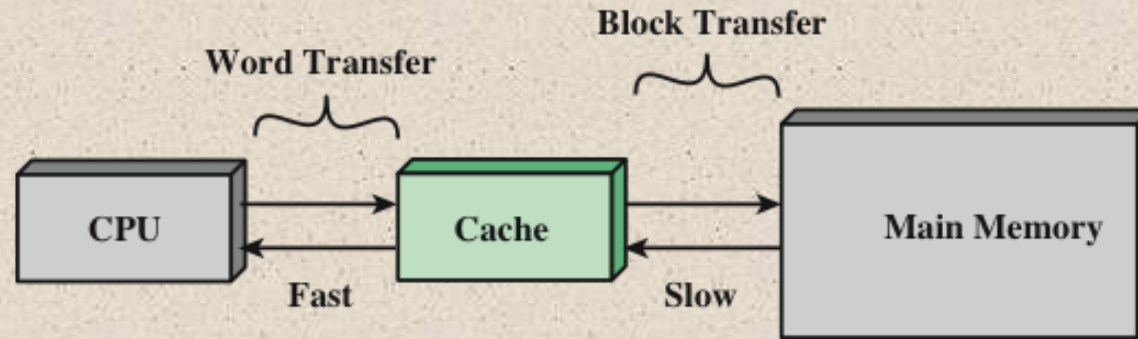+

William Stallings
Computer Organization
and Architecture
9th Edition
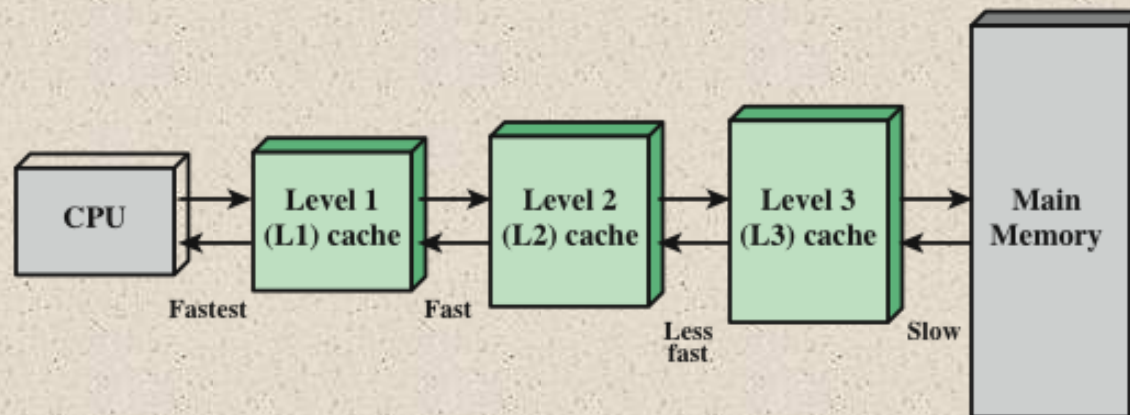
# + Chapter 4

