+

William Stallings
Computer Organization
and Architecture
9th Edition

# + Chapter 4

Cache Memory

# Memory

- The most common forms are:
  - Semiconductor memory
  - Magnetic surface memory
  - Optical
  - Magneto-optical

- Several physical characteristics of data storage are important:
  - Volatile memory
    - Information decays naturally or is lost when electrical power is switched off
  - Nonvolatile memory
    - Once recorded, information remains without deterioration until deliberately changed
    - No electrical power is needed to retain information
  - Magnetic-surface memories
    - Are nonvolatile
  - Semiconductor memory
    - May be either volatile or nonvolatile
  - Nonerasable memory
    - Cannot be altered, except by destroying the storage unit
    - Semiconductor memory of this type is known as read-only memory (ROM)

- For random-access memory the organization is a key design issue
  - Organization refers to the physical arrangement of bits to form words

# + Memory Hierarchy

- Design constraints on a computer's memory can be summed up by three questions:
  - How much, how fast, how expensive

- There is a trade-off among capacity, access time, and cost
  - Faster access time, greater cost per bit
  - Greater capacity, smaller cost per bit
  - Greater capacity, slower access time

- The way out of the memory dilemma is not to rely on a single memory component or technology, but to employ a memory hierarchy
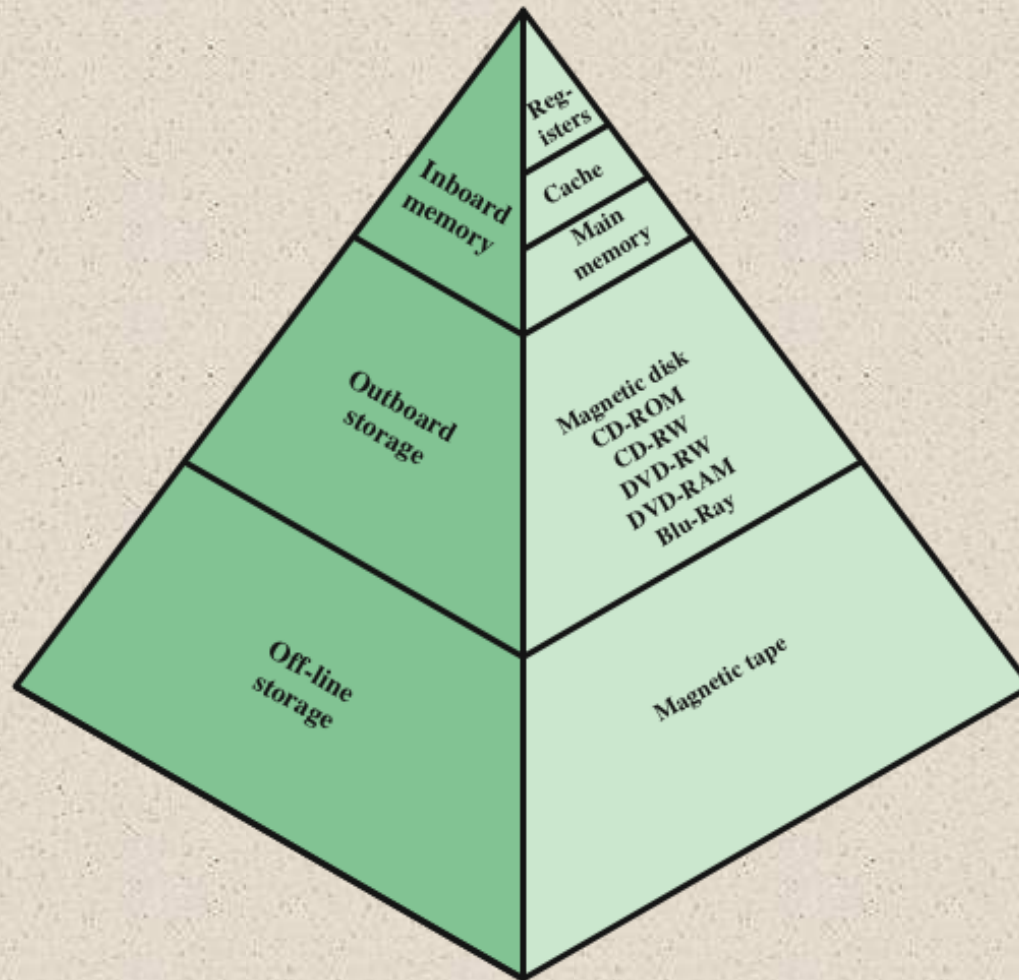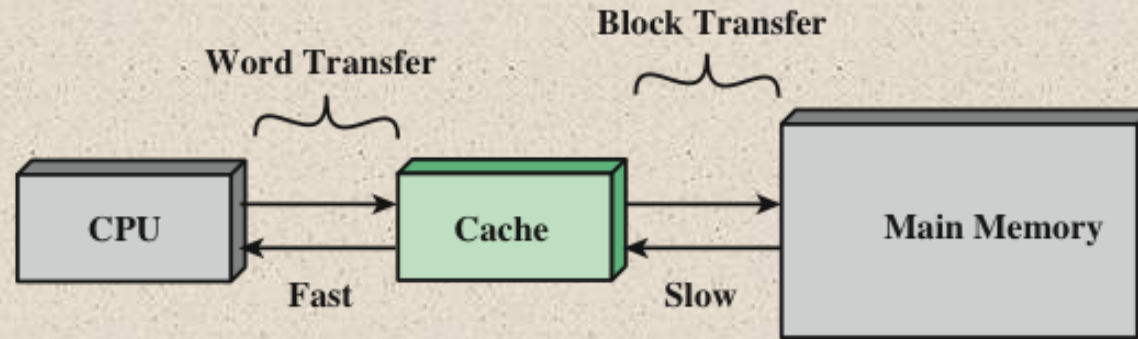
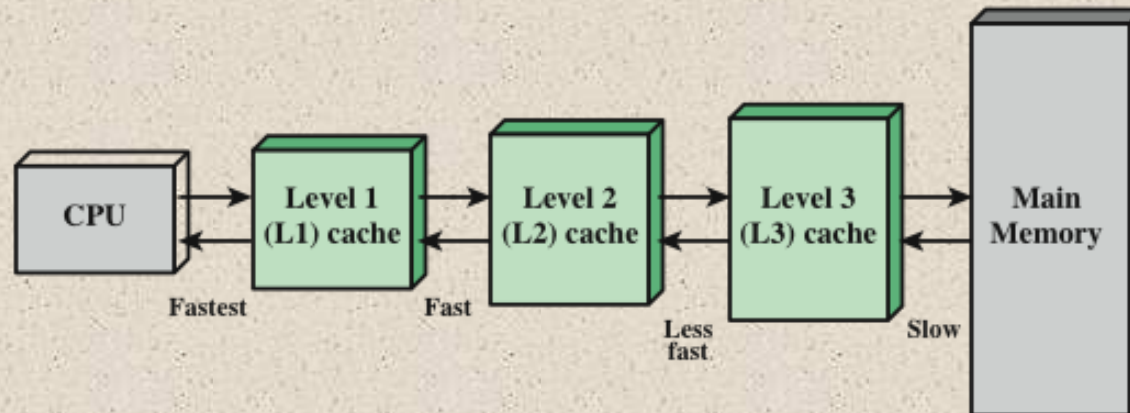# Memory Hierarchy - Diagram



**Figure 4.1   The Memory Hierarchy**

# Locality of Reference

- Two or more levels of memory can be used to produce average access time approaching the highest level

- The reason that this works well is called "**locality of reference**"
  - **Spatial locality** refers to the tendency of execution to involve a number of memory locations that are clustered.
  - **Temporal locality** refers to the tendency for a processor to access memory locations that have been used recently.

