

# Lecture No.9

## Lecture Outlines

### 4.3 Elements of Cache Design

- Cache Addresses
- Cache Size
- Mapping Function

### 4.3 ELEMENTS OF CACHE DESIGN

This section provides an overview of cache design parameters and reports some typical results. We occasionally refer to the use of caches in **high-performance computing (HPC)**. HPC deals with supercomputers and their software, especially for scientific applications that involve large amounts of data, vector and matrix computation, and the use of parallel algorithms. Cache design for HPC is quite different than for other hardware platforms and applications. Indeed, many researchers have found that HPC applications perform poorly on computer architectures that employ caches. Other researchers have since shown that a cache hierarchy can be useful in improving performance if the application software is tuned to exploit the cache.

Although there are a large number of cache implementations, there are a few basic design elements that serve to classify and differentiate cache architectures. Table 4.2 lists key elements.

#### Cache Addresses

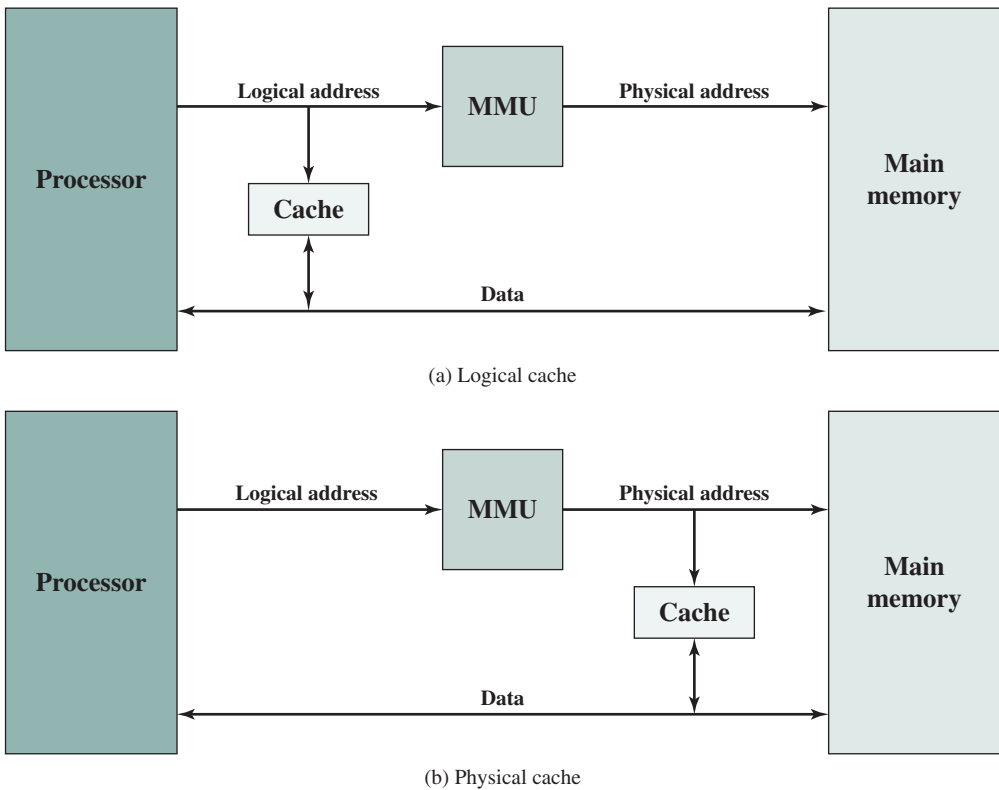
Almost all nonembedded processors, and many embedded processors, support virtual memory, a concept discussed in Chapter 8. In essence, virtual memory is a facility that allows programs to address memory from a logical point of view, without regard to the amount of main memory physically available. When virtual memory is 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.

When virtual addresses are used, the system designer may choose to place the cache between the processor and the MMU or between the MMU and main memory (Figure 4.7). A **logical cache**, also known as a **virtual cache**, stores data using

**Table 4.2** Elements of Cache Design

<b>Cache Addresses</b>	<b>Write Policy</b>
Logical	Write through
Physical	Write back
<b>Cache Size</b>	<b>Line Size</b>
<b>Mapping Function</b>	<b>Number of Caches</b>
Direct	Single or two level
Associative	Unified or split
Set associative	
<b>Replacement Algorithm</b>	
Least recently used (LRU)	
First in first out (FIFO)	
Least frequently used (LFU)	
Random	

**Figure 4.7** Logical and Physical Caches

**virtual addresses.** The processor accesses the cache directly, without going through the MMU. A **physical cache** stores data using main memory **physical addresses**.

One obvious advantage of the logical cache is that cache access speed is faster than for a physical cache, because the cache can respond before the MMU performs an address translation. The disadvantage has to do with the fact that most virtual memory systems supply each application with the same virtual memory address space. That is, each application sees a virtual memory that starts at address 0. Thus, identify which virtual address space this address refers to.

The subject of logical versus physical cache is a complex one, and beyond the scope of this book. For a more in-depth discussion, see [CEKL97] and [JACO08].

