0.50	1.00	1.00	0.50

■ Nearly 2x speedup for Int *, FP +, FP *

- Reason: Breaks sequential dependency

```
x = x OP (d[i] OP d[i+1]);
```

2 func. units for FP *
2 func. units for load

4 func. units for int +
2 func. units for load

- Why is that? (next slide)

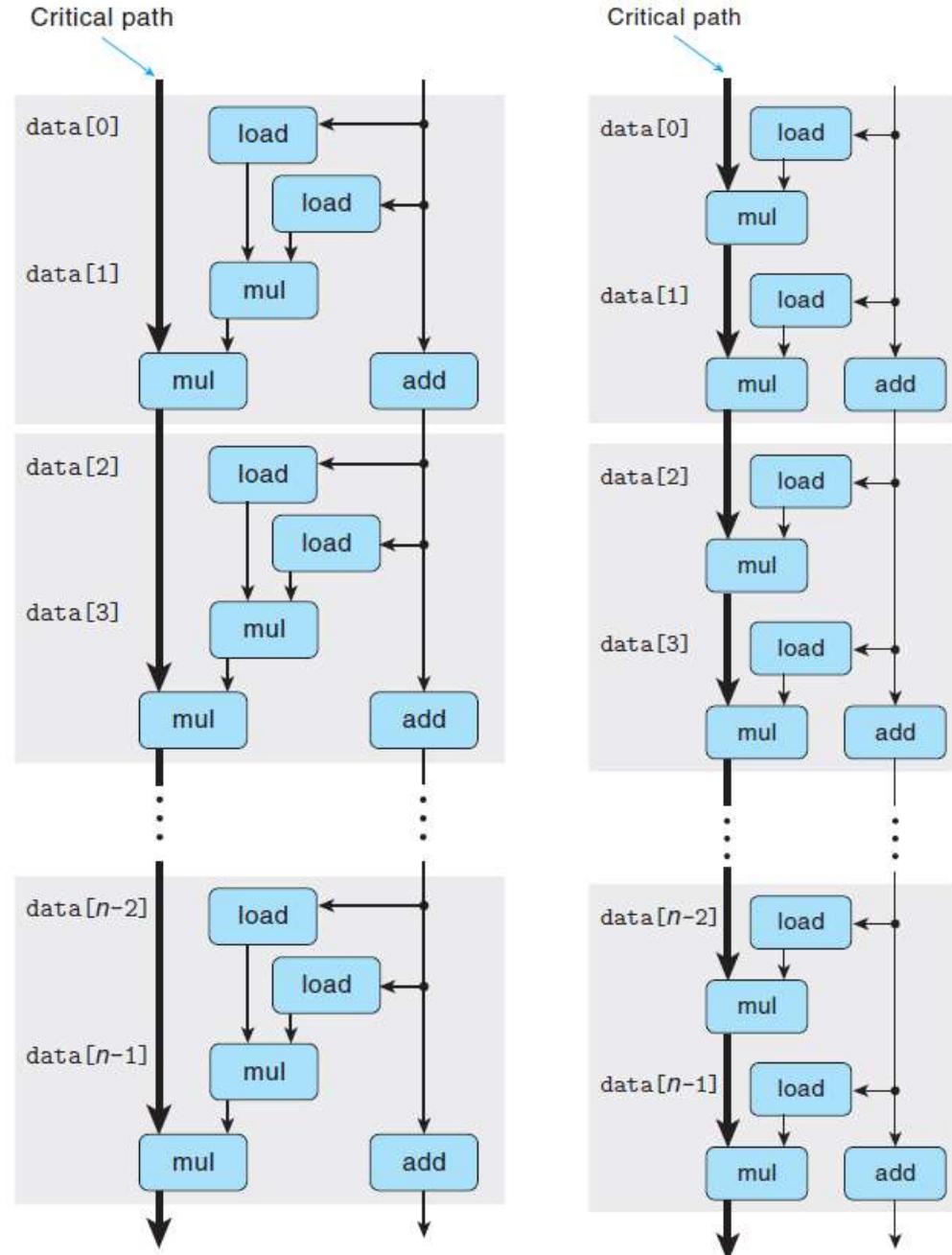
Reassociated Computation

■ What changed:

- Ops in the next iteration can be started early (no dependency)

■ Overall Performance

- N elements, D cycles latency/op
- $(N/2+1)*D$ cycles:
CPE = D/2



Loop Unrolling with Separate Accumulators (2x2)

```
void unroll2a_combine(vec_ptr v, data_t *dest)
{
    long length = vec_length(v);
    long limit = length-1;
    data_t *d = get_vec_start(v);
    data_t x0 = IDENT;
    data_t x1 = IDENT;
    long i;
    /* Combine 2 elements at a time */
    for (i = 0; i < limit; i+=2) {
        x0 = x0 OP d[i];
        x1 = x1 OP d[i+1];
    }
    /* Finish any remaining elements */
    for (; i < length; i++) {
        x0 = x0 OP d[i];
    }
    *dest = x0 OP x1;
}
```

■ Different form of reassociation

Effect of Separate Accumulators

Method	Integer		Double FP	
Operation	Add	Mult	Add	Mult
Combine4	1.27	3.01	3.01	5.01
Unroll 2x1	1.01	3.01	3.01	5.01
Unroll 2x1a	1.01	1.51	1.51	2.51
Unroll 2x2	0.81	1.51	1.51	2.51
Latency Bound	1.00	3.00	3.00	5.00
Throughput Bound	0.50	1.00	1.00	0.50

- Int + makes use of two load units

```
x0 = x0 OP d[i];
x1 = x1 OP d[i+1];
```

- 2x speedup (over unroll2) for Int *, FP +, FP *

Separate Accumulators

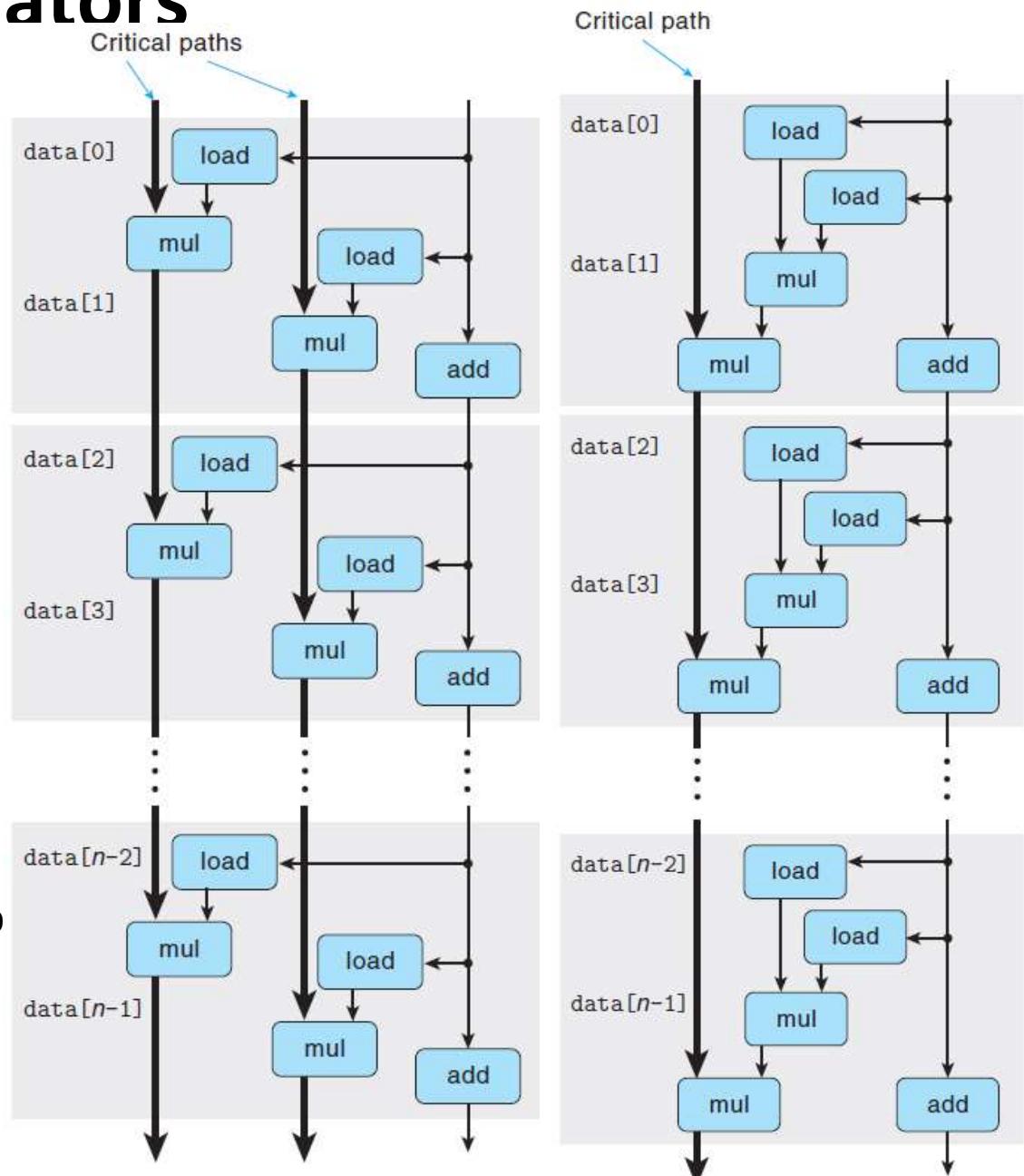
What Now?

What changed:

- Two independent “streams” of operations

Overall Performance

- N elements, D cycles latency/op
- Should be $(N/2+1)*D$ cycles:
CPE = D/2
- CPE matches prediction!



Unrolling & Accumulating

■ Idea

- Can unroll to any degree L
- Can accumulate K results in parallel
- L must be multiple of K

■ Limitations

- Diminishing returns
 - Cannot go beyond throughput limitations of execution units
- Large overhead for short lengths
 - Finish off iterations sequentially

Achievable Performance

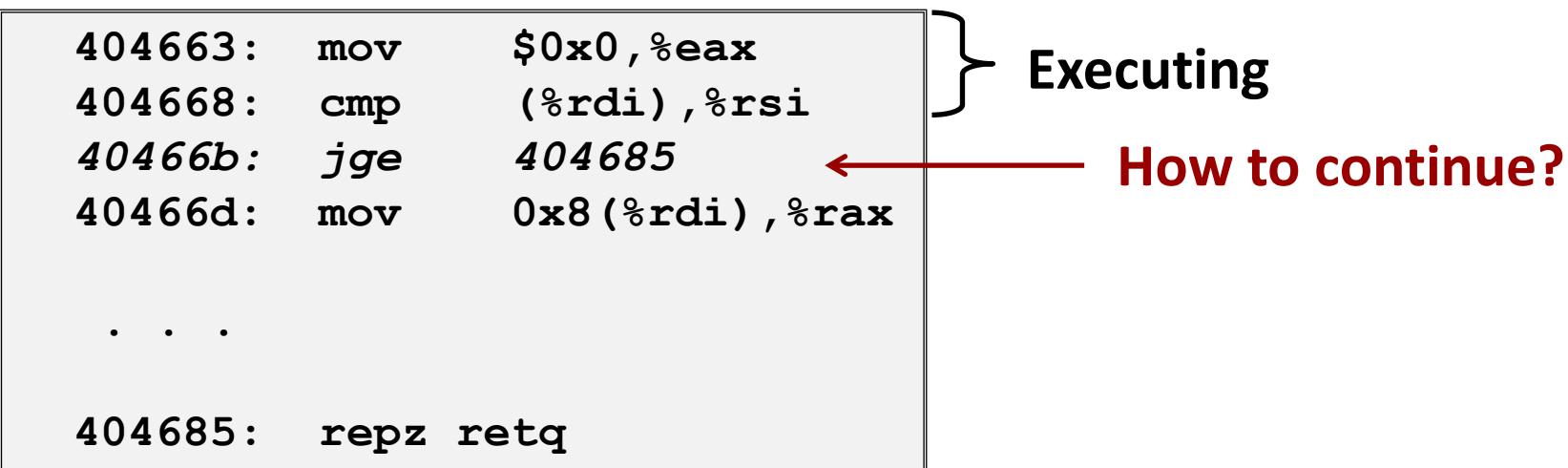
Method	Integer		Double FP	
Operation	Add	Mult	Add	Mult
Best	0.54	1.01	1.01	0.52
Latency Bound	1.00	3.00	3.00	5.00
Throughput Bound	0.50	1.00	1.00	0.50

- Limited only by throughput of functional units
- Up to 42X improvement over original, unoptimized code

What About Branches?

