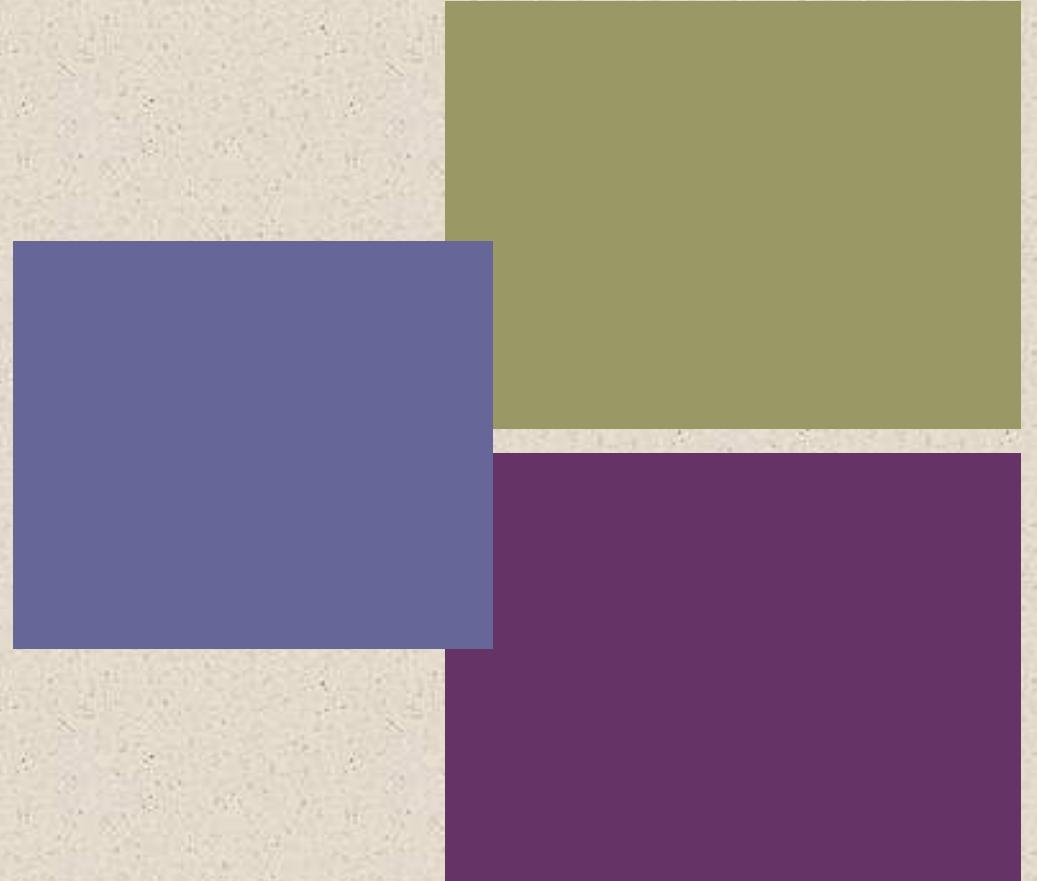


Cache Memory



+

Location

Internal (e.g. processor registers, cache, main memory)

External (e.g. optical disks, magnetic disks, tapes)

Capacity

Number of words

Number of bytes

Unit of Transfer

Word

Block

Access Method

Sequential

Direct

Random

Associative

Performance

Access time

Cycle time

Transfer rate

Physical Type

Semiconductor

Magnetic

Optical

Magneto-optical

Physical Characteristics

Volatile/nonvolatile

Erasable/nonerasable

Organization

Memory modules

Characteristics of Memory Systems

- Location
 - Refers to whether memory is internal and external to the computer
 - Internal memory is often equated with main memory
 - Processor requires its own local memory, in the form of registers
 - Cache is another form of internal memory
 - External memory consists of peripheral storage devices that are accessible to the processor via I/O controllers
- Capacity
 - Memory is typically expressed in terms of bytes
- Unit of transfer
 - For internal memory the unit of transfer is equal to the number of electrical lines into and out of the memory module

Method of Accessing Units of Data

Sequential access

Memory is organized into units of data called records

Access must be made in a specific linear sequence

Access time is variable

Tape

✓ Direct access

Involves a shared read-write mechanism

Individual blocks or records have a unique address based on physical location

Access time is variable

Random access

Each addressable location in memory has a unique, physically wired-in addressing mechanism

The time to access a given location is independent of the sequence of prior accesses and is constant

Any location can be selected at random and directly addressed and accessed

Main memory and some cache systems are random access

✓ Associative

A word is retrieved based on a portion of its contents rather than its address

Each location has its own addressing mechanism and retrieval time is constant independent of location or prior access patterns

~~Cache memories may employ associative access~~

Capacity and Performance:

The two most important characteristics of memory

Three performance parameters are used:

Access time (latency)

- For random-access memory it is the time it takes to perform a read or write operation
- For non-random-access memory it is the time it takes to position the read-write mechanism at the desired location

Memory cycle time

- Access time plus any additional time required before second access can commence
- Additional time may be required for transients to die out on signal lines or to regenerate data if they are read destructively
- Concerned with the system bus, not the processor

Transfer rate

- The rate at which data can be transferred into or out of a memory unit
- For random-access memory it is equal to $1/(cycle\ time)$

Memory

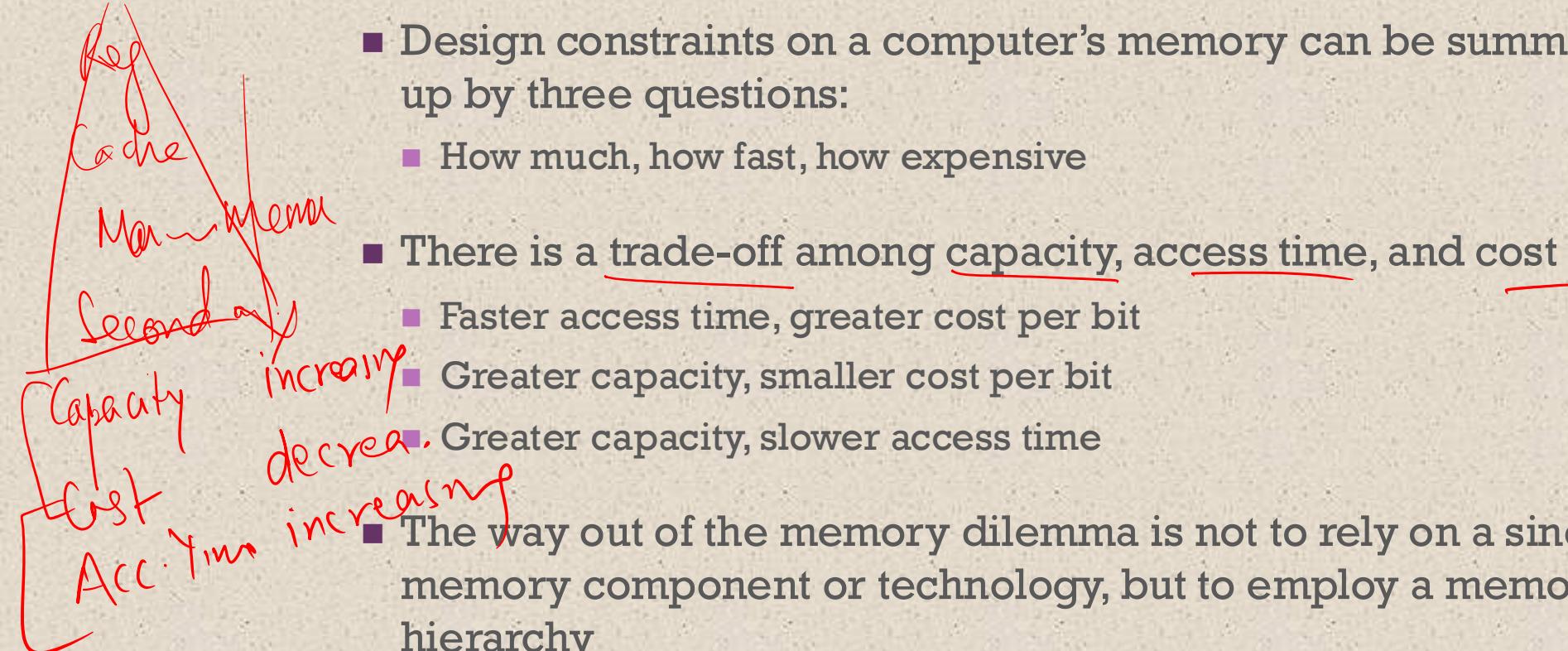


- The most common forms are:
 - Semiconductor memory
 - Magnetic surface memory
 - Optical
 - Magneto-optical
- Several physical characteristics of data storage are important:
 - Volatile memory
 - Information decays naturally or is lost when electrical power is switched off
 - Nonvolatile memory
 - 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



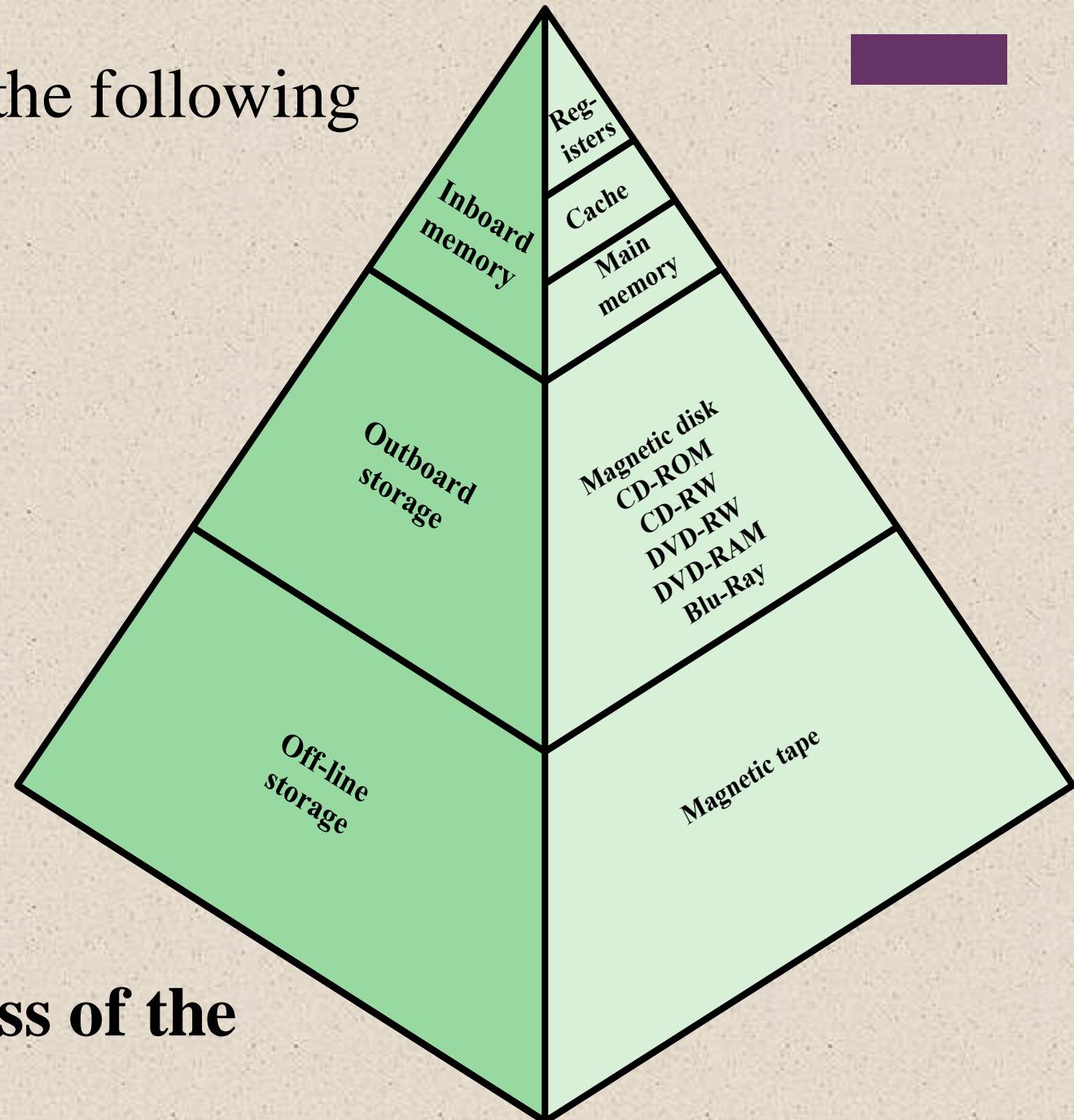
As one goes down the hierarchy, the following occur:

a. Decreasing cost per bit

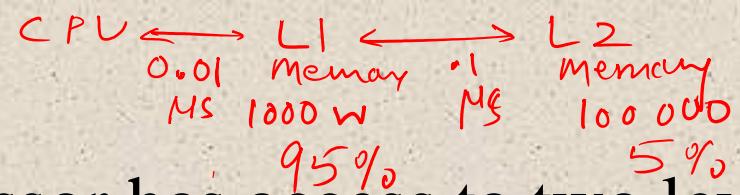
b. Increasing capacity

c. Increasing access time

d. Decreasing frequency of access of the memory by the processor



Example



Suppose that the processor has access to two levels of memory. Level 1 contains 1000 words and has an access time of 0.01 μ s; level 2 contains 100,000 words and has an access time of 0.1 μ s. Assume that if a word to be accessed is in level 1, then the processor accesses it directly. If it is in level 2, then the word is first transferred to level 1 and then accessed by the processor.