## Cache Size

The second item in Table 4.2, cache size, has already been discussed. We would like the size of the cache to be small enough so that the overall average cost per bit is close to that of main memory alone and large enough so that the overall average access time is close to that of the cache alone. There are several other motivations for minimizing cache size. The larger the cache, the larger the number of gates involved in addressing the cache. The result is that large caches tend to be slightly slower than small ones—even when built with the same integrated circuit technology and put in the same place on chip and circuit board. The available chip and board area also limits cache size. Because the performance of the cache is very sensitive to the nature of the workload, it is impossible to arrive at a single “optimum” cache size.

## Mapping Function

Because there are fewer cache lines than main memory blocks, an algorithm is needed for mapping main memory blocks into cache lines. Further, a means is needed for determining which main memory block currently occupies a cache line. The choice of the mapping function dictates how the cache is organized. Three techniques can be used: direct, associative, and set-associative. We examine each of these in turn. In each case, we look at the general structure and then a specific example.

**EXAMPLE 4.2** For all three cases, the example includes the following elements:

- The cache can hold 64 kB.
- Data are transferred between main memory and the cache in blocks of 4 bytes each. This means that the cache is organized as  $16\text{K} = 2^{14}$  lines of 4 bytes each.
- The main memory consists of 16 MB, with each byte directly addressable by a 24-bit address ( $2^{24} = 16\text{M}$ ). Thus, for mapping purposes, we can consider main memory to consist of 4M blocks of 4 bytes each.

**DIRECT MAPPING** The simplest technique, known as direct mapping, maps each block of main memory into only one possible cache line. The mapping is expressed as

$$i = j \text{ modulo } m$$

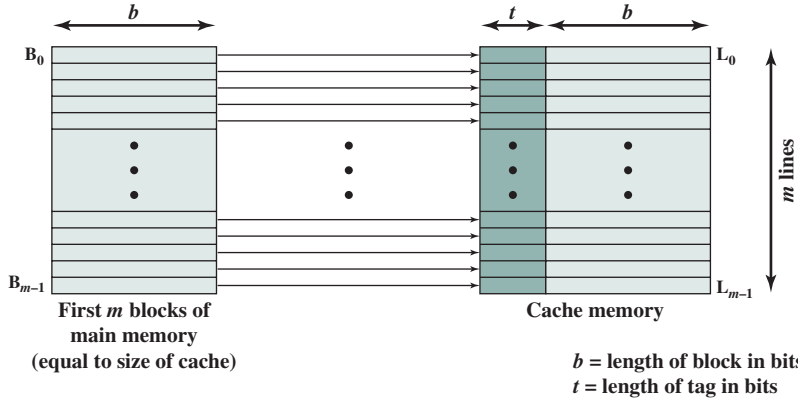
where

$i$  = cache line number

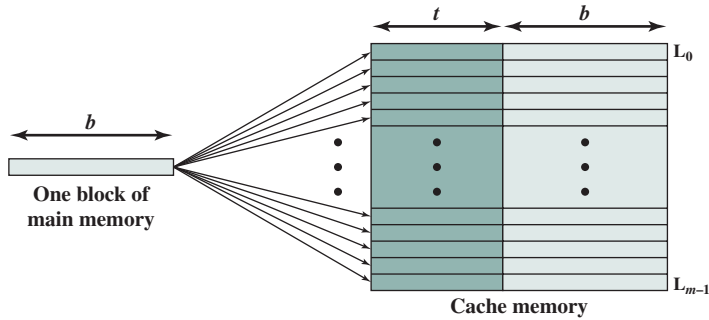
$j$  = main memory block number

$m$  = number of lines in the cache

Figure 4.8a shows the mapping for the first  $m$  blocks of main memory. Each block of main memory maps into one unique line of the cache. The next  $m$  blocks



(a) Direct mapping



(b) Associative mapping

**Figure 4.8** Mapping from Main Memory to Cache: Direct and Associative

of main memory map into the cache in the same fashion; that is, block  $B_m$  of main memory maps into line  $L_0$  of cache, block  $B_{m+1}$  maps into line  $L_1$ , and so on.

The mapping function is easily implemented using the main memory address. Figure 4.9 illustrates the general mechanism. For purposes of cache access, each main memory address can be viewed as consisting of three fields. The least significant  $w$  bits identify a unique word or byte within a block of main memory; in most contemporary machines, the address is at the byte level. The remaining  $s$  bits specify one of the  $2^s$  blocks of main memory. The cache logic interprets these  $s$  bits as a tag of  $s - r$  bits (most significant portion) and a line field of  $r$  bits. This latter field identifies one of the  $m = 2^r$  lines of the cache. To summarize,

- 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 =  $\frac{2^{s+w}}{2^w} = 2^s$
- Number of lines in cache =  $m = 2^r$
- Size of cache =  $2^{r+w}$  words or bytes
- Size of tag =  $(s - r)$  bits