■ Challenge

- Instruction Control Unit must work well ahead of Execution Unit to generate enough operations to keep EU busy



- When encounters conditional branch, cannot reliably determine where to continue fetching

Branch Outcomes

- When encounter conditional branch, cannot determine where to continue fetching
 - Branch Taken: Transfer control to branch target
 - Branch Not-Taken: Continue with next instruction in sequence
- Cannot resolve until outcome determined by branch/integer unit

```
404663:    mov      $0x0,%eax
404668:    cmp      (%rdi),%rsi
40466b:    jge      404685
40466d:    mov      0x8(%rdi),%rax
.
.
.
404685:    repz    retq
```

Branch Not-Taken

Branch Taken

Branch Prediction

■ Idea

- Guess which way branch will go
- Begin executing instructions at predicted position
 - But don't actually modify register or memory data

```
404663:  mov    $0x0,%eax
404668:  cmp    (%rdi),%rsi
40466b:  jge    404685
40466d:  mov    0x8(%rdi),%rax
```

Predict Taken

. . .

```
404685:  repz  retq
```

} Begin
Execution

Branch Prediction Through Loop

```

401029:  vmulsd (%rdx), %xmm0, %xmm0
40102d:  add    $0x8, %rdx
401031:  cmp    %rax, %rdx
401034:  jne    401029
          i = 98
  
```

Assume
vector length = 100

```

401029:  vmulsd (%rdx), %xmm0, %xmm0
40102d:  add    $0x8, %rdx
401031:  cmp    %rax, %rdx
401034:  jne    401029
          i = 99
  
```

Predict Taken (OK)

```

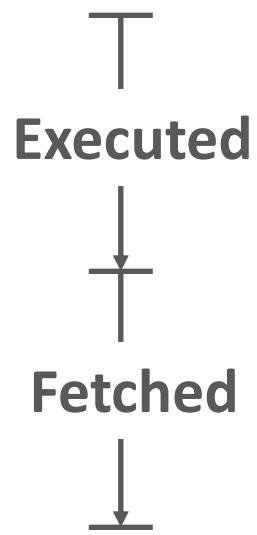
401029:  vmulsd (%rdx), %xmm0, %xmm0
40102d:  add    $0x8, %rdx
401031:  cmp    %rax, %rdx
401034:  jne    401029
          i = 100
  
```

Predict Taken
(Oops)

```

401029:  vmulsd (%rdx), %xmm0, %xmm0
40102d:  add    $0x8, %rdx
401031:  cmp    %rax, %rdx
401034:  jne    401029
          i = 101
  
```

Read
invalid
location



Branch Misprediction Invalidation

```
401029:  vmulsd (%rdx),%xmm0,%xmm0  
40102d:  add    $0x8,%rdx  
401031:  cmp    %rax,%rdx  
401034:  jne    401029      i = 98
```

Assume
vector length = 100

```
401029:  vmulsd (%rdx),%xmm0,%xmm0  
40102d:  add    $0x8,%rdx  
401031:  cmp    %rax,%rdx  
401034:  jne    401029      i = 99
```

Predict Taken (OK)

```
401029:  vmulsd (%rdx),%xmm0,%xmm0  
40102d:  add    $0x8,%rdx  
401031:  cmp    %rax,%rdx  
401034:  jne    401029      i = 100
```

Predict Taken
(Oops)

```
401029:  vmulsd (%rdx),%xmm0,%xmm0  
40102d:  add    $0x8,%rdx  
401031:  cmp    %rax,%rdx  
401034:  jne    401029      i = 101
```

Invalidate

Branch Misprediction Recovery

```
401029:  vmulsd (%rdx), %xmm0, %xmm0  
40102d:  add    $0x8,%rdx  
401031:  cmp    %rax,%rdx  
401034:  jne    401029  
401036:  jmp    401040  
...  
401040:  vmovsd %xmm0, (%r12)
```

i = 99

Definitely not taken

Reload
Pipeline

■ Performance Cost

- Multiple clock cycles on modern processor
- Can be a major performance limiter

Getting High Performance

- **Good compiler and flags**
- **Don't do anything stupid**
 - Watch out for hidden algorithmic inefficiencies
 - Write compiler-friendly code
 - Watch out for optimization blockers:
procedure calls & memory references
 - Look carefully at innermost loops (where most work is done)
- **Tune code for machine**
 - Exploit instruction-level parallelism
 - Avoid unpredictable branches
 - Make code cache friendly

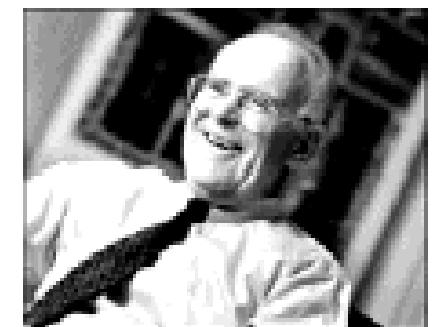
Introduction to Parallel Processing Architecture

Introduction

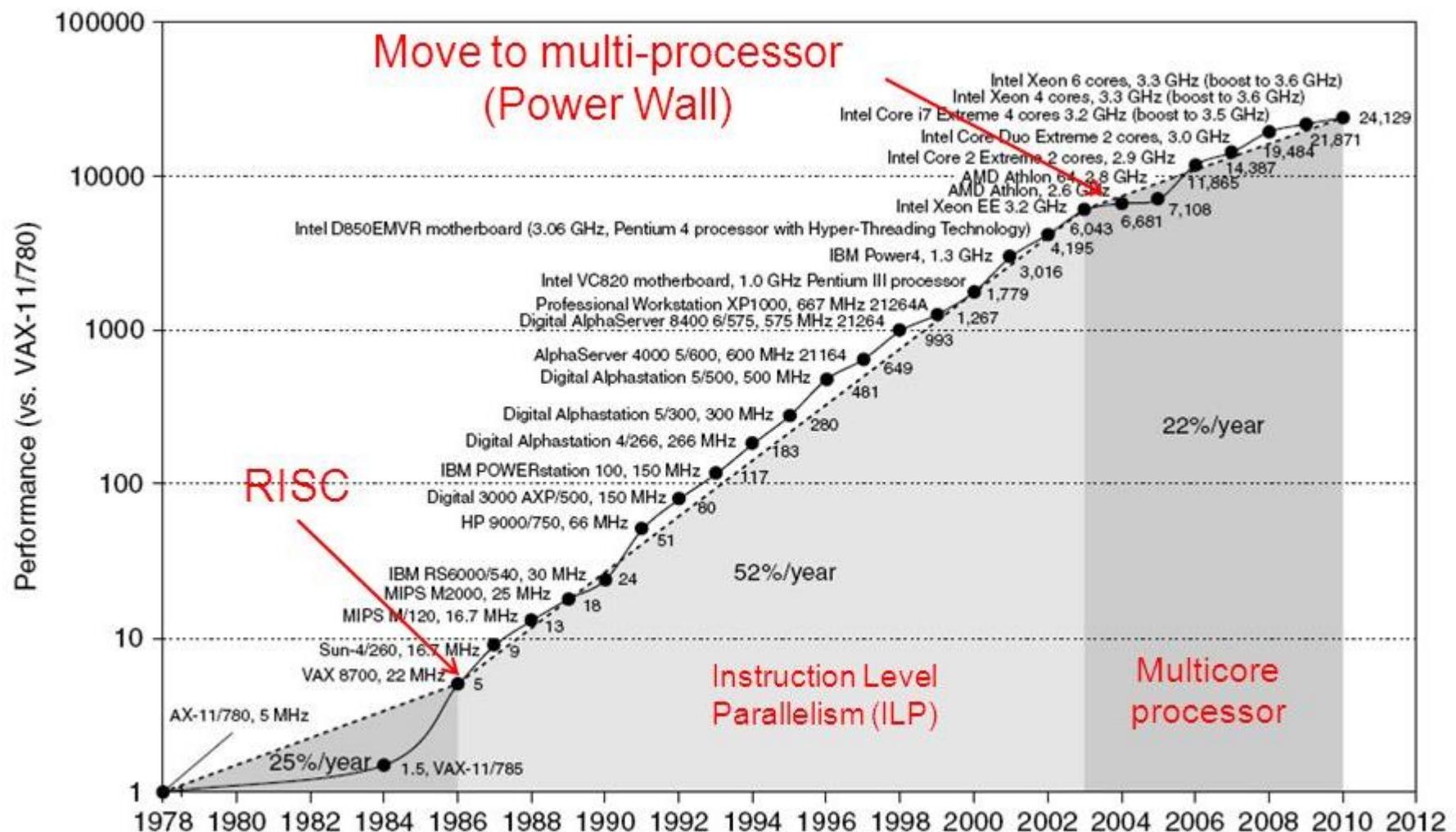
- Performance improvements:
- Improvements in semiconductor technology
 - Feature size, clock speed
- Improvements in computer architectures
 - Enabled by HLL compilers, UNIX
 - Lead to RISC architectures
- Together have enabled:
 - Lightweight computers
 - Productivity-based managed/interpreted programming languages
 - SaaS, Virtualization, Cloud
- Applications evolution:
 - Speech, sound, images, video, “augmented/extended reality”, “big data”

Moore's Law

- **Gordon Moore (co-founder of Intel) predicted in 1965 that the transistor density of semiconductor chips would double roughly every 18 months.**



Single Processor Performance



- So, we need Parallel Computing!

• **Current Trends in Architecture**

- Cannot continue to leverage Instruction-Level parallelism (ILP)
- Single processor performance improvement ended in 2003
- New model for performance: **Parallelism**
- These require explicit restructuring of the application

Level of Parallelism

- Bit level parallelism: 1970 to ~1985
 - 4 bits, 8 bit, 16 bit, 32 bit microprocessors
- Instruction level parallelism: ~1985 - today
 - Pipelining
 - Superscalar
 - VLIW
 - Out-of-order execution / Dynamic Instr. Scheduling
- Process level or thread level parallelism
 - Servers are parallel
 - Desktop dual processor PCs
 - Multicore architectures (CPUs, GPUs)

Parallelism

- Classes of parallelism in applications:
- Data-Level Parallelism (DLP)
- Task-Level Parallelism (TLP)
- Methods of exploiting architectural parallelism by hardware:
- Instruction-Level Parallelism (ILP)
- Vector architectures/Graphic Processor Units (GPUs)
- Thread-Level Parallelism
- Request-Level Parallelism

CPU Design

Reduced instruction set computer (RISC) architecture

- RISC-CPUs developed for high-performance computers and now used broadly.
- It increases the arithmetic speed of the CPU by decreasing the number of instructions the CPU.

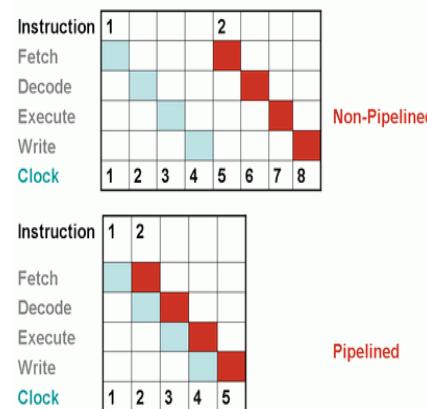
Complex instruction set computer (CISC), architecture

- Pipelining: concurrent preparation of several instructions that are then executed successively.
- High level compilers, to improve performance.