Cache Memory

# Cache and Main Memory
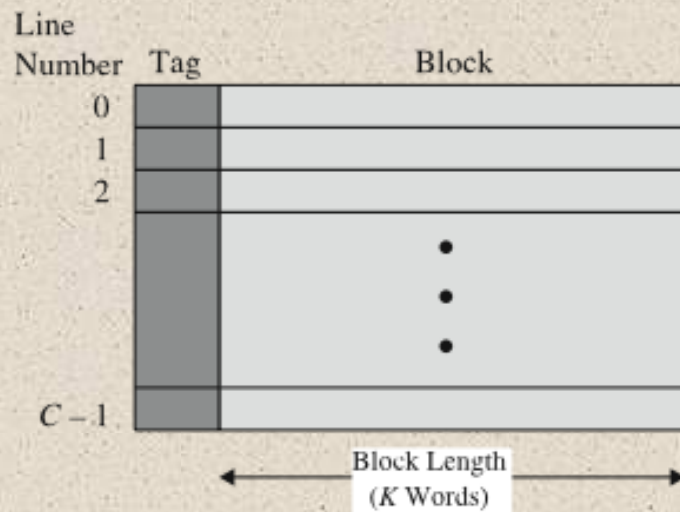


(a) Single cache

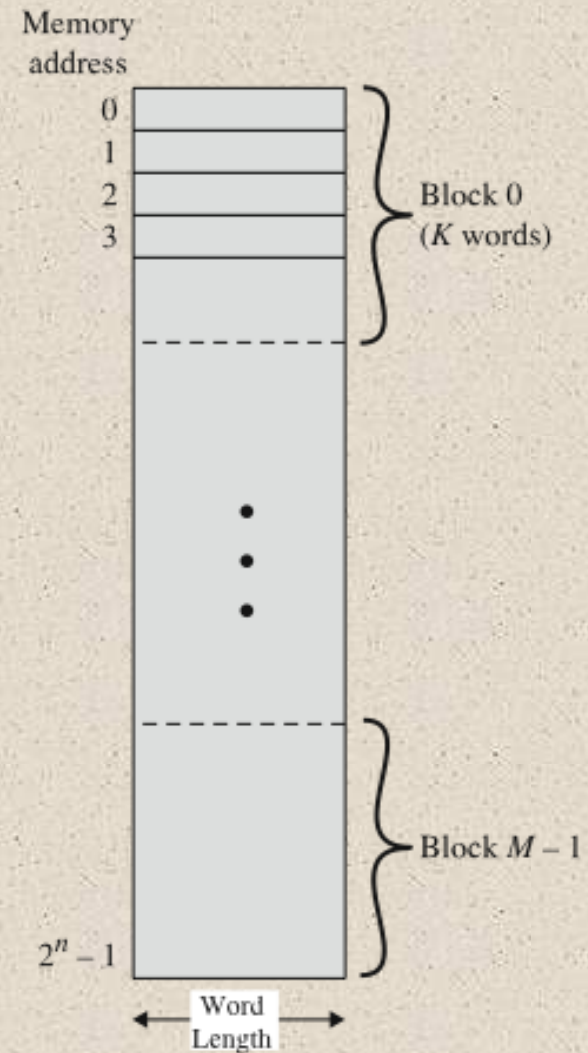(b) Three-level cache organization

Figure 4.3  Cache and Main Memory