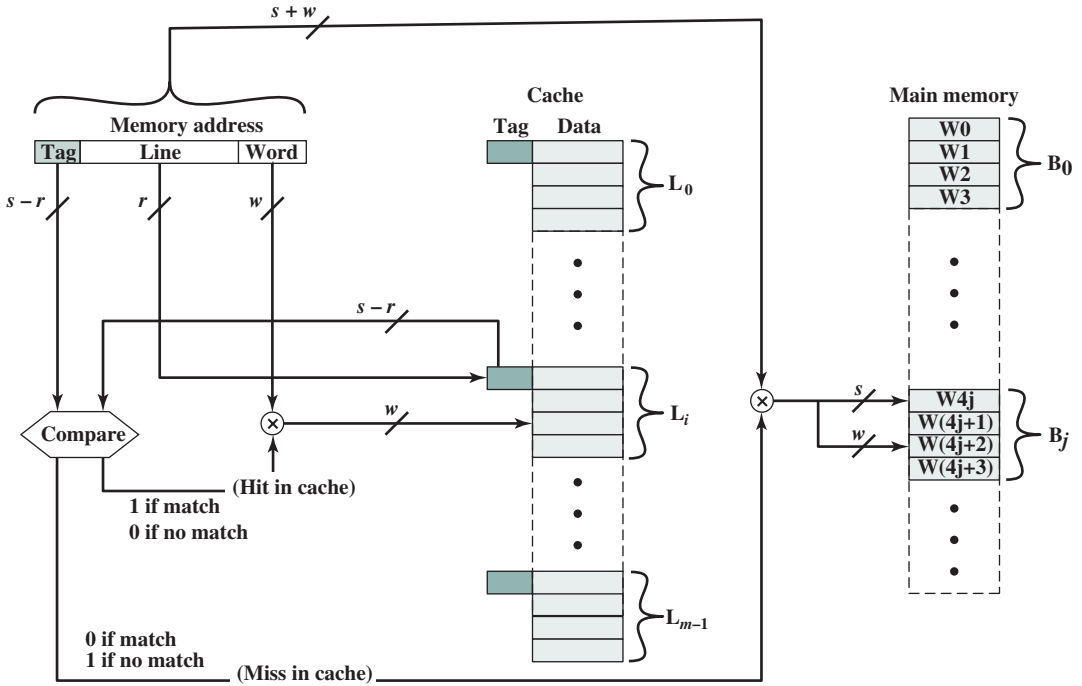


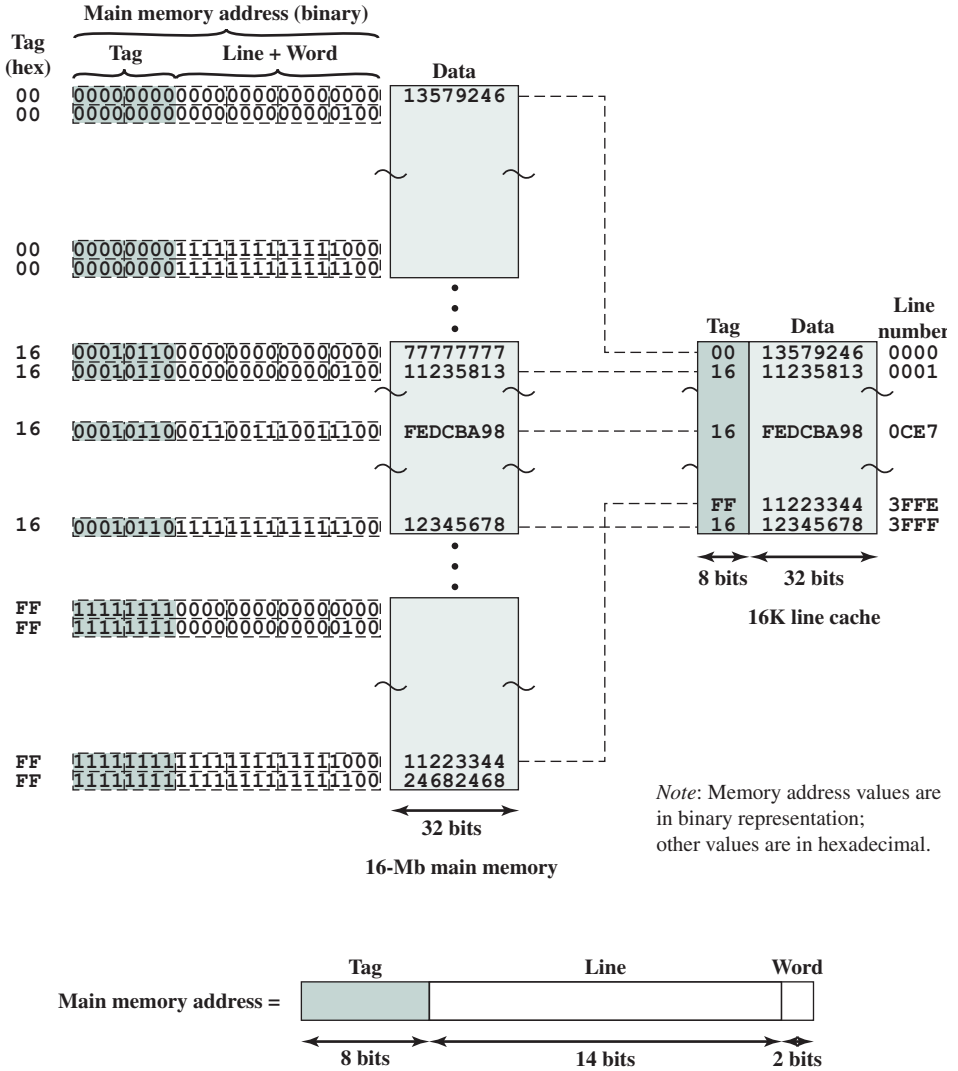
Figure 4.9 Direct-Mapping Cache Organization