CPU Design

Reduced instruction set computer (RISC) architecture

- **Pipelining**
- An instruction can be executed while the next one is being decoded and the next one is being fetched.
- It overlaps various stages of instruction execution to achieve performance
- The speed of a pipeline is eventually **limited by the slowest stage**.



Multiple pipelines.

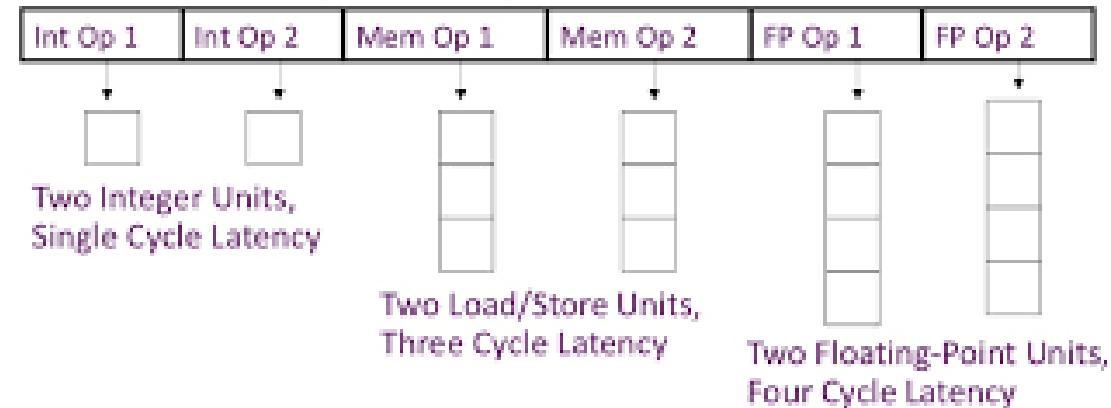
CPU Design

Very Long Instruction Word (VLIW) Processors

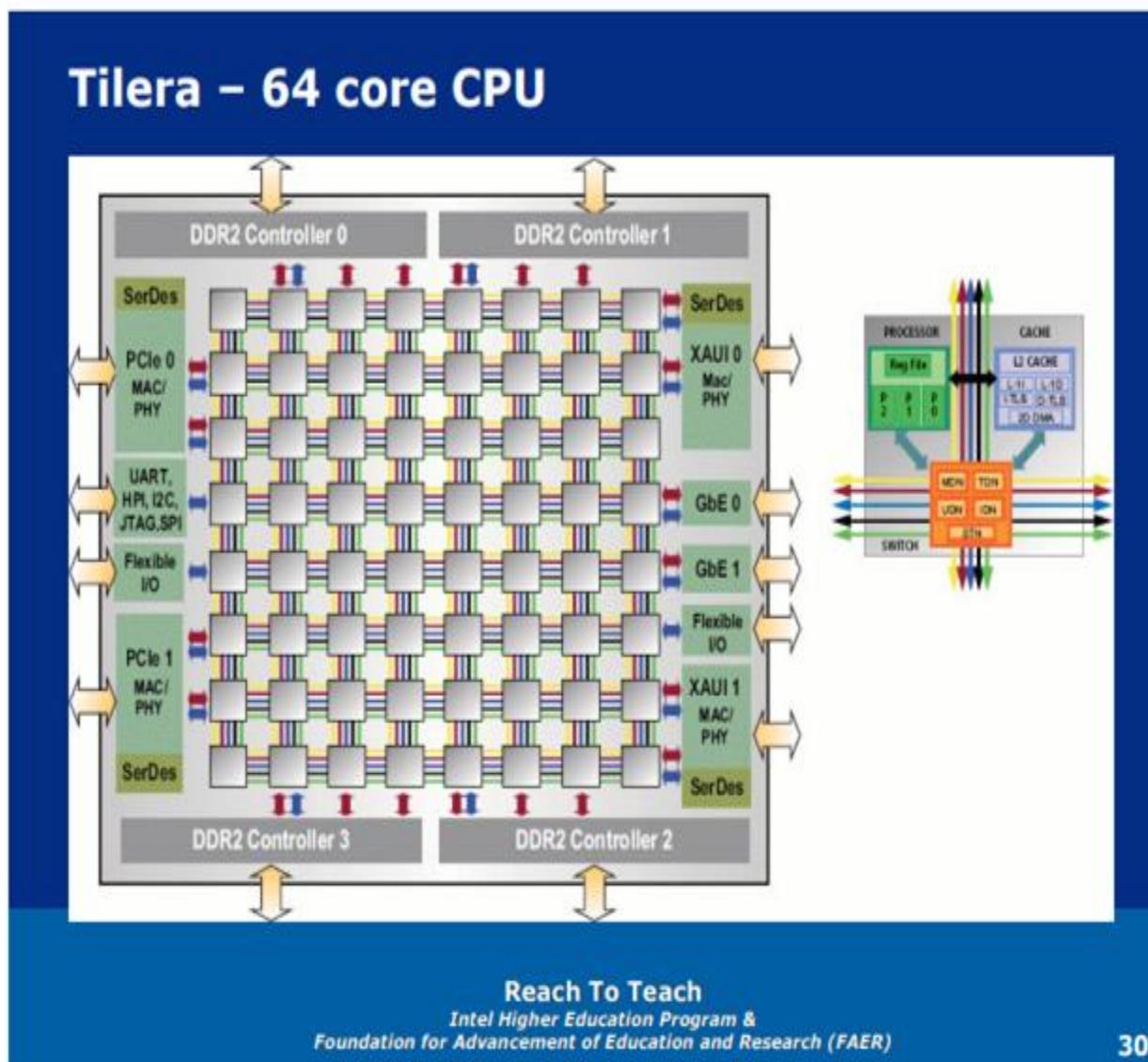
- VLIW processors rely on compile time analysis to identify and bundle together instructions that can be executed concurrently.
- These instructions are packed and dispatched together, and thus the name very long instruction word

Intel IA64 processors.

VLIW: Very Long Instruction Word



Multicore Processor



- A multi-core processor implements multiprocessing in a single physical package.
- Cores in a multi-core device may be coupled together tightly or loosely.
- For example, cores may or may not share caches, and they may implement message passing or shared memory inter-core communication methods.
- Common network topologies to interconnect cores include: bus, ring, 2-dimensional mesh, and crossbar.
- All cores are identical in *symmetric* multi-core systems and they are not identical in *asymmetric* multi-core systems.
- Just as with single-processor systems, cores in multi-core systems may implement architectures such as superscalar, vector processing, or multithreading.

- Multi-core processors are widely used across many application domains including: general-purpose, embedded, network, digital signal processing, and graphics.
- The amount of performance gained by the use of a multi-core processor is strongly dependent on the software algorithms and implementation.
- Multi-core processing is a growing industry trend as single core processors rapidly reach the physical limits of possible complexity and speed.
- Companies that have produced or are working on multi-core products include AMD, ARM, Broadcom, Intel, and VIA.

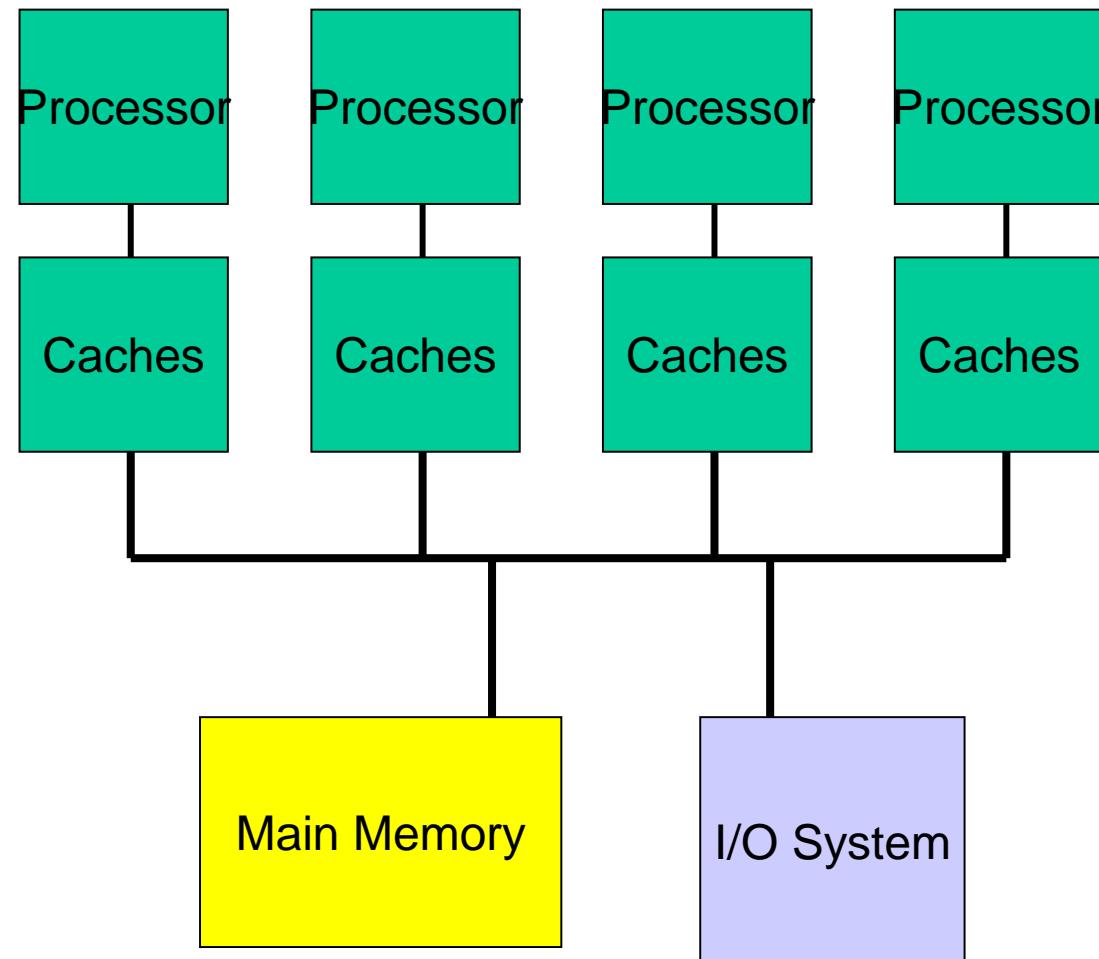
Flynn's Taxonomy

- Single instruction stream, single data stream (SISD)
 - uniprocessor
- Single instruction stream, multiple data streams (SIMD)
 - Vector architectures
 - Multimedia extensions
 - Graphics processor units
- Multiple instruction streams, single data stream (MISD)
 - No commercial implementation
 - Imagine data going through a pipeline of execution engines
- Multiple instruction streams, multiple data streams (MIMD)
 - Tightly-coupled MIMD
 - Loosely-coupled MIMD
 - Most multiprocessors today: easy to construct with off-the-shelf computers, most flexibility

Multiple-instruction, multiple-data (MIMD)

- MIMD machines are broadly categorized based on the way Processing Elements (PE) are coupled to the main memory.
- **Shared-memory MIMD** (**Types: Centralized and Distributed**)
 - Tightly coupled multiprocessor systems
 - Connected to a single global memory and they all have access to it.
 - Symmetric Multi-Processing.
 - Less likely to scale
- **Distributed-memory MIMD**
 - Loosely coupled multiprocessor systems
 - All PEs have a local memory
 - Communication between PEs in this model takes place through the interconnection network.