- During the course of the execution of a program, memory references tend to cluster
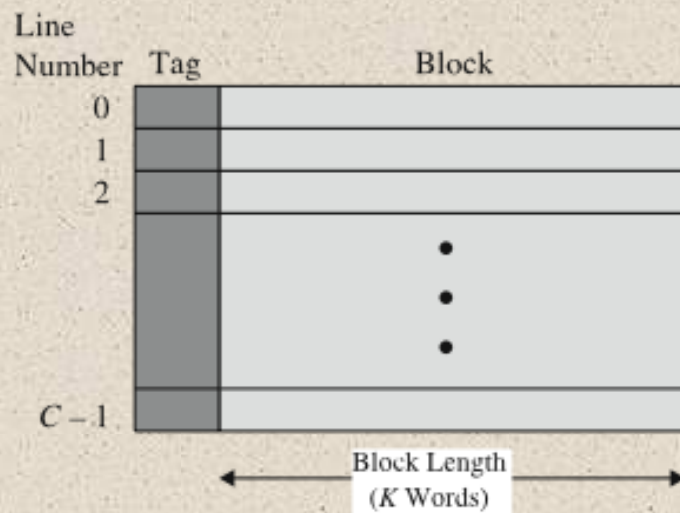  - e.g. loops
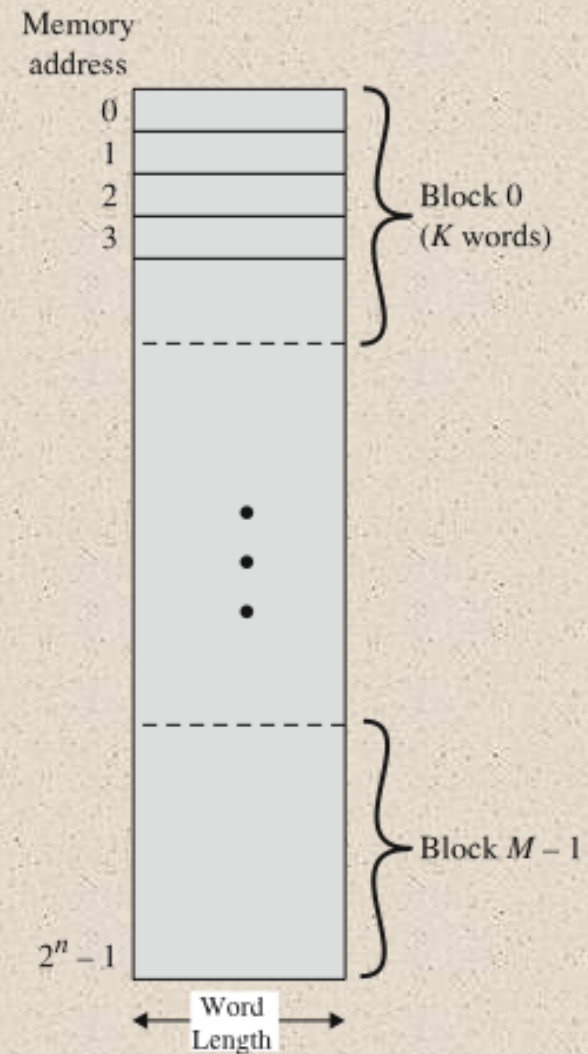
# Cache and Main Memory



(a) Single cache

(b) Three-level cache organization

Figure 4.3 Cache and Main Memory

# Cache/Main Memory Structure



Figure 4.4 Cache/Main-Memory Structure

# Cache Read Operation

**START**

Receive address RA from CPU

Is block containing RA in cache?

