

Memory

Eri Prasetyo Wibowo

<http://eri.staff.gunadarma.ac.id>

Sem icon

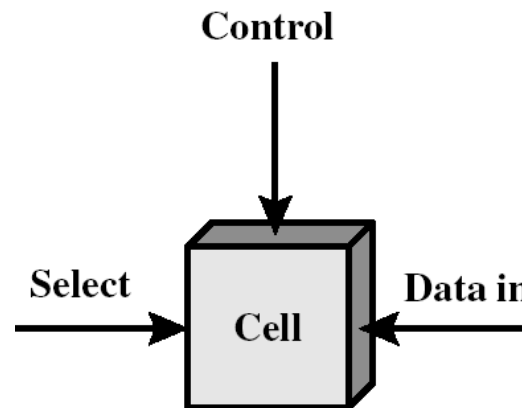
Table 5.1 Semiconductor Memory Types

Memory Type	Category	Erase	Write Mechanism	Volatility
Random-access memory (RAM)	Read-write memory	Electrically, byte-level	Electrically	Volatile
Read-only memory (ROM)	Read-only memory	Not possible	Masks	
Programmable ROM (PROM)			Read-mostly memory	UV light, chip-level
Erasable PROM (EPROM)				
Electrically Erasable PROM (EEPROM)	Electrically, byte-level			
Flash memory		Electrically, block-level		

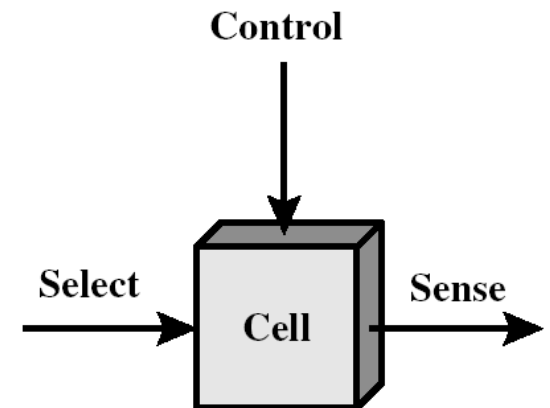
Semiconductor Memory

RAM

- Misnamed as all semiconductor memory is random access
- Read/ Write
- Volatile
- Temporary storage
- Static or dynamic



(a) Write



(b) Read

Two basic types devised:

- Static RAM (SRAM)
- Dynamic RAM (DRAM)

Static RAM faster operation than dynamic RAM (around 10 times faster) but requires more internal components (see next).

Hence DRAM has more memory cells in chip (around 4-8 times capacity).

Usually main memory is dynamic and cache memory usually static.

Random-Access Memory (RAM)

Key features

- RAM adalah sebagai kemasan chip.
- Dasar penyimpanan unit adalah sel (satu bit per sel).
- Beberapa chip RAM membentuk memori

Static RAM (SRAM)

- Setiap sel menyimpan bit dengan enam transistor-sirkuit.
- Nilai tetap tentu, selama ini disimpan daya.
- Relatif kebal untuk gangguan listrik seperti kebisingan.
- Lebih cepat dan lebih mahal daripada DRAM.

Dynamic RAM (DRAM)

- Setiap sel menyimpan bit dengan kapasitor dan transistor.
- Nilai harus refresh setiap 10-100 ms.
- Sensitif terhadap gangguan.
- Lambat dan lebih murah dibandingkan SRAM.

SRAM vs DRAM summary

	Tran. per bit	Access time	Persist?	Sensitive?	Cost	Applications
SRAM	6	1X	Yes	No	100x	cache memories
DRAM	1	10X	No	Yes	1X	Main memories, frame buffers

Dyn

ami

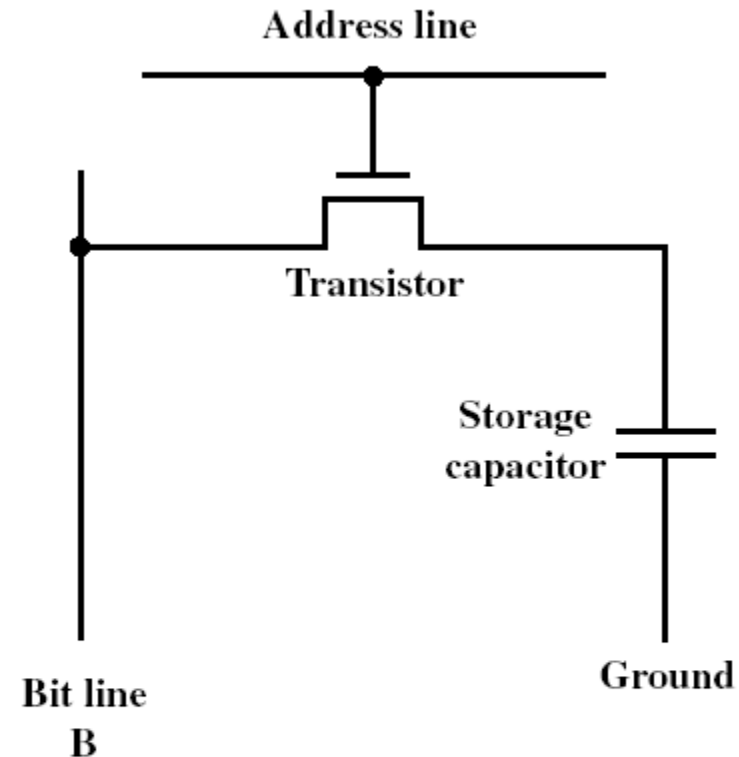
C

RA

M

- Bits stored as charge in capacitors
- Charges leak
- Need refreshing even when powered
- Simpler construction
- Smaller per bit
- Less expensive
- Need refresh circuits
- Slower
- Main memory
- Essentially analogue

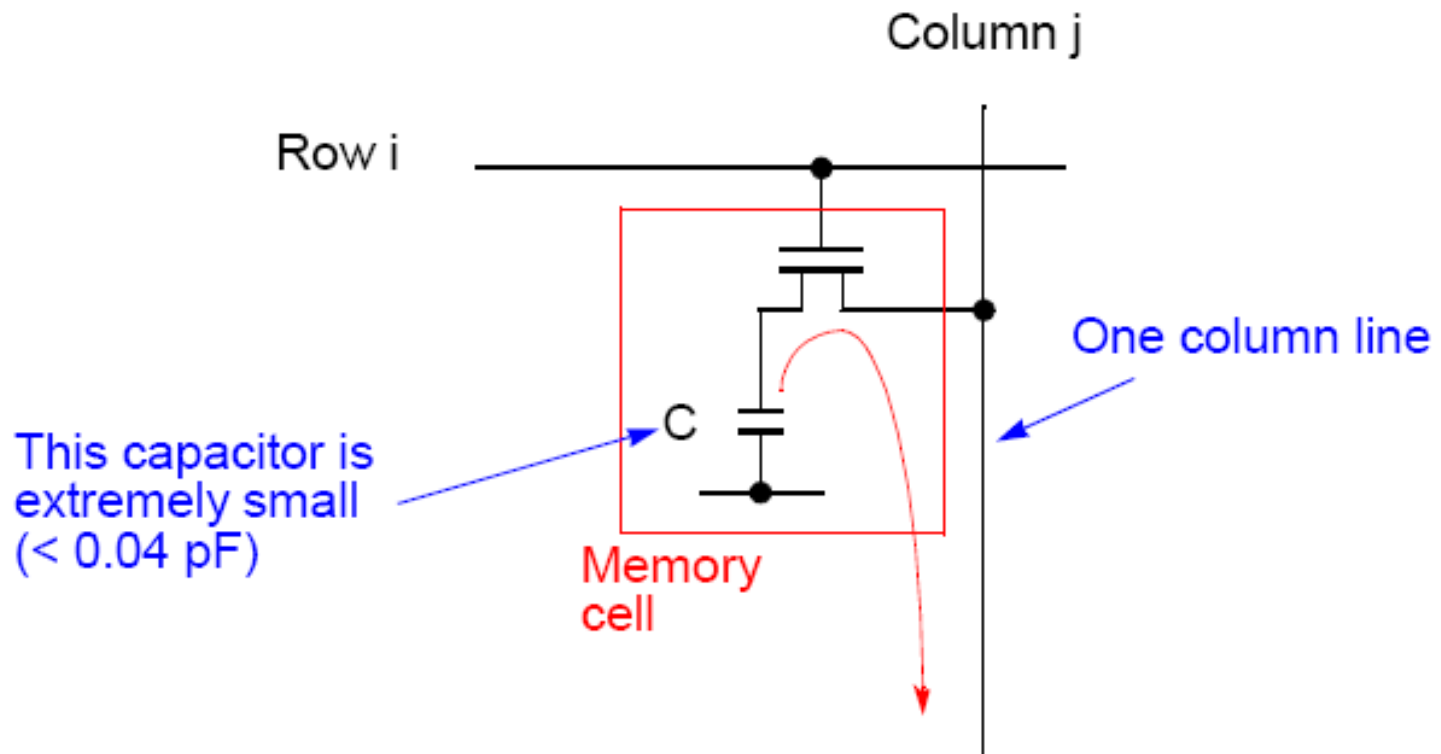
- Level of charge determines value



(a) Dynamic RAM (DRAM) cell

Dynamic RAM cell circuit

Each memory cell uses a capacitor to store charge. Two states, charged and not charged represent the two binary values 0 and 1. Less components, higher memory density. Usual design:



Because charge will decay, memory cell must be refreshed periodically, typically every 2-4 ms. Refresh consists of reading data and rewriting it. Usually done automatically within chip.

DR

AM

Operati

on

Address line active when bit read or written

- Transistor switch closed (current flows)

Write

- Voltage to bit line
 - **High for 1 low for 0**
- Then signal address line
 - **Transfers charge to capacitor**

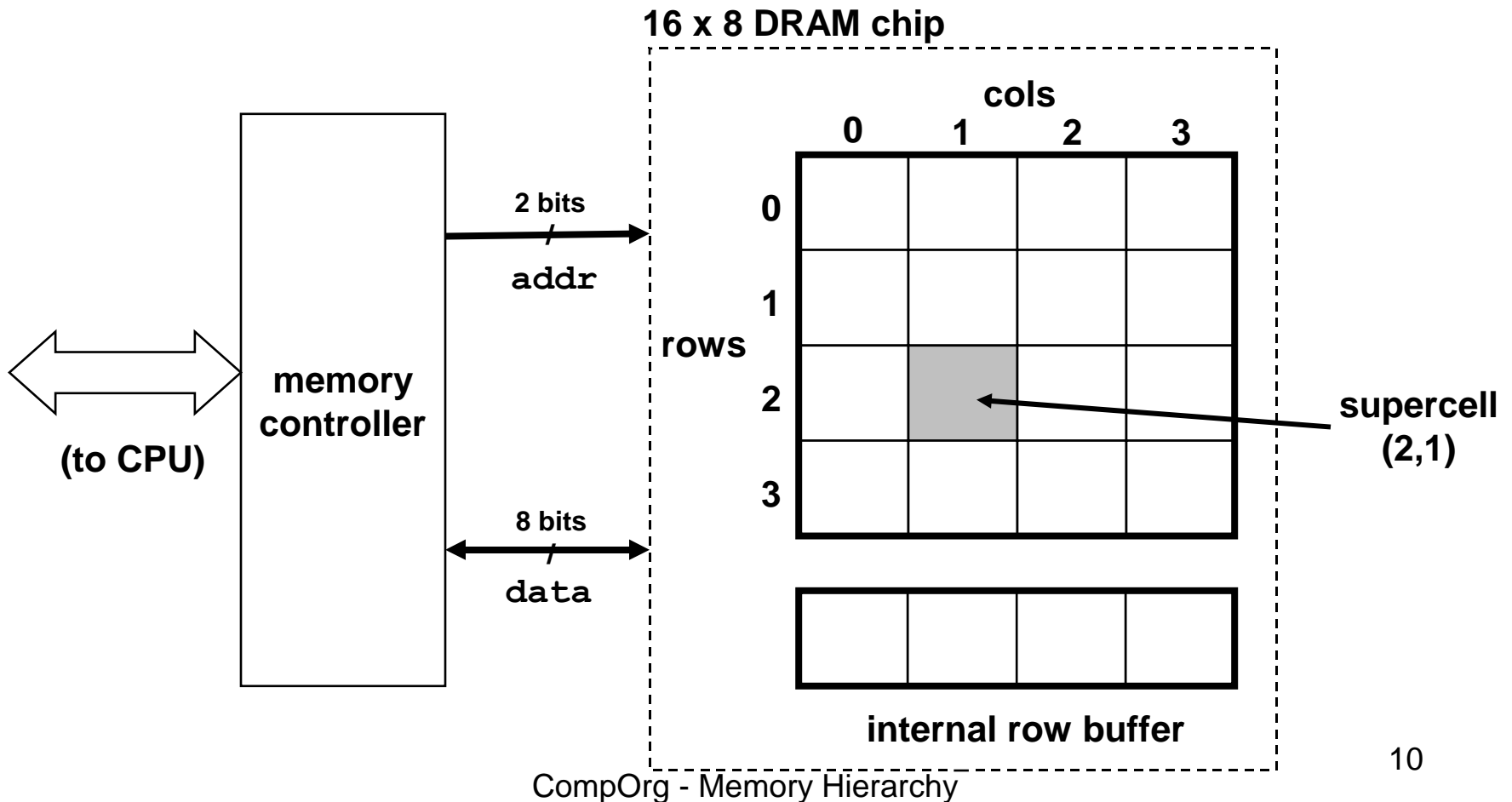
Read

- Address line selected
 - **transistor turns on**
- Charge from capacitor fed via bit line to sense amplifier
 - **Compares with reference value to determine 0 or 1**
- Capacitor charge must be restored

Conventional DRAM organization

$d \times w$ DRAM:

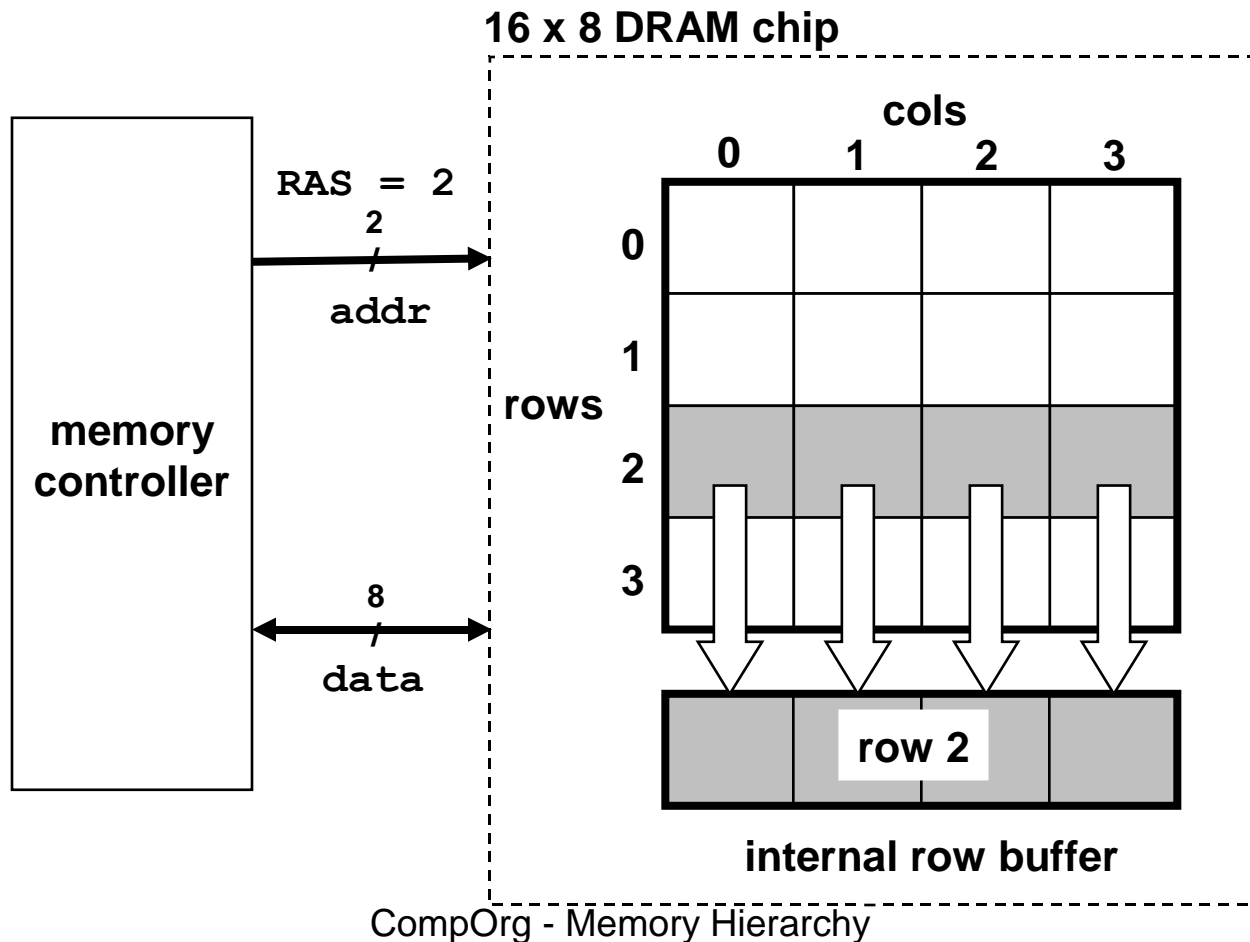
- dw total bits organized as d supercells of size w bits



Reading DRAM supercell

Step 1(a): Row access strobe (RAS) selects row 2.