**EXAMPLE 4.2a** Figure 4.10 shows our example system using direct mapping. In the example,  $m = 16K = 2^{14}$  and  $i = j$  modulo  $2^{14}$ . The mapping becomes

Cache Line	Starting Memory Address of Block
0	000000, 010000, ..., FF0000
1	000004, 010004, ..., FF0004
$\vdots$	$\vdots$
$2^{14} - 1$	00FFFC, 01FFFC, ..., FFFFFC

Note that no two blocks that map into the same line number have the same tag number. Thus, blocks with starting addresses 000000, 010000, ..., FF0000 have tag numbers 00, 01, ..., FF, respectively.

Referring back to Figure 4.5, a read operation works as follows. The cache system is presented with a 24-bit address. The 14-bit line number is used as an index into the cache to access a particular line. If the 8-bit tag number matches the tag number currently stored in that line, then the 2-bit word number is used to select one of the 4 bytes in that line. Otherwise, the 22-bit tag-plus-line field is used to fetch a block from main memory. The actual address that is used for the fetch is the 22-bit tag-plus-line concatenated with two 0 bits, so that 4 bytes are fetched starting on a block boundary.



**Figure 4.10** Direct Mapping Example

The effect of this mapping is that blocks of main memory are assigned to lines of the cache as follows:

Cache line	Main memory blocks assigned
0	$0, m, 2m, \dots, 2^s - m$
1	$1, m + 1, 2m + 1, \dots, 2^s - m + 1$
⋮	⋮
$m - 1$	$m - 1, 2m - 1, 3m - 1, \dots, 2^s - 1$

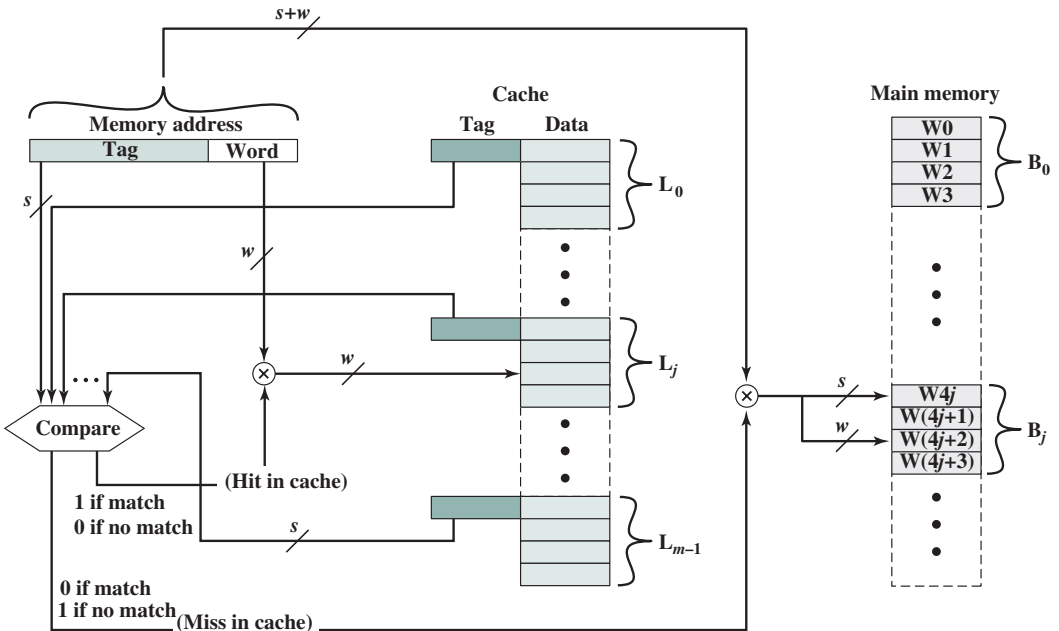
Thus, the use of a portion of the address as a line number provides a unique mapping of each block of main memory into the cache. When a block is actually

read into its assigned line, it is necessary to tag the data to distinguish it from other blocks that can fit into that line. The most significant  $s - r$  bits serve this purpose.

The direct mapping technique is simple and inexpensive to implement. Its main disadvantage is that there is a fixed cache location for any given block. Thus, if a program happens to reference words repeatedly from two different blocks that map into the same line, then the blocks will be continually swapped in the cache, and the hit ratio will be low (a phenomenon known as *thrashing*).

One approach to lower the **miss** penalty is to remember what was discarded in case it is needed again. Since the discarded data has already been fetched, it can be used again at a small cost. Such recycling is possible using a victim cache. Victim cache was originally proposed as an approach to reduce the conflict misses of direct mapped caches without affecting its fast access time. Victim cache is a fully associative cache, whose size is typically 4 to 16 cache lines, residing between a direct mapped L1 cache and the next level of memory.