Centralized Shared-Memory



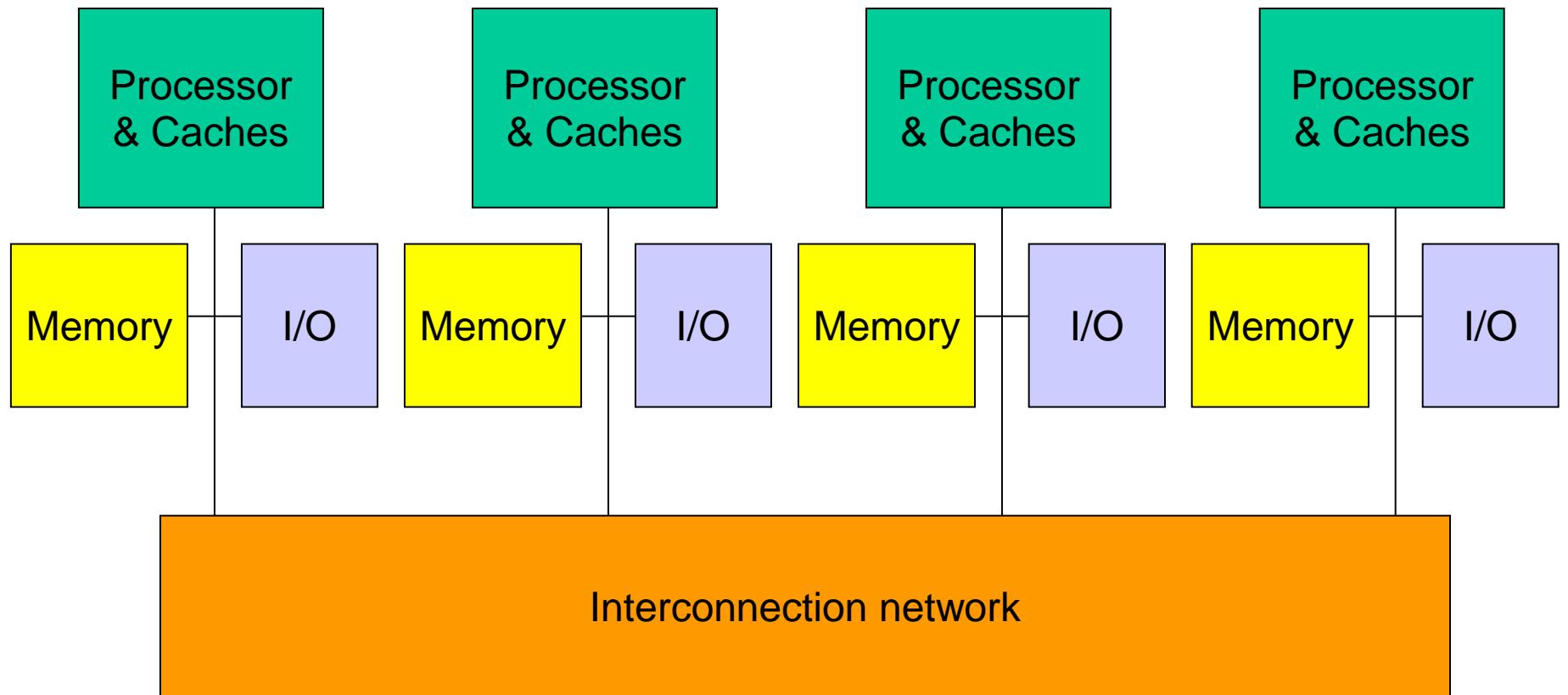
• **Advantages**

- Global address space provides a user-friendly programming perspective to memory
- Data sharing between tasks is both fast and uniform due to the proximity of memory to CPUs

• **Disadvantages**

- Primary disadvantage is the lack of scalability between memory and CPUs. Adding more CPUs can geometrically increase traffic on the shared memory-CPU path, and for cache coherent systems, geometrically increase traffic associated with cache/memory management.
- Programmer responsibility for synchronization constructs that ensure "correct" access of global memory.
- Expense: it becomes increasingly difficult and expensive to design and produce shared memory machines with ever increasing numbers of processors.

Distributed Memory Multiprocessors



Distributed Memory Multiprocessors

- Require a communication network to connect inter-processor memory.
- Processors have their own local memory. Memory addresses in one processor do not map to another processor
- No concept of global address space across all processors.
- Changes to local memory have no effect on the memory of other processors.

• **Advantages**

- Memory is scalable with number of processors. Increase the number of processors and the size of memory increases proportionately.
- Each processor can rapidly access its own memory without interference and without the overhead incurred with trying to maintain cache coherency.
- Cost effectiveness: can use commodity, off-the-shelf processors and networking.

• **Disadvantages**

- The programmer is responsible for many of the details associated with data communication between processors.
- It may be difficult to map existing data structures, based on global memory, to this memory organization.
- Non-uniform memory access (NUMA) times

Interconnection Network

- Both shared memory and message passing architectures can use
 - **static interconnection network**
 - can not be modified without a physical re-designing of a system
 - static networks usually in message passing computers
 - **dynamic interconnection network**
 - dynamic networks usually in shared memory computers
 - enable changing (reconfiguring) of the connection structure in a system.

Interconnection Network

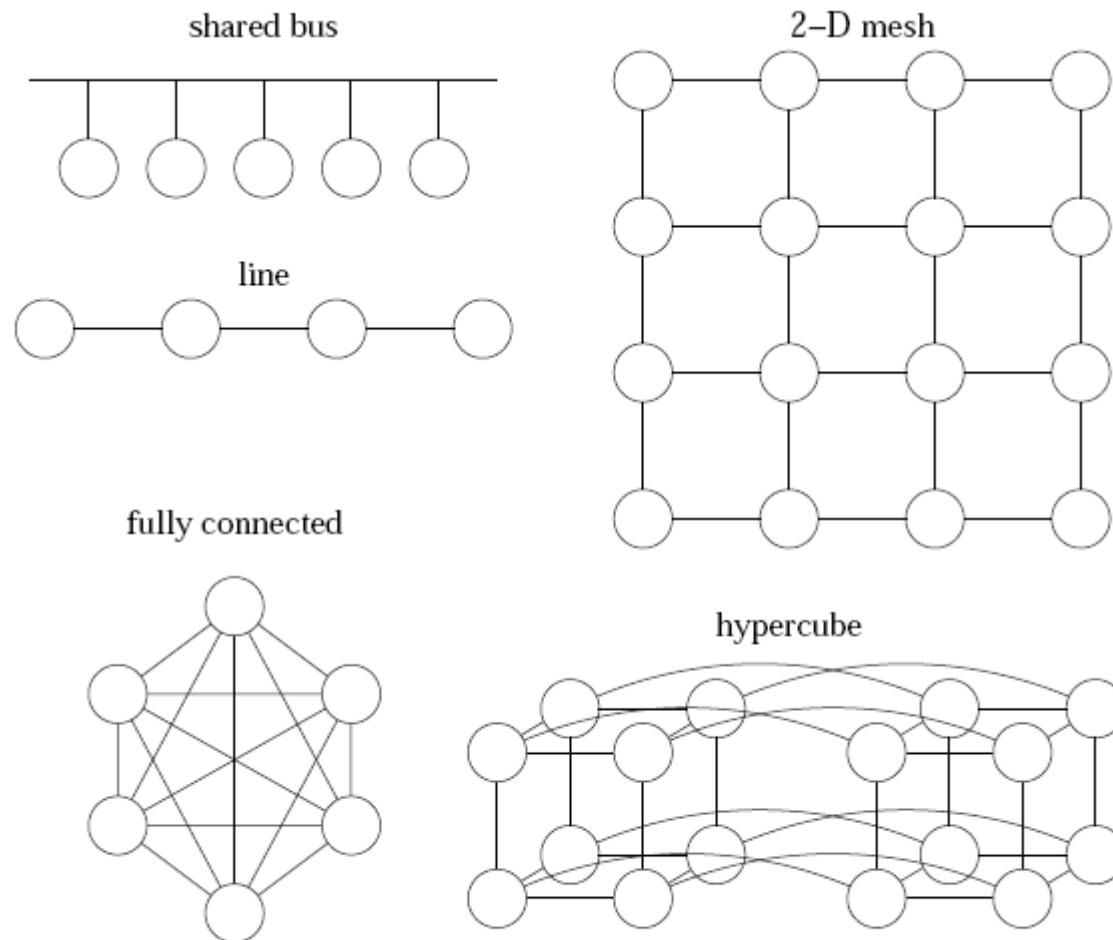
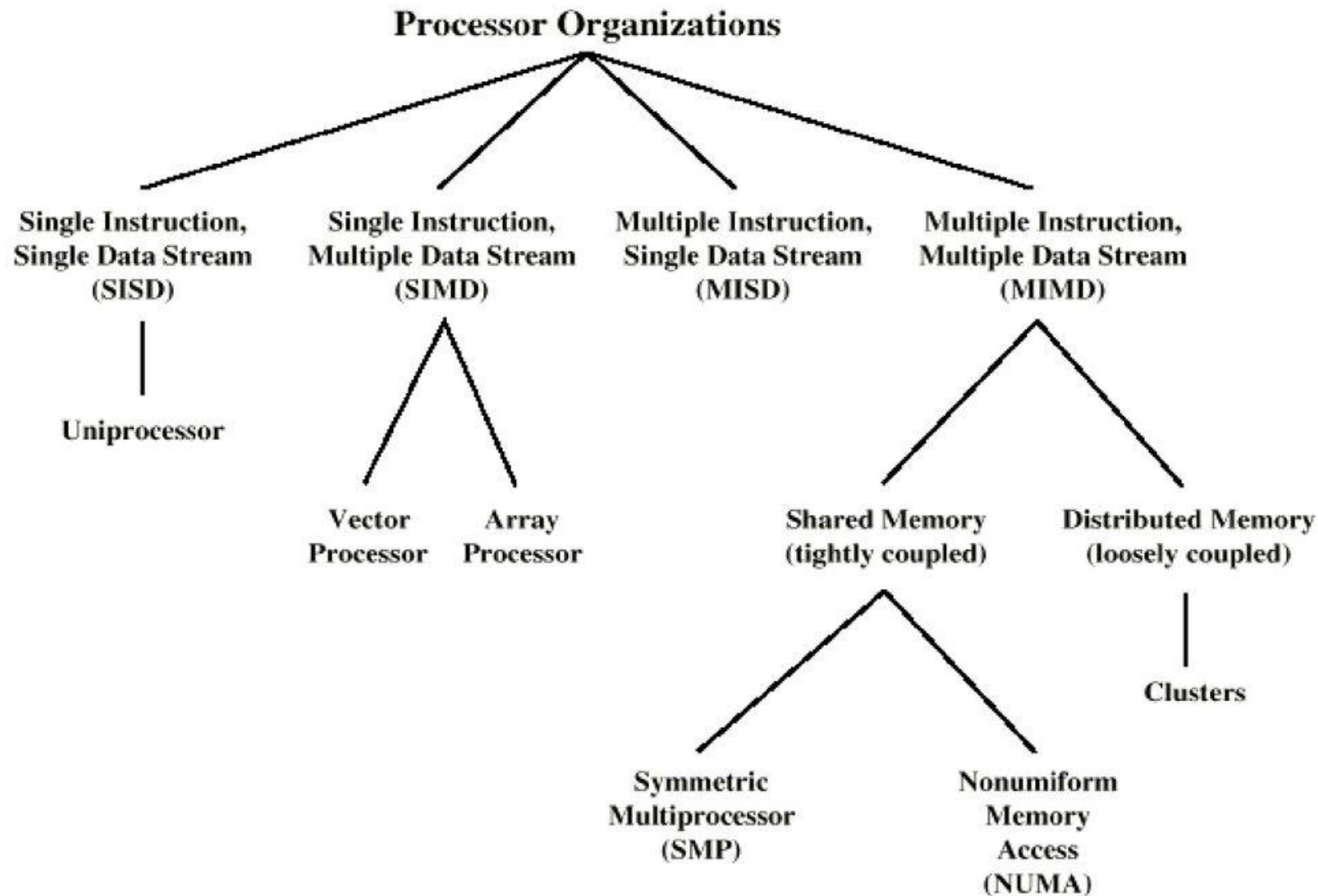


Figure 5.6: Five important interconnect network topologies.