Suppose 95% of the memory accesses are found in level 1. Then the average time to access a word can be expressed as

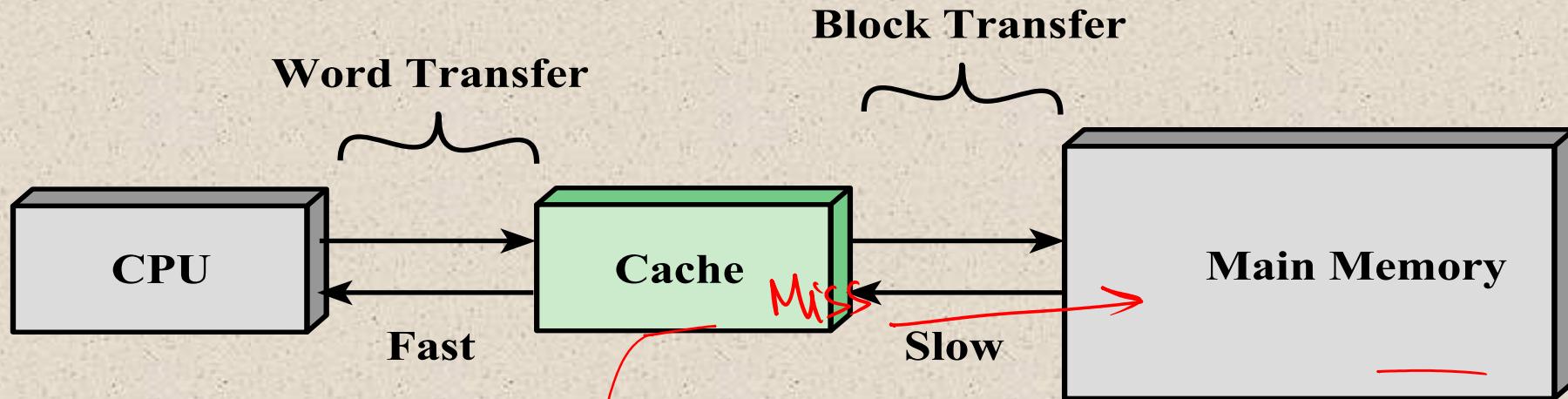
$$(0.95)(0.01 \mu\text{s}) + (0.05)(0.01 \mu\text{s} + 0.1 \mu\text{s}) = 0.0095 + 0.0055 = 0.015 \mu\text{s}$$

The average access time is much closer to 0.01 μ s than to 0.1 μ s, as desired.

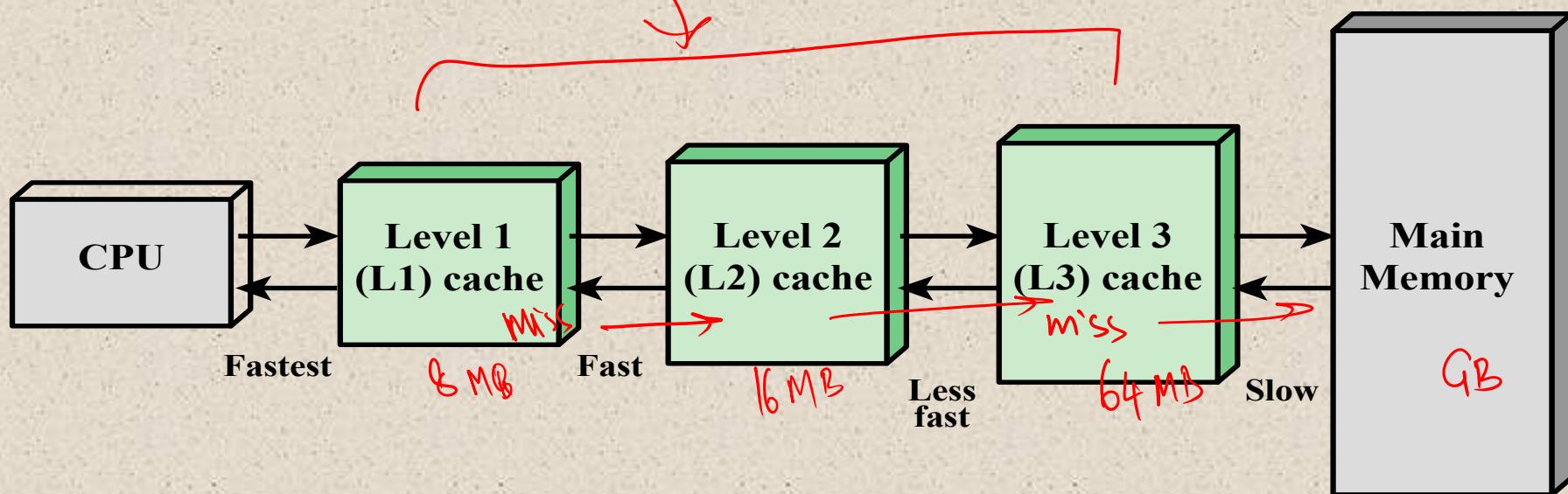
$$.05 \times .01 \mu\text{s} + .95 (.11 \mu\text{s}) = 0.10455 \mu\text{s}$$

+ Memory

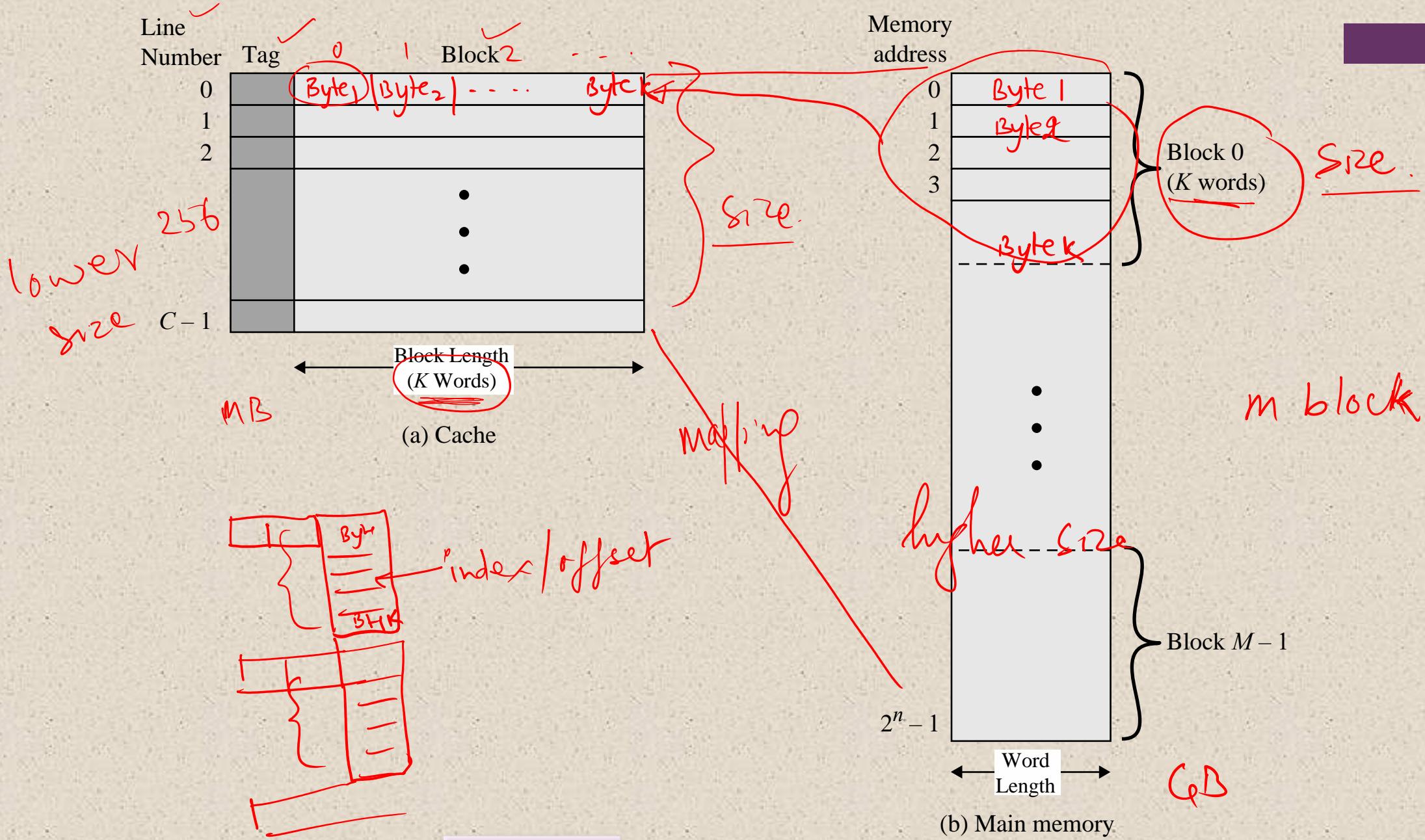
- The use of three levels exploits the fact that semiconductor memory comes in a variety of types which differ in speed and cost
- Data are stored more permanently on external mass storage devices
- External, nonvolatile memory is also referred to as **secondary memory** or **auxiliary memory**
- Disk cache
 - A portion of main memory can be used as a buffer to hold data temporarily that is to be read out to disk
 - A few large transfers of data can be used instead of many small transfers of data
 - Data can be retrieved rapidly from the software cache rather than slowly from the disk