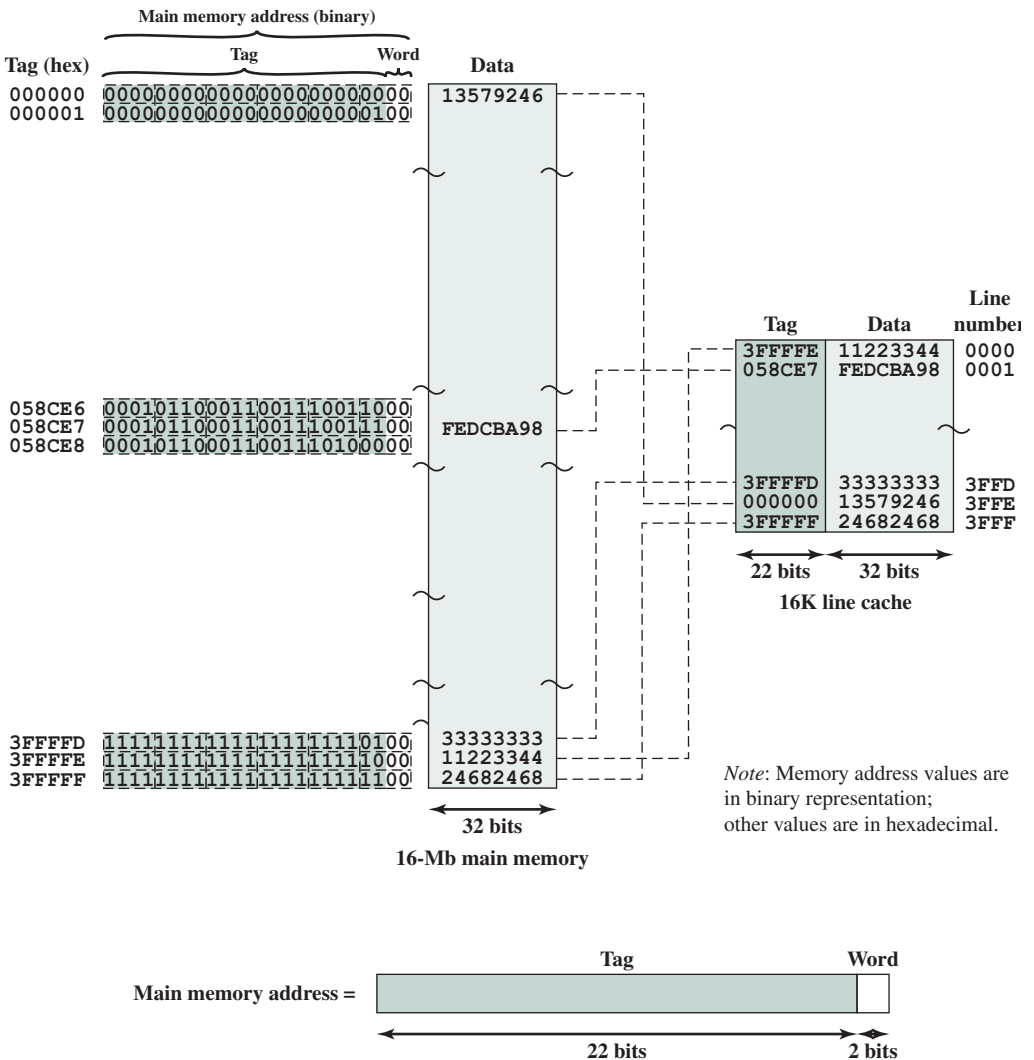
**ASSOCIATIVE MAPPING** Associative mapping overcomes the disadvantage of direct mapping by permitting each main memory block to be loaded into any line of the cache (Figure 4.8b). In this case, the cache control logic interprets a memory address simply as a Tag and a Word field. The Tag field uniquely identifies a block of main memory. To determine whether a block is in the cache, the cache control logic must simultaneously examine every line's tag for a match. Figure 4.11 illustrates the logic.



**Figure 4.11** Fully Associative Cache Organization

**EXAMPLE 4.2b** Figure 4.12 shows our example using associative mapping. A main memory address consists of a 22-bit tag and a 2-bit byte number. The 22-bit tag must be stored with the 32-bit block of data for each line in the cache. Note that it is the leftmost (most significant) 22 bits of the address that form the tag. Thus, the 24-bit hexadecimal address 16339C has the 22-bit tag 058CE7. This is easily seen in binary notation:

Memory address	0001	0110	0011	0011	1001	1100	(binary)
	1	6	3	3	9	C	(hex)
Tag (leftmost 22 bits)	00	0101	1000	1100	1110	0111	(binary)
	0	5	8	C	E	7	(hex)



**Figure 4.12** Associative Mapping Example



Note that no field in the address corresponds to the line number, so that the number of lines in the cache is not determined by the address format. To summarize,

- 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 =  $\frac{2^{s+w}}{2^w} = 2^s$
- Number of lines in cache = undetermined
- Size of tag =  $s$  bits

With associative mapping, there is flexibility as to which block to replace when a new block is read into the cache. Replacement algorithms, discussed later in this section, are designed to maximize the hit ratio. The principal disadvantage of associative mapping is the complex circuitry required to examine the tags of all cache lines in parallel.