# Cache/Main Memory Structure



Figure 4.4  Cache/Main-Memory Structure

# Cache Read Operation



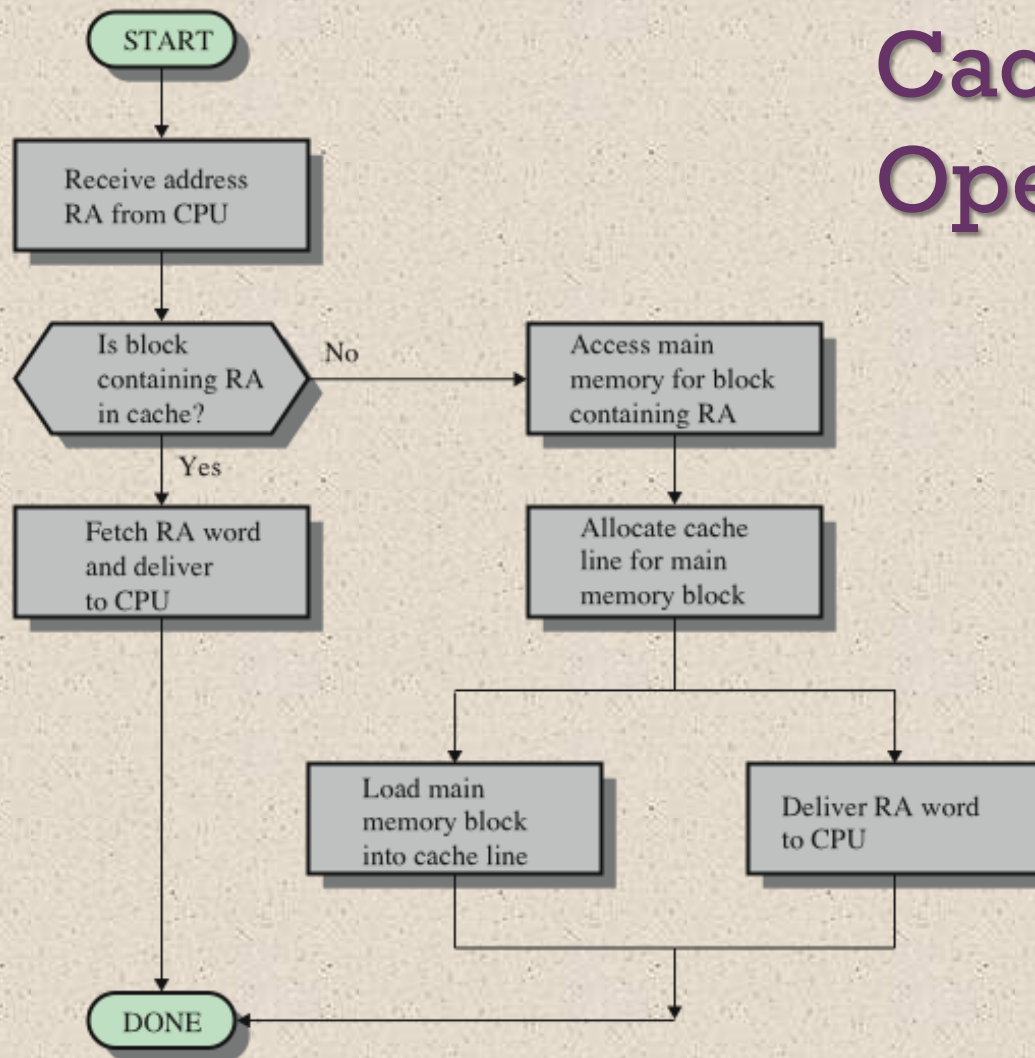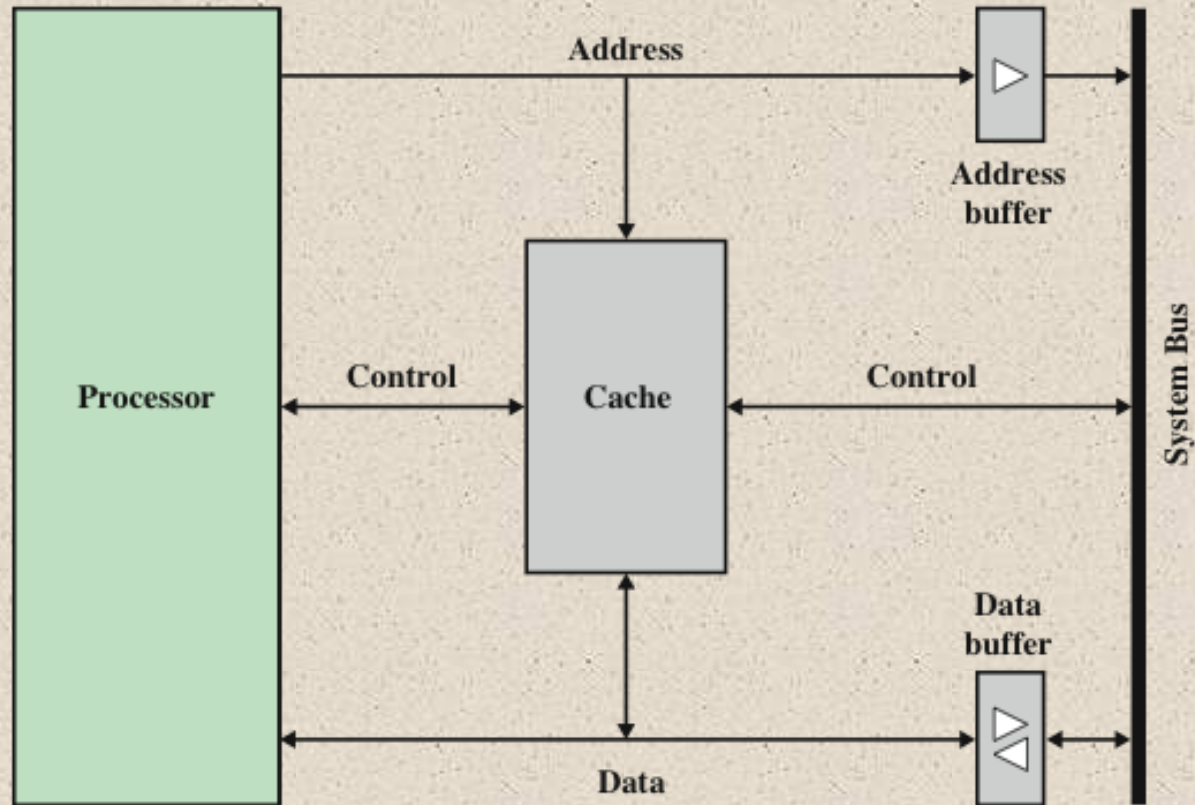**Figure 4.5   Cache Read Operation**