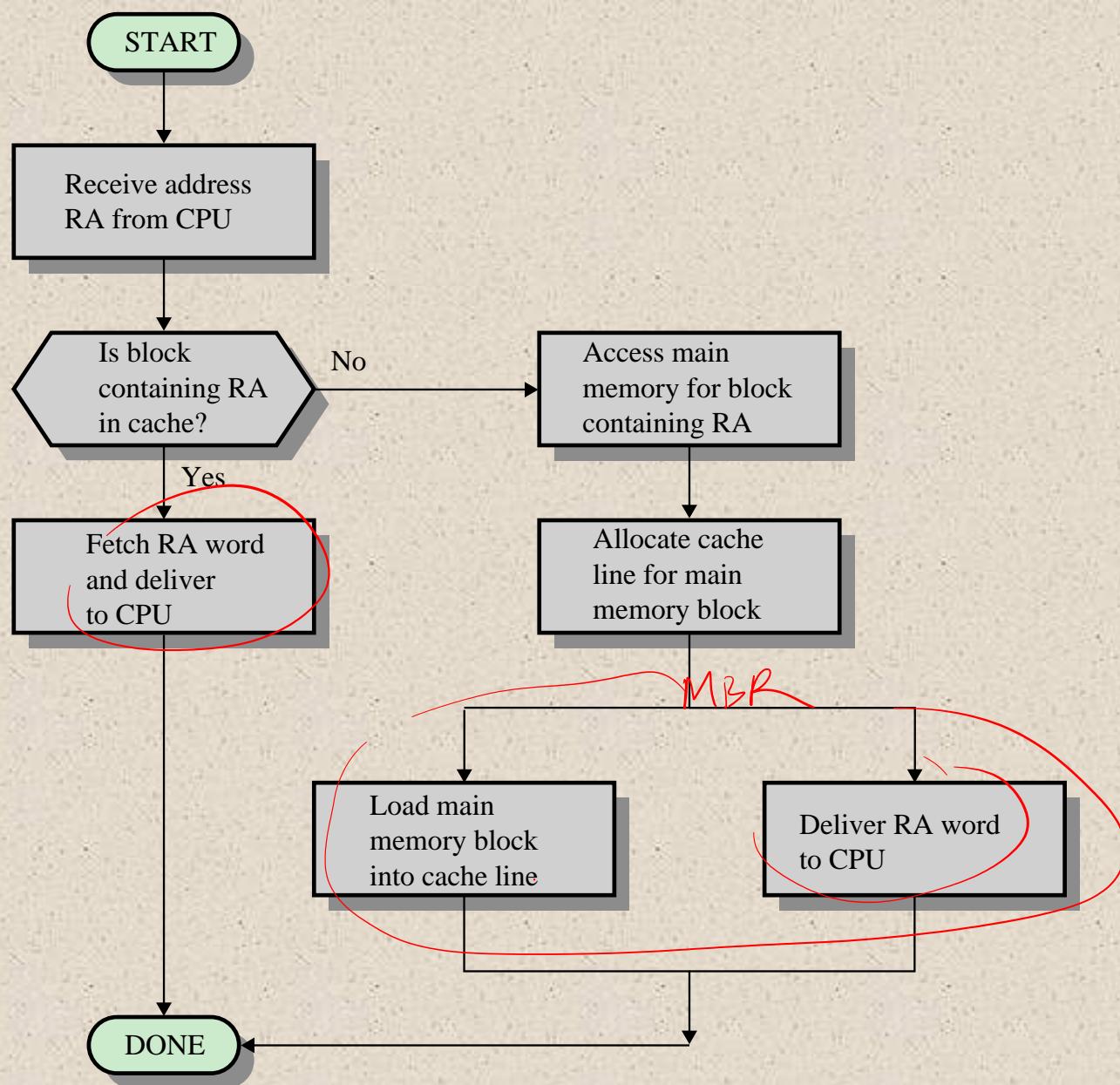
(a) Single cache



(b) Three-level cache organization



Cache/Main-Memory Structure



- Read address (RA) of a word to be read

Figure 4.5 Cache Read Operation

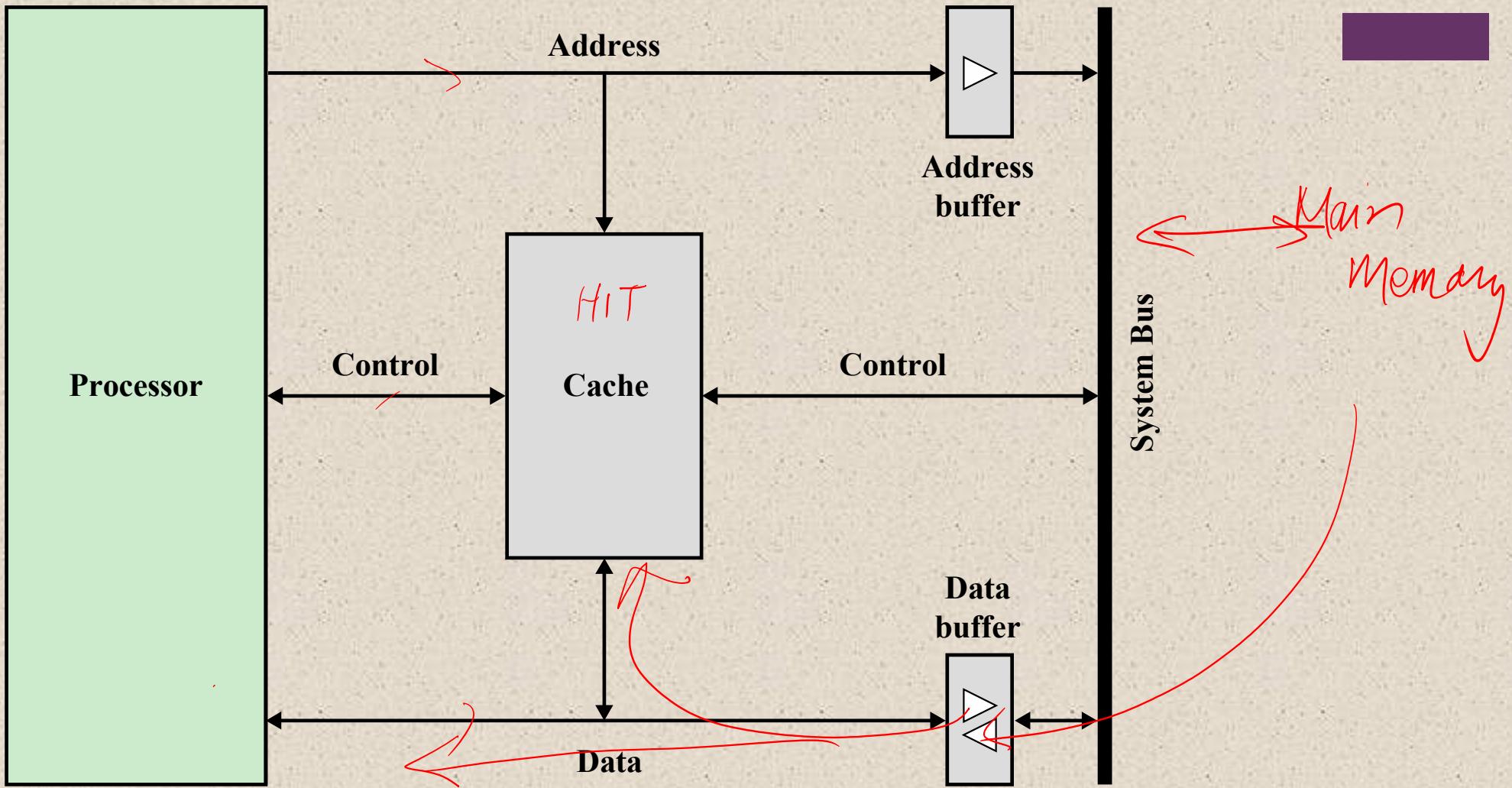


Figure 4.6 Typical Cache Organization

Elements of Cache Design

Cache Addresses

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

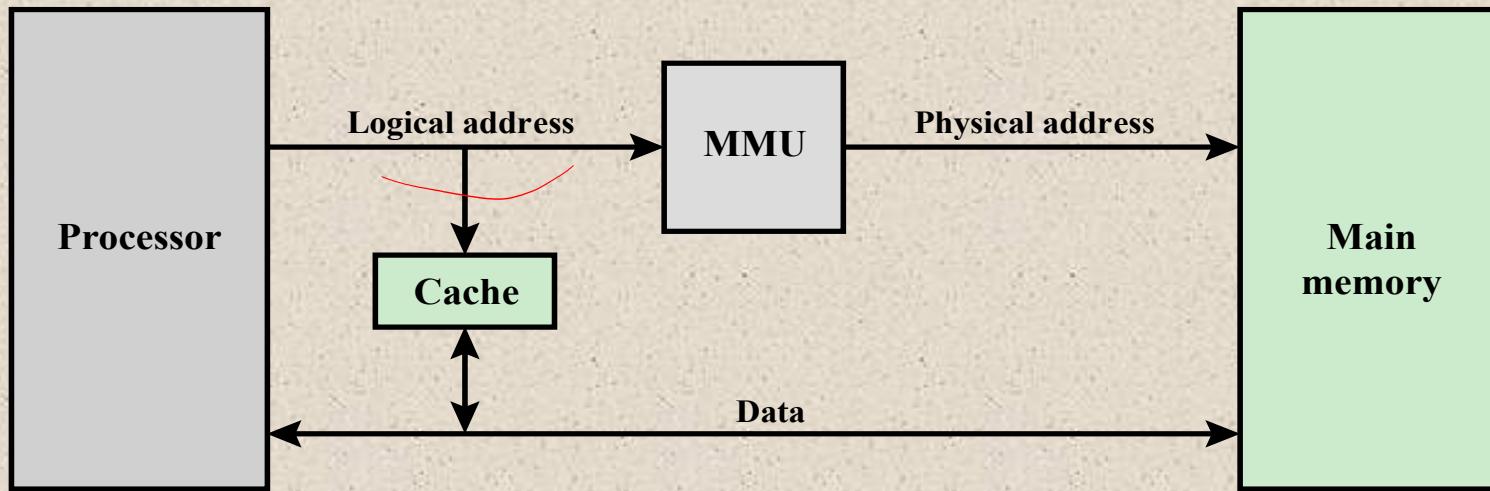
Unified or split

+ Cache Addresses

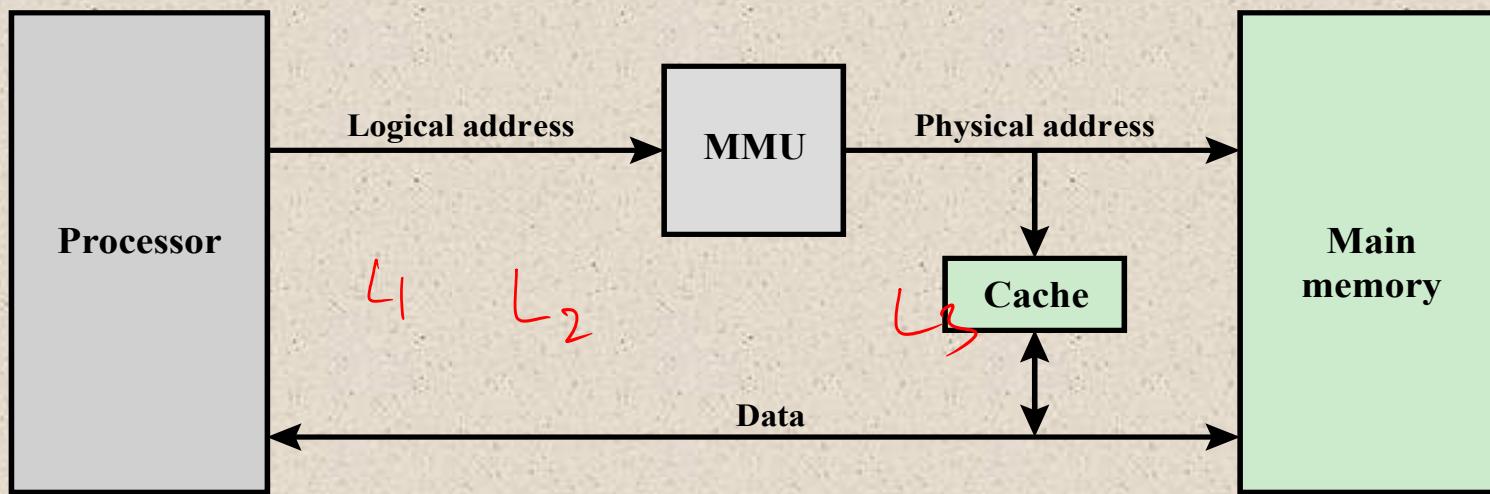
Virtual Memory

■ 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



(a) Logical Cache



(b) Physical Cache

Cache Sizes of Some Processors

Processor	Type	Year of Introduction	L1 Cache	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 MTab	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	Workstation/ 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.

(Table can be found on page 134 in the textbook.)