**SET-ASSOCIATIVE MAPPING** Set-associative mapping is a compromise that exhibits the strengths of both the direct and associative approaches while reducing their disadvantages.

In this case, the cache consists of number sets, each of which consists of a number of lines. The relationships are

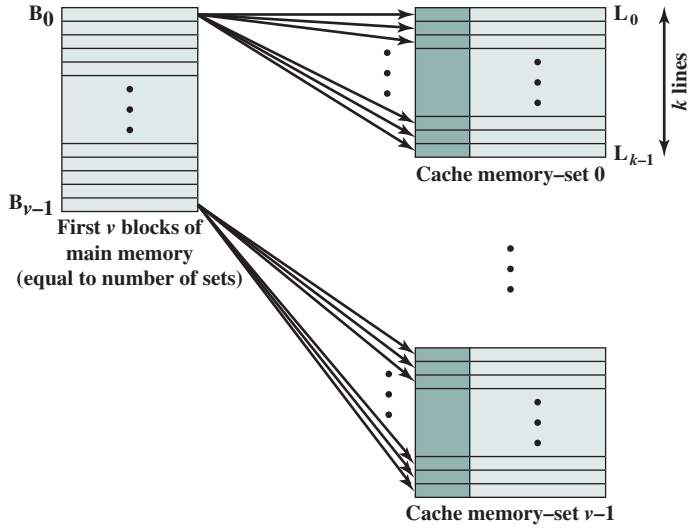
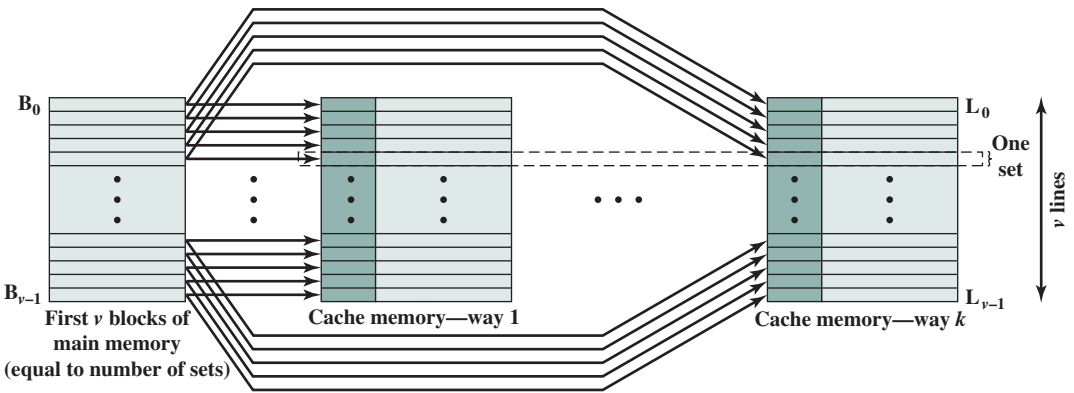
$$m = v \times k$$

$$i = j \text{ modulo } v$$

where

- $i$  = cache set number
- $j$  = main memory block number
- $m$  = number of lines in the cache
- $v$  = number of sets
- $k$  = number of lines in each set