Taxonomy of Parallel Processor Architectures



- Parallel architectures can be classified in several ways
- Flynn's taxonomy
- Control mechanism
- Address-space organization
- Interconnection network
- Granularity of processors

Granularity of Processors

- Granularity is the ratio of computation time to communication time.
- Computation time – Time required to perform computation of a task
- Communication time is the time required to exchange data between processors
- Fine-grained, medium-grained and coarse-grained parallelism

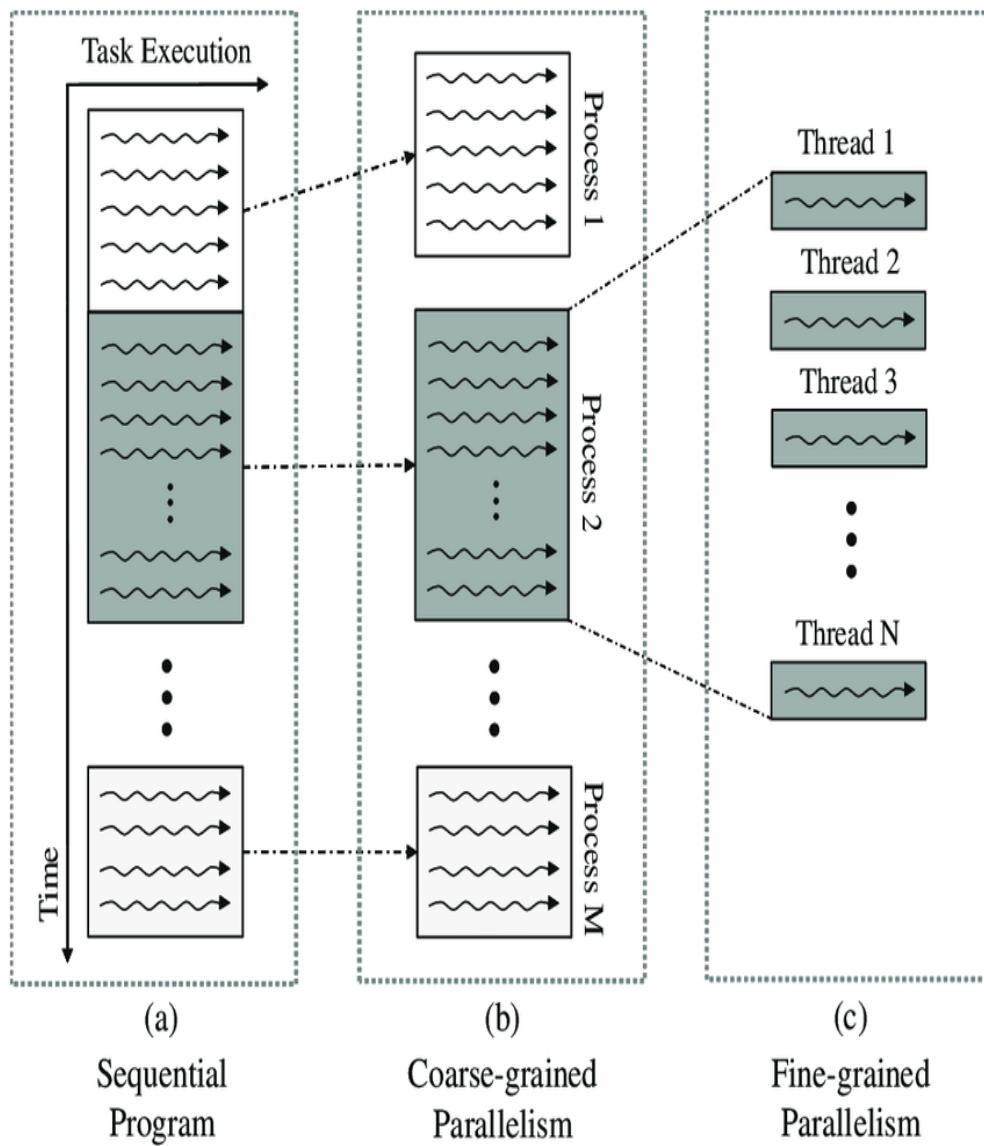
- **Fine-grained parallelism**
 - a program is broken down to a large number of small **tasks**.
 - **Tasks** are assigned individually to many processors.
 - The work is evenly distributed among the processors.
- Increases the communication and synchronization overhead

Fine-grained parallelism is best exploited in architectures which support fast communication

Shared memory architecture are most suitable for fine-grained parallelism.

- **Coarse-grained parallelism**

- a program is split into **large tasks**. a large amount of computation takes place in processors.
- Load imbalance
- Certain tasks process the bulk of the data while others might be idle.
- Low communication and synchronization overhead.



- **Medium-grained parallelism**

- Medium-grained parallelism is used relatively to fine-grained and coarse-grained parallelism.
- Task size and communication time greater than fine-grained parallelism and lower than coarse-grained parallelism
- General-purpose parallel computers fall in this category

- Example :
- Consider a 10×10 image that needs to be processed, given that, processing of the 100 pixels is independent of each other. Assume there are
- **Fine-grained parallelism:**
- 100 processors that are responsible for processing the 10×10 image.
- 100 processors can process the 10×10 image in 1 clock cycle
- Each processor is working on 1 pixel of the image
- **Medium-grained parallelism:**
- 25 processors processing the 10×10 image.
- The processing of the image will now take 4 clock cycles.
- **Coarse-grained parallelism**
- we reduce the processors to 2, processing will take 50 clock cycles

- Relationship between levels of parallelism, grain size and degree of parallelism

<u>Levels</u>	Grain Size	Parallelism
– Instruction level	-- Fine	-- Highest
– Loop level	-- Fine	-- Moderate
– Sub-routine level	-- Medium	-- Moderate
– Program level	-- Coarse	-- Least

• **Impact of granularity on performance**

- Fine grains or small tasks results in more parallelism.
- Synchronization overhead, scheduling strategies
- Coarse grained tasks have less communication overhead but they often cause load imbalance.

Finding the best grain size, depends on a number of factors and varies greatly from problem-to-problem.

References

- [https://en.wikipedia.org/wiki/Granularity_\(parallel_computing\)](https://en.wikipedia.org/wiki/Granularity_(parallel_computing))
- Hwang, Kai (1992). Advanced Computer Architecture: Parallelism,Scalability,Programmability (1st ed.). McGraw-Hill Higher Education. ISBN 978-0070316225.

Linking

Instructor:

R. Shathanaa

Example C Program

```
int sum(int *a, int n);  
  
int array[2] = {1, 2};  
  
int main()  
{  
    int val = sum(array, 2);  
    return val;  
}
```

main.c

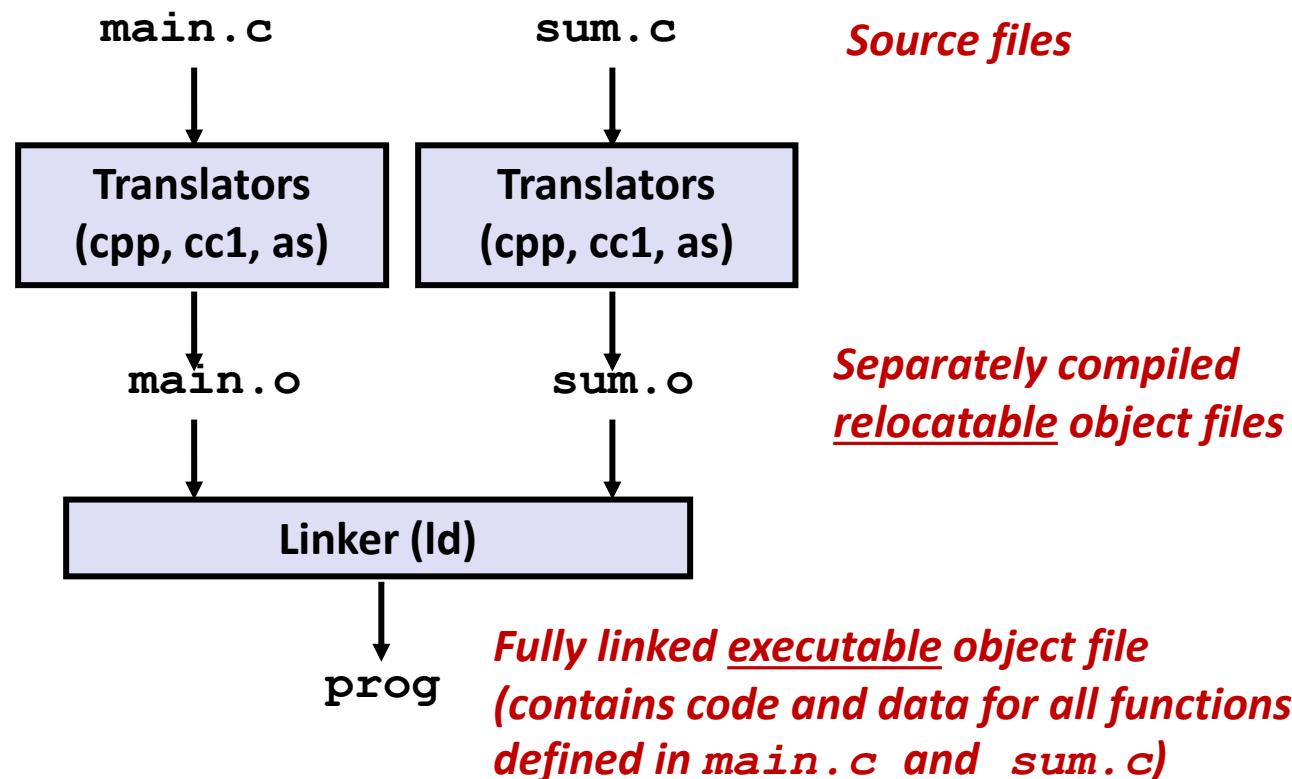
```
int sum(int *a, int n)  
{  
    int i, s = 0;  
  
    for (i = 0; i < n; i++) {  
        s += a[i];  
    }  
    return s;  
}
```

sum.c

Static Linking

- Programs are translated and linked using a *compiler driver*:

- linux> `gcc -Og -o prog main.c sum.c`
- linux> `./prog`



Why Linkers?

■ Reason 1: Modularity

- Program can be written as a collection of smaller source files, rather than one monolithic mass.
- Can build libraries of common functions (more on this later)
 - e.g., Math library, standard C library

Why Linkers? (cont)

■ Reason 2: Efficiency

- Time: Separate compilation
 - Change one source file, compile, and then relink.
 - No need to recompile other source files.
- Space: Libraries
 - Common functions can be aggregated into a single file...
 - Yet executable files and running memory images contain only code for the functions they actually use.

What Do Linkers Do?

■ Step 1: Symbol resolution

- Programs define and reference *symbols* (global variables and functions):
 - `void swap() { ... } /* define symbol swap */`
 - `swap(); /* reference symbol swap */`
 - `int *xp = &x; /* define symbol xp, reference x */`
- Symbol definitions are stored in object file (by assembler) in *symbol table*.
 - Symbol table is an array of structs
 - Each entry includes name, size, and location of symbol.
- **During symbol resolution step, the linker associates each symbol reference with exactly one symbol definition.**

What Do Linkers Do? (cont)

■ Step 2: Relocation

- Merges separate code and data sections into single sections
- Relocates symbols from their relative locations in the .o files to their final absolute memory locations in the executable.
- Updates all references to these symbols to reflect their new positions.

Let's look at these two steps in more detail....

Three Kinds of Object Files (Modules)

■ Relocatable object file (.o file)

- Contains code and data in a form that can be combined with other relocatable object files to form executable object file.
 - Each .o file is produced from exactly one source (.c) file

■ Executable object file (a .out file)

- Contains code and data in a form that can be copied directly into memory and then executed.