No → Access main memory for block containing RA

Yes ↓

Fetch RA word and deliver to CPU

Allocate cache line for main memory block

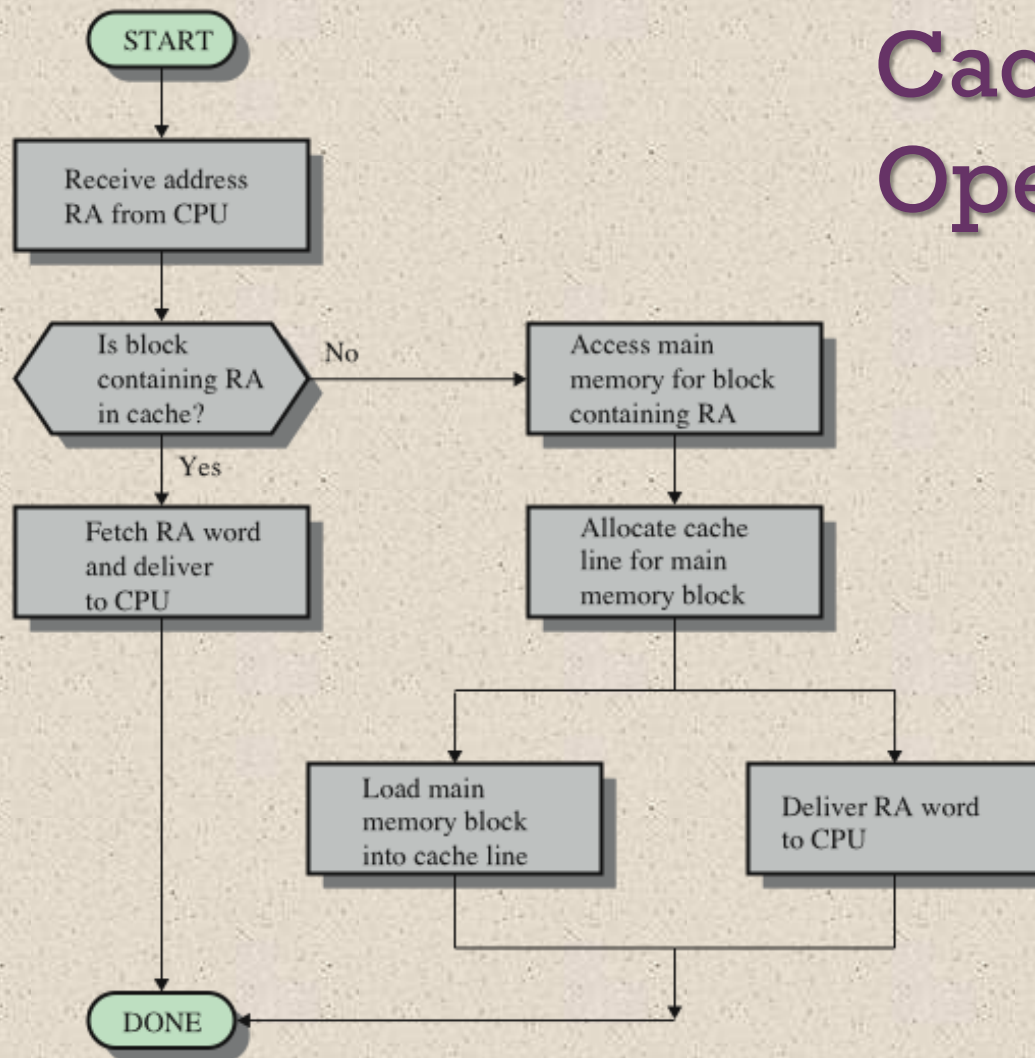Load main memory block into cache line

Deliver RA word to CPU

**DONE**

**Figure 4.5   Cache Read Operation**

# Typical Cache Organization



Figure 4.6   Typical Cache Organization

# Elements of Cache Design

| | |
|---|---|
| **Cache Addresses** | **Write Policy** |
| Logical | Write through |
| Physical | Write back |
| **Cache Size** | **Line Size** |
| **Mapping Function** | **Number of caches** |
| Direct | Single or two level |
| Associative | Unified or split |
| Set Associative | |
| **Replacement Algorithm** | |
| Least recently used (LRU) | |
| First in first out (FIFO) | |
| Least frequently used (LFU) | |
| Random | |