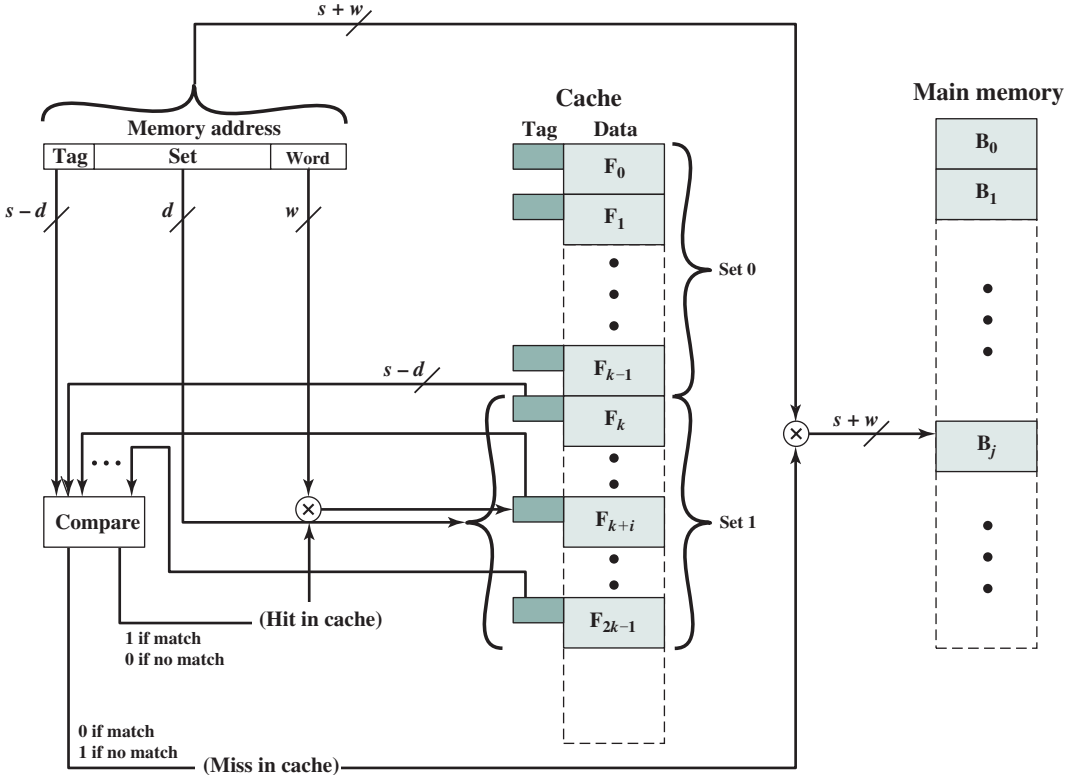
This is referred to as  $k$ -way set-associative mapping. With set-associative mapping, block  $B_j$  can be mapped into any of the lines of set  $j$ . Figure 4.13a illustrates this mapping for the first  $v$  blocks of main memory. As with associative mapping, each word maps into multiple cache lines. For set-associative mapping, each word maps into all the cache lines in a specific set, so that main memory block  $B_0$  maps into set 0, and so on. Thus, the set-associative cache can be physically implemented as  $v$  associative caches. It is also possible to implement the set-associative cache as  $k$  direct mapping caches, as shown in Figure 4.13b. Each direct-mapped cache is referred to as a *way*, consisting of  $v$  lines. The first  $v$  lines of main memory are direct mapped into the  $v$  lines of each way; the next group of  $v$  lines of main memory are similarly mapped, and so on. The direct-mapped implementation is typically used

(a)  $v$  associative-mapped caches(b)  $k$  direct-mapped caches**Figure 4.13** Mapping from Main Memory to Cache:  $k$ -WaySet Associative

for small degrees of associativity (small values of  $k$ ) while the associative-mapped implementation is typically used for higher degrees of associativity [JACO08].

For set-associative mapping, the cache control logic interprets a memory address as three fields: Tag, Set, and Word. The  $d$  set bits specify one of  $v = 2^d$  sets. The  $s$  bits of the Tag and Set fields specify one of the  $2^s$  blocks of main memory. Figure 4.14 illustrates the cache control logic. With fully associative mapping, the tag in a memory address is quite large and must be compared to the tag of every line in the cache. With  $k$ -way set-associative mapping, the tag in a memory address is much smaller and is only compared to the  $k$  tags within a single set. To summarize,

- Address length =  $(s + w)$  bits
- Number of addressable units =  $2^{s+w}$  words or bytes



**Figure 4.14**  $k$ -Way Set-Associative Cache Organization

- Block size = line size =  $2^w$  words or bytes
- Number of blocks in main memory =  $\frac{2^{s+w}}{2^w} = 2^s$
- Number of lines in set =  $k$
- Number of sets =  $v = 2^d$
- Number of lines in cache =  $m = kv = k \times 2^d$
- Size of cache =  $k \times 2^{d+w}$  words or bytes
- Size of tag =  $(s - d)$  bits

**EXAMPLE 4.2c** Figure 4.15 shows our example using set-associative mapping with two lines in each set, referred to as two-way set-associative. The 13-bit set number identifies a unique set of two lines within the cache. It also gives the number of the block in main memory, modulo  $2^{13}$ . This determines the mapping of blocks into lines. Thus, blocks 000000, 008000, ..., FF8000 of main memory map into cache set 0. Any of those blocks can be loaded into either of the two lines in the set. Note that no two blocks that map into the same cache set have the same tag number. For a read operation, the 13-bit set number is used to determine which set of two lines is to be examined. Both lines in the set are examined for a match with the tag number of the address to be accessed.