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

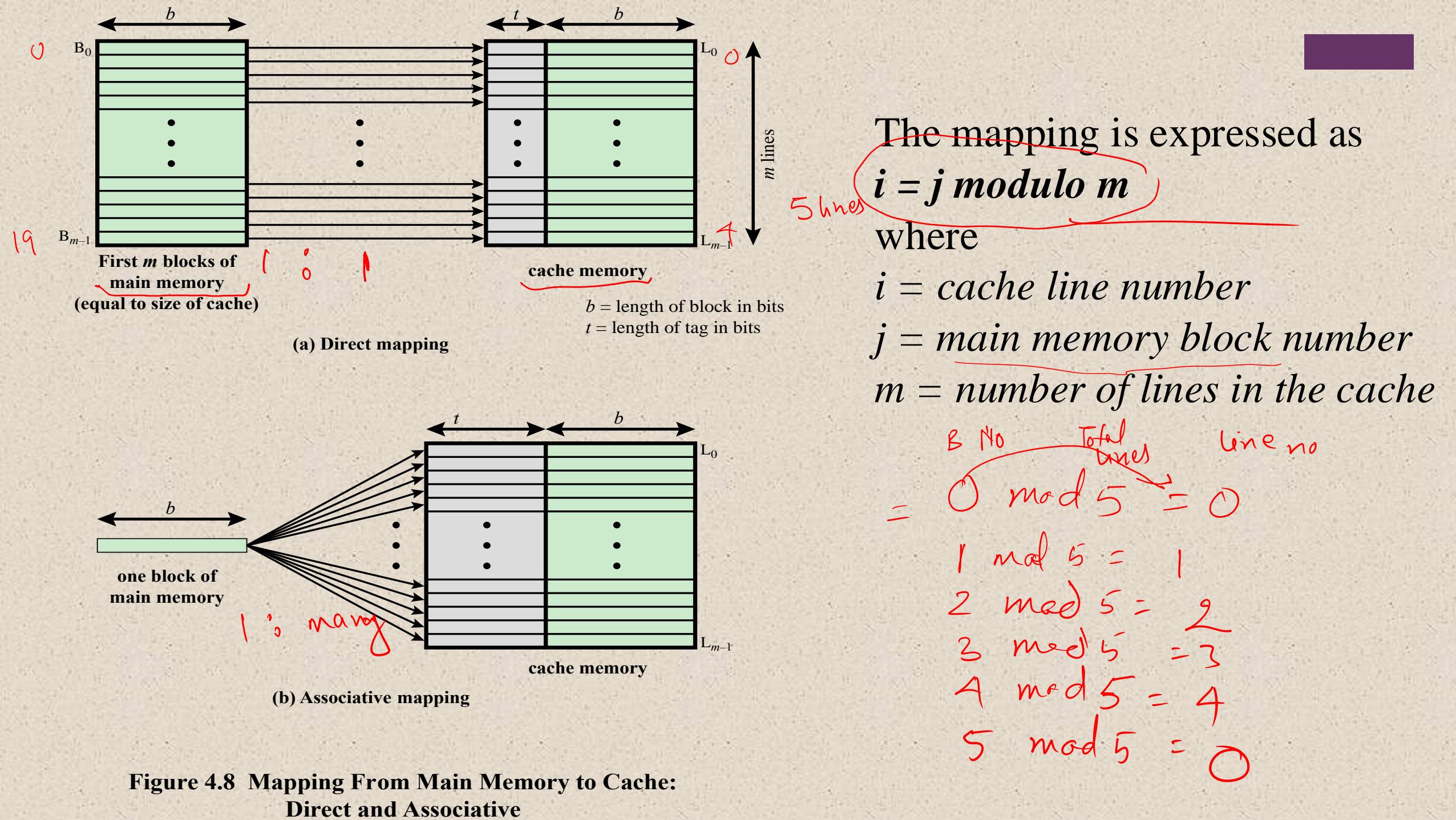


Figure 4.8 Mapping From Main Memory to Cache:
Direct and Associative

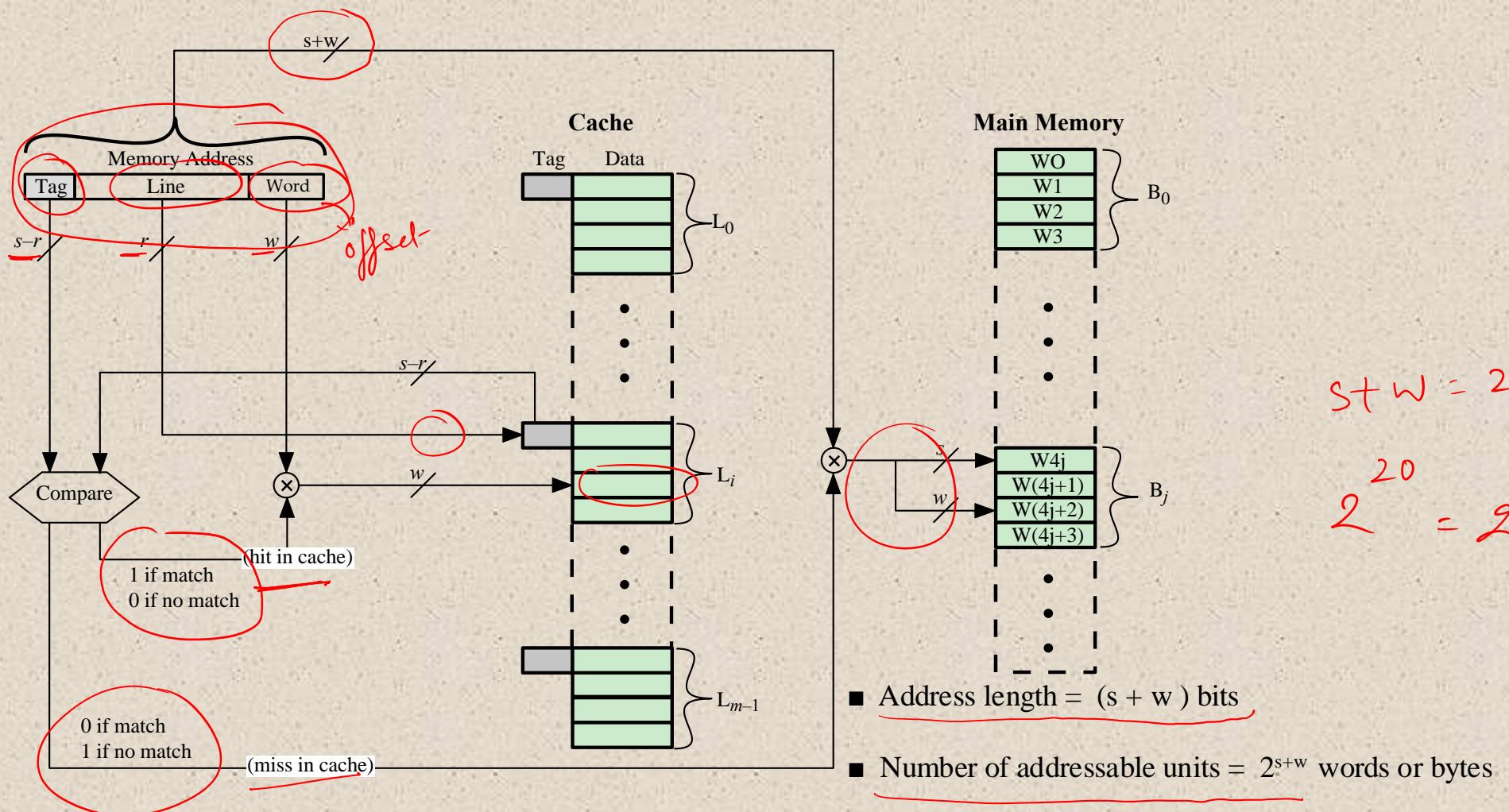


Figure 4.9 Direct-Mapping Cache Organization

- 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 = $m = 2^r$

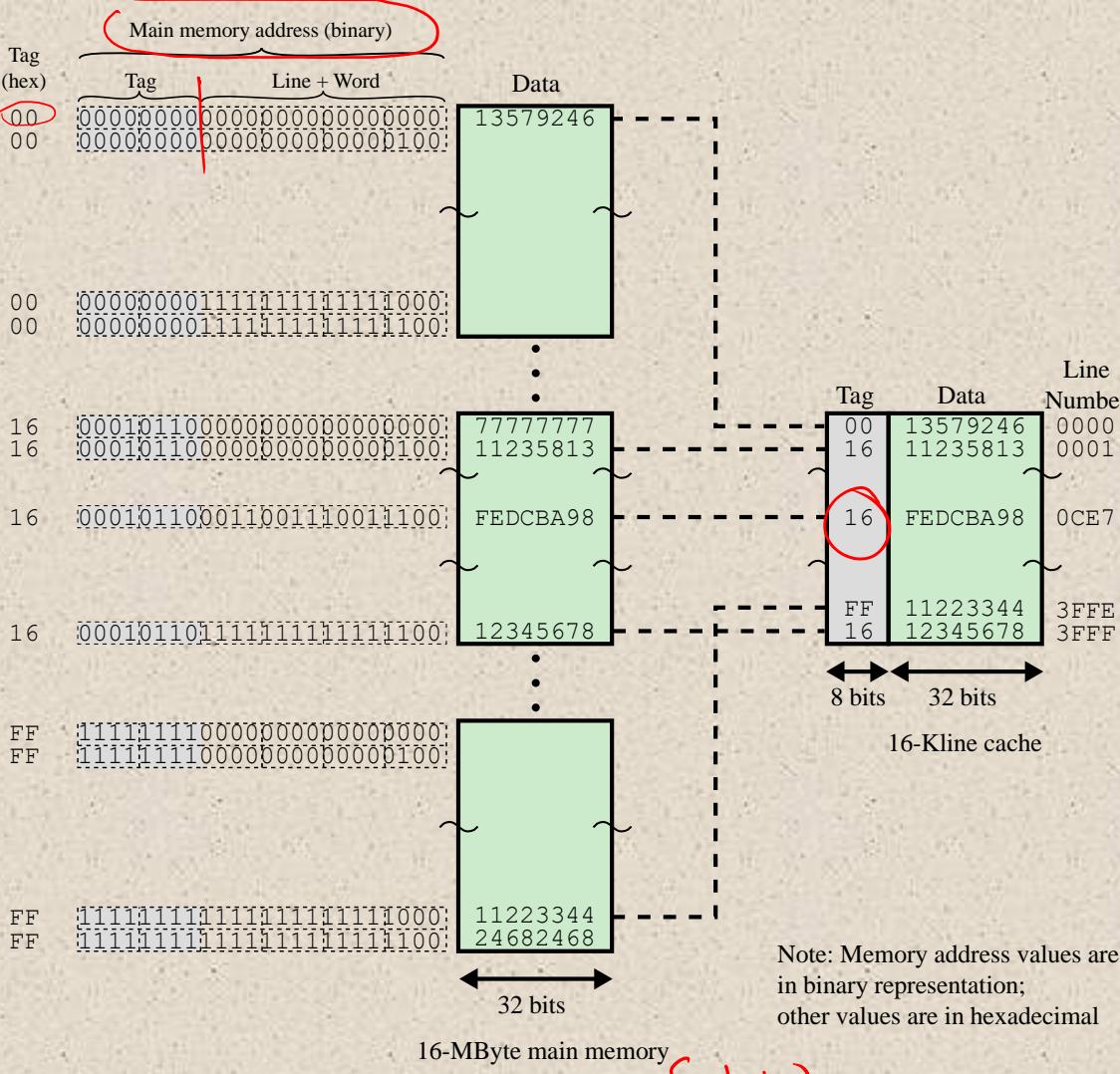
- Size of cache = 2^{r+w} words or bytes

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

Add.

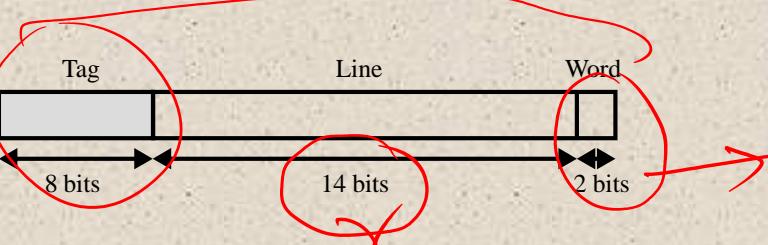
$$s + w = 20 \text{ bits}$$

$$2^{20} = 2^{s+w}$$



16-MByte main memory

Main memory address =



$$S - R = tag$$

Figure 4.10 Direct Mapping Example

+

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*).

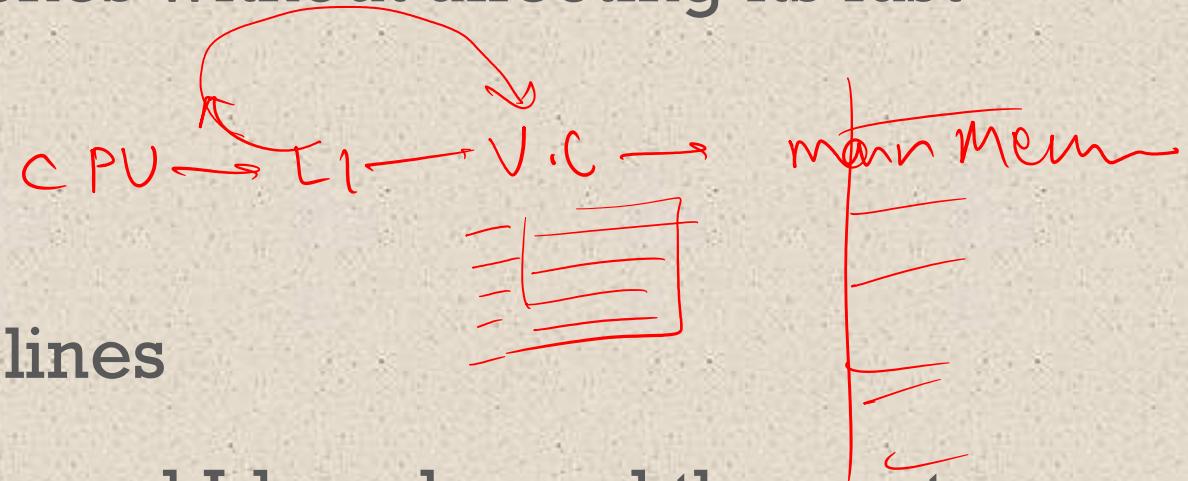




Victim Cache



- Originally proposed as an approach to reduce the conflict misses of direct mapped caches without affecting its fast access time
- Fully associative cache
- Typical size is 4 to 16 cache lines
- Residing between direct mapped L1 cache and the next level of memory



Associative mapping overcomes the disadvantage of direct mapping by permitting each main memory block to be loaded into any line of the cache