■ Shared object file (.so file)

- Special type of relocatable object file that can be loaded into memory and linked dynamically, at either load time or run-time.
- Called *Dynamic Link Libraries* (DLLs) by Windows

Executable and Linkable Format (ELF)

- Standard binary format for object files
- One unified format for
 - Relocatable object files (.o),
 - Executable object files (a.out)
 - Shared object files (.so)
- Generic name: ELF binaries

ELF Object File Format

■ Elf header

- Word size, byte ordering, file type (.o, exec, .so), machine type, etc.

■ Segment header table

- Page size, virtual addresses memory segments (sections), segment sizes.

■ .text section

- Code

■ .rodata section

- Read only data: jump tables, ...

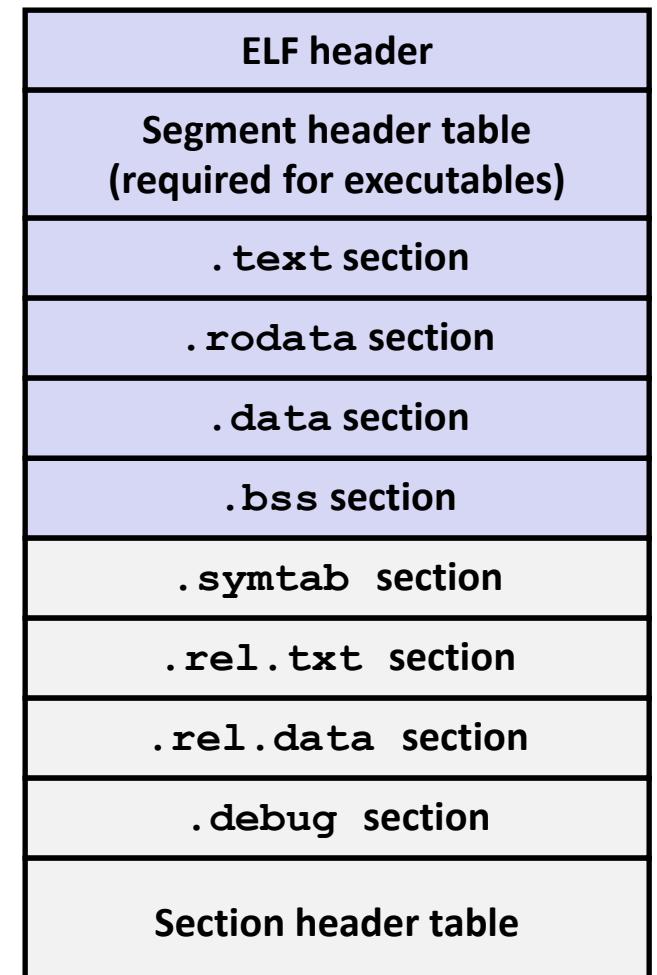
■ .data section

- Initialized global variables

■ .bss section

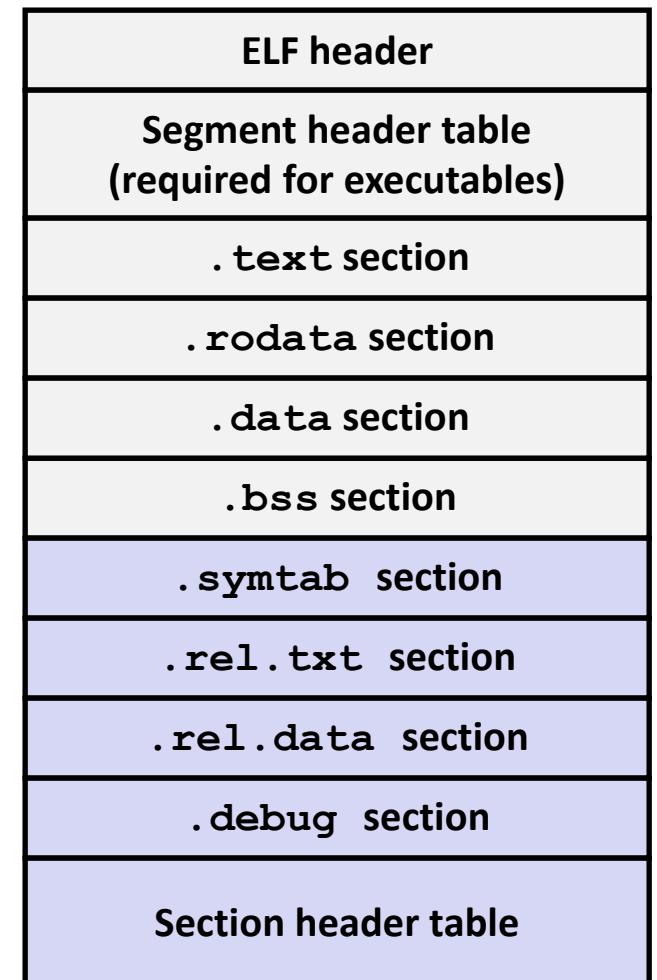
- Uninitialized global variables
- “Block Started by Symbol”
- “Better Save Space”

- Has section header but occupies no space



ELF Object File Format (cont.)

- **.syntab section**
 - Symbol table
 - Procedure and static variable names
 - Section names and locations
- **.rel.text section**
 - Relocation info for **.text** section
 - Addresses of instructions that will need to be modified in the executable
 - Instructions for modifying.
- **.rel.data section**
 - Relocation info for **.data** section
 - Addresses of pointer data that will need to be modified in the merged executable
- **.debug section**
 - Info for symbolic debugging (**gcc -g**)
- **Section header table**
 - Offsets and sizes of each section



Linker Symbols

■ Global symbols

- Symbols defined by module m that can be referenced by other modules.
- E.g.: non-**static** C functions and non-**static** global variables.

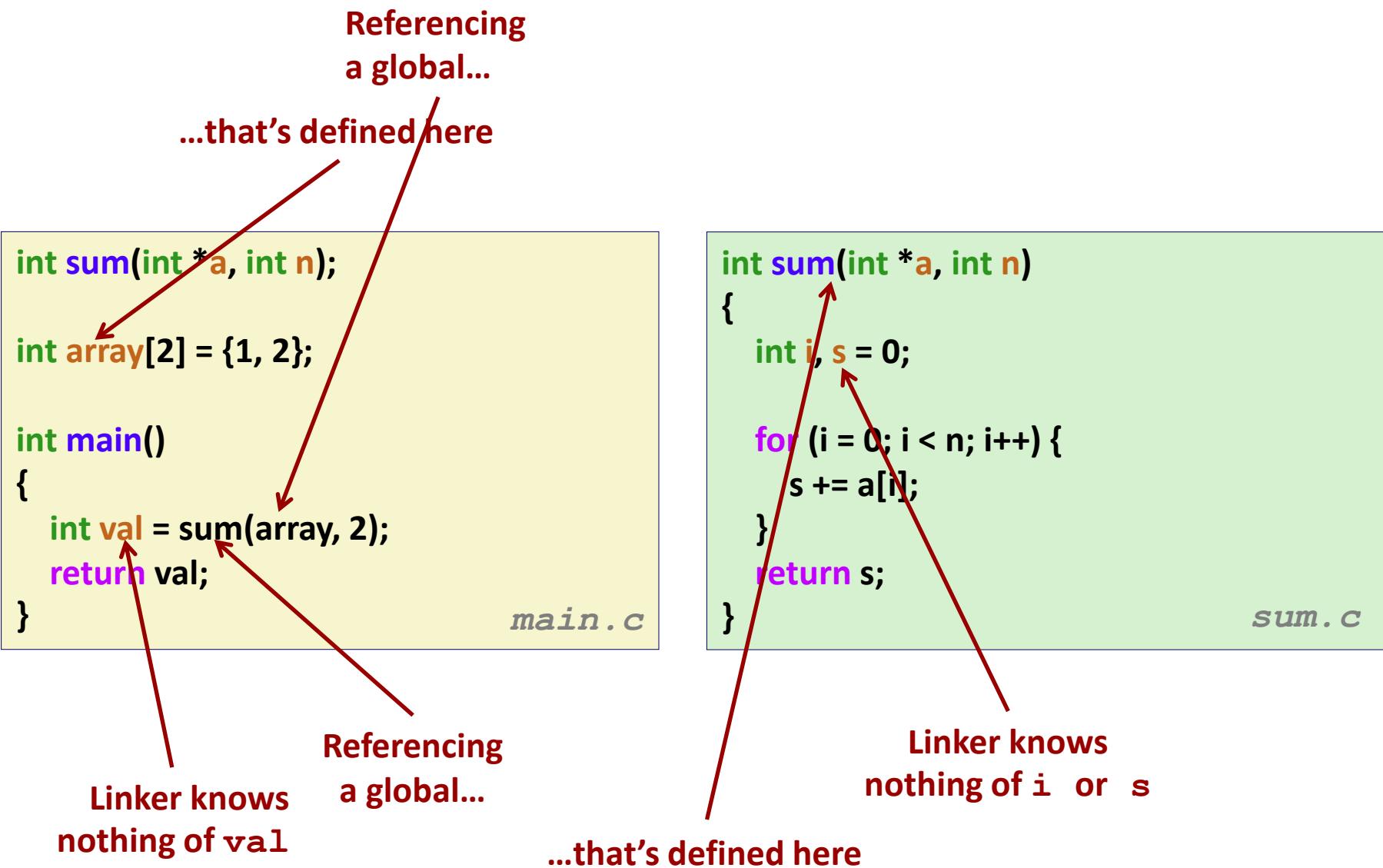
■ External symbols

- Global symbols that are referenced by module m but defined by some other module.

■ Local symbols

- Symbols that are defined and referenced exclusively by module m .
- E.g.: C functions and global variables defined with the **static** attribute.
- **Local linker symbols are *not* local program variables**

Step 1: Symbol Resolution



Local Symbols

■ Local non-static C variables vs. local static C variables

- local non-static C variables: stored on the stack
- local static C variables: stored in either .bss, or .data

```
int f()
{
    static int x = 0;
    return x;
}

int g()
{
    static int x = 1;
    return x;
}
```

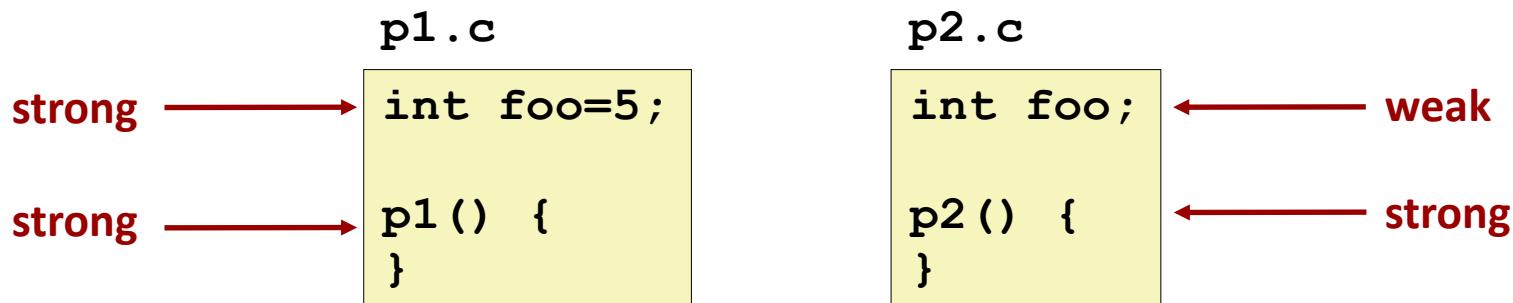
Compiler allocates space in .data for each definition of x

Creates local symbols in the symbol table with unique names, e.g., x . 1 and x . 2.