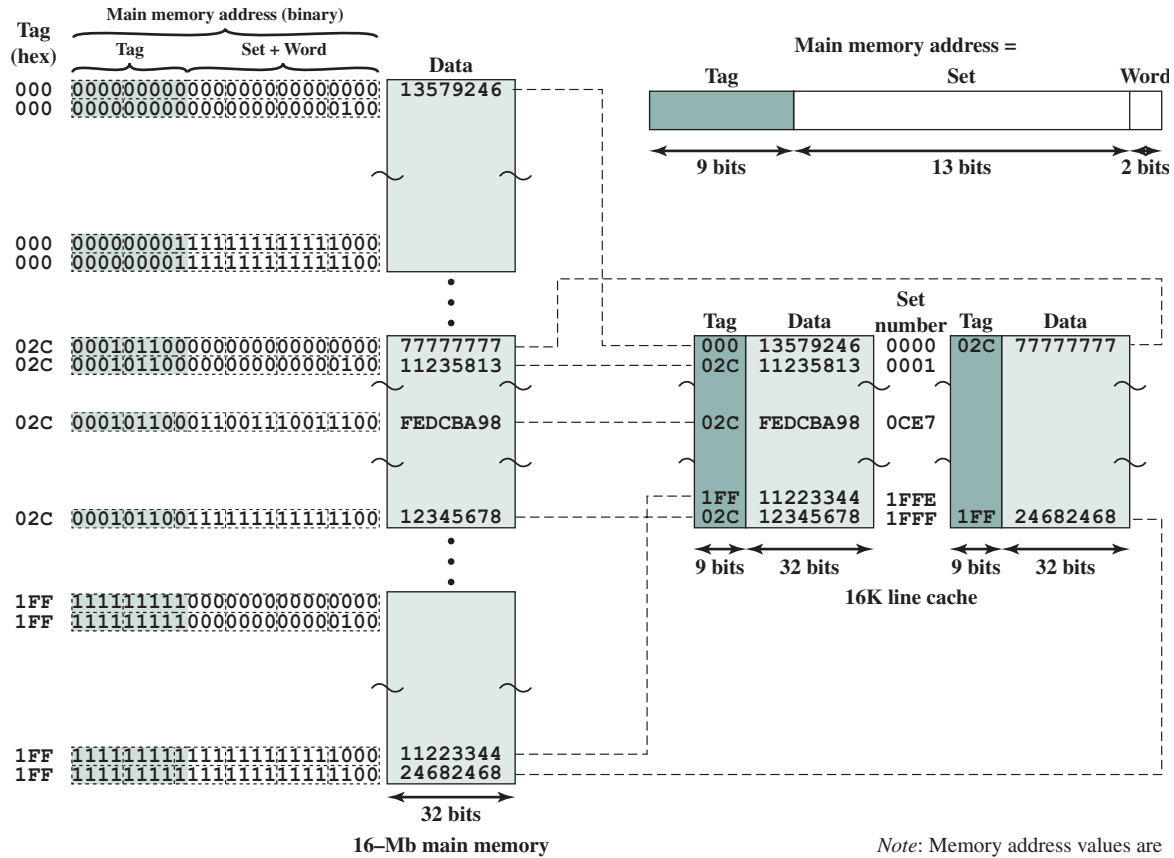
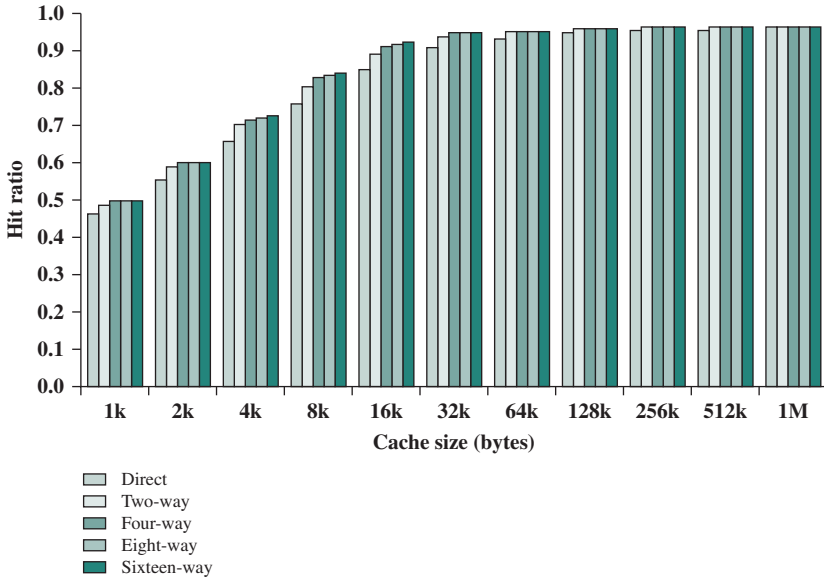


Figure 4.15 Two-Way Set-Associative Mapping Example



**Figure 4.16** Varying Associativity over Cache Size

In the extreme case of  $v = m, k = 1$ , the set-associative technique reduces to direct mapping, and for  $v = 1, k = m$ , it reduces to associative mapping. The use of two lines per set ( $v = m/2, k = 2$ ) is the most common set-associative organization. It significantly improves the hit ratio over direct mapping. Four-way set associative ( $v = m/4, k = 4$ ) makes a modest additional improvement for a relatively small additional cost [MAYB84, HILL89]. Further increases in the number of lines per set have little effect.

Figure 4.16 shows the results of one simulation study of set-associative cache performance as a function of cache size [GENU04]. The difference in performance between direct and two-way set associative is significant up to at least a cache size of 64 kB. Note also that the difference between two-way and four-way at 4 kB is much less than the difference in going from for 4 kB to 8 kB in cache size. The complexity of the cache increases in proportion to the associativity, and in this case would not be justifiable against increasing cache size to 8 or even 16 kB. A final point to note is that beyond about 32 kB, increase in cache size brings no significant increase in performance.

The results of Figure 4.16 are based on simulating the execution of a GCC compiler. Different applications may yield different results. For example, [CANT01] reports on the results for cache performance using many of the CPU2000 SPEC benchmarks. The results of [CANT01] in comparing hit ratio to cache size follow the same pattern as Figure 4.16, but the specific values are somewhat different.