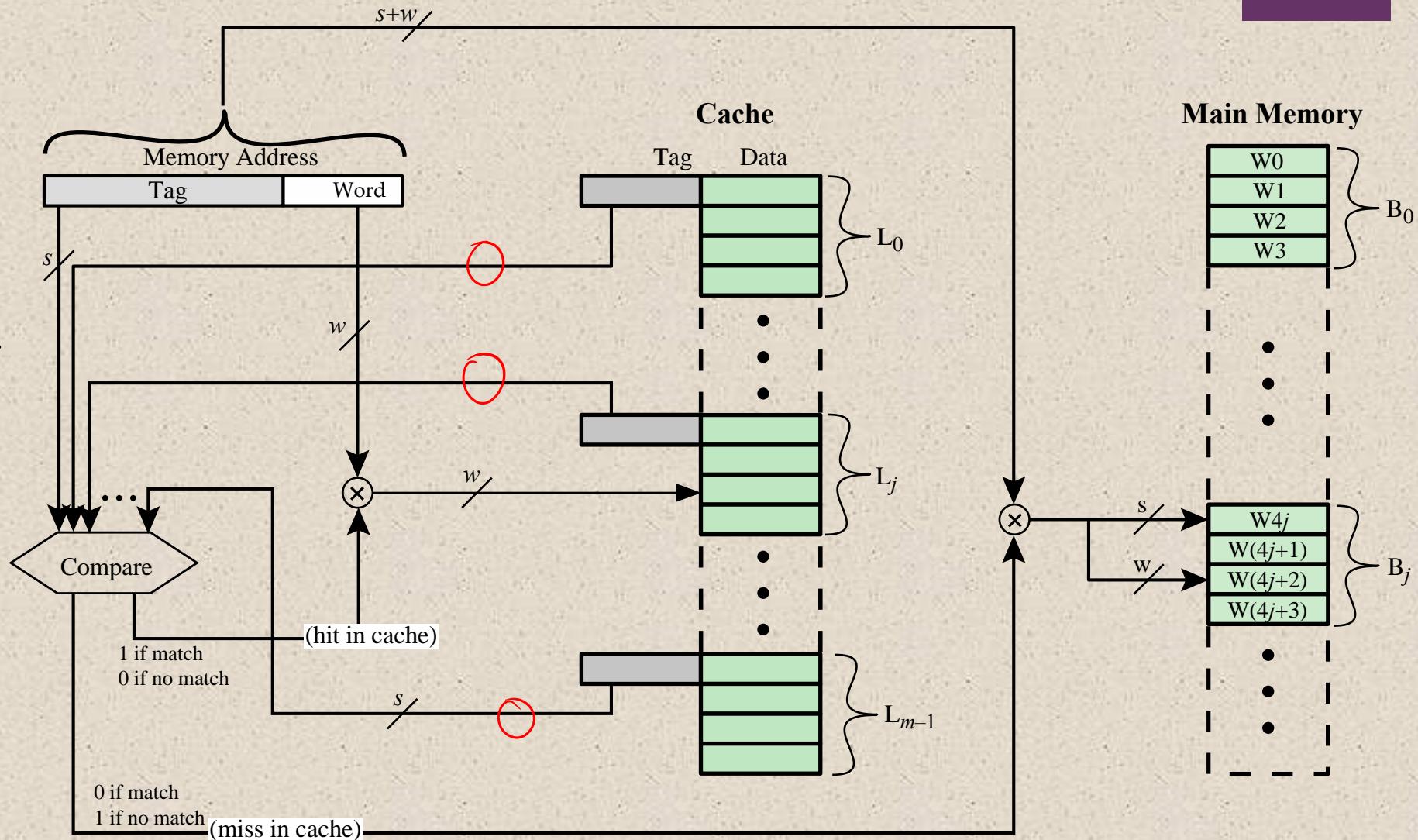


Figure 4.11 Fully Associative Cache Organization

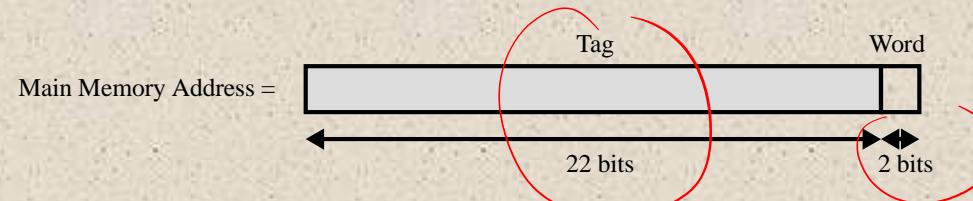
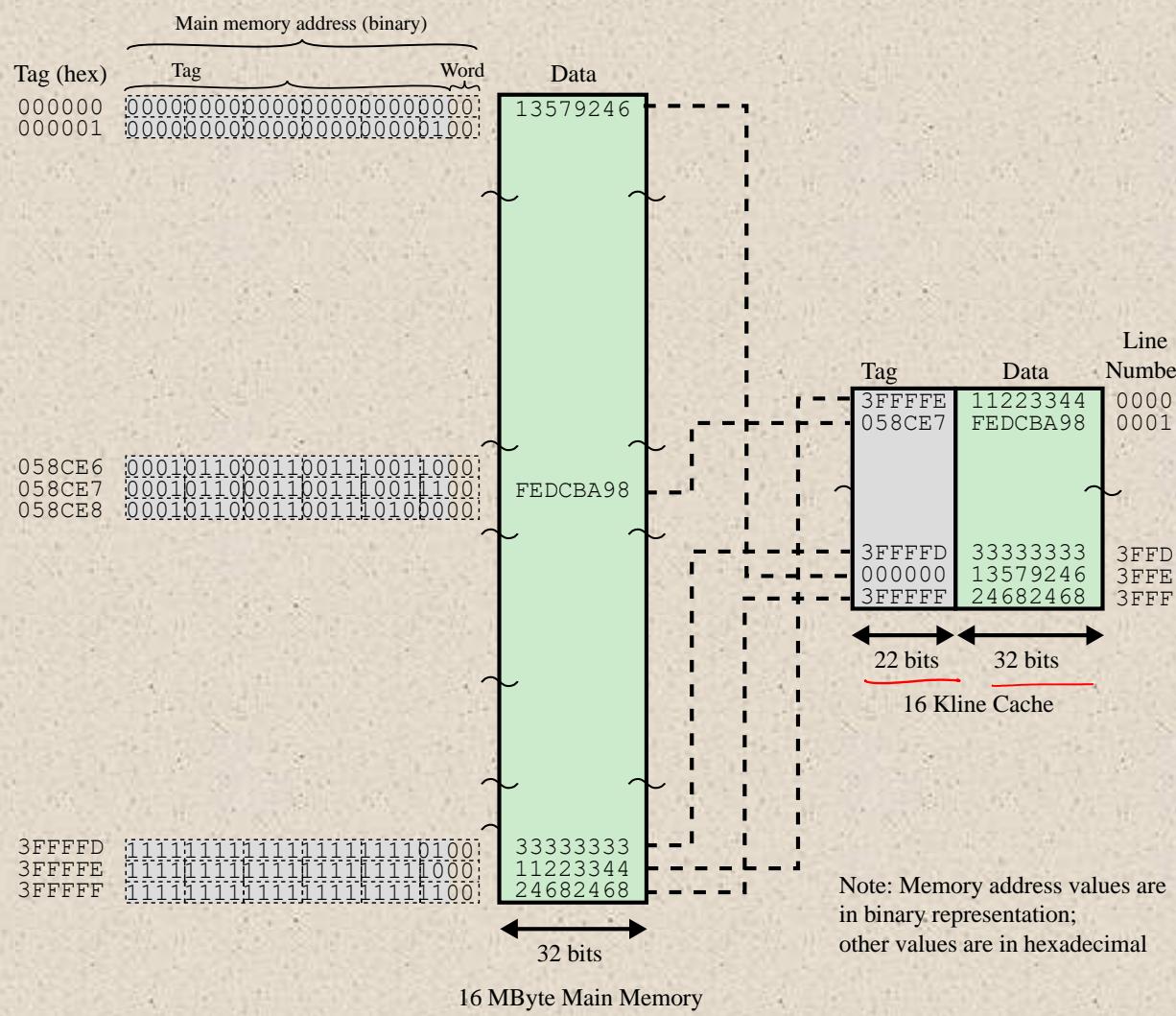


Figure 4.12 Associative Mapping Example

+ Associative 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 = undetermined
- Size of tag = s bits



Set Associative Mapping

- is a Compromise that exhibits the strengths of both the direct and associative approaches while reducing their disadvantages

- Cache consists of a number of sets

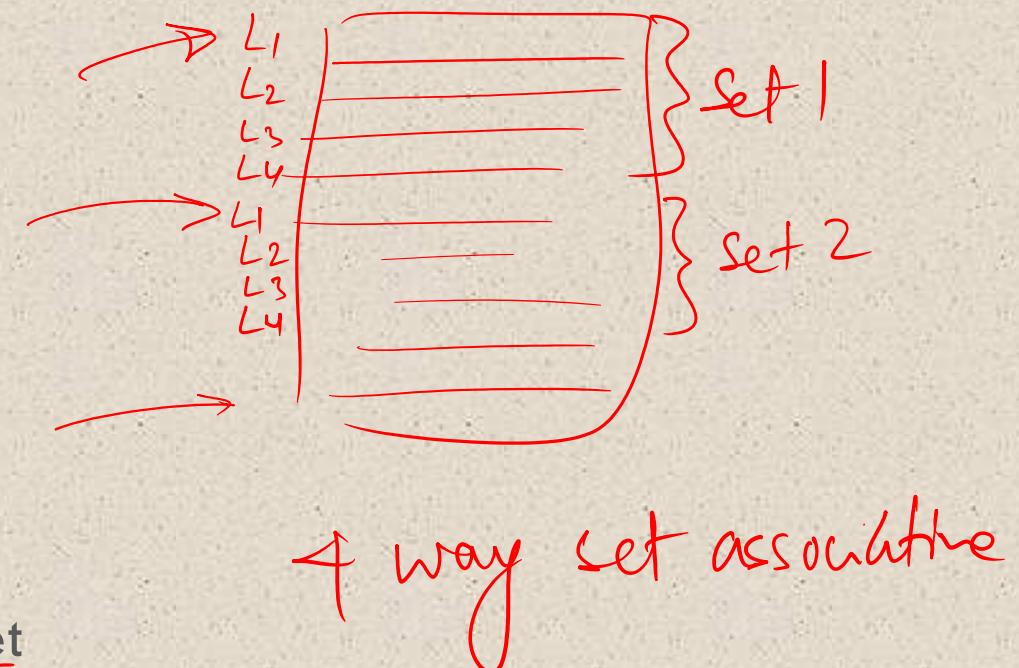
- Each set contains a number of lines

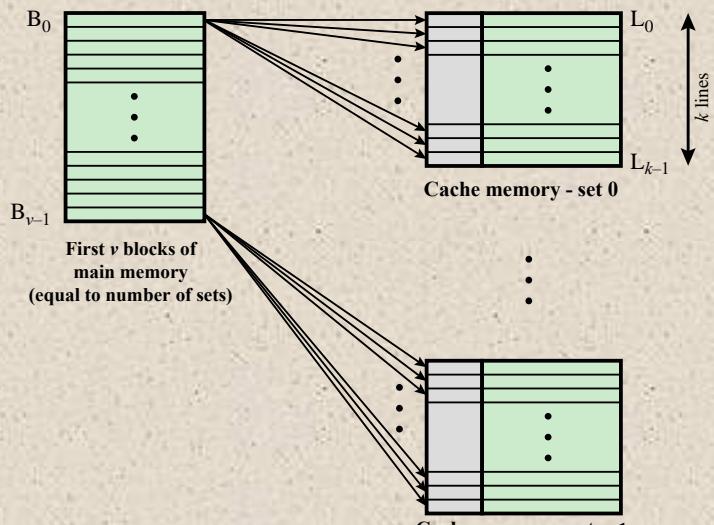
- A given block maps to any line in a given set

- e.g. 2 lines per set

- 2 way associative mapping

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





(a) v associative-mapped caches

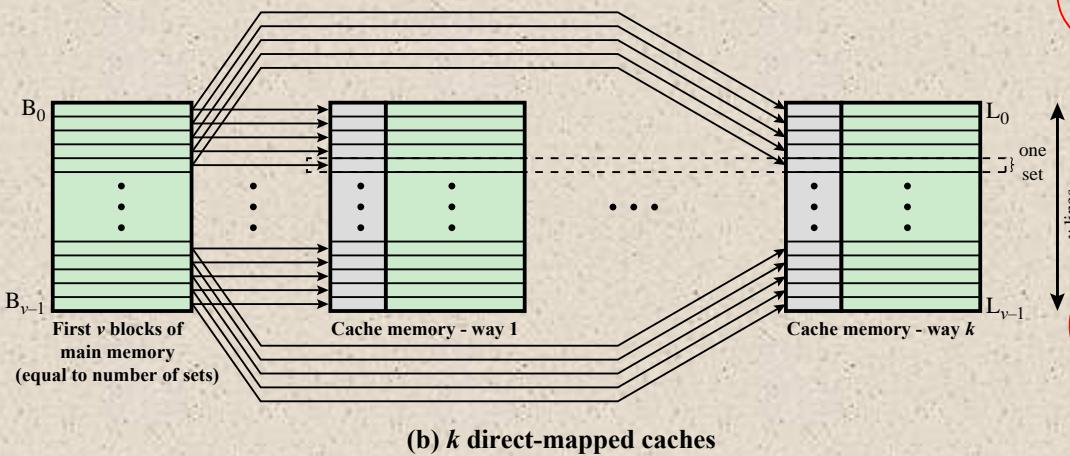


Figure 4.13 Mapping From Main Memory to Cache:
 k -way Set Associative

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

$$m = v * k \\ i = j \text{ modulo } v$$

where

$i = \text{cache set number}$

$j = \text{main memory block number}$

$m = \text{number of lines in the cache}$

$v = \text{number of sets}$

$k = \text{number of lines in each set}$

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 .

