

Figure 4.4 Cache/Main-Memory Structure

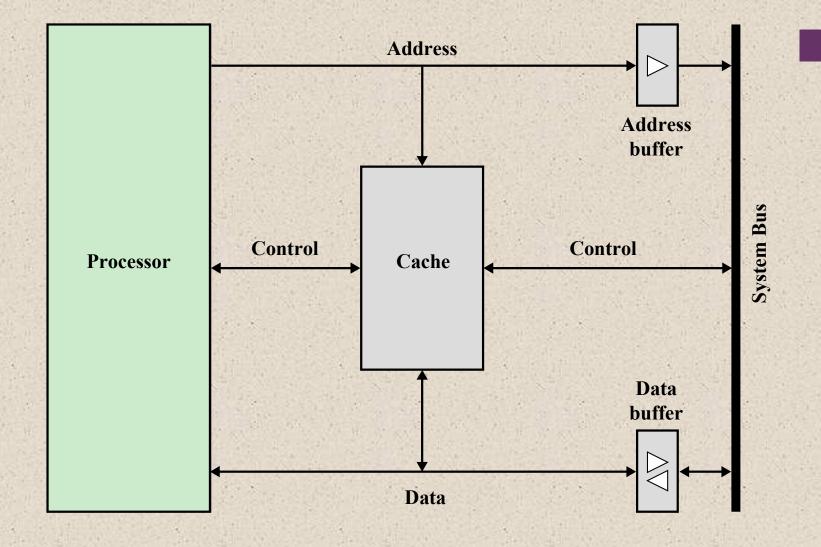


Figure 4.6 Typical Cache Organization



Logical

Physical

Cache Size

Mapping Function

Direct

Associative

Set Associative

Replacement Algorithm

Least recently used (LRU)

First in first out (FIFO)

Least frequently used (LFU)