How Linker Resolves Duplicate Symbol Definitions

■ Program symbols are either *strong* or *weak*

- *Strong*: procedures and initialized globals
- *Weak*: uninitialized globals



Linker's Symbol Rules

- **Rule 1: Multiple strong symbols are not allowed**
 - Each item can be defined only once
 - Otherwise: Linker error
- **Rule 2: Given a strong symbol and multiple weak symbols, choose the strong symbol**
 - References to the weak symbol resolve to the strong symbol
- **Rule 3: If there are multiple weak symbols, pick an arbitrary one**
 - Can override this with `gcc -fno-common`

Linker Puzzles

```
int x;  
p1() {}
```

```
p1() {}
```

Link time error: two strong symbols (**p1**)

```
int x;  
p1() {}
```

```
int x;  
p2() {}
```

References to **x** will refer to the same uninitialized int. Is this what you really want?

```
int x;  
int y;  
p1() {}
```

```
double x;  
p2() {}
```

Writes to **x** in **p2** might overwrite **y**!
Evil!

```
int x=7;  
int y=5;  
p1() {}
```

```
double x;  
p2() {}
```

Writes to **x** in **p2** will overwrite **y**!
Nasty!

```
int x=7;  
p1() {}
```

```
int x;  
p2() {}
```

References to **x** will refer to the same initialized variable.

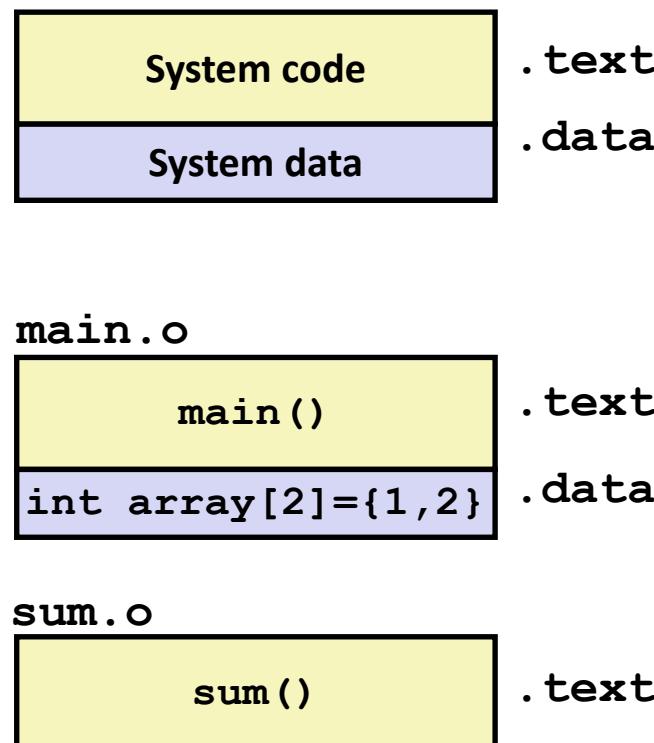
Nightmare scenario: two identical weak structs, compiled by different compilers with different alignment rules.

Global Variables

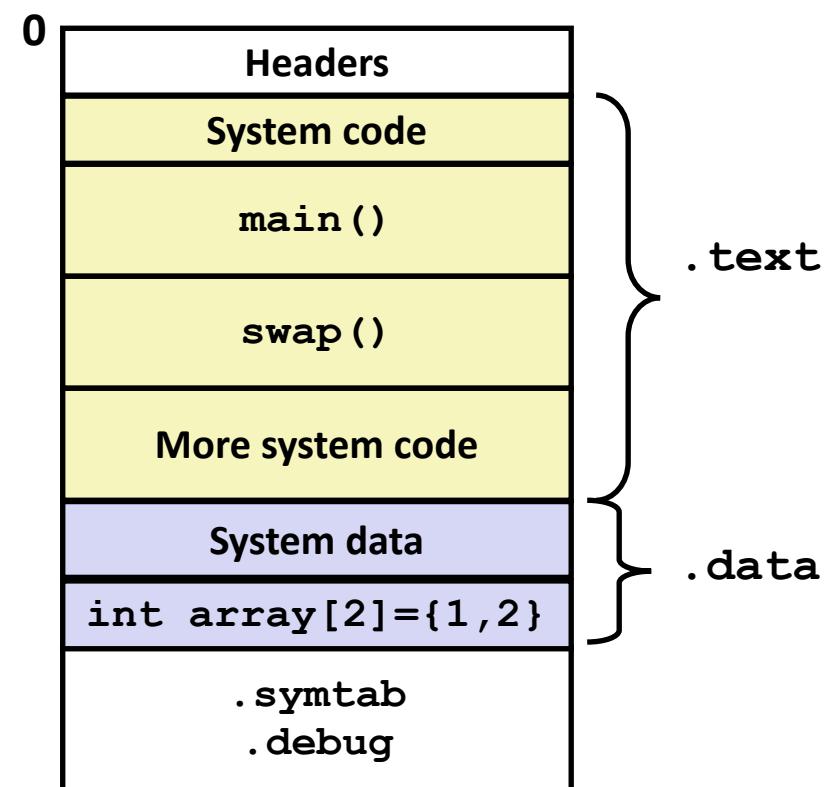
- **Avoid if you can**
- **Otherwise**
 - Use **static** if you can
 - Initialize if you define a global variable
 - Use **extern** if you reference an external global variable

Step 2: Relocation

Relocatable Object Files



Executable Object File



Relocation Entries

```
int array[2] = {1, 2};
```

```
int main()
{
    int val = sum(array, 2);
    return val;
}
```

main.c

```
0000000000000000 <main>:
0: 48 83 ec 08      sub  $0x8,%rsp
4: be 02 00 00 00    mov   $0x2,%esi
9: bf 00 00 00 00    mov   $0x0,%edi  # %edi = &array
    a: R_X86_64_32 array    # Relocation entry

e: e8 00 00 00 00    callq 13 <main+0x13> # sum()
    f: R_X86_64_PC32 sum-0x4  # Relocation entry
13: 48 83 c4 08      add   $0x8,%rsp
17: c3                retq
```

main.o

Relocated .text section

00000000004004d0 <main>:

```
4004d0: 48 83 ec 08    sub $0x8,%rsp
4004d4: be 02 00 00 00  mov $0x2,%esi
4004d9: bf 18 10 60 00  mov $0x601018,%edi # %edi = &array
4004de: e8 05 00 00 00  callq 4004e8 <sum> # sum()
4004e3: 48 83 c4 08  add $0x8,%rsp
4004e7: c3           retq
```

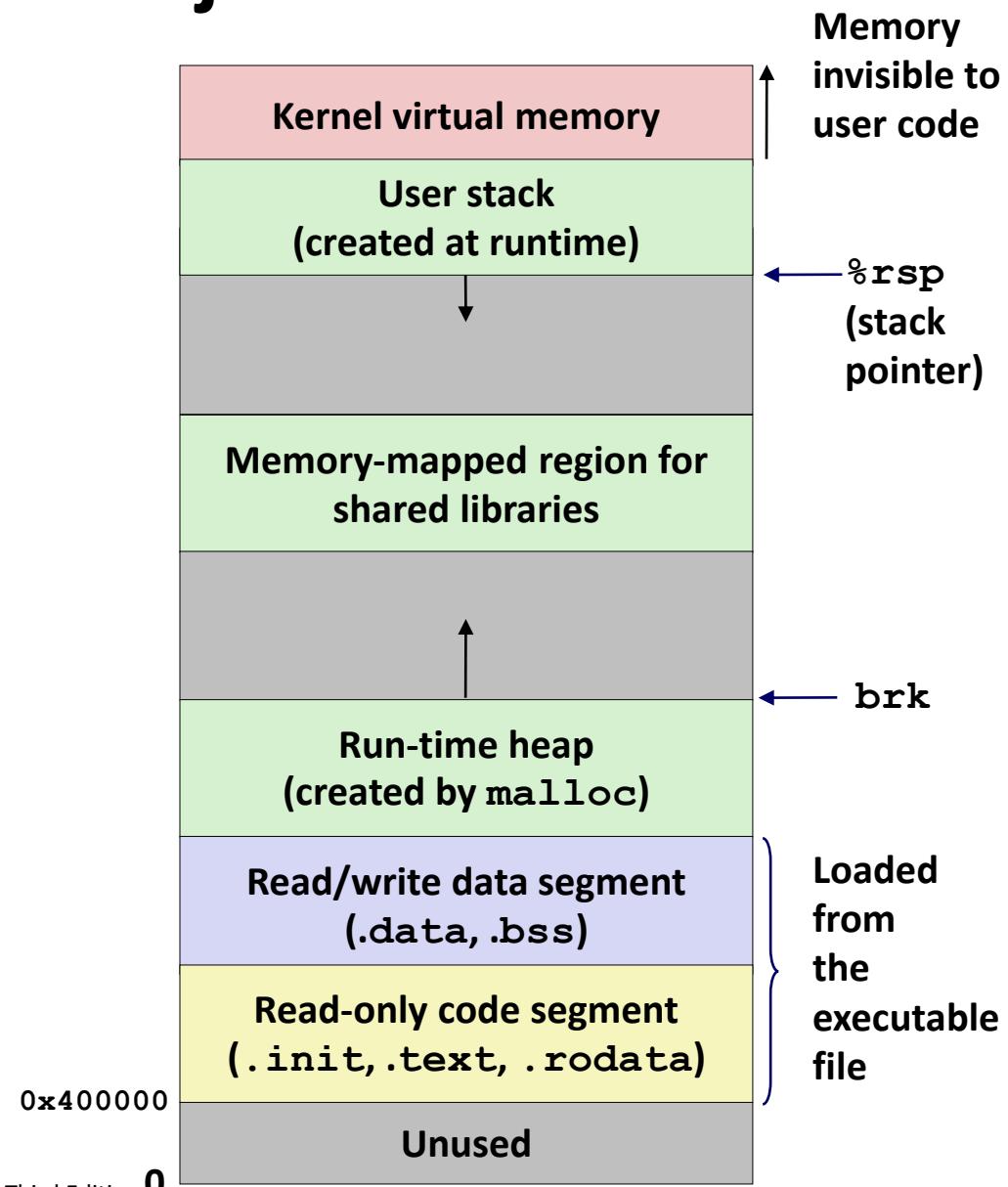
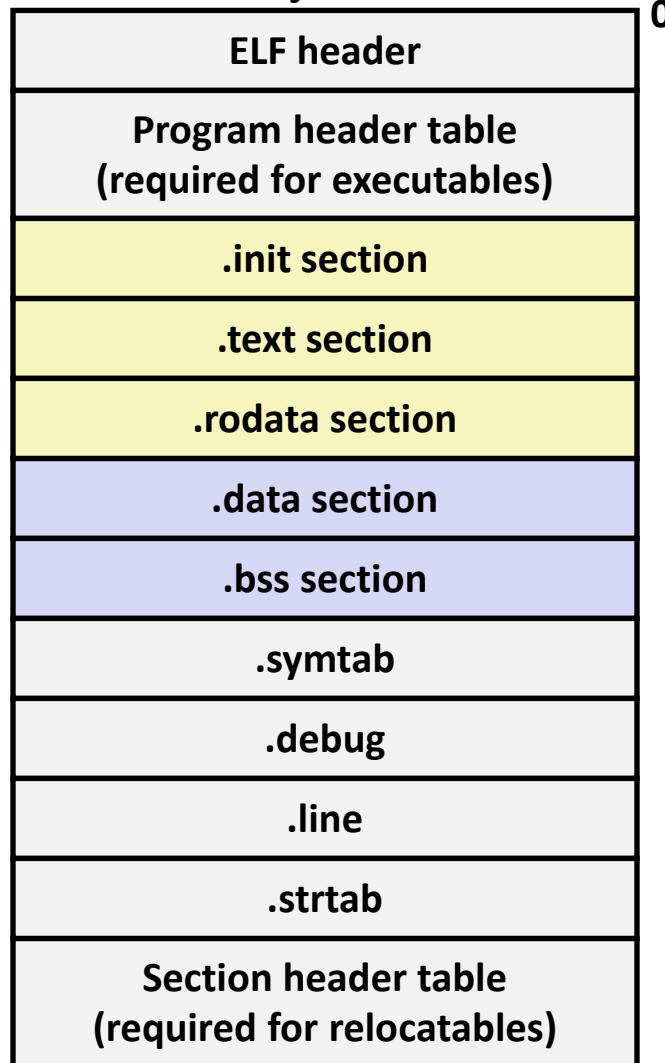
00000000004004e8 <sum>:

```
4004e8: b8 00 00 00 00  mov $0x0,%eax
4004ed: ba 00 00 00 00  mov $0x0,%edx
4004f2: eb 09          jmp 4004fd <sum+0x15>
4004f4: 48 63 ca        movslq %edx,%rcx
4004f7: 03 04 8f        add (%rdi,%rcx,4),%eax
4004fa: 83 c2 01        add $0x1,%edx
4004fd: 39 f2          cmp %esi,%edx
4004ff: 7c f3          jl 4004f4 <sum+0xc>
400501: f3 c3          repz retq
```