# Typical Cache Organization



Figure 4.6  Typical Cache Organization

# Elements of Cache Design

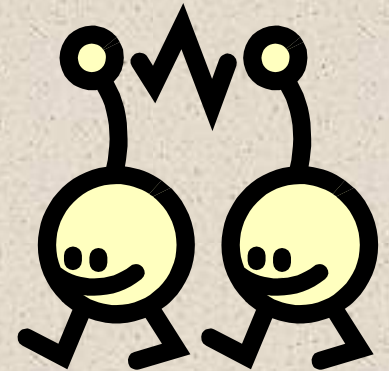| | |
|---|---|
| **Cache Addresses** | **Write Policy** |
|     Logical |     Write through |
|     Physical |     Write back |
| **Cache Size** | **Line Size** |
| **Mapping Function** | **Number of caches** |
|     Direct |     Single or two level |
|     Associative |     Unified or split |
|     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
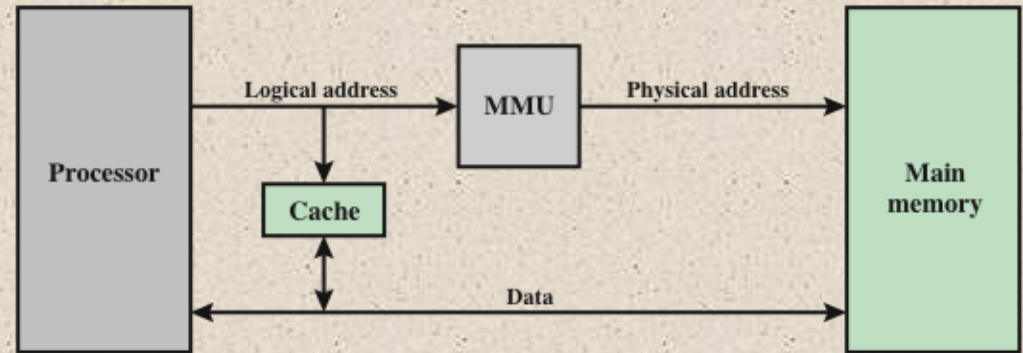
# Cache Addresses

## 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
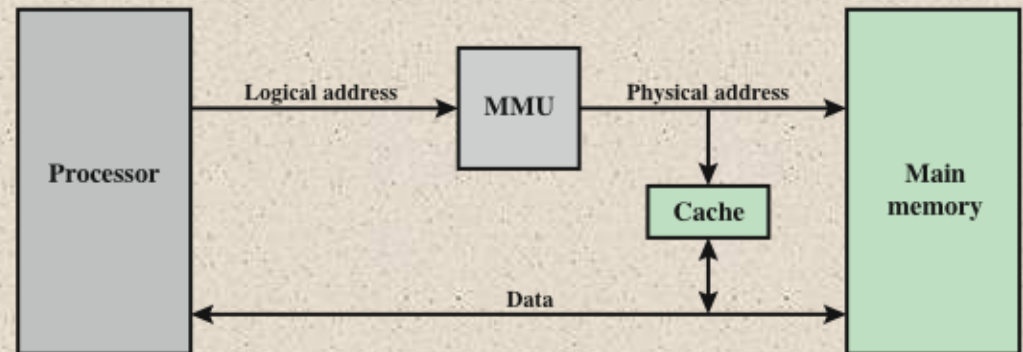
# Logical and Physical Caches