Random

Write Policy

Write through

Write back

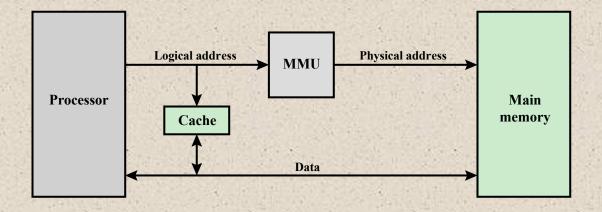
Line Size

Number of caches

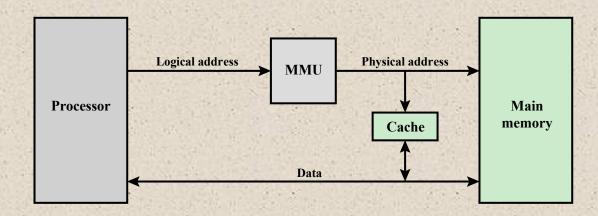
Single or two level

Unified or split

Table 4.2 Elements of Cache Design



(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

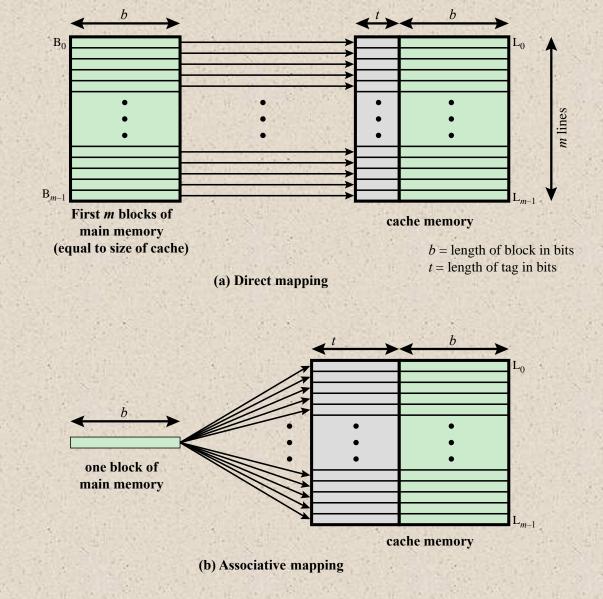


Figure 4.8 Mapping From Main Memory to Cache: Direct and Associative

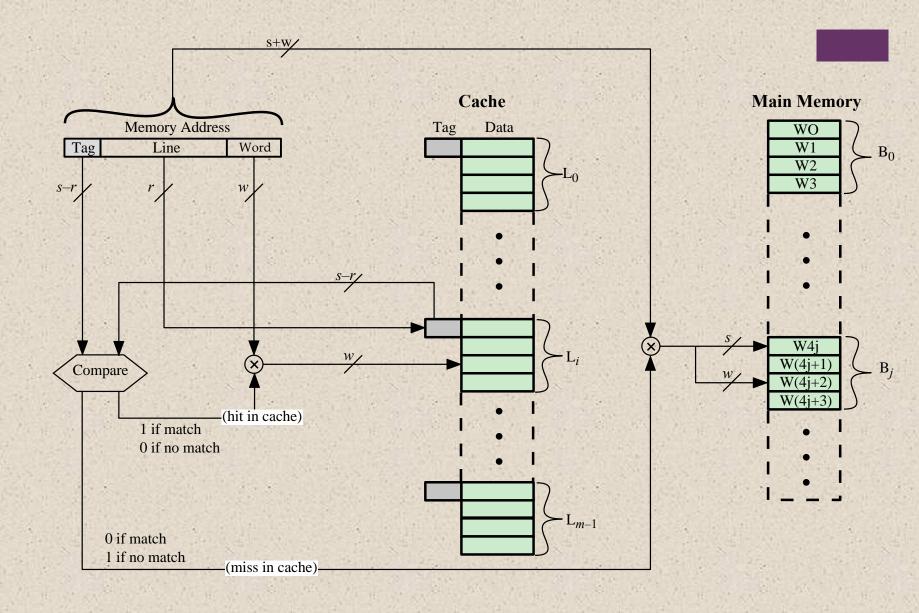


Figure 4.9 Direct-Mapping Cache Organization

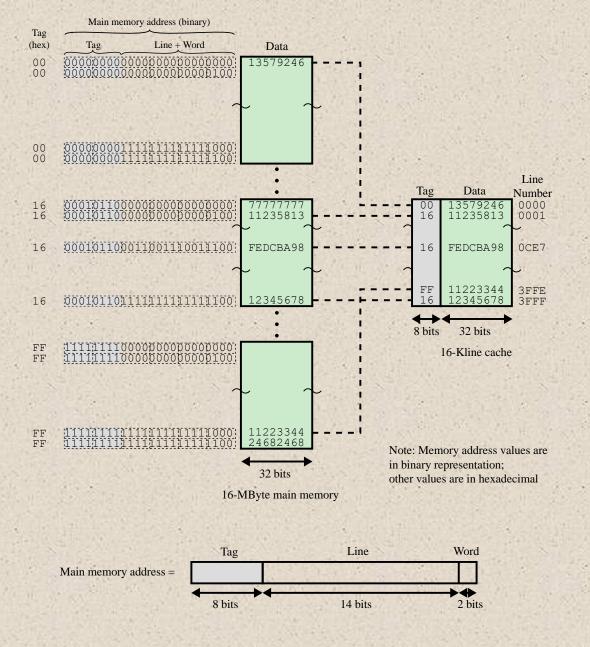
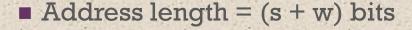


Figure 4.10 Direct Mapping Example

Direct Mapping Summary



- 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



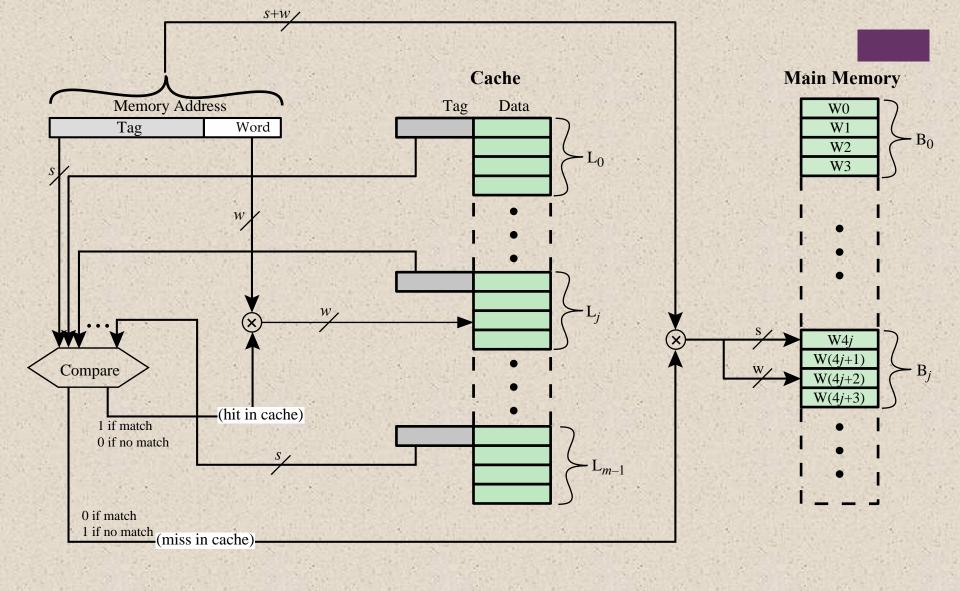


Figure 4.11 Fully Associative Cache Organization

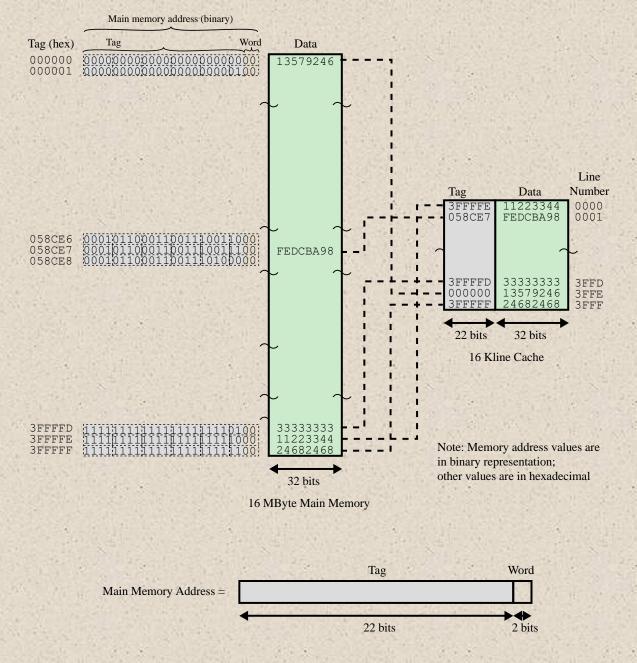
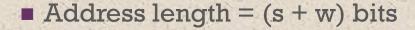


Figure 4.12 Associative Mapping Example

Associative Mapping Summary



- 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