Using PC-relative addressing for sum(): **0x4004e8 = 0x4004e3 + 0x5**

Loading Executable Object Files

Executable Object File



Packaging Commonly Used Functions

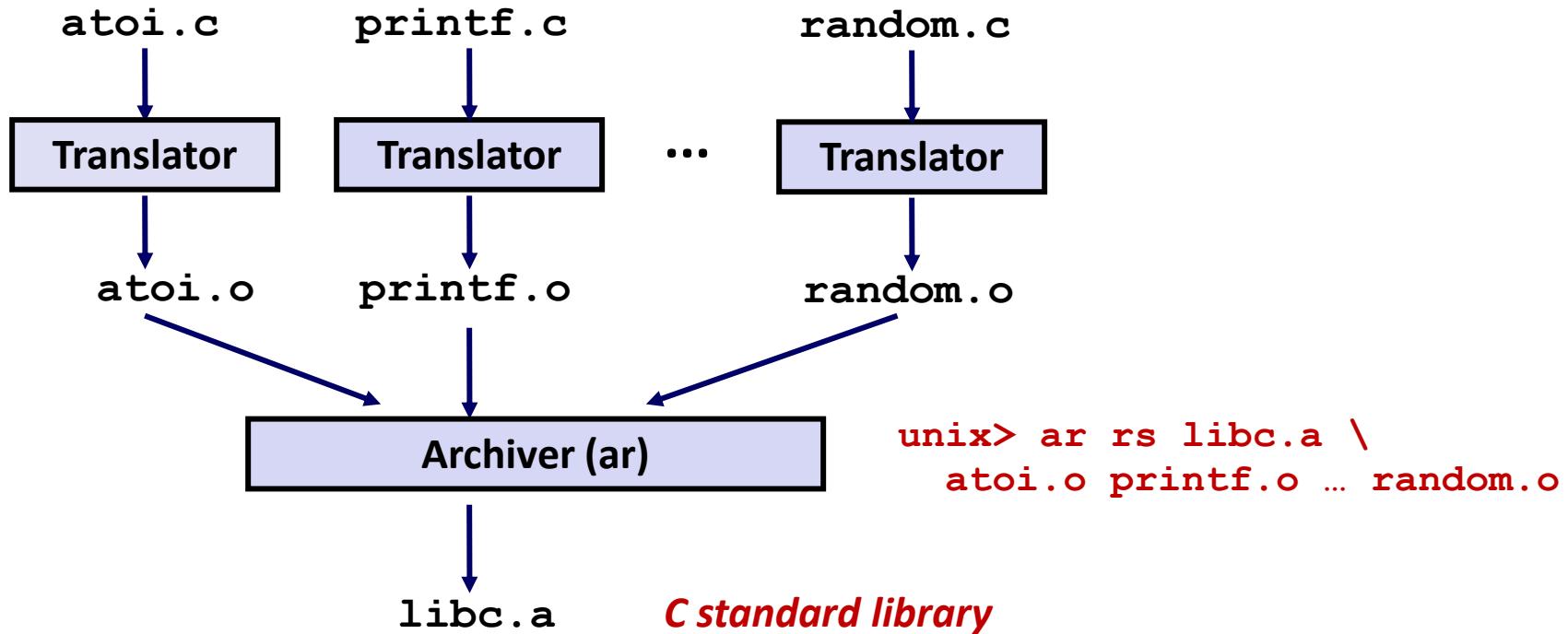
- **How to package functions commonly used by programmers?**
 - Math, I/O, memory management, string manipulation, etc.
- **Awkward, given the linker framework so far:**
 - **Option 1:** Put all functions into a single source file
 - Programmers link big object file into their programs
 - Space and time inefficient
 - **Option 2:** Put each function in a separate source file
 - Programmers explicitly link appropriate binaries into their programs
 - More efficient, but burdensome on the programmer

Old-fashioned Solution: Static Libraries

■ **Static libraries (.a archive files)**

- Concatenate related relocatable object files into a single file with an index (called an *archive*).
- Enhance linker so that it tries to resolve unresolved external references by looking for the symbols in one or more archives.
- If an archive member file resolves reference, link it into the executable.

Creating Static Libraries



- Archiver allows incremental updates
- Recompile function that changes and replace .o file in archive.

Commonly Used Libraries

libc.a (the C standard library)

- 4.6 MB archive of 1496 object files.
- I/O, memory allocation, signal handling, string handling, data and time, random numbers, integer math

libm.a (the C math library)

- 2 MB archive of 444 object files.
- floating point math (sin, cos, tan, log, exp, sqrt, ...)

```
% ar -t libc.a | sort
...
fork.o
...
fprintf.o
fpu_control.o
fputc.o
freopen.o
fscanf.o
fseek.o
fstab.o
...
```

```
% ar -t libm.a | sort
...
e_acos.o
e_acosf.o
e_acosh.o
e_acoshf.o
e_acoshl.o
e_acosl.o
e_asin.o
e_asinf.o
e_asinl.o
...
```

Linking with Static Libraries

```
#include <stdio.h>
#include "vector.h"

int x[2] = {1, 2};
int y[2] = {3, 4};
int z[2];

int main()
{
    addvec(x, y, z, 2);
    printf("z = [%d %d]\n",
           z[0], z[1]);
    return 0;
}
```

main2.c

libvector.a

```
void addvec(int *x, int *y,
            int *z, int n) {
    int i;

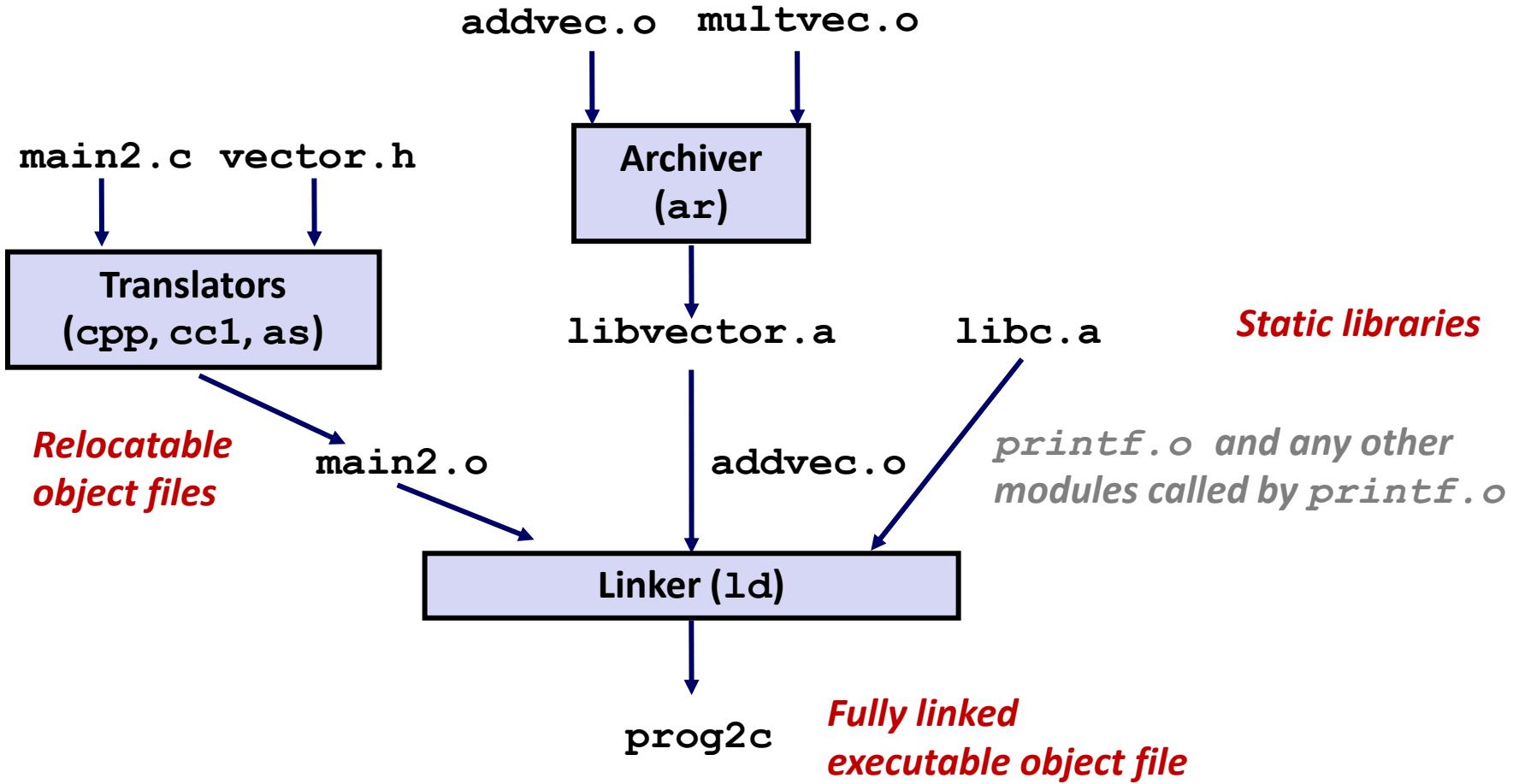
    for (i = 0; i < n; i++)
        z[i] = x[i] + y[i];
}
```

```
void multvec(int *x, int *y,
             int *z, int n)
{
    int i;

    for (i = 0; i < n; i++)
        z[i] = x[i] * y[i];
}
```

addvec.c *multvec.c*

Linking with Static Libraries



Using Static Libraries

■ Linker's algorithm for resolving external references:

- Scan **.o** files and **.a** files in the command line order.
- During the scan, keep a list of the current unresolved references.
- As each new **.o** or **.a** file, *obj*, is encountered, try to resolve each unresolved reference in the list against the symbols defined in *obj*.
- If any entries in the unresolved list at end of scan, then error.

■ Problem:

- Command line order matters!
- Moral: put libraries at the end of the command line.

```
unix> gcc -L. libtest.o -lmine
unix> gcc -L. -lmine libtest.o
libtest.o: In function `main':
libtest.o(.text+0x4): undefined reference to `libfun'
```

Modern Solution: Shared Libraries

■ Static libraries have the following disadvantages:

- Duplication in the stored executables (every function needs libc)
- Duplication in the running executables
- Minor bug fixes of system libraries require each application to explicitly relink

■ Modern solution: Shared Libraries

- Object files that contain code and data that are loaded and linked into an application *dynamically*, at either *load-time* or *run-time*
- Also called: dynamic link libraries, DLLs, .so files

Shared Libraries (cont.)

- **Dynamic linking can occur when executable is first loaded and run (load-time linking).**
 - Common case for Linux, handled automatically by the dynamic linker (**ld-linux.so**) .
 - Standard C library (**libc.so**) usually dynamically linked.
- **Dynamic linking can also occur after program has begun (run-time linking).**
 - In Linux, this is done by calls to the **dlopen()** interface .
 - Distributing software.
 - High-performance web servers.
 - Runtime library interpositioning.