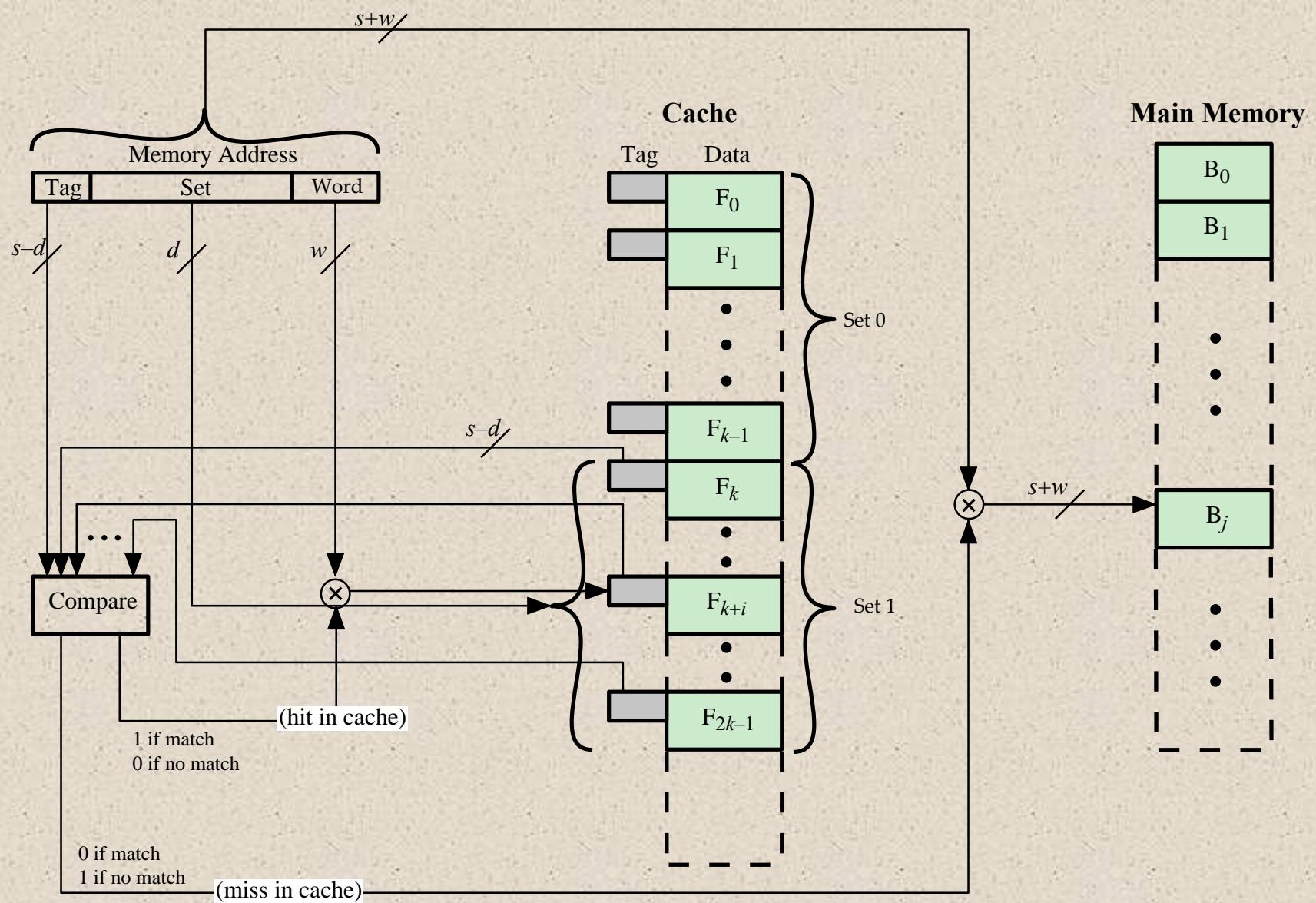


Figure 4.14 *k*-Way Set Associative Cache Organization

Set Associative 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 set = k
- Number of sets = $v = 2^d$
- Number of lines in cache = $m=kv = k * 2^d$
- Size of cache = $k * 2^{d+w}$ words or bytes
- Size of tag = $(s - d)$ bits



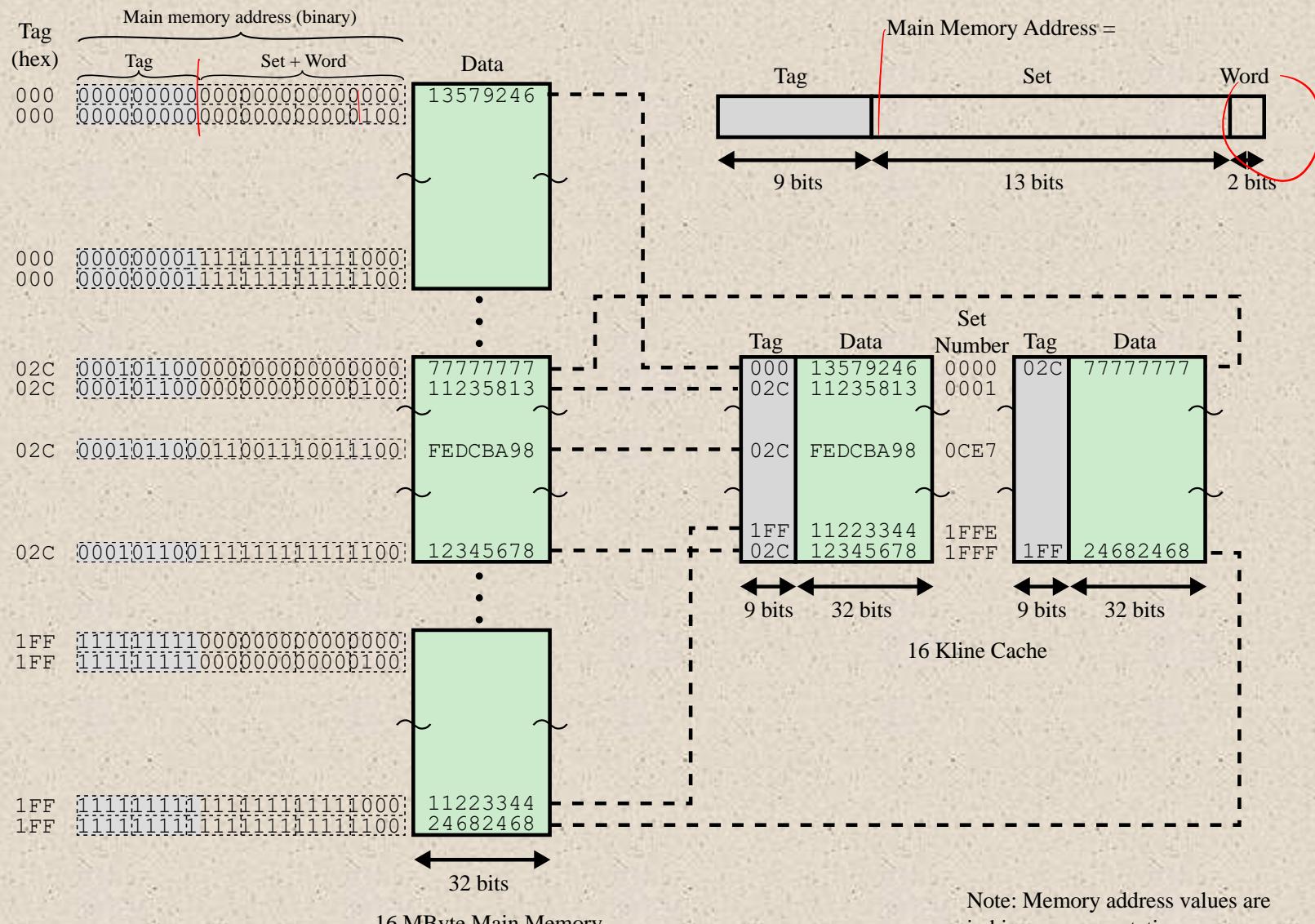


Figure 4.15 Two-Way Set Associative Mapping Example

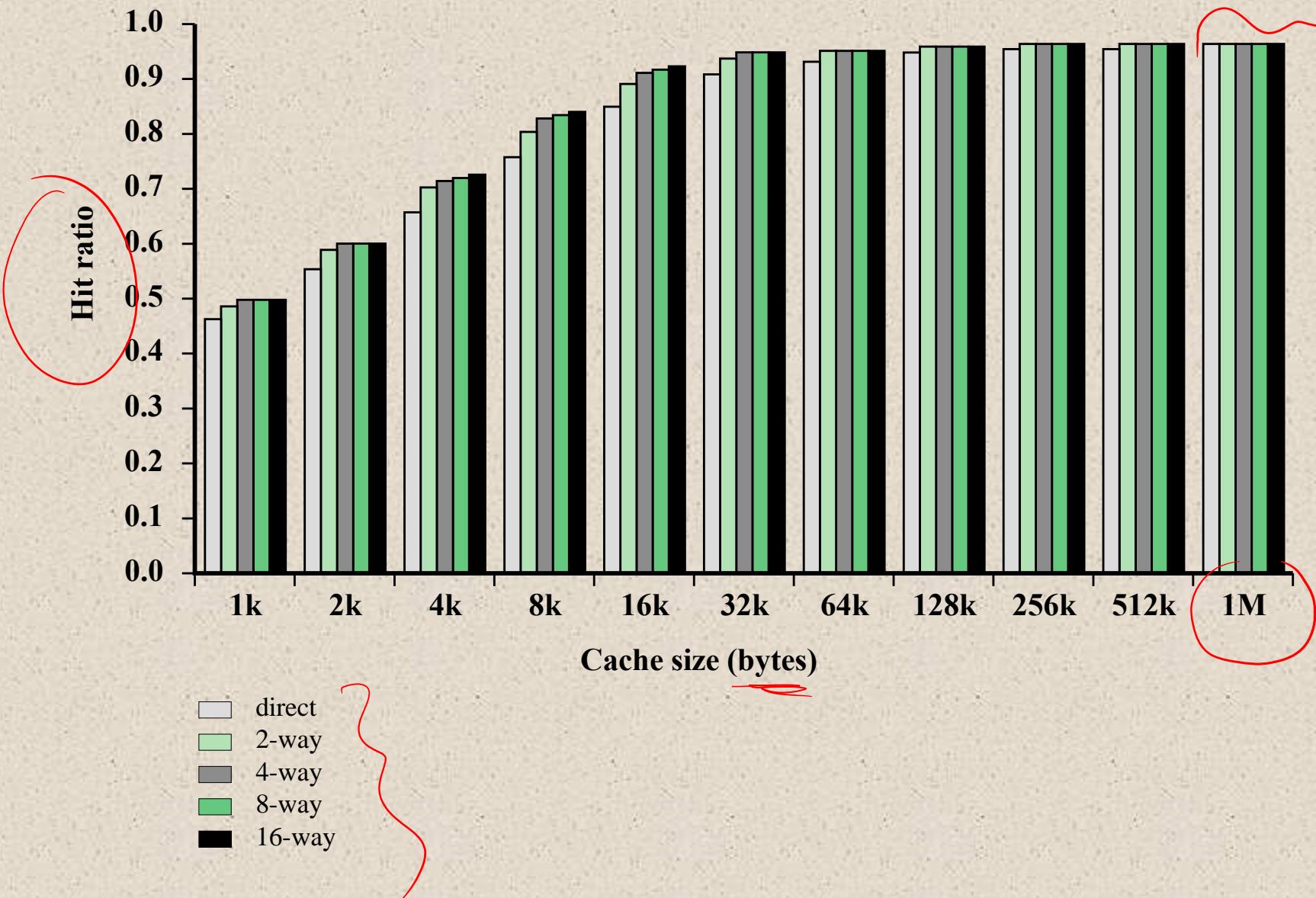
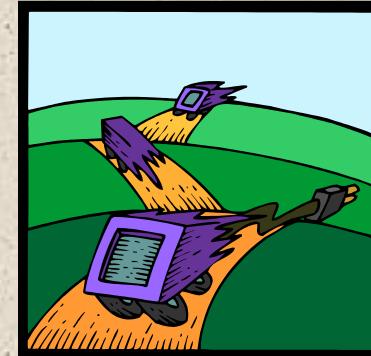


Figure 4.16 Varying Associativity over Cache Size

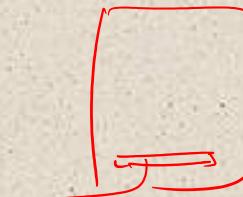
Replacement Algorithms



- Once the cache has been filled, when a new block is brought into the cache, one of the existing blocks must be replaced
- For direct mapping there is only one possible line for any particular block and no choice is possible
- For the associative and set-associative techniques a replacement algorithm is needed
- To achieve high speed, an algorithm must be implemented in hardware

The most common replacement algorithms are:

- Least recently used (LRU)
 - Most effective
 - Replace that block in the set that has been in the cache longest with no reference to it
 - Because of its simplicity of implementation, LRU is the most popular replacement algorithm
- First-in-first-out (FIFO)
 - Replace that block in the set that has been in the cache longest
 - Easily implemented as a round-robin or circular buffer technique
- Least frequently used (LFU)
 - Replace that block in the set that has experienced the fewest references
 - Could be implemented by associating a counter with each line



Write Policy

When a block that is resident in the cache is to be replaced there are two cases to consider:

If the old block in the cache has not been altered then it may be overwritten with a new block without first writing out the old block

If at least one write operation has been performed on a word in that line of the cache then main memory must be updated by writing the line of cache out to the block of memory before bringing in the new block

There are two problems to contend with:

More than one device may have access to main memory

A more complex problem occurs when multiple processors are attached to the same bus and each processor has its own local cache - if a word is altered in one cache it could conceivably invalidate a word in other caches



Write Through and Write Back

- Write through
 - Simplest technique
 - All write operations are made to main memory as well as to the cache
 - The main disadvantage of this technique is that it generates substantial memory traffic and may create a bottleneck

- Write back
 - Minimizes memory writes
 - Updates are made only in the cache
 - Portions of main memory are invalid and hence accesses by I/O modules can be allowed only through the cache
 - This makes for complex circuitry and a potential bottleneck

Line Size

When a block of data is retrieved and placed in the cache not only the desired word but also some number of adjacent words are retrieved

As the block size increases more useful data are brought into the cache

Two specific effects come into play:

- Larger blocks reduce the number of blocks that fit into a cache
- As a block becomes larger each additional word is farther from the requested word



Multilevel Caches

- As logic density has increased it has become possible to have a cache on the same chip as the processor
- The on-chip cache reduces the processor's external bus activity and speeds up execution time and increases overall system performance
 - When the requested instruction or data is found in the on-chip cache, the bus access is eliminated
 - On-chip cache accesses will complete appreciably faster than would even zero-wait state bus cycles
 - During this period the bus is free to support other transfers
- Two-level cache:
 - Internal cache designated as level 1 (L1)
 - External cache designated as level 2 (L2)
- Potential savings due to the use of an L2 cache depends on the hit rates in both the L1 and L2 caches
- The use of multilevel caches complicates all of the design issues related to caches, including size, replacement algorithm, and write policy

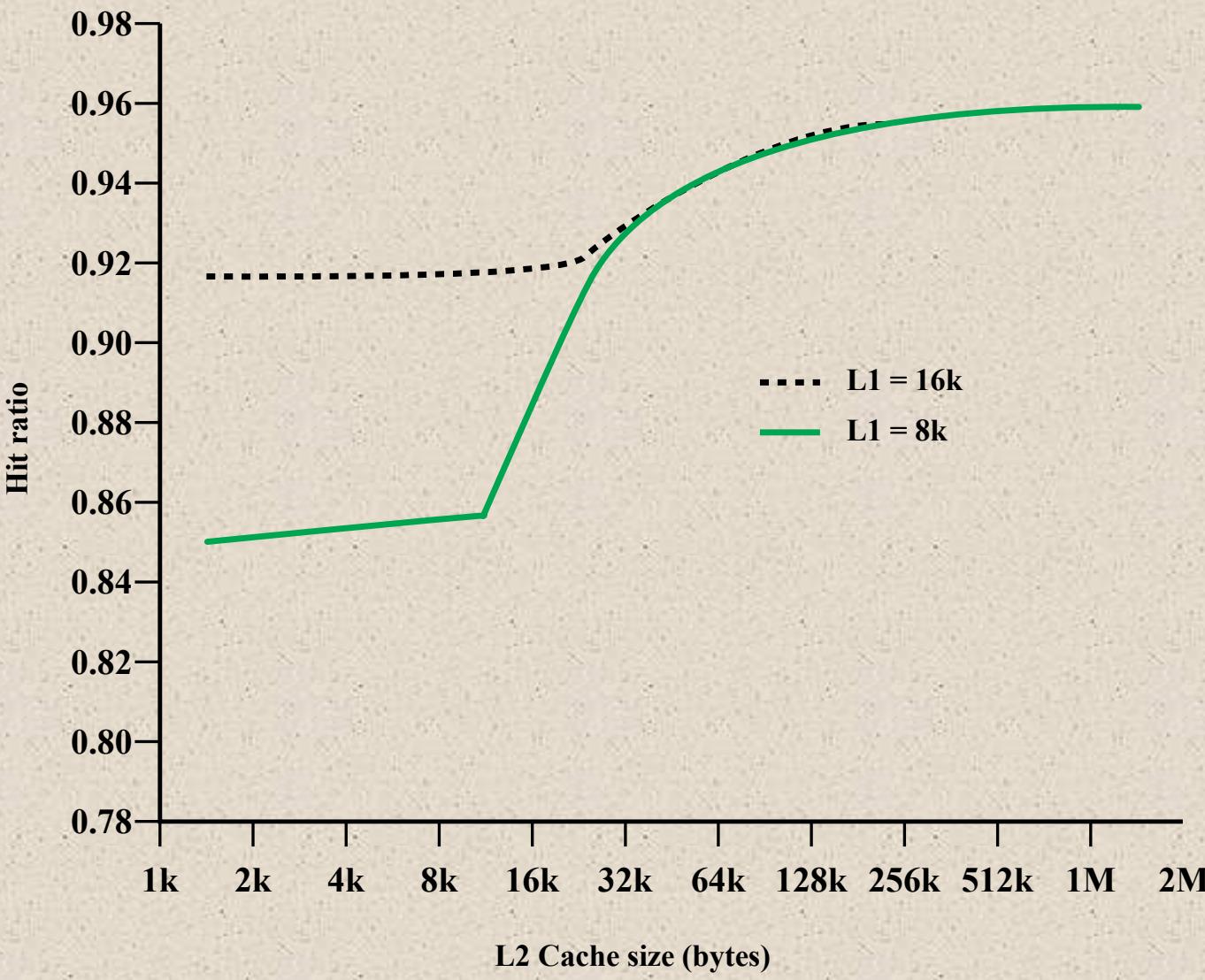


Figure 4.17 Total Hit Ratio (L1 and L2) for 8 Kbyte and 16 Kbyte L1



Unified Versus Split Caches

- Has become common to split cache:
 - One dedicated to instructions
 - One dedicated to data
 - Both exist at the same level, typically as two L1 caches
- Advantages of unified cache:
 - Higher hit rate
 - Balances load of instruction and data fetches automatically
 - Only one cache needs to be designed and implemented
- Trend is toward split caches at the L1 and unified caches for higher levels
- Advantages of split cache:
 - Eliminates cache contention between instruction fetch/decode unit and execution unit
 - Important in pipelining

Table 4.4

**Intel
Cache
Evolution**

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

(Table is on page 150 in the textbook.)

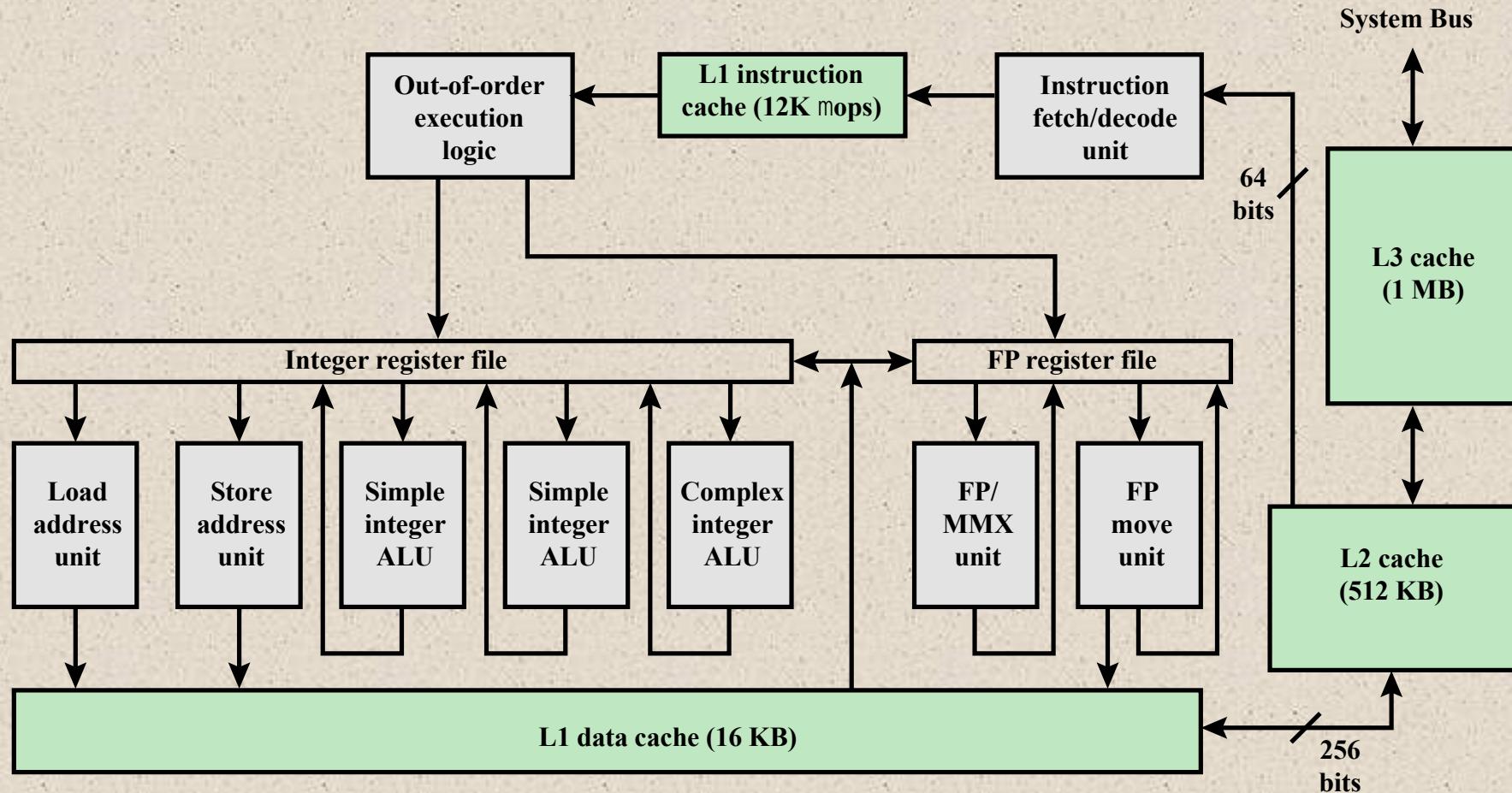
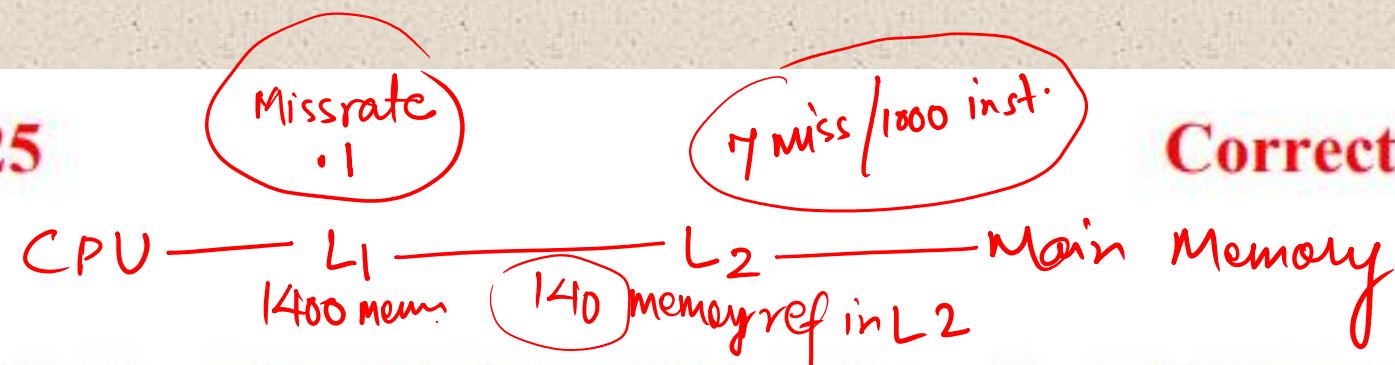


Figure 4.18 Pentium 4 Block Diagram

Table 4.5 Pentium 4 Cache Operating Modes

Control Bits		Operating Mode		
CD	NW	Cache Fills	Write Throughs	Invalidates
0	0	Enabled	Enabled	Enabled
1	0	Disabled	Enabled	Enabled
1	1	Disabled	Disabled	Disabled

Note: CD = 0; NW = 1 is an invalid combination.