(a) Logical Cache

(b) Physical Cache

Figure 4.7   Logical and Physical Caches

# 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
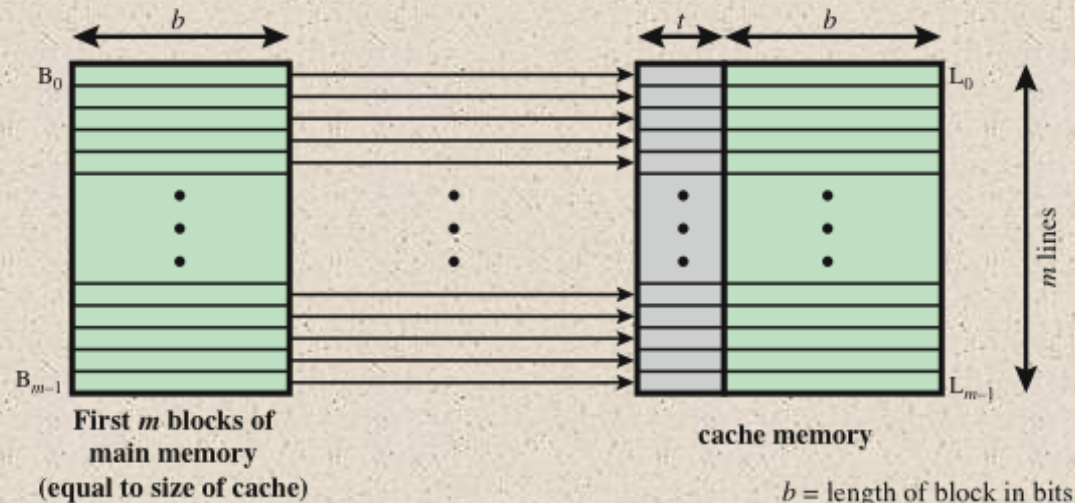
## Associative

- 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
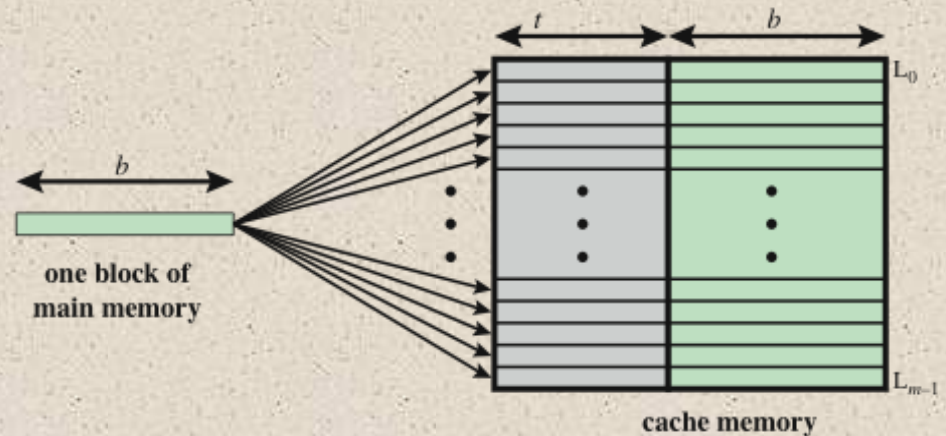
## Set Associative