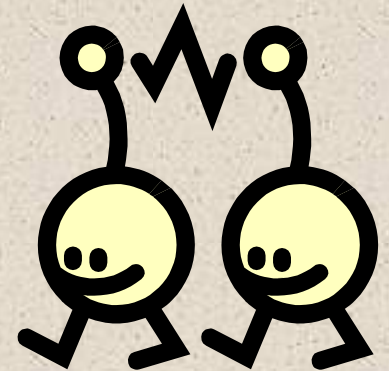Table 4.2    Elements of Cache Design
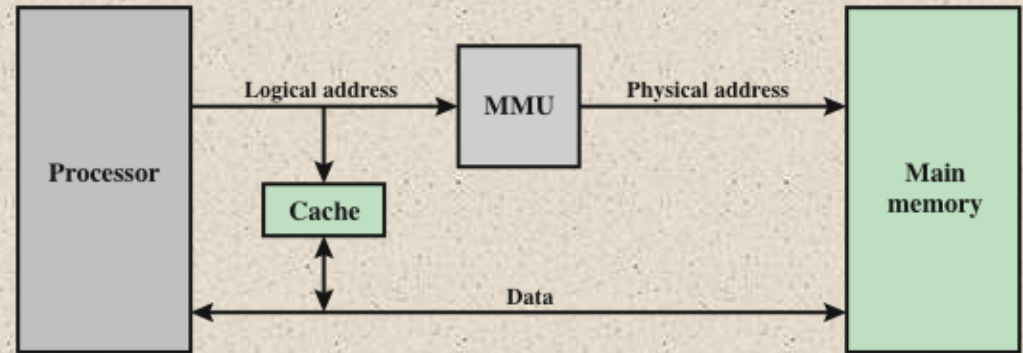
# Cache Addresses

## Virtual Memory

- Virtual memory
  - Facility that allows programs to address memory from a logical point of view, without regard to the amount of main memory physically available
  - When used, the address fields of machine instructions contain virtual addresses
  - For reads to and writes from main memory, a hardware memory management unit (MMU) translates each virtual address into a physical address in main memory
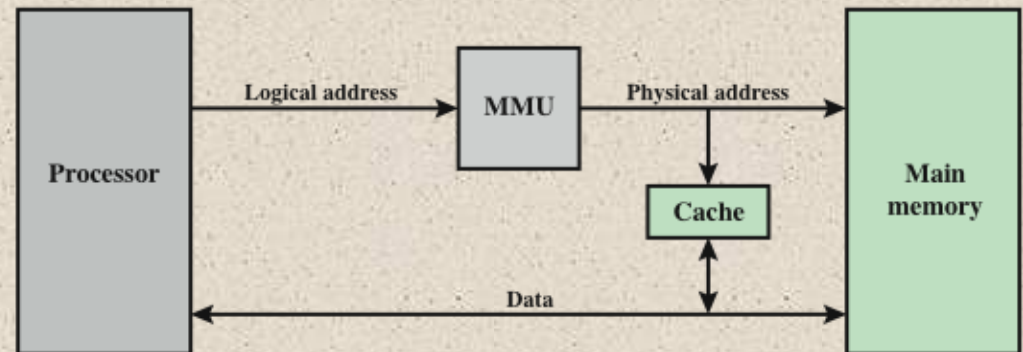
# Logical and Physical Caches



(a) Logical Cache

(b) Physical Cache

**Figure 4.7   Logical and Physical Caches**

**Table 4.3**

**Cache Sizes of Some Processors**

| 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 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 MTA[b] | Supercomputer | 2000 | 8 kB | 2 MB | — |
| Itanium | PC/server | 2001 | 16 kB/16 kB | 96 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 | 36 MB |
| CRAY XD-1 | Supercomputer | 2004 | 64 kB/64 kB | 1MB | — |
| IBM POWER6 | PC/server | 2007 | 64 kB/64 kB | 4 MB | 32 MB |
| IBM z10 | Mainframe | 2008 | 64 kB/128 kB | 3 MB | 24-48 MB |
| Intel Core i7 EE 990 | Workstaton/ server | 2011 | 6 × 32 kB/32 kB | 1.5 MB | 12 MB |
| IBM zEnterprise 196 | Mainframe/ Server | 2011 | 24 × 64 kB/ 128 kB | 24 × 1.5 MB | 24 MB L3 192 MB L4 |

[a] Two values separated by a slash refer to instruction and data caches.

[b] Both caches are instruction only; no data caches.

# Mapping Function

- Because there are fewer cache lines than main memory blocks, an algorithm is needed for mapping main memory blocks into cache lines

- Three techniques can be used:

## Direct

- The simplest technique
- Maps each block of main memory into only one possible cache line

## Associative

- Permits each main memory block to be loaded into any line of the cache
- The cache control logic interprets a memory address simply as a Tag and a Word field
- To determine whether a block is in the cache, the cache control logic must simultaneously examine every line's Tag for a match

## Set Associative

- A compromise that exhibits the strengths of both the direct and associative approaches while reducing their disadvantages