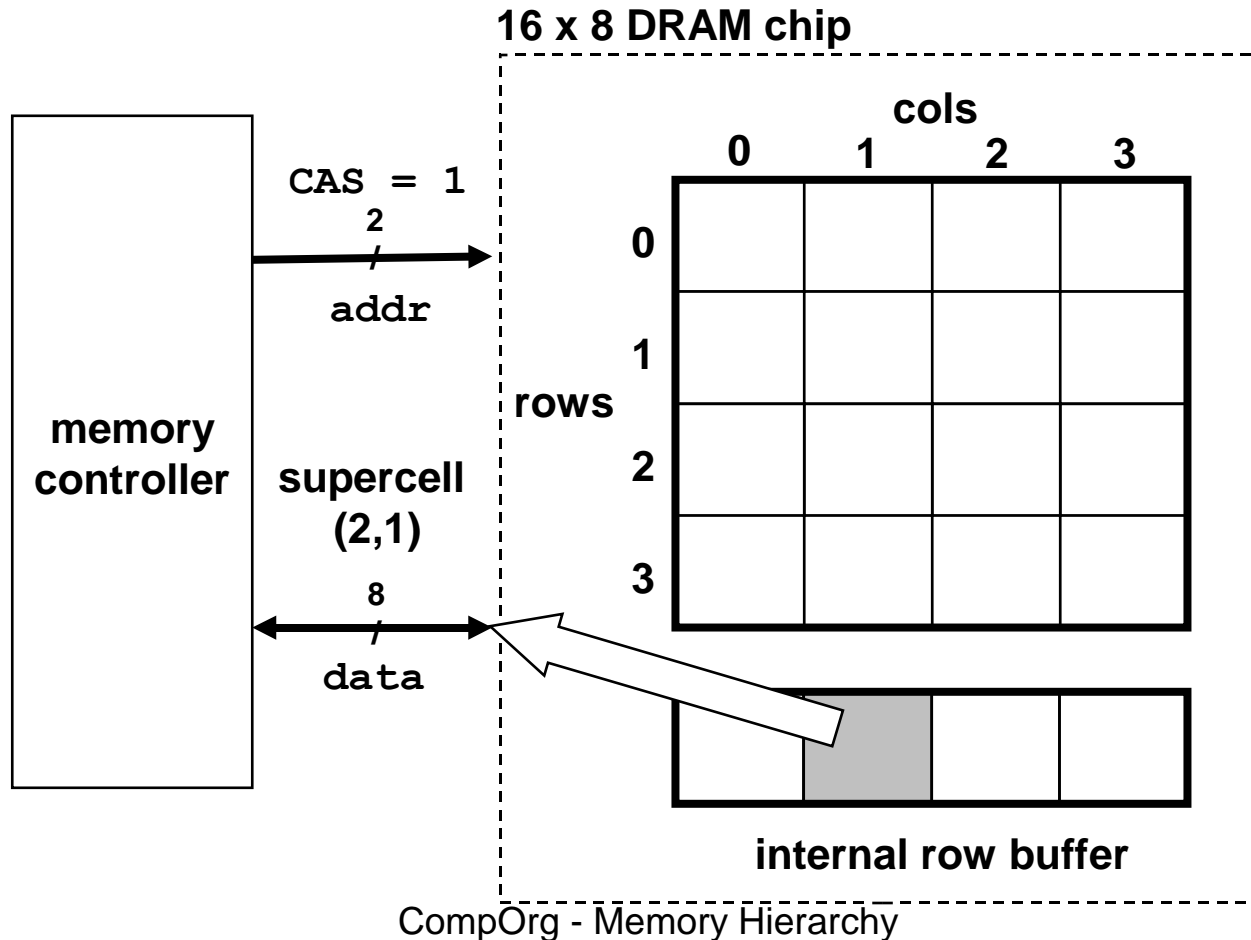
Step 1(b): Row 2 copied from DRAM array to row buffer.



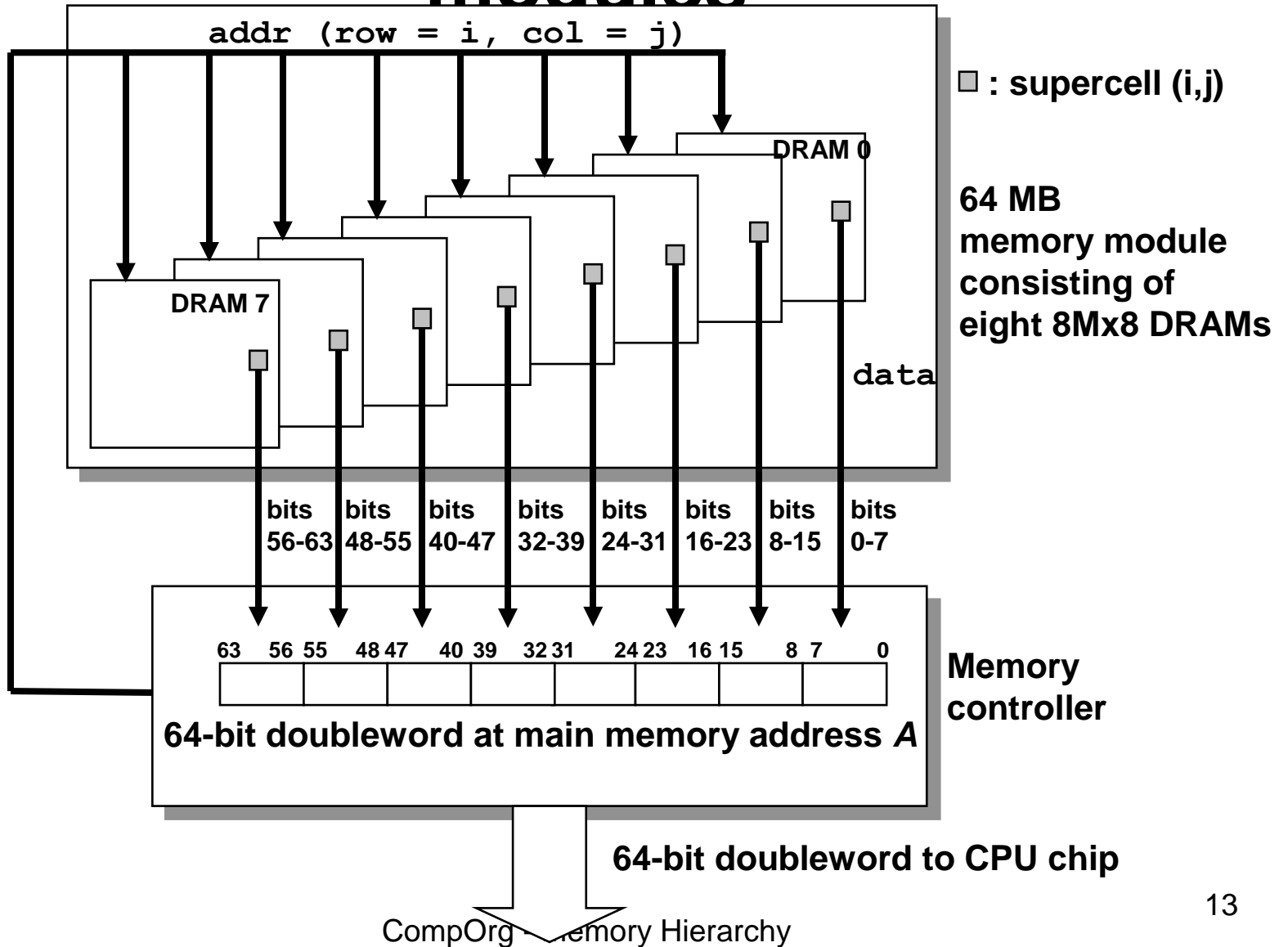
Reading DRAM supercell

Step 2(a): Column access **(2,1)** (CAS) selects column 1.

Step 2(b): Supercell (2,1) copied from buffer to data lines, and eventually back to the CPU.



Memory modules



Enhanced

DRAMs

All enhanced DRAMs are built around the conventional DRAM core.

- **Fast page mode DRAM (FPM DRAM)**
 - Access contents of row with [RAS, CAS, CAS, CAS, CAS] instead of [(RAS,CAS), (RAS,CAS), (RAS,CAS), (RAS,CAS)].
- **Extended data out DRAM (EDO DRAM)**
 - Enhanced FPM DRAM with more closely spaced CAS signals.
- **Synchronous DRAM (SDRAM)**
 - Driven with rising clock edge instead of asynchronous control signals.
- **Double data-rate synchronous DRAM (DDR SDRAM)**
 - Enhancement of SDRAM that uses both clock edges as control signals.
- **Video RAM (VRAM)**
 - Like FPM DRAM, but output is produced by shifting row buffer
 - Dual ported (allows concurrent reads and writes)

Stat

Bits stored as on/off switches

No charges to leak

No refreshing needed when powered

More complex construction

Larger per bit

More expensive

Does not need refresh circuits

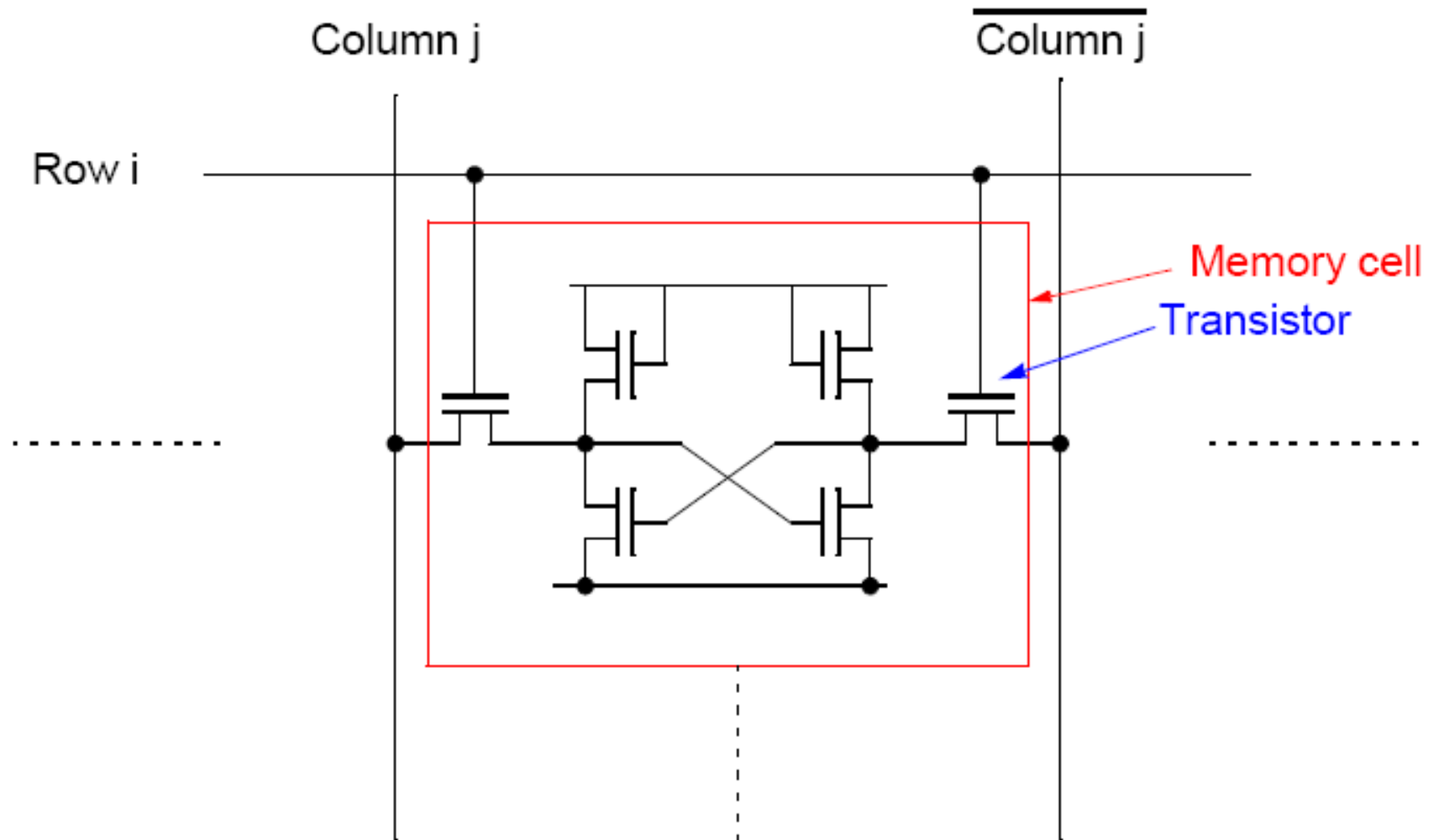
Faster

Cache

Digital

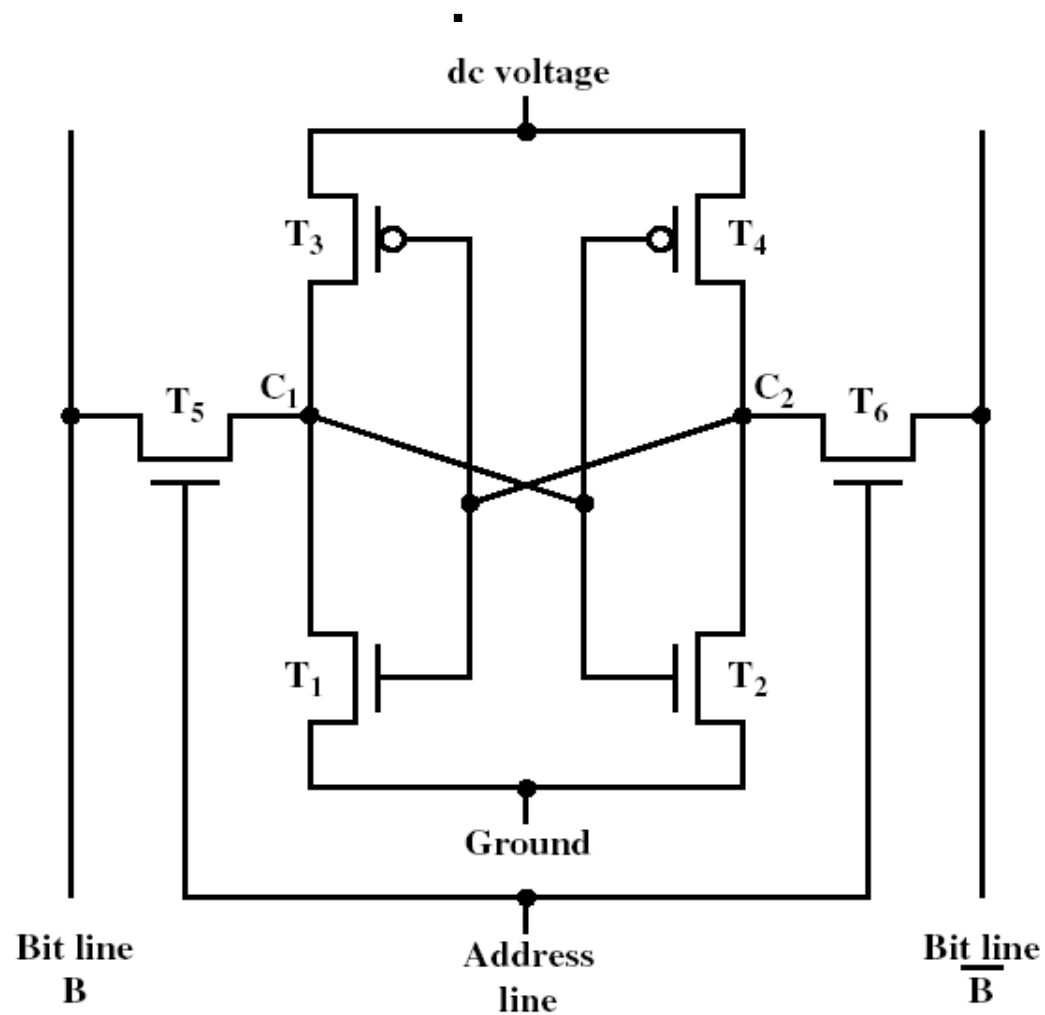
- Uses flip-flops

Static RAM Memory Cell Circuit



Design requires 4 or 6 transistors and two columns select lines per column.
Do not worry about its operation. Key point is number of components.

Stat



(b) Static RAM (SRAM) cell

Stat

ic
RA

M

Ope
rati
on

Transistor arrangement gives stable logic state

State 1

- C_1 high, C_2 low
- T_1 T_4 off, T_2 T_3 on

State 0

- C_2 high, C_1 low
- T_2 T_3 off, T_1 T_4 on

Address line transistors T_5 T_6 is switch

Write – apply value to B & compliment to B

Read – value is on line B

—

SRA

M

vs

DR

AM

Both volatile

- Power needed to preserve data

Dynamic cell

- Simpler to build, smaller
- More dense
- Less expensive
- Needs refresh
- Larger memory units

Static

- Faster
- Cache

Nonvolatile

DRAM and SRAM are volatile memories

- Lose information if powered off.

Nonvolatile memories retain value even if powered off.

- Generic name is read-only memory (ROM).
- Misleading because some ROMs can be read and modified.

Types of ROMs

- Programmable ROM (PROM)
- Erasable programmable ROM (EPROM)
- Electrically erasable PROM (EEPROM)
- Flash memory

Firmware

- Program stored in a ROM
 - Boot time code, BIOS (basic input/output system)
 - graphics cards, disk controllers.

Permanent storage

- Nonvolatile

Microprogramming (see later)

Library subroutines

Systems programs (BIOS)

Function tables

Read
Only
Memory
(ROM)

Types

Written during manufacture

- Very expensive for small runs

Programmable (once)

- PROM
- Needs special equipment to program

Read “mostly”

- Erasable Programmable (EPROM)
 - **Erased by UV**
- Electrically Erasable (EEPROM)
 - **Takes much longer to write than read**
- Flash memory
 - **Erase whole memory electrically**

Org

A 16Mbit chip can be organised as 1M of 16 bit words

A bit per chip system has 6 bits of 1Mbit chip with bit 1 of each word in chip 1 and so on

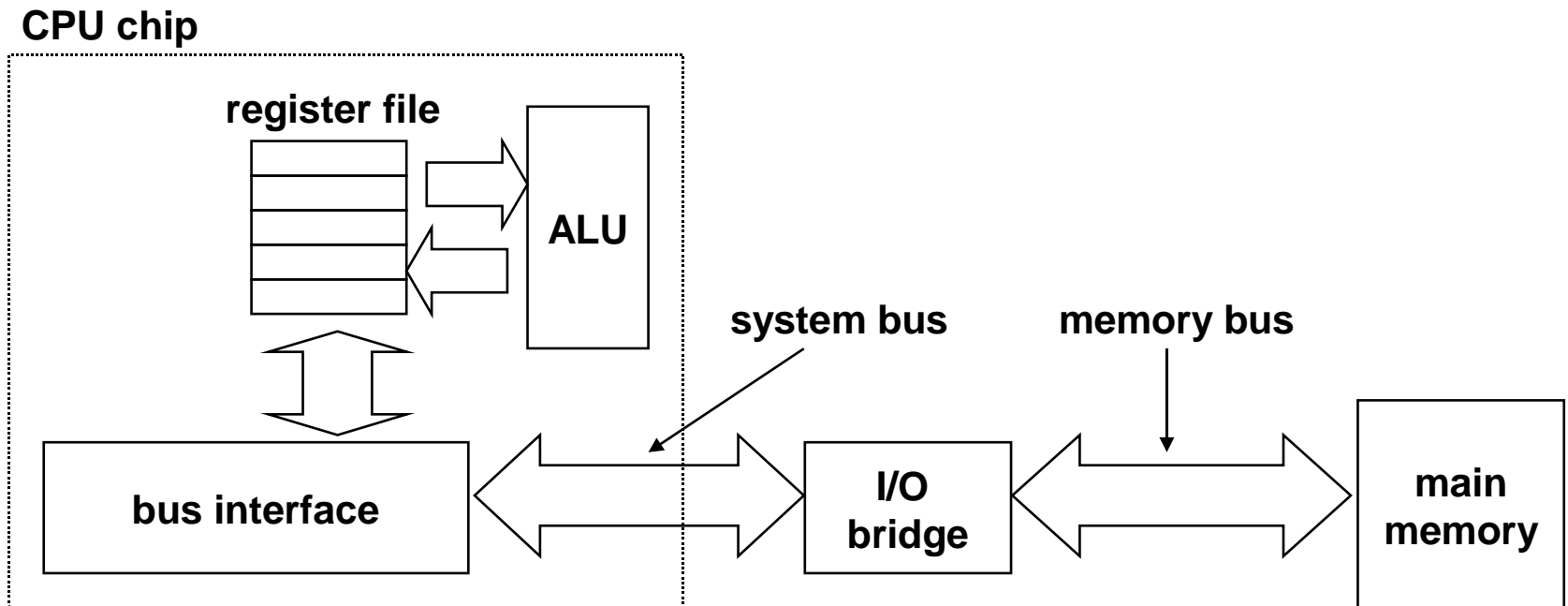
A 16Mbit chip can be organised as a 2048 x 2048 x 4bit array

- Reduces number of address pins
 - Multiplex row address and column address
 - 11 pins to address ($2^{11}=2048$)
 - Adding one more pin doubles range of values so x4 capacity (2^{12} x4 Capacity with 2^{11})

Bus structure connecting CPU and memory

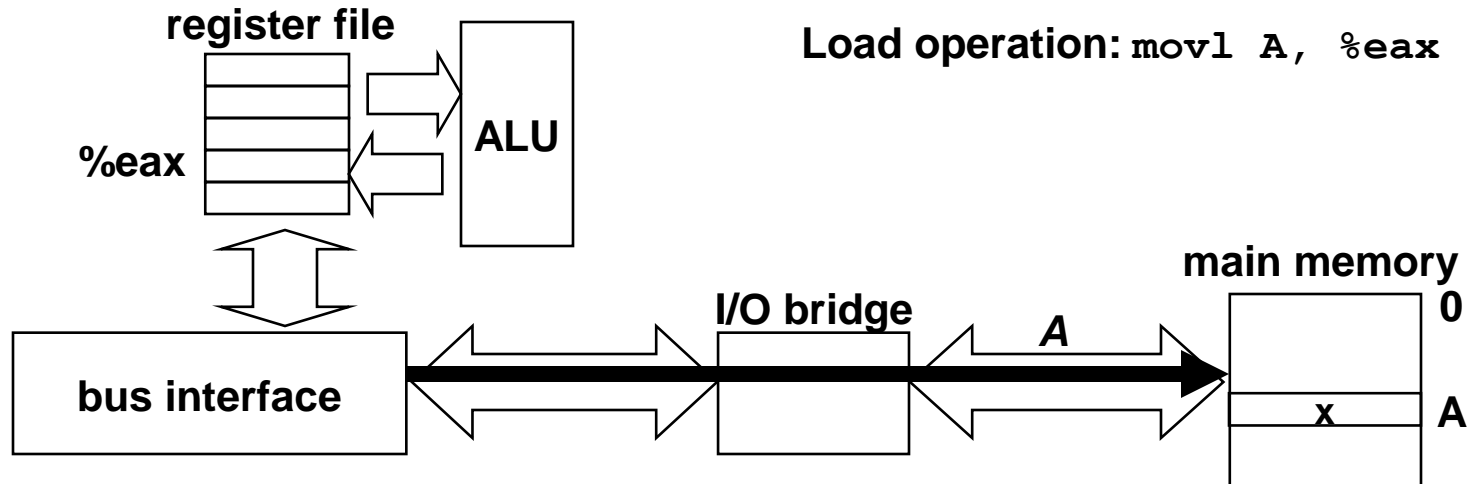
A *bus* is a collection of parallel wires that carry address, data, and control signals.

Buses are typically shared by multiple devices.



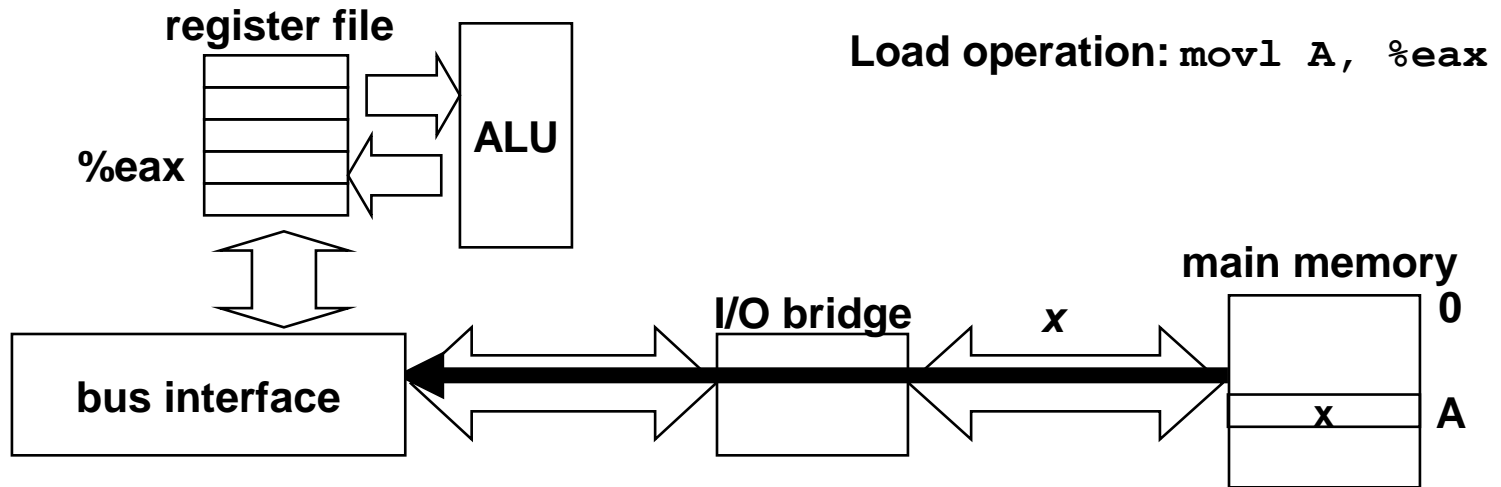
Memory read transaction

CPU places address **A** on the memory bus. **(1)**



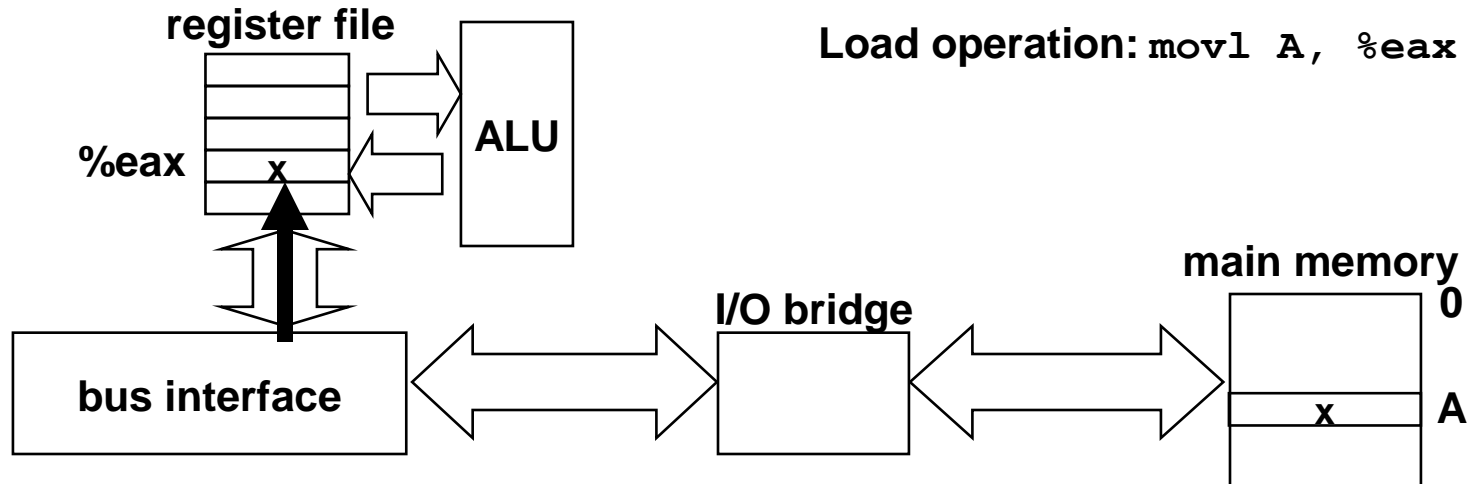
Memory read transaction

(2)
Main memory reads A from the memory bus, retrieves word x , and places it on the bus.



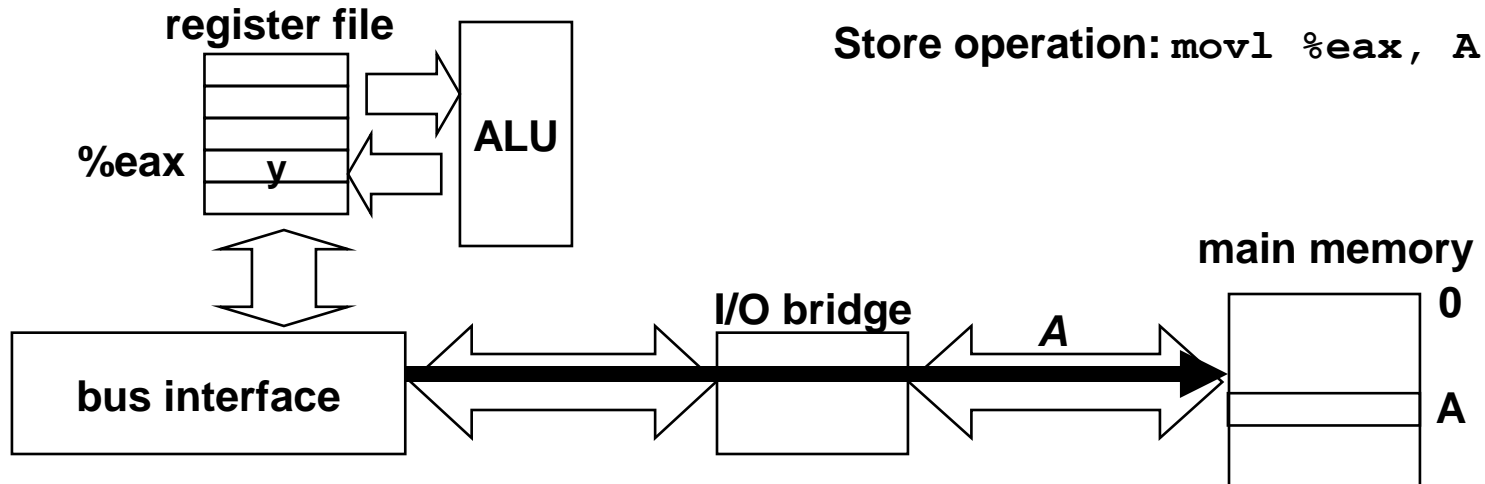
Memory read transaction

CPU read word x from the **(3)** bus and copies it into register `%eax`.



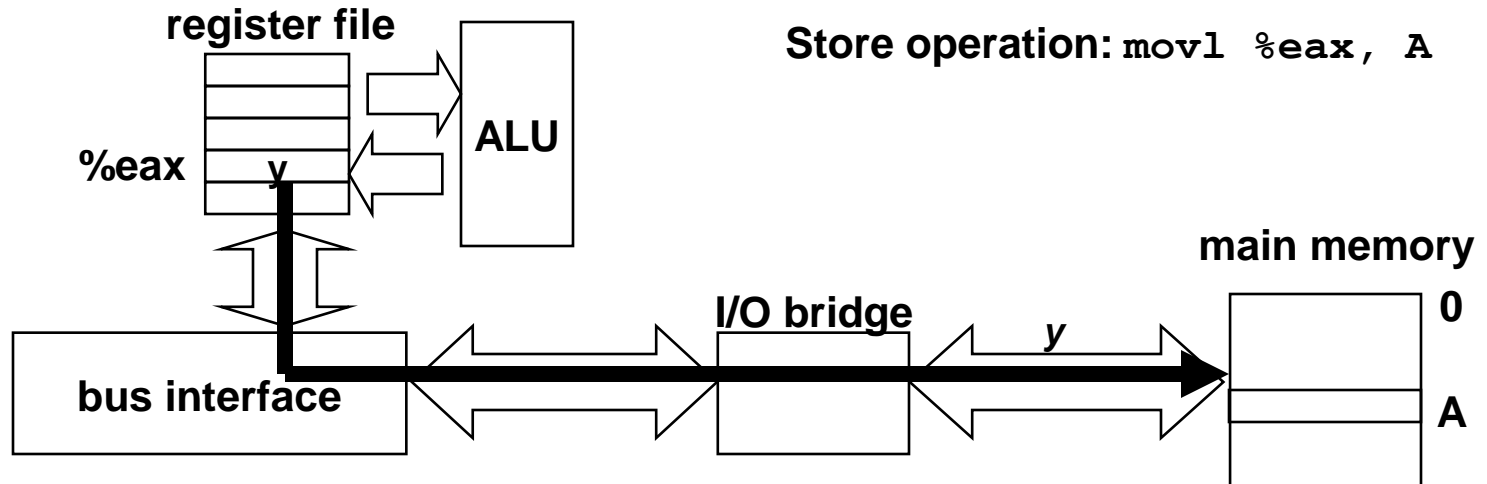
Memory write transaction

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



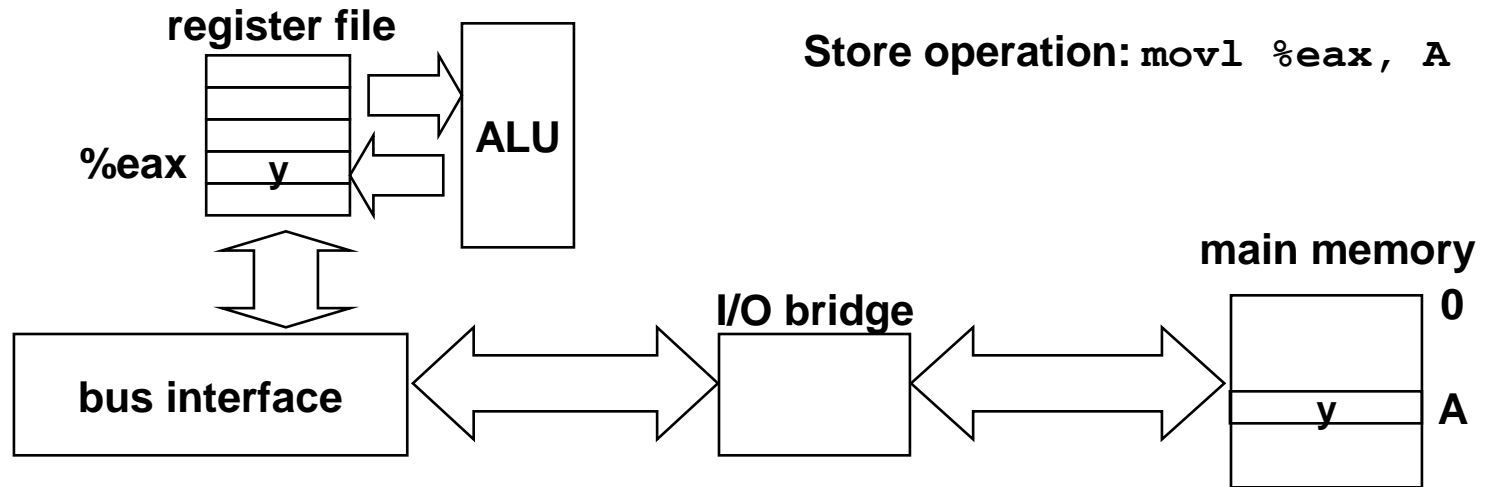
Memory write transaction

CPU places data word y on the bus. **(2)**



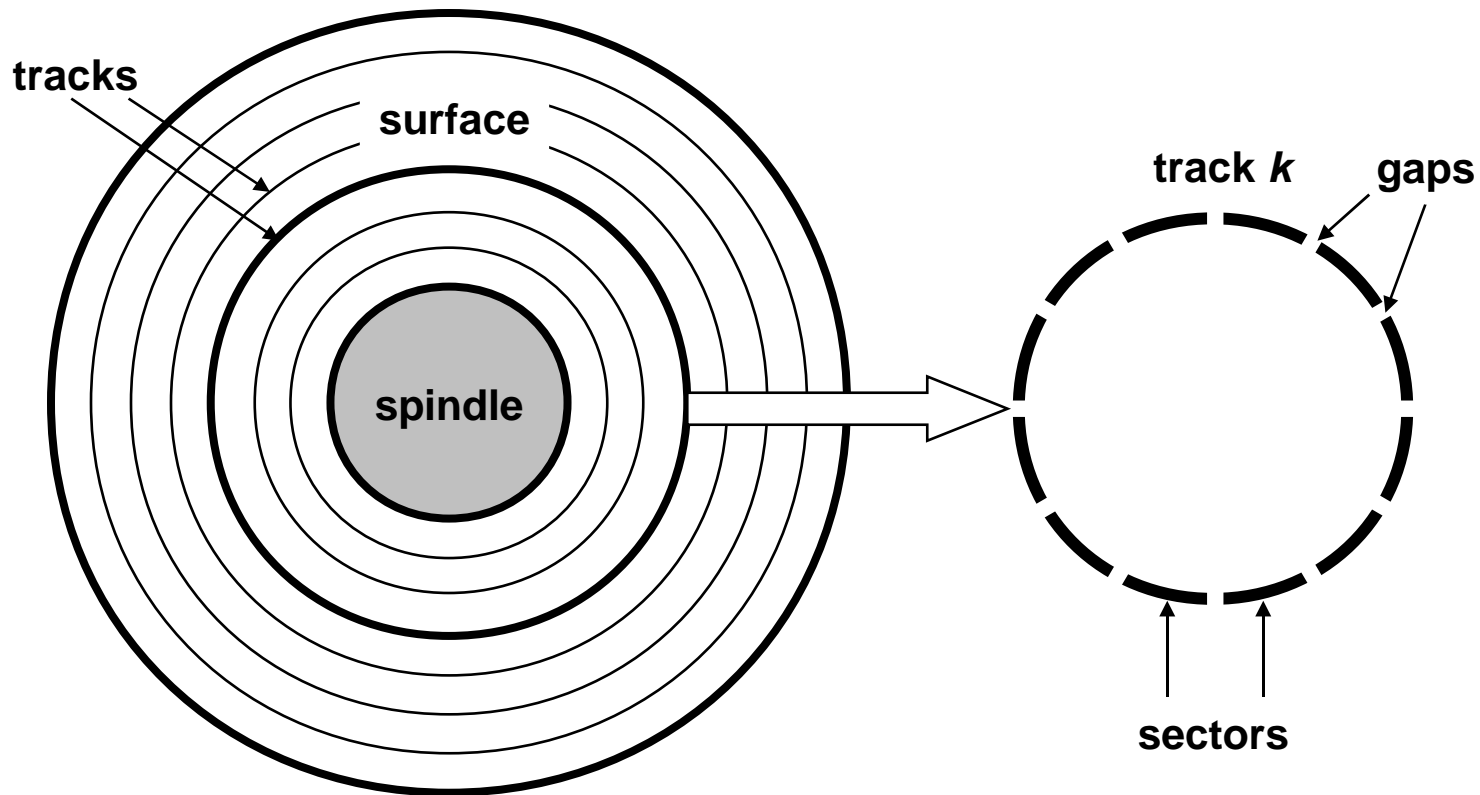
Memory write transaction

(3)
Main memory read data word y from the bus and stores it at address A .



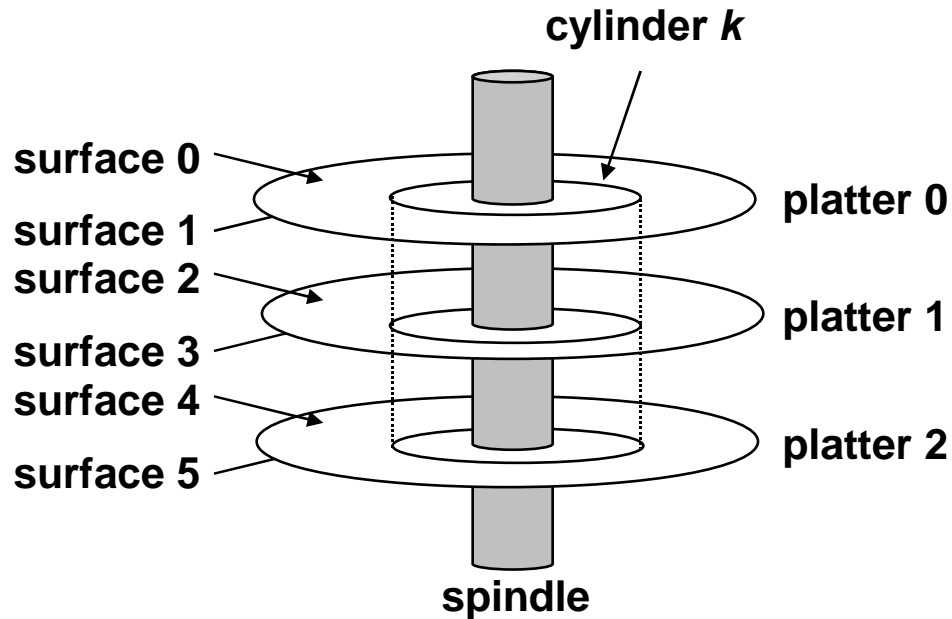
Disk

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: maximum number of bytes that can be stored.

- Vendors express capacity in units of gigabytes (GB), where $1 \text{ GB} = 10^6$.

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.

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

Computing disk capacity

$$\text{Capacity} = (\# \text{ bytes/sector}) \times (\text{avg. } \# \text{ sectors/track}) \times (\# \text{ tracks/surface}) \times (\# \text{ surfaces/platter}) \times (\# \text{ platters/disk})$$

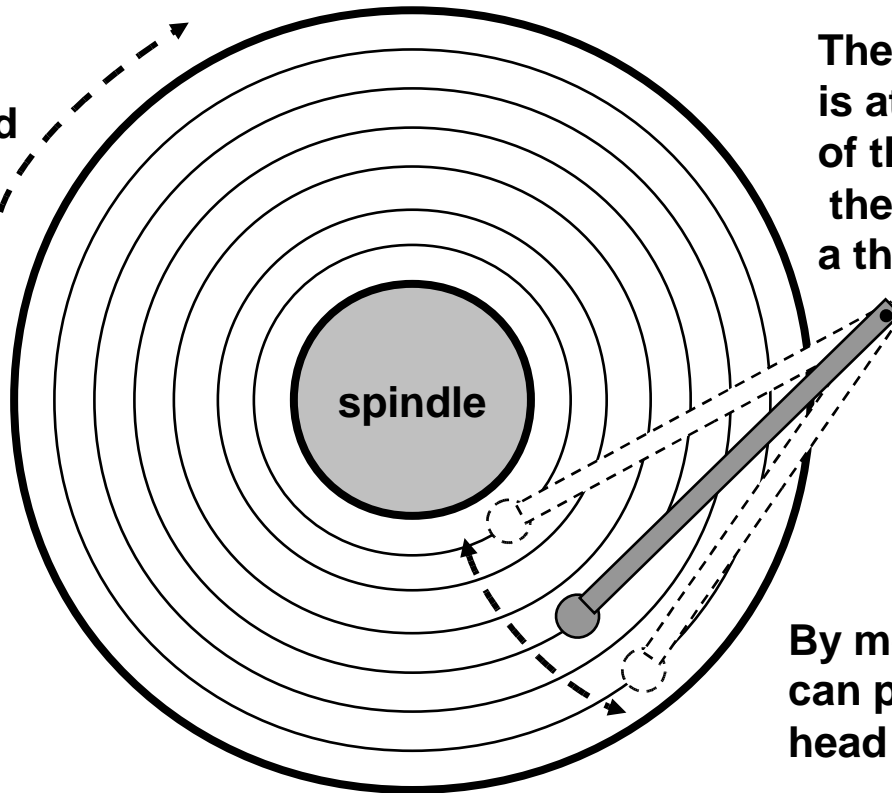
Example:

- 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)

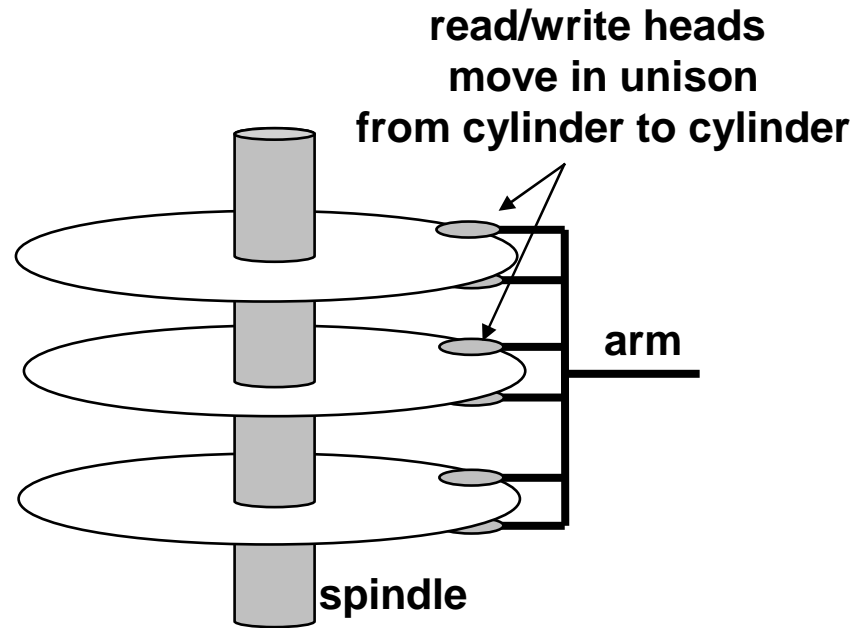
The disk surface spins at a fixed rotational rate



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 access

Average time to access ~~some~~ **time** target sector approximated by :

- $T_{\text{access}} = T_{\text{avg seek}} + T_{\text{avg rotation}} + T_{\text{avg transfer}}$

Seek time

- Time to position heads over cylinder containing target sector.
- Typical $T_{\text{avg seek}} = 9 \text{ ms}$

Rotational latency

- Time waiting for first bit of target sector to pass under r/w head.
- $T_{\text{avg rotation}} = 1/2 \times 1/\text{RPMs} \times 60 \text{ sec}/1 \text{ min}$

Transfer time

- Time to read the bits in the target sector.
- $T_{\text{avg transfer}} = 1/\text{RPM} \times 1/(\text{avg \# sectors/track}) \times 60 \text{ secs}/1 \text{ min.}$

Disk access time example

Given:

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

Derived:

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

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 4ns/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, ...)

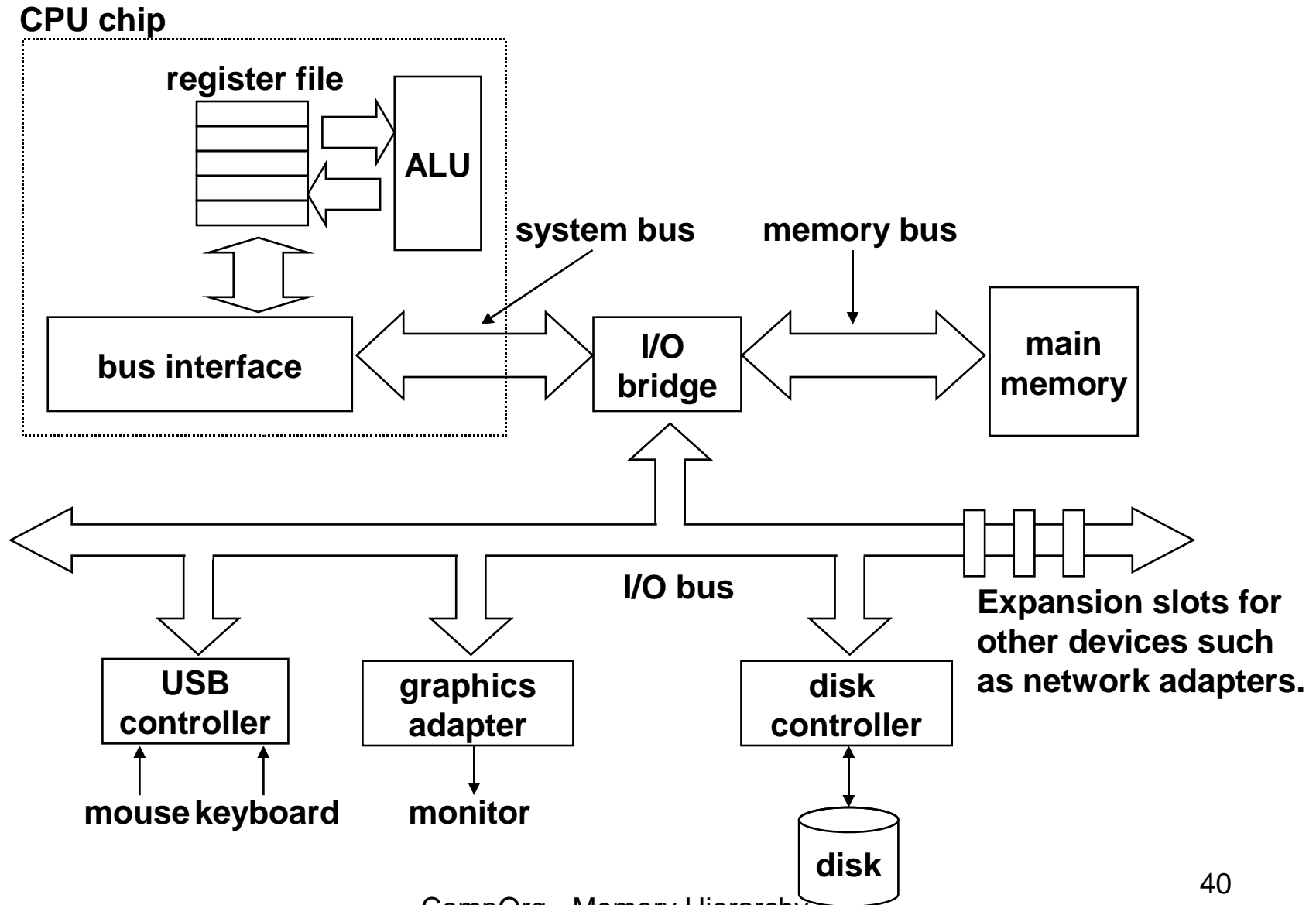
Mapping between logical blocks and actual (physical) sectors

- Maintained by hardware/firmware device called disk controller.
- Converts requests for logical blocks into (surface, track, sector) triples.

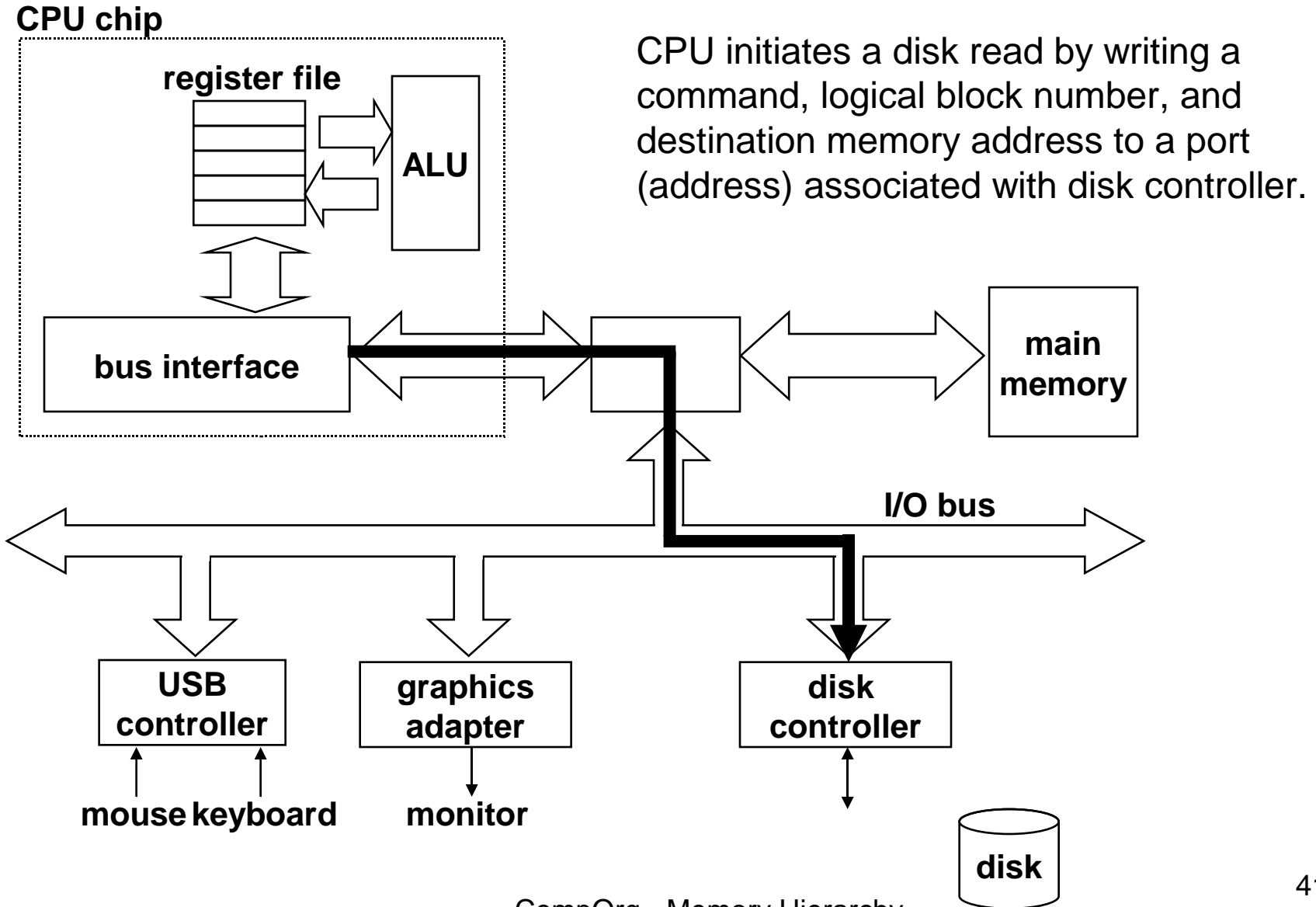
Allows controller to set aside spare cylinders for each zone.

- Accounts for the difference in “formatted capacity” and “maximum capacity”.

Bus structure connecting I/O and CPU



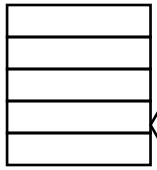
Reading a disk sector (1)



Reading a disk sector (2)

CPU chip

register file



ALU

bus interface

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

main memory

I/O bus

USB

controller

mouse keyboard

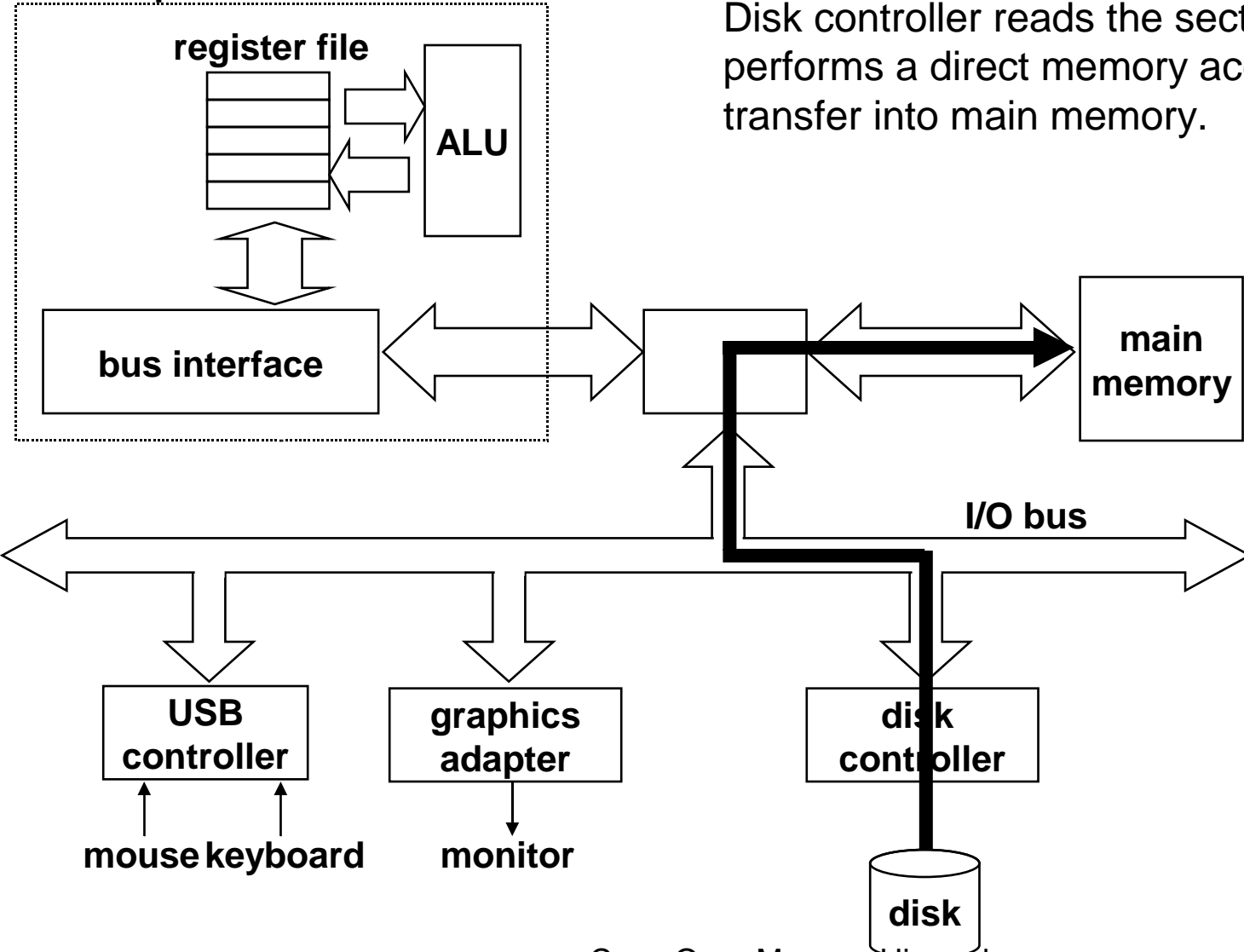
graphics

adapter

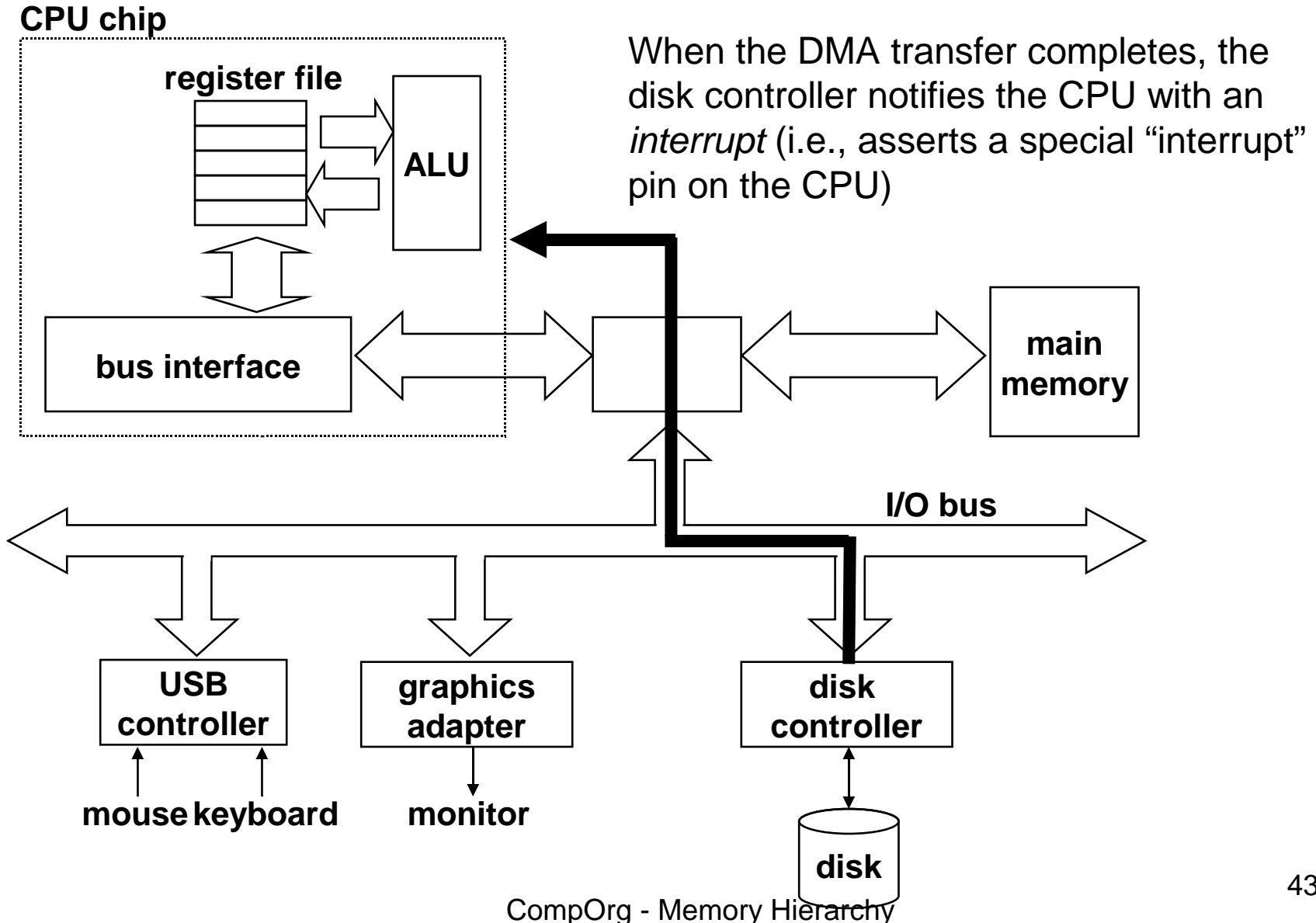
monitor

disk
controller

disk



Reading a disk sector (3)



Storage trends

SRAM

metric	1980	1985	1990	1995	2000	2000:1980
\$/MB	19,200	2,900	320	256	100	190
access (ns)	300	150	35	15	2	100

DRAM

metric	1980	1985	1990	1995	2000	2000:1980
\$/MB	8,000	880	100	30	1	8,000
access (ns)	375	200	100	70	60	6
typical size(MB)	0.064	0.256	4	16	64	1,000

Disk

metric	1980	1985	1990	1995	2000	2000:1980
\$/MB	500	100	8	0.30	0.05	10,000
access (ms)	87	75	28	10	8	11
typical size(MB)	1	10	160	1,000	9,000	9,000

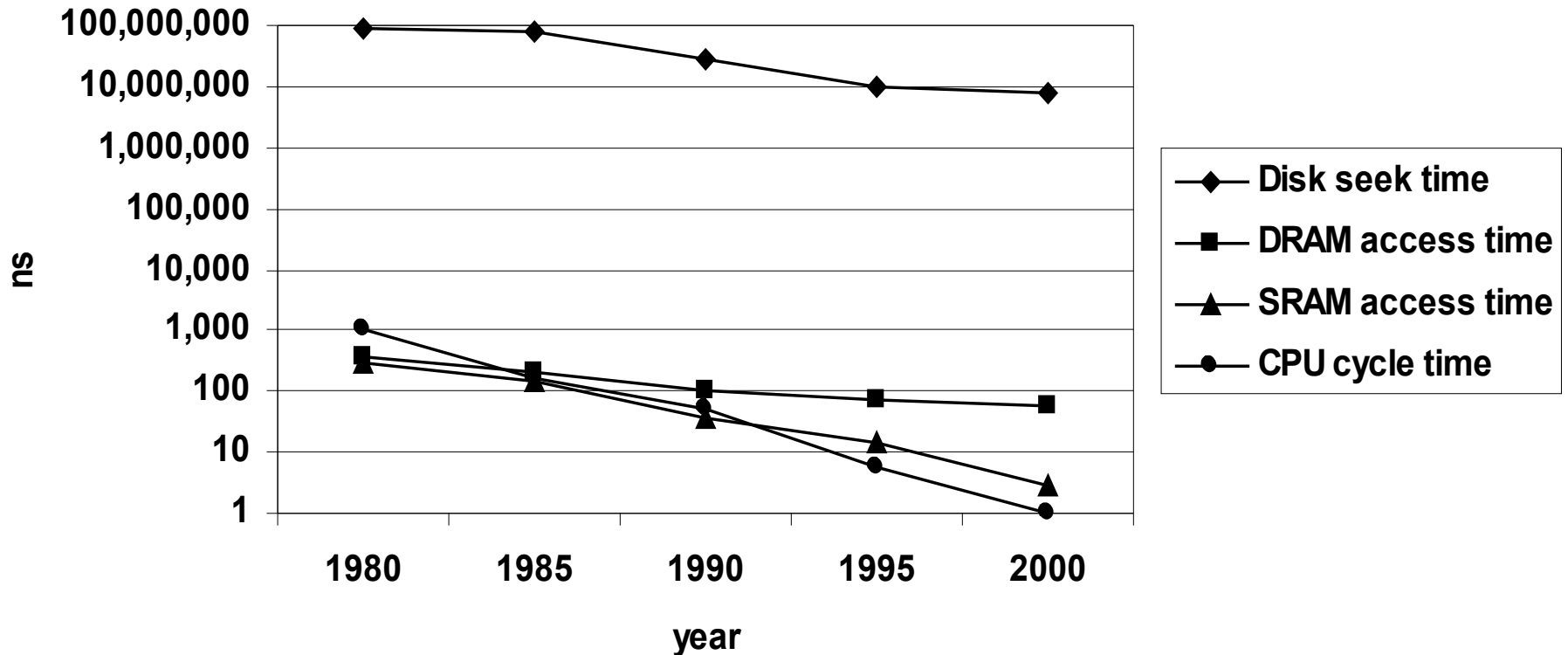
(Culled from back issues of Byte and PC Magazine)⁴⁴

CPU clock rates

	1980	1985	1990	1995	2000	<i>2000:1980</i>
processor	8080	286	386	Pent	P-III	
clock rate(MHz)	1	6	20	150	750	750
cycle time(ns)	1,000	166	50	6	1.6	750

The CPU-Memory

The increasing gap between **Gap** DRAM, disk, and CPU speeds.



Memory

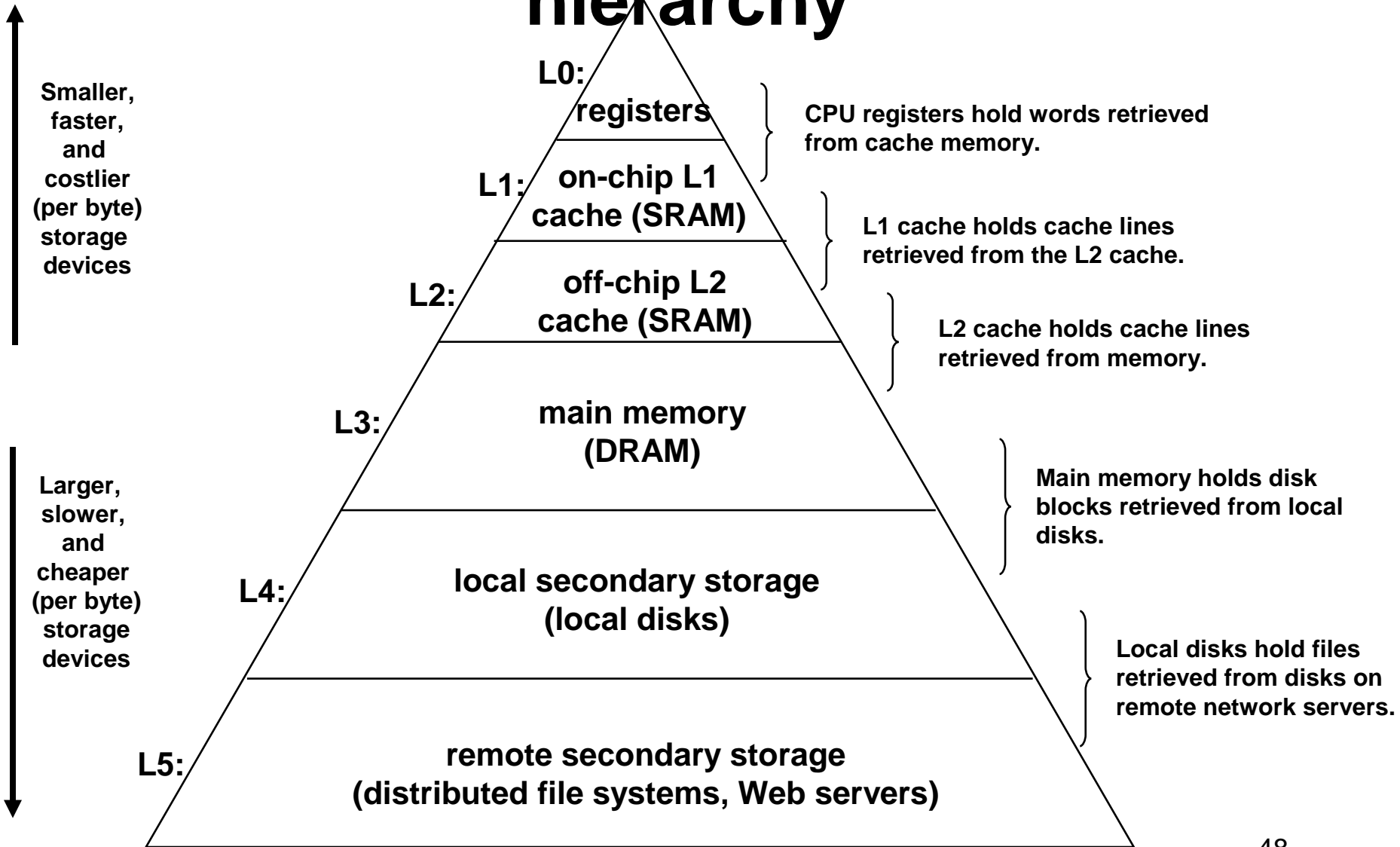
Some fundamental and enduring properties of hardware and software:

- **Fast storage technologies cost more per byte and have less capacity.**
- **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.

Suggest an approach for organizing memory and storage systems known as a “memory hierarchy”.

An example memory hierarchy



Cache

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?

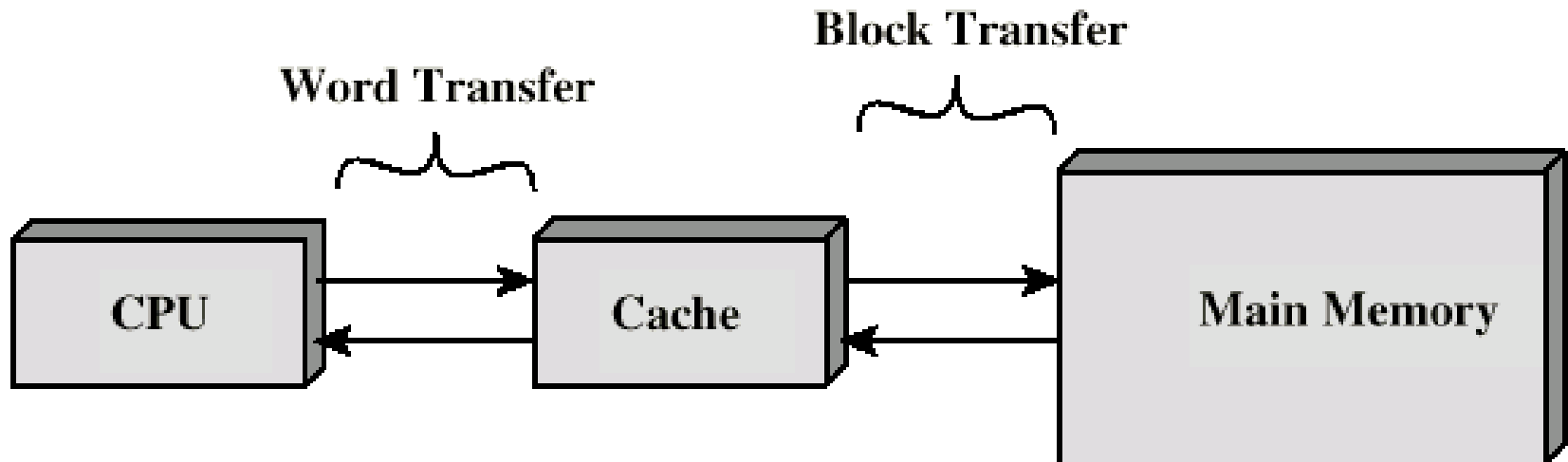
- Programs tend to access the data at level k more often than they access the data at level $k+1$.
- Thus, the storage at level $k+1$ can be slower, and thus larger and cheaper per bit.
- Net effect: A large pool of memory that costs as much as the cheap storage near the bottom, but that serves data to programs at the rate of the fast storage near the top.

Cache

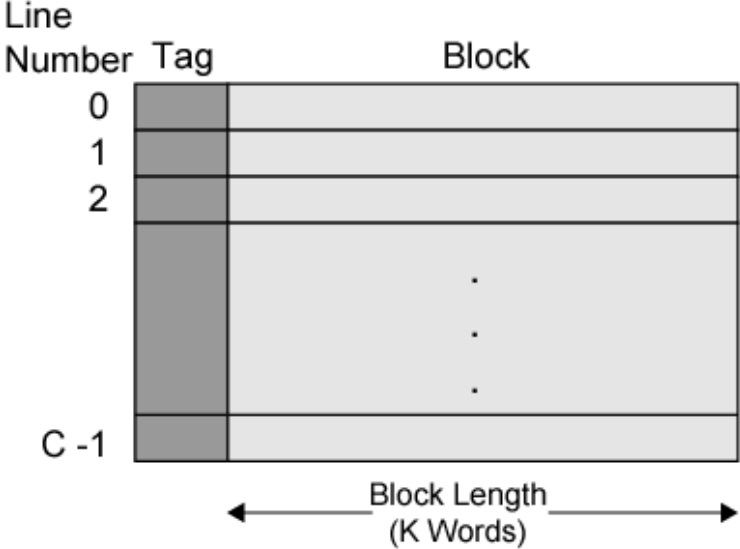
Small amount of fast memory

Sits between normal main memory and CPU

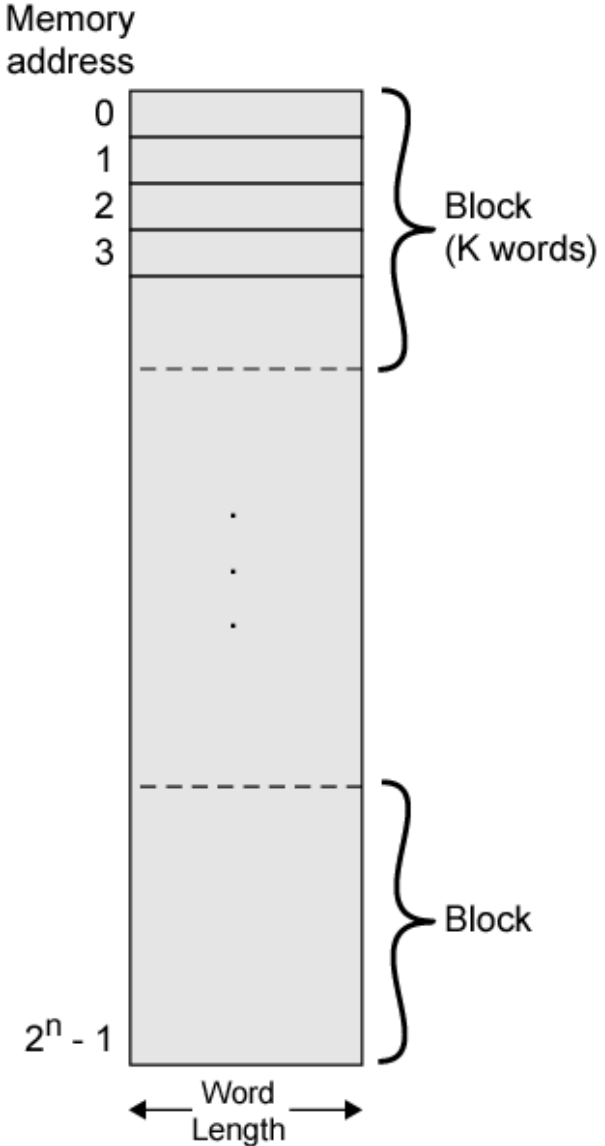
May be located on CPU chip or module



Cache



(a) Cache



(b) Main memory

Cac

CPU requests contents of memory location

Check cache for this data

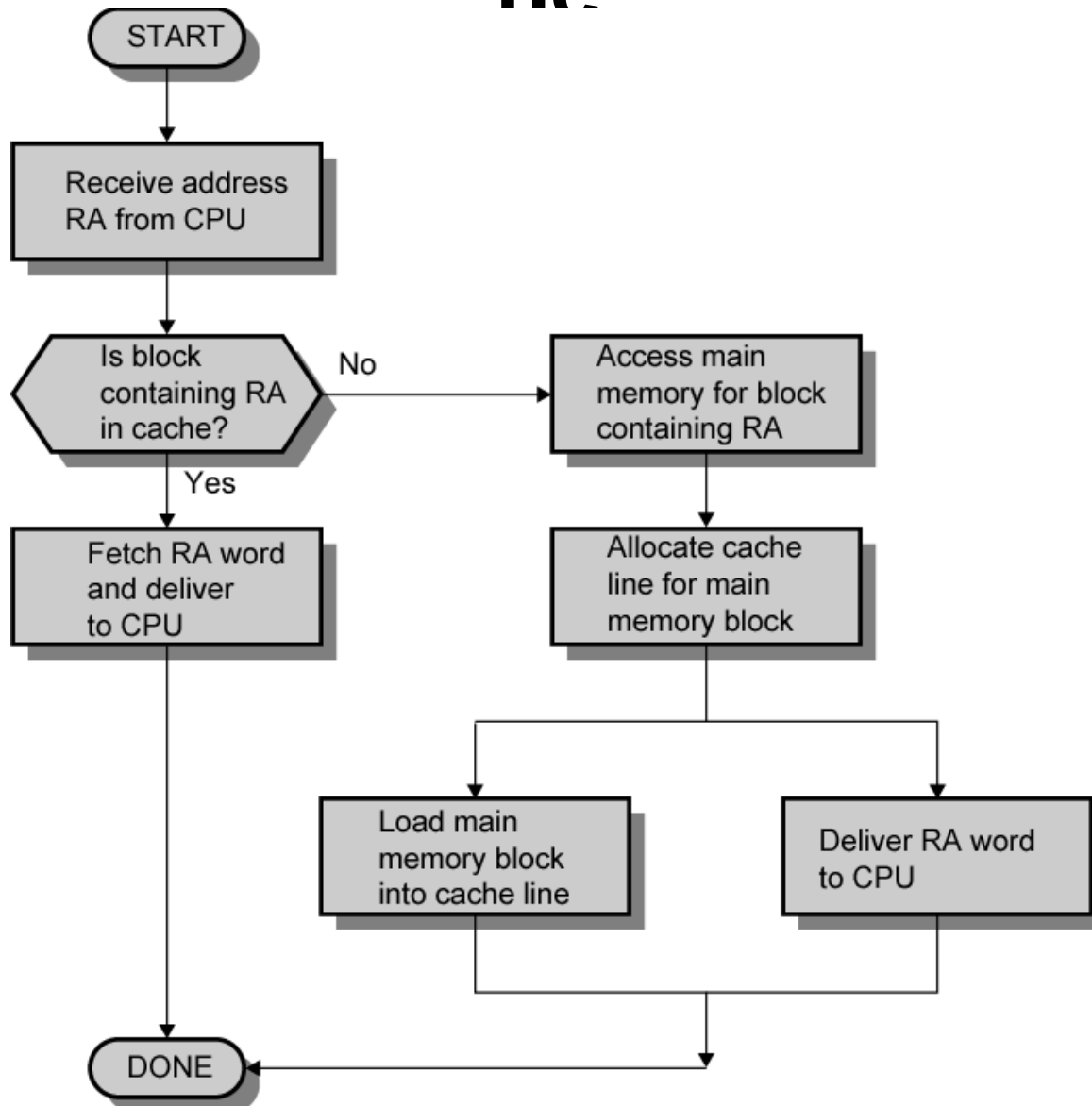
If present, get from cache (fast)

If not present, read required block from main memory to cache

Then deliver from cache to CPU

Cache includes tags to identify which block of main memory is in each cache slot

Cache



Cache Design

Size

Mapping Function

Replacement Algorithm

Write Policy

Block Size

Number of Caches

**Size
does
not
matter**

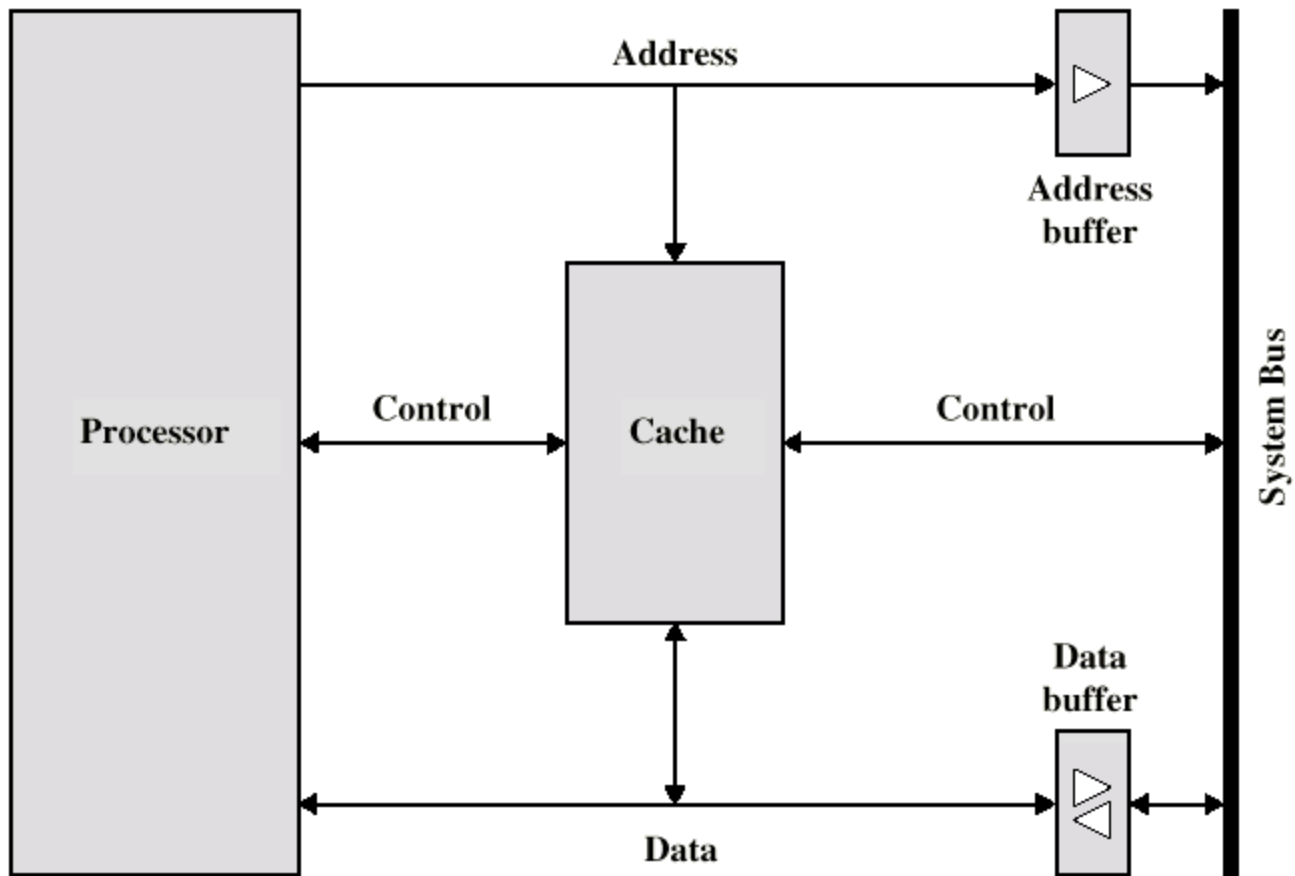
Cost

- More cache is expensive

Speed

- More cache is faster (up to a point)
- Checking cache for data takes time

Typical Cache



Comparison of Cache Sizes

Processor	Type	Year of Introduction	L1 cache ^a	L2 cache	L3 cache
IBM 360/85	Mainframe	1968	16 to 32 KB	—	—
PDP-11/70	Minicomputer	1975	1 KB	—	—
VAX 11/780	Minicomputer	1978	16 KB	—	—
IBM 3033	Mainframe	1978	64 KB	—	—
IBM 3090	Mainframe	1985	128 to 256 KB	—	—
Intel 80486	PC	1989	8 KB	—	—
Pentium	PC	1993	8 KB/8 KB	256 to 512 KB	—
PowerPC 601	PC	1993	32 KB	—	—
PowerPC 620	PC	1996	32 KB/32 KB	—	—
PowerPC G4	PC/server	1999	32 KB/32 KB	256 KB to 1 MB	2 MB
IBM S/390 G4	Mainframe	1997	32 KB	256 KB	2 MB
IBM S/390 G6	Mainframe	1999	256 KB	8 MB	—
Pentium 4	PC/server	2000	8 KB/8 KB	256 KB	—
IBM SP	High-end server/ supercomputer	2000	64 KB/32 KB	8 MB	—
CRAY MTAb	Supercomputer	2000	8 KB	2 MB	—
Itanium	PC/server	2001	16 KB/16 KB	96 KB	4 MB
SGI Origin 2001	High-end server	2001	32 KB/32 KB	4 MB	—
Itanium 2	PC/server	2002	32 KB	256 KB	6 MB
IBM POWER5	High-end server	2003	64 KB	1.9 MB	57 MB
CRAY XD-1	Supercomputer	2004	64 KB/64 KB	1MB	—

Map pin g Func tio n

Cache of 64kByte

Cache block of 4 bytes

- i.e. cache is 16k (2^{14}) lines of 4 bytes

16MBytes main memory

24 bit address

- ($2^{24}=16M$)

Direct Mapping

Each block of main memory maps to only one cache line

- i.e. if a block is in cache, it must be in one specific place

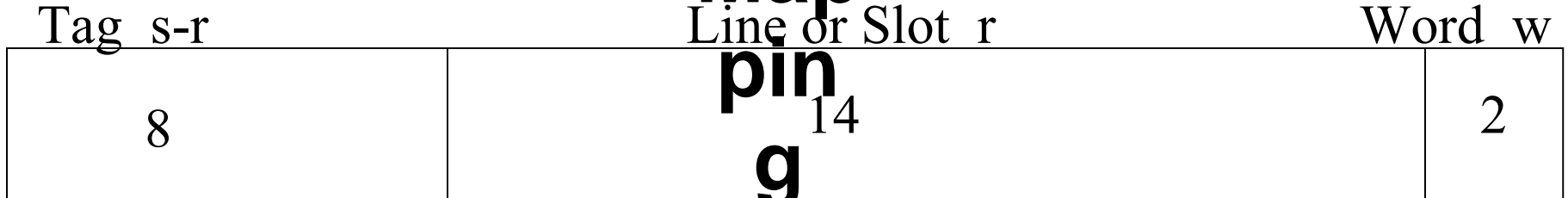
Address is in two parts

Least Significant w bits identify unique word

Most Significant s bits specify one memory block

The MSBs are split into a cache line field r and a tag of $s-r$ (most significant)

Direct Mapping Address Structure



24 bit address

2 bit word identifier (4 byte block)

22 bit block identifier

- 8 bit tag (=22-14)
- 14 bit slot or line

No two blocks in the same line have the same Tag field

Check contents of cache by finding line and checking Tag

Direct

Cache line

0

1

m-1

Main Memory blocks held

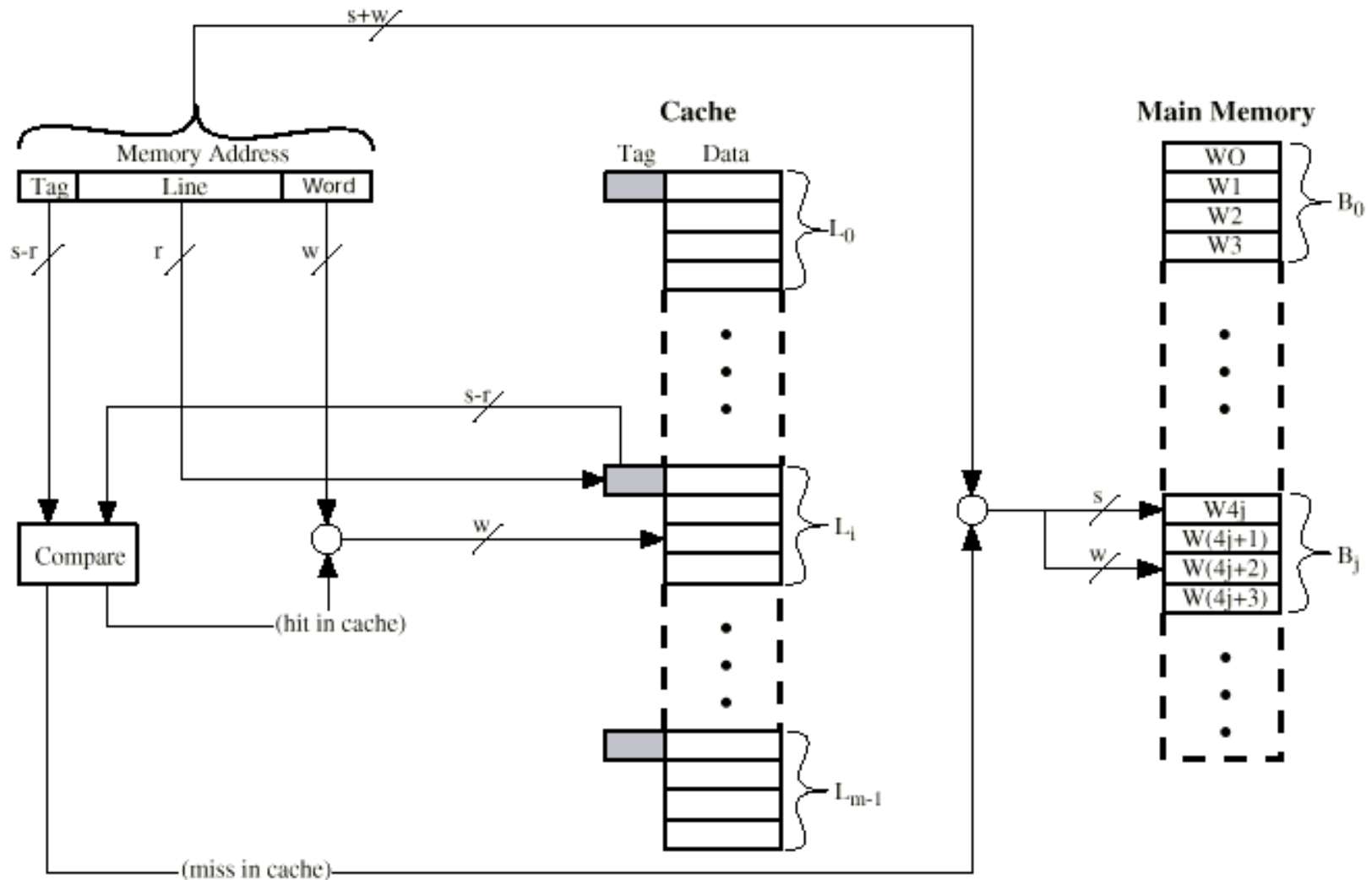
0, m, 2m, 3m, ..., 2s-m

1, m+1, 2m+1, ..., 2s-m+1

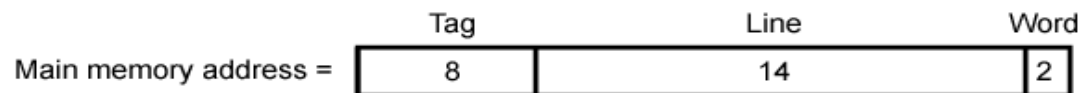
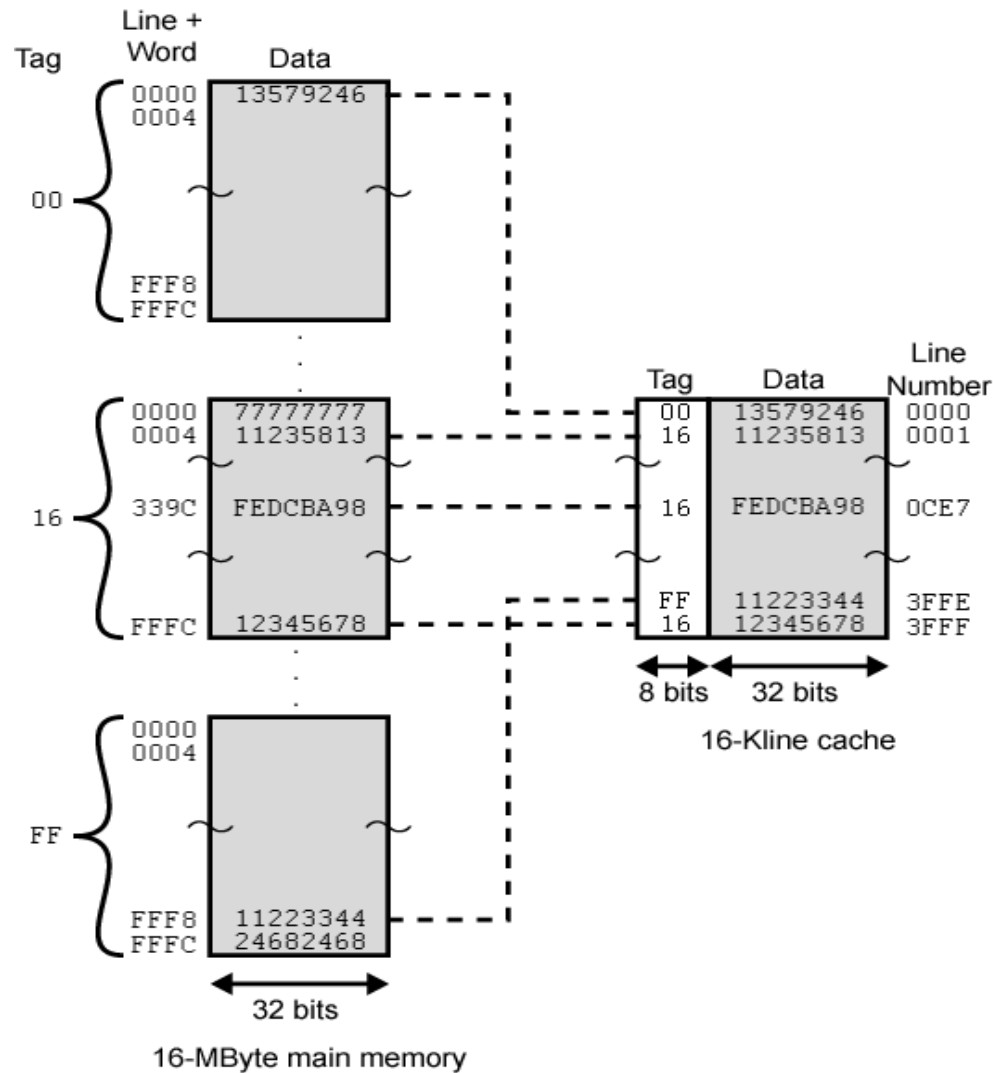
m-1, 2m-1, 3m-1, ..., 2s-1

Mapping Cache Line Table

Direct Man



Direct Mapping Example



Direct Mapping

Address length = $(s + w)$ bits

Number of addressable units = 2^{s+w} words or bytes

Block size = line size = 2^w words or bytes

Number of blocks in main memory = $2^{s+w} / 2^w = 2^s$

Number of lines in cache = $m = 2^r$

Size of tag = $(s - r)$ bits

Summary

Direct Mapping

Simple

Inexpensive

Fixed location for given block

- If a program accesses 2 blocks that map to the same line repeatedly, cache misses are very high

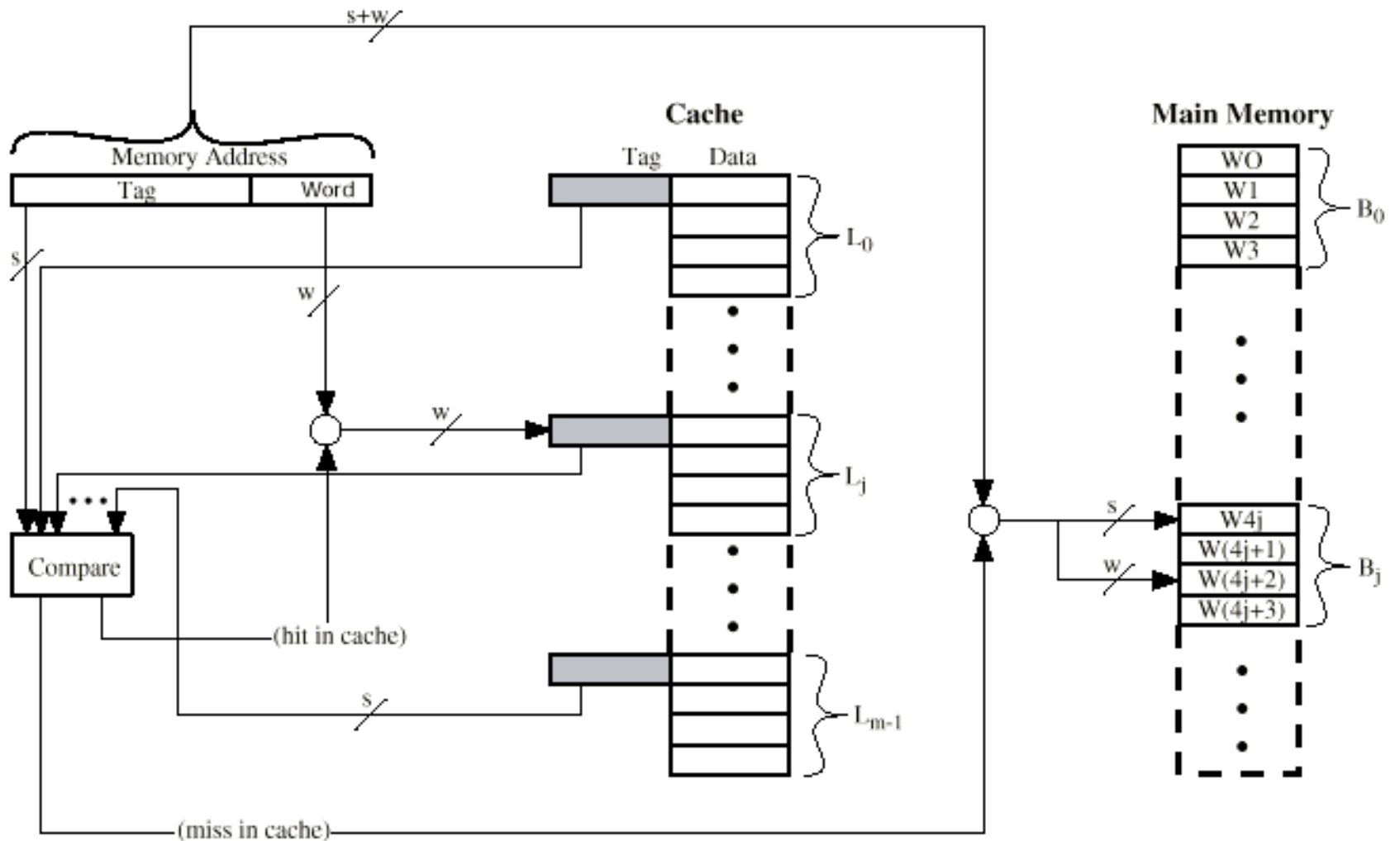
pros & cons

Ass

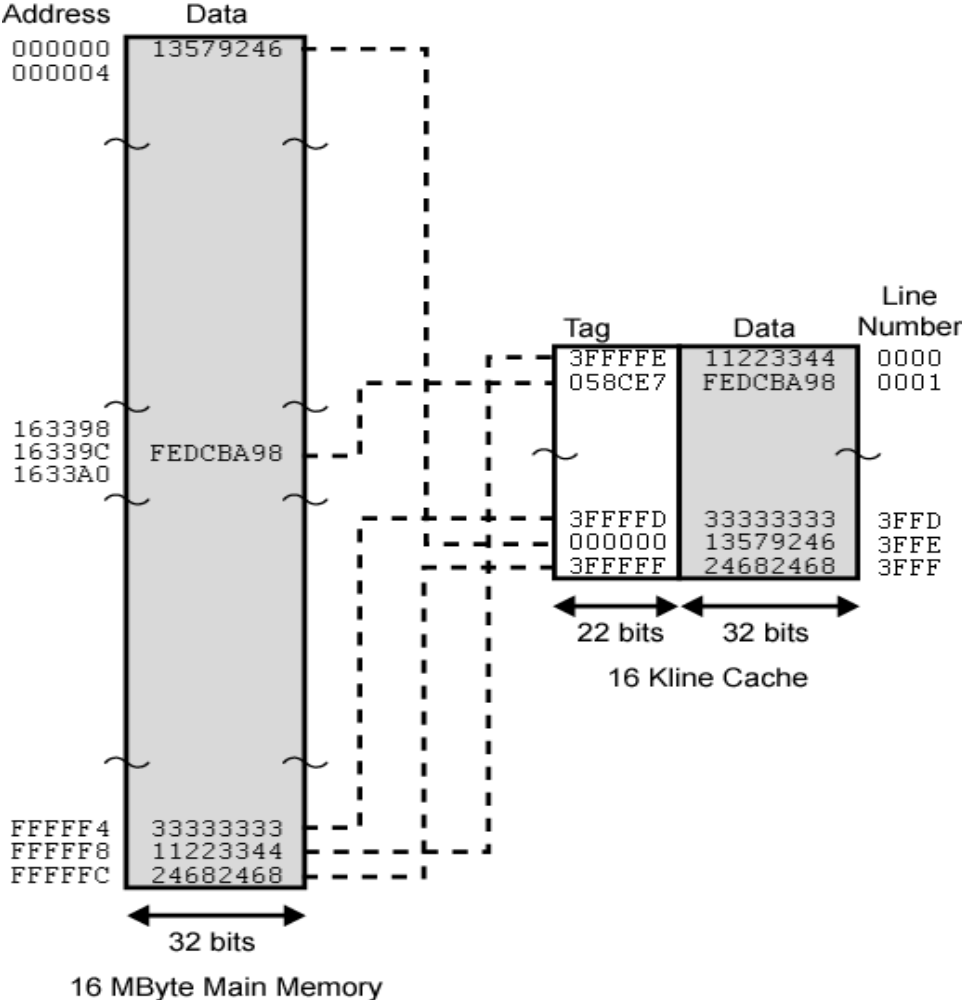
A main memory block can load into any line of cache
Memory address is interpreted as tag and word
Tag uniquely identifies block of memory
Every line's tag is examined for a match
Cache searching gets expensive

Associative Mapping

Full y Ass



Associative Mapping Example



Associative Map

Tag 22 bit

Word
2 bit

22 bit tag stored with each 32 bit block of data

Compare tag field with tag entry in cache to check for hit

Least significant 2 bits of address identify which 16 bit word is required from 32 bit data block

e.g.

- | | Address | Tag | Data | Cache line |
|---|---------|---------|--------------|------------|
| • | FFFFFFC | FFFFFFC | 246824683FFF | |

Structure

e

Ass

Address length = $(s + w)$ bits

Number of addressable units = 2^{s+w} words or bytes

Block size = line size = 2^w words or bytes

Number of blocks in main memory = $2^{s+w} / 2^w = 2^s$

Number of lines in cache = undetermined

Size of tag = s bits

Associative Mapping Summary

Set

Cache is divided into a number of sets

Each set contains a number of lines

A given block maps to any line in a given set

- e.g. Block B can be in any line of set i

e.g. 2 lines per set

- 2 way associative mapping
- A given block can be in one of 2 lines in only one set

Set

13 bit set number

Ass

Block number in main memory is modulo 2^{13}

000000, 00A000, 00B000, 00C000 ... map to same set

Map

pin

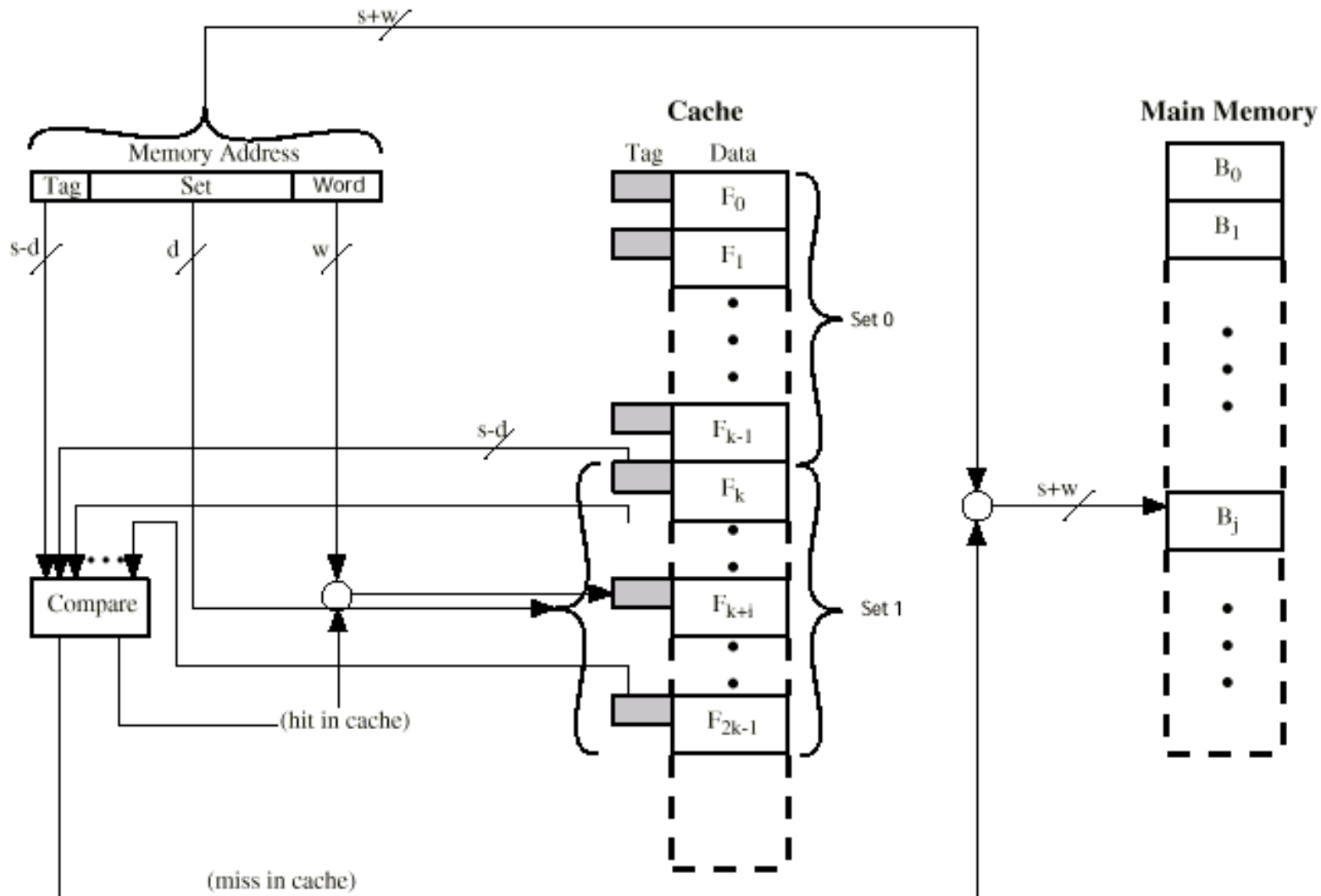
g

Exa

mpl

e

Two Way Set Associative Cache Organization



Set Associative Mapping

Tag 9 bit	Set 13 bit	Word 2 bit
-----------	------------	------------

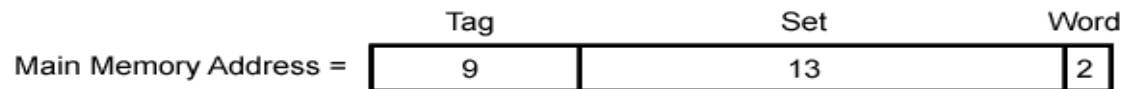
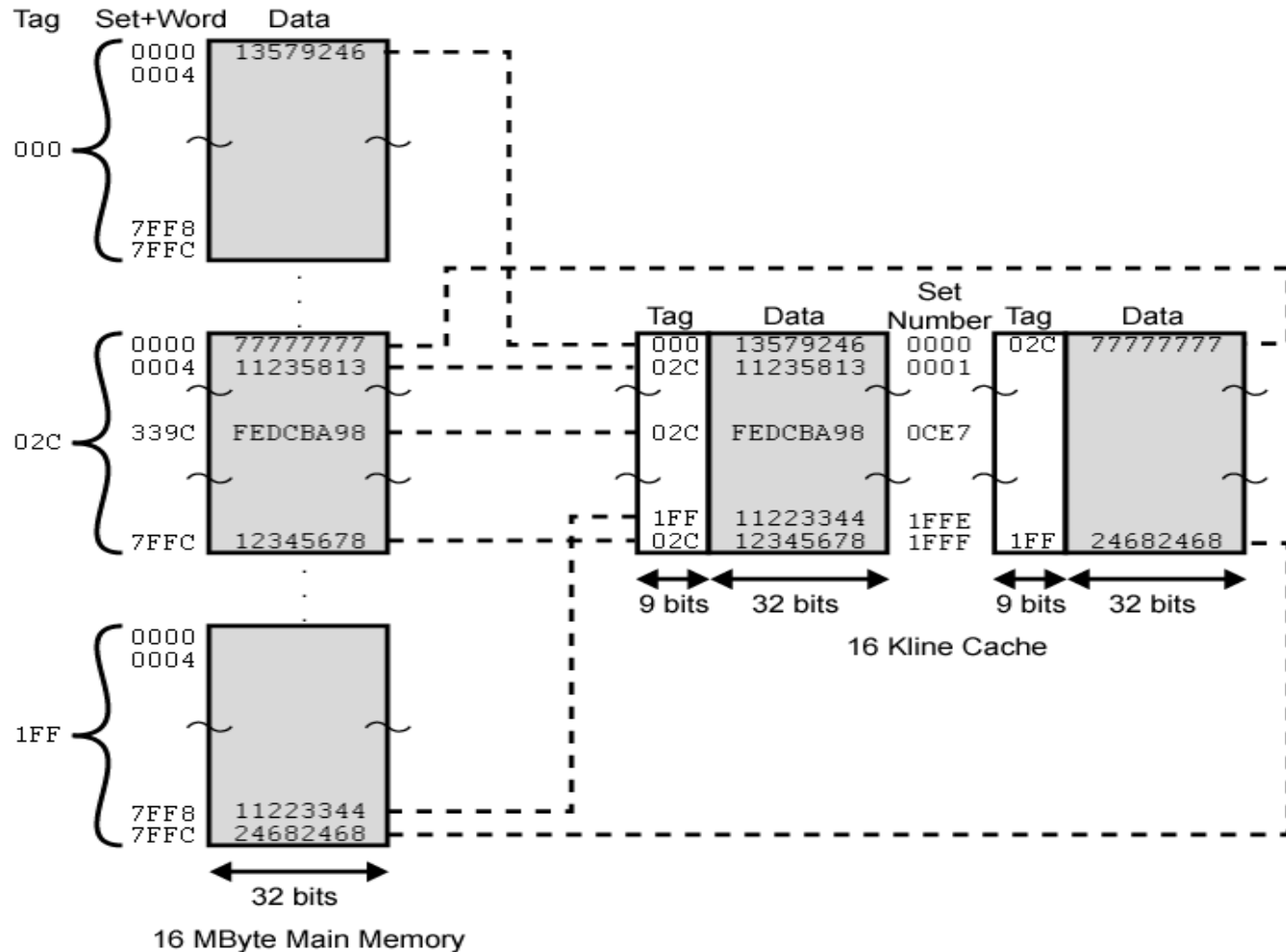
Use set field to determine cache set to look in
Compare tag field to see if we have a hit

e.g

- | Address | Tag | Data | Set number |
|----------------|----------|------|------------|
| • 1FF 7FFC 1FF | 12345678 | | 1FFF |
| • 001 7FFC 001 | 11223344 | | 1FFF |

Structur

Two Way set Associative Mapping Example



Set

Address length = $(s + w)$ bits

Number of addressable units = 2^{s+w} words or bytes

Block size = line size = 2^w words or bytes

Number of blocks in main memory = 2^d

Number of lines in set = k

Number of sets = $v = 2^d$

Number of lines in cache = $kv = k * 2^d$

Size of tag = $(s - d)$ bits

Associative Mapping

Summary

No choice

Each block only maps to one line

Replace that line

Rep

lace

ment

t

Alg

orit

hms

(1)

Dire

ct

map

pin

d

Replacement Algorithms (2)

Associative & Set Associative

Hardware implemented algorithm (speed)

Least Recently used (LRU)

e.g. in 2 way set associative

- Which of the 2 block is lru?

First in first out (FIFO)

- replace block that has been in cache longest

Least frequently used

- replace block which has had fewest hits

Random

Write

Must not overwrite a cache block unless main memory is up to date

Policy

Multiple CPUs may have individual caches

I/O may address main memory directly

Write

All writes go to main memory as well as cache

Multiple CPUs can monitor main memory traffic to keep local (to CPU) cache up to date

Lots of traffic

Slows down writes

Remember bogus write through caches!

Write

Updates initially made in cache only

Update bit for cache slot **is set** when update occurs

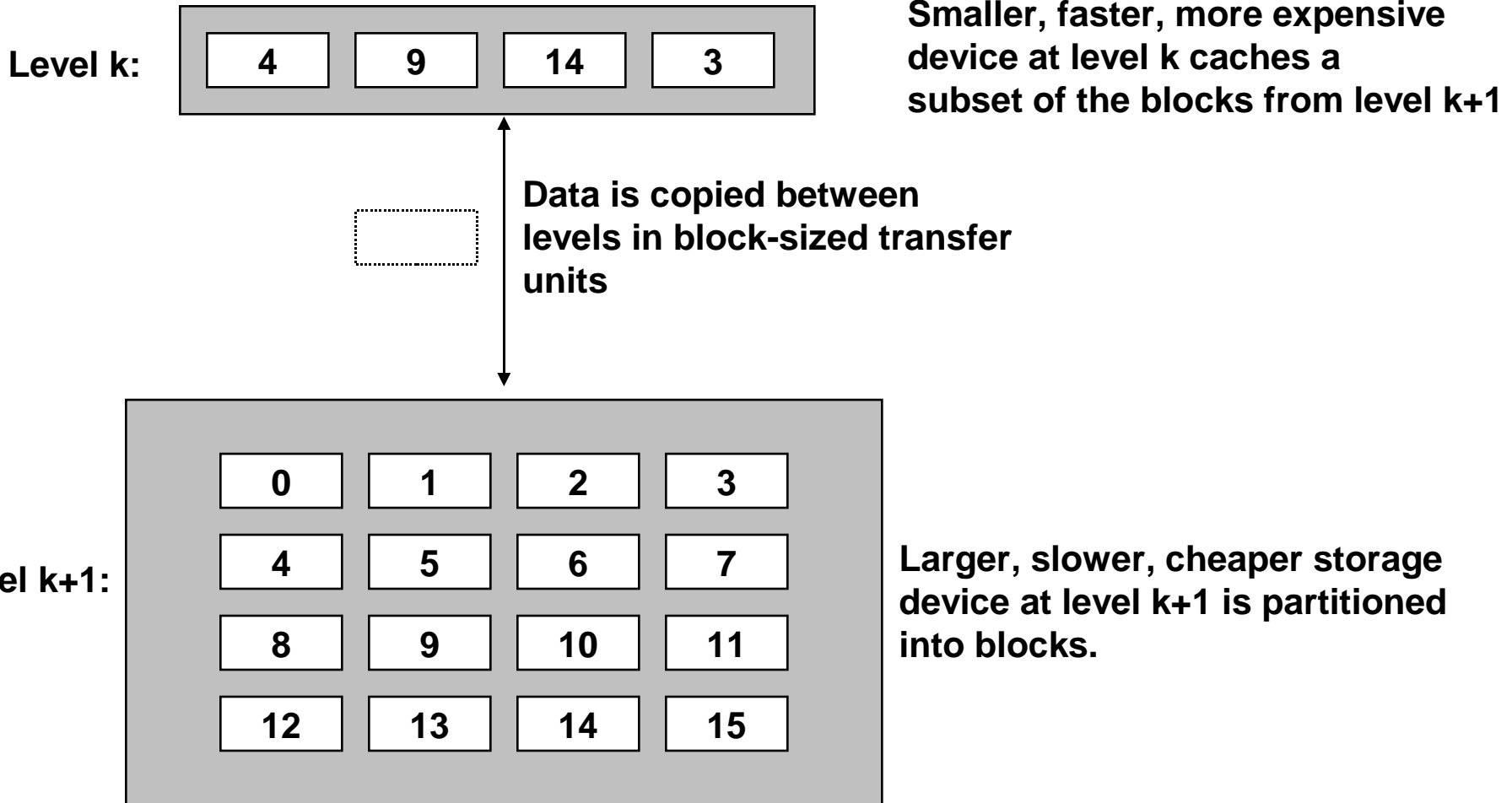
If block is to be replaced, **write** to main memory only if update bit is set

Other caches get out of sync

I/O must access main memory through cache

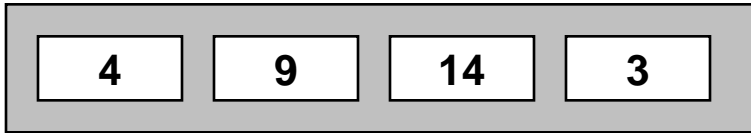
N.B. 15% of memory references are writes

Caching in a memory hierarchy

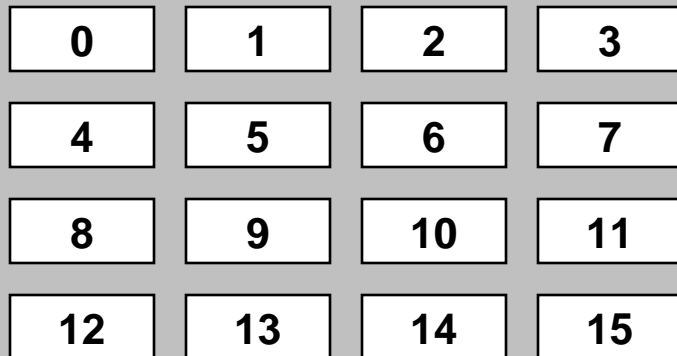


General caching concepts

Level k:



Level k+1:



Program needs object **d**, which is stored in some block **b**.

Cache hit

- Program finds **b** in the cache at level **k**. E.g. block 14.

Cache miss

- **b** is not at level **k**, so level **k** cache must fetch it from level **k+1**. E.g. block 12.
- If level **k** cache is full, then some current block must be *replaced (evicted)*.
- Which one? Determined by replacement policy. E.g. evict **least recently used block**.

General caching concepts

Types of cache misses:

- **Cold (compulsary) 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 $k+1$.
 - 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 hierarchy

Cache Type	What Cached	Where Cached	Latency (cycles)	Managed By
Registers	4-byte word	CPU registers	0	Compiler
TLB	Address translations	On-Chip TLB	0	Hardware
L1 cache	32-byte block	On-Chip L1	1	Hardware
L2 cache	32-byte block	Off-Chip L2	10	Hardware
Virtual Memory	4-KB page	Main memory	100	Hardware+ OS
Buffer cache	Parts of files	Main memory	100	OS
Network buffer cache	Parts of files	Local disk	10,000,000	AFS/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

Pen

tiu

m 4

Cache

he

80386 – no on chip cache

80486 – 8k using 16 byte lines and four way set associative organization

Pentium (all versions) – two on chip L1 caches

- Data & instructions

Pentium III – L3 cache added off chip

Pentium 4

- **L1 caches**
 - 8k bytes
 - 64 byte lines
 - four way set associative
- **L2 cache**
 - Feeding both L1 caches
 - 256k
 - 128 byte lines
 - 8 way set associative
- **L3 cache on chip**

Intel

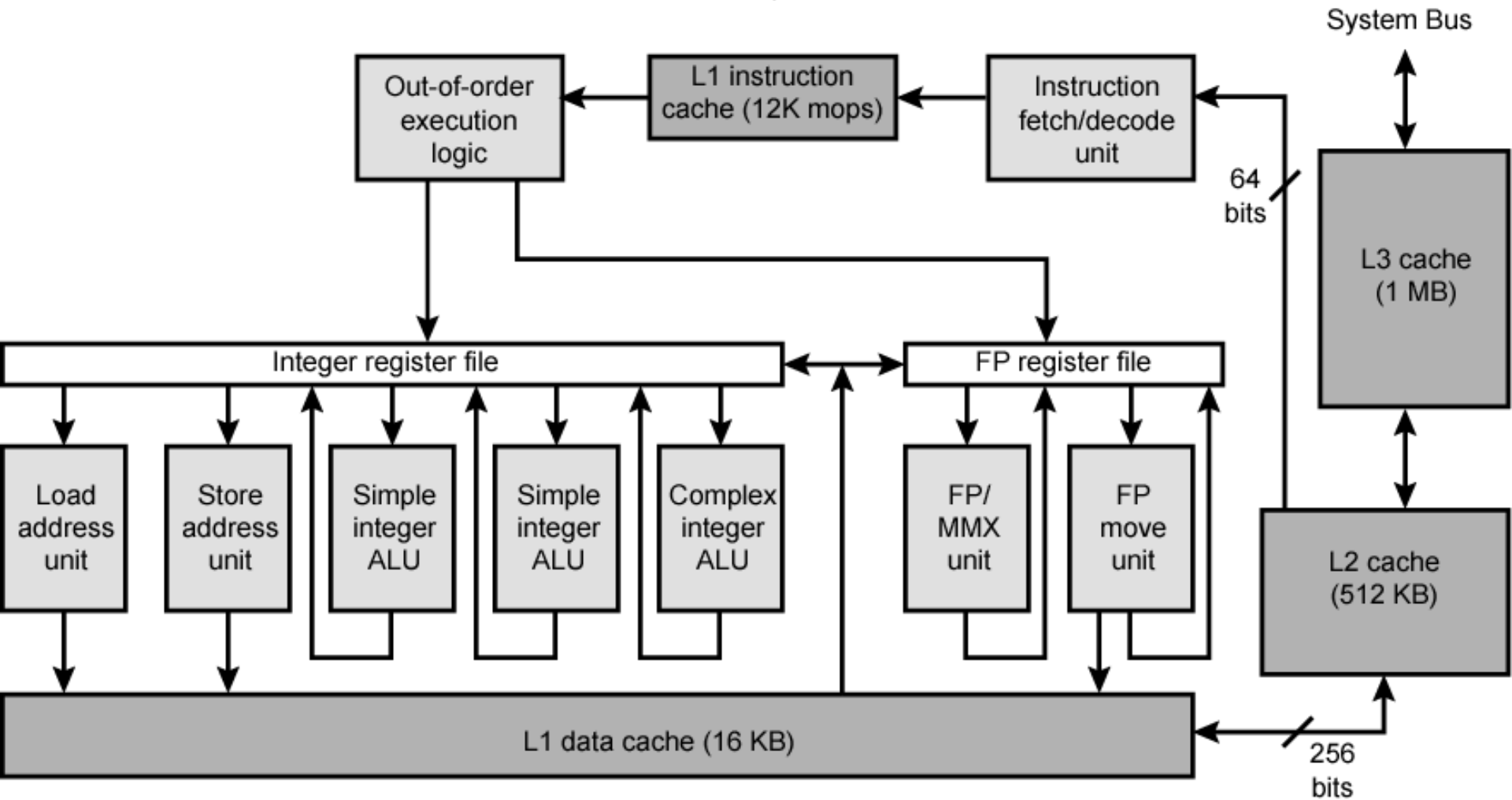
Cache

Evolution

on

Problem	Solution	Processor on which feature first appears
External memory slower than the system bus.	Add external cache using faster memory technology.	386
Increased processor speed results in external bus becoming a bottleneck for cache access.	Move external cache on-chip, operating at the same speed as the processor.	486
Internal cache is rather small, due to limited space on chip	Add external L2 cache using faster technology than main memory	486
Contention occurs when both the Instruction Prefetcher and the Execution Unit simultaneously require access to the cache. In that case, the Prefetcher is stalled while the Execution Unit's data access takes place.	Create separate data and instruction caches.	Pentium
Increased processor speed results in external bus becoming a bottleneck for L2 cache access.	Create separate back-side bus that runs at higher speed than the main (front-side) external bus. The BSB is dedicated to the L2 cache.	Pentium Pro
	Move L2 cache on to the processor chip.	Pentium II
Some applications deal with massive databases and must have rapid access to large amounts of data. The on-chip caches are too small.	Add external L3 cache.	Pentium III
	Move L3 cache on-chip.	Pentium 4

Pen tiu



Pen tiu m 4 Cor e Pro ces sor

Fetch/Decode Unit

- Fetches instructions from L2 cache
- Decode into micro-ops
- Store micro-ops in L1 cache

Out of order execution logic

- Schedules micro-ops
- Based on data dependence and resources
- May speculatively execute

Execution units

- Execute micro-ops
- Data from L1 cache
- Results in registers

Memory subsystem

- L2 cache and systems bus

Pen

Decodes instructions into RISC like micro-ops before L1 cache

Micro-ops fixed length

- Superscalar pipelining and scheduling

Pentium instructions long & complex

Performance improved by separating decoding from scheduling & pipelining

- (More later – ch14)

Data cache is write back

- Can be configured to write through

L1 cache controlled by 2 bits in register

- CD = cache disable
- NW = not write through
- 2 instructions to invalidate (flush) cache and write back then invalidate

L2 and L3 8-way set-associative

- Line size 128 bytes

ti m 4 Des ign Rea soni ng

PowerPC

601 – single 32kb 8 way set associative

603 – 16kb (2 x 8kb) two way set associative

604 – 32kb

620 – 64kb

G3 & G4

- 64kb L1 cache
 - 8 way set associative
- 256k, 512k or 1M L2 cache
 - two way set associative

G5

- 32kB instruction cache
- 64kB data cache

Cache Organization

Pow erD