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

# Direct Mapping



Figure 4.8  Mapping From Main Memory to Cache: Direct and Associative
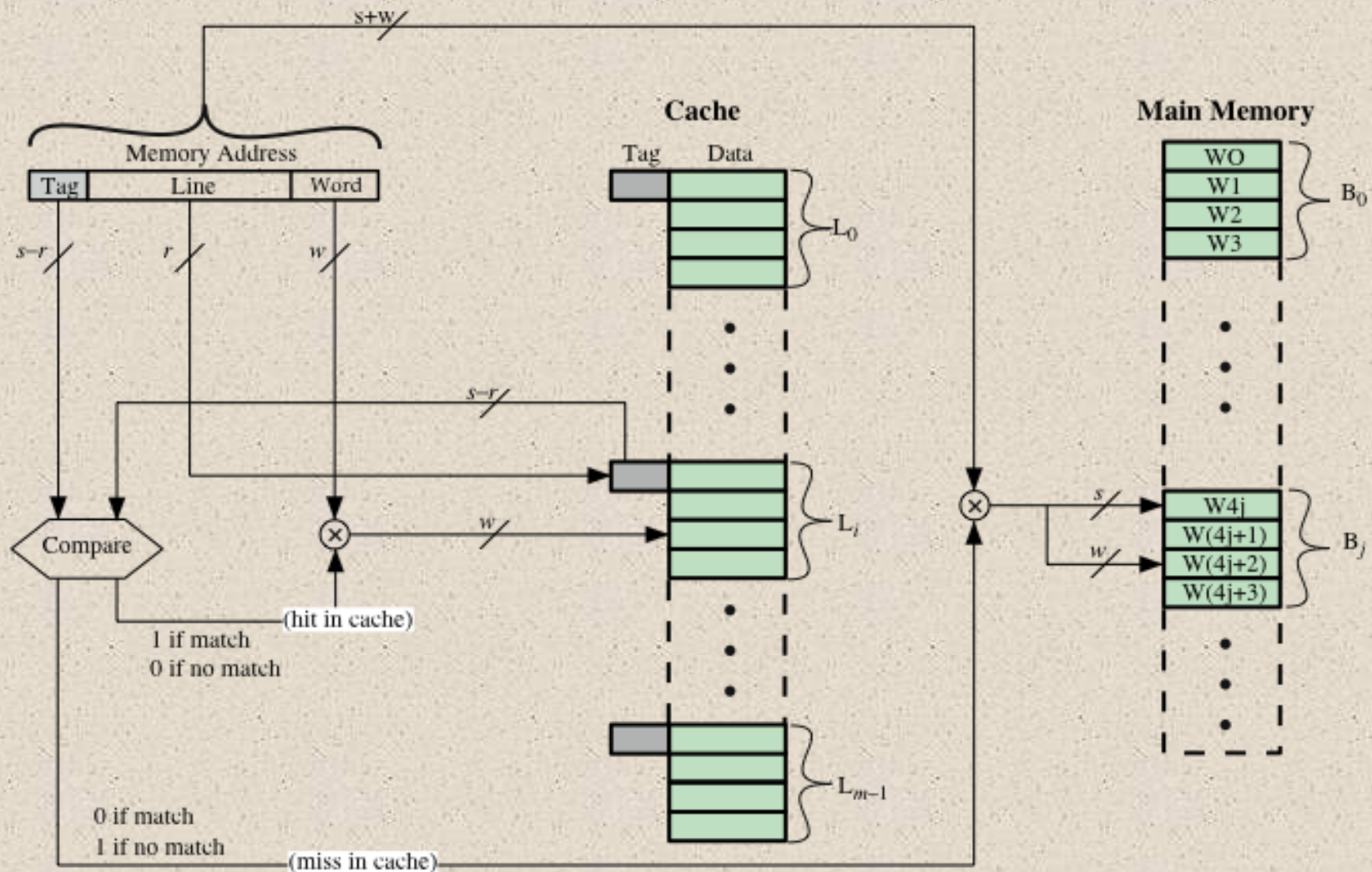
# Direct Mapping Cache Organization