Question Number : 25**Correct : 1 Wrong : 0**

Consider a two-level cache hierarchy with L1 and L2 caches. An application incurs 1.4 memory accesses per instruction on average. For this application, the miss rate of L1 cache is 0.1; the L2 cache experiences, on average, 7 misses per 1000 instructions. The miss rate of L2 expressed correct to two decimal places is _____.

Given = 1.4 memory access/instrⁿ on avg.

$$\text{Total memory access w.r.t 1000 inst} = 1.4 \times 1000 = 1400$$

$$\text{Total No. of miss in L1} \Rightarrow 1400 \times 0.1 = 140$$

$$\frac{140}{1000} \times 7 = \frac{98}{100} = .98$$

$$\text{Total no. of sets} = \frac{256}{2} = 128 \text{ sets.}$$

Question Number : 51

Correct : 2 Wrong : 0

Consider a 2-way set associative cache with 256 blocks and uses LRU replacement. Initially the cache is empty. Conflict misses are those misses which occur due to contention of multiple blocks for the same cache set. Compulsory misses occur due to first time access to the block. The following sequence of accesses to memory blocks

(0, 128, 256, 128, 0, 128, 256, 128, 1, 129, 257, 129, 1, 129, 257, 129) } $\times 10$

is repeated 10 times. The number of *conflict misses* experienced by the cache is _____.

set0	(0, 128)	compulsory miss + (0, -)	Compulsory Miss (1, -)
set1	()	"	"
set2	()	"	"
	HIT	(256, 128)	(1, 129)
	Conflict Miss	(256, 128)	(257, 129)
	HIT	(0, 128)	(1, 129)
	Conflict Miss	(256, 128)	(1, 129)
→ HIT	(256, 128)	Conflict Miss (257, 129)	(257, 129)
		→ HIT	(257, 129)

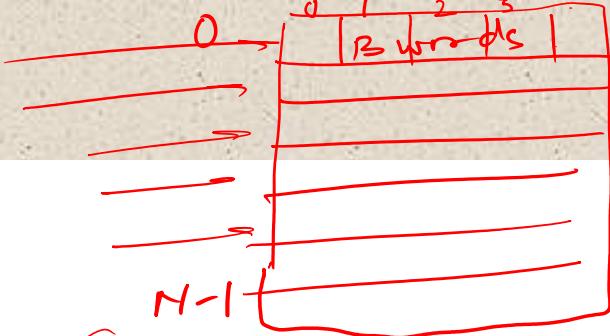
Set 0 (2nd iteration)
(2)

Set 1 (2nd iteration)



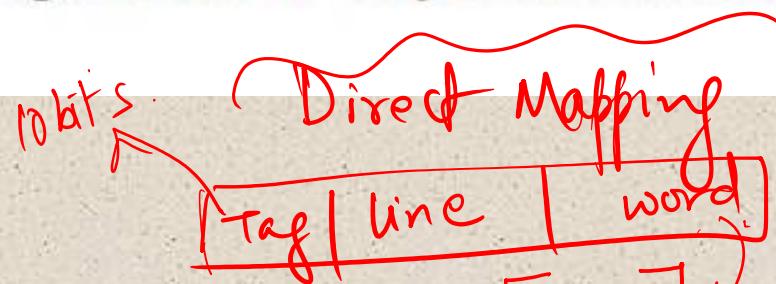


Question Number : 54



Correct : 2 Wrong : 0

A cache memory unit with capacity of N words and block size of B words is to be designed. If it is designed as a direct mapped cache, the length of the TAG field is 10 bits. If the cache unit is now designed as a 16-way set-associative cache, the length of the TAG field is _____ bits.



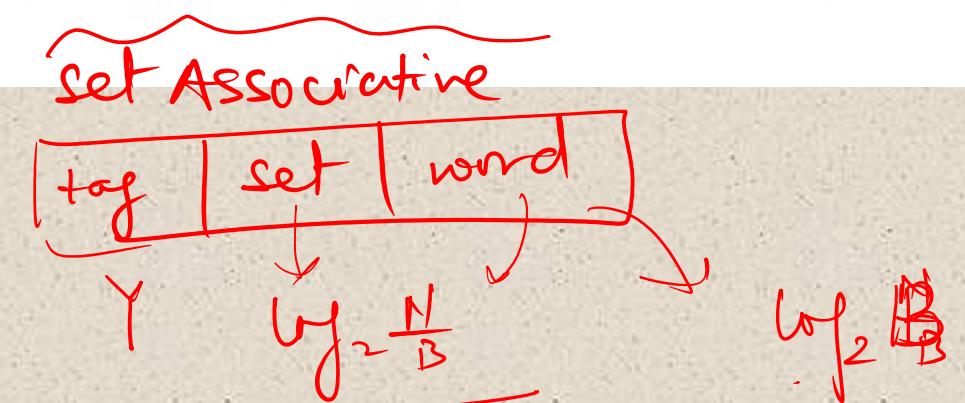
$$\text{offset for blocks} = \lceil \log_2 B \rceil$$

$$\text{No. of blocks in Cache} = \frac{N}{B}$$

$$\text{No. of bits required to represent the blocks} = \left\lceil \log_2 \left(\frac{N}{B} \right) \right\rceil / 16$$

$$10 + \cancel{\log \frac{N}{B}} + \cancel{\log B} = Y + \cancel{\frac{\log \frac{N}{B}}{16}} + \cancel{\frac{\log B}{16}}$$

$$10 = Y \neq 1, Y = 10 + 4 = 14 \text{ bits}$$



$$\frac{\log \frac{N}{B}}{16} - \log^{16}$$

$$\frac{\log \frac{N}{B}}{16} + \cancel{\log B}$$

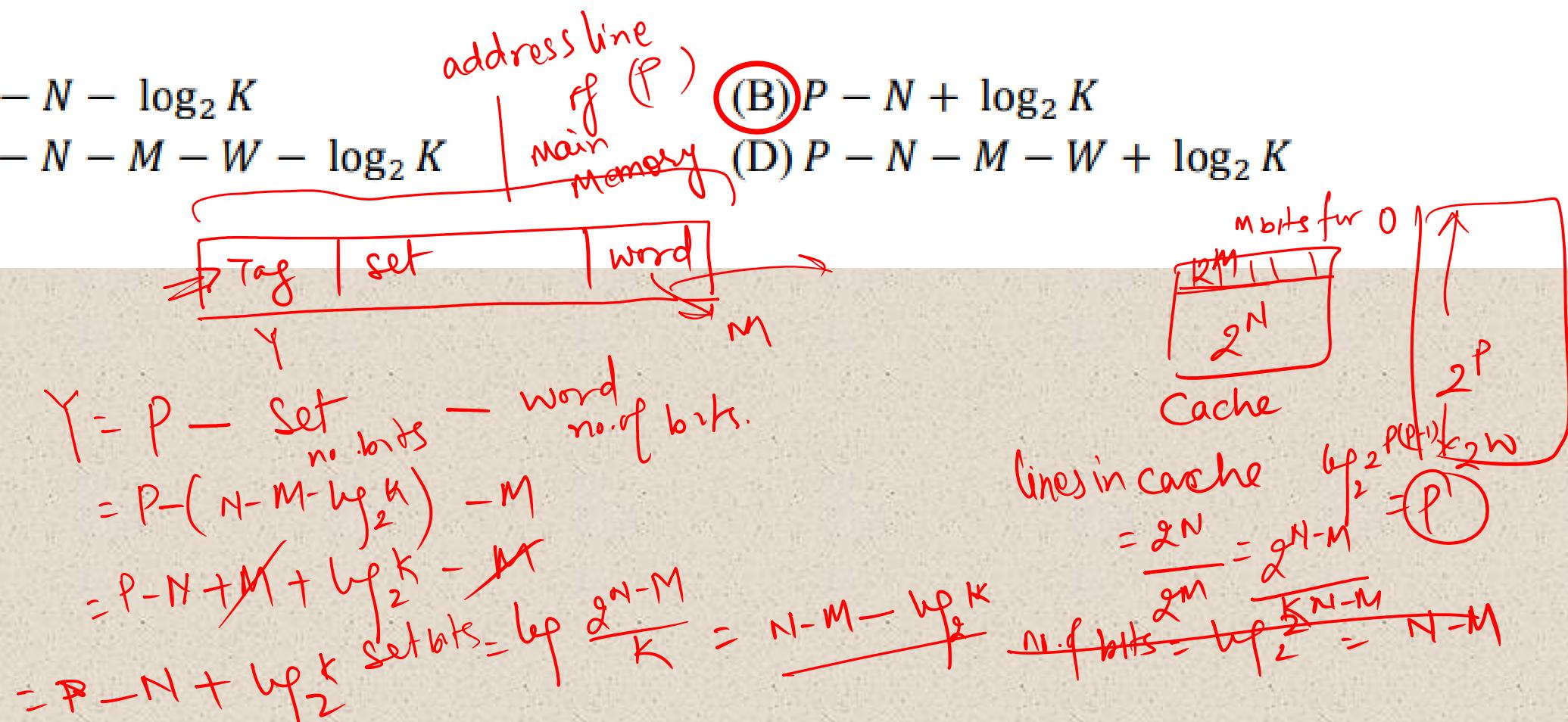
- Q.34 The size of the physical address space of a processor is 2^P bytes. The word length is 2^W bytes. The capacity of cache memory is 2^N bytes. The size of each cache block is 2^M words. For a K -way set-associative cache memory, the length (in number of bits) of the tag field is

(A) $P - N - \log_2 K$

(C) $P - N - M - W - \log_2 K$

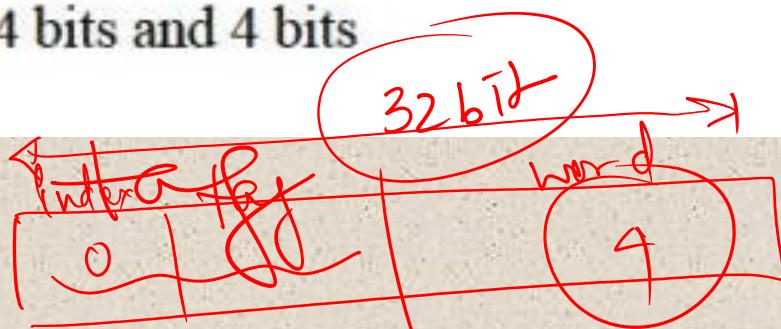
(B) $P - N + \log_2 K$

(D) $P - N - M - W + \log_2 K$



Q.1 A certain processor uses a fully associative cache of size 16 kB. The cache block size is 16 bytes. Assume that the main memory is byte addressable and uses a 32-bit address. How many bits are required for the *Tag* and the *Index* fields respectively in the addresses generated by the processor?

- (A) 24 bits and 0 bits
- (C) 24 bits and 4 bits



$$16 \text{ kB} = 2^4$$

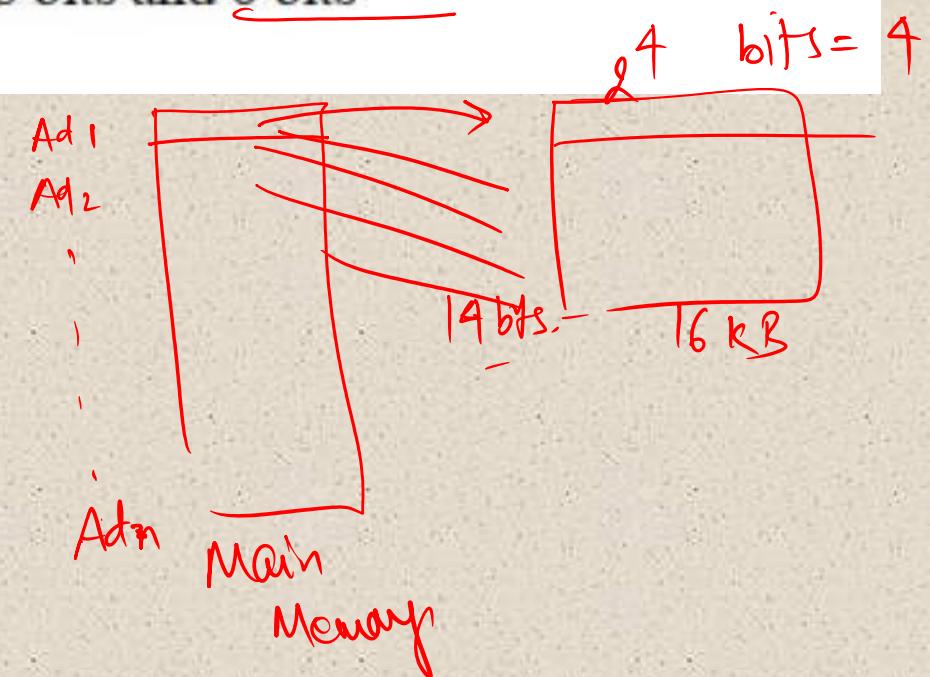
$$\lg(16 \text{ kB}) = 14 \text{ bits}$$

~~$32 - 4 = 28$~~

$$32 - 4 = 28 \text{ bits}$$

- (B) 28 bits and 4 bits
- (