**Figure 4.9  Direct-Mapping Cache Organization**

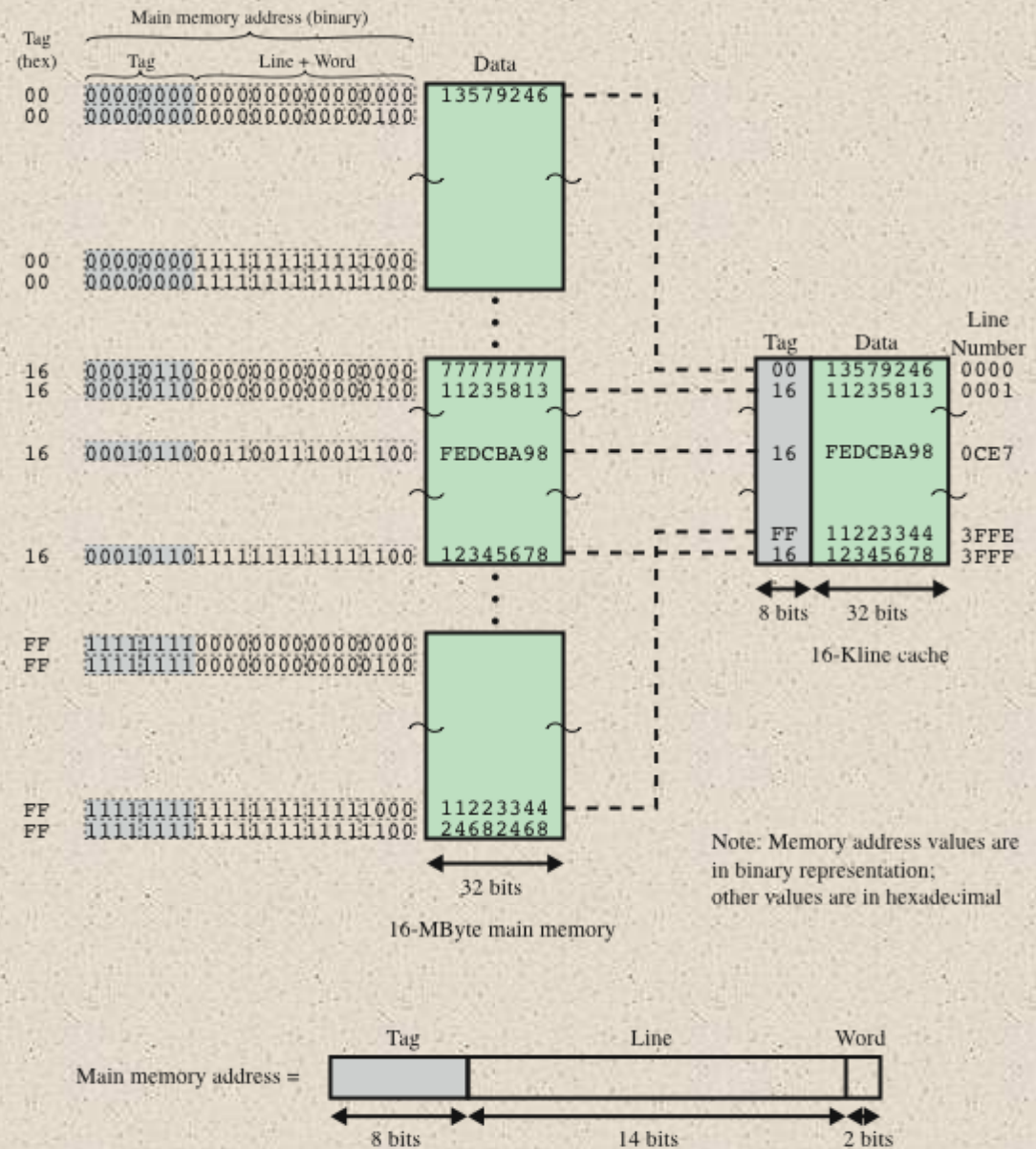# Direct Mapping Example



**Figure 4.10  Direct Mapping Example**

# Direct Mapping Summary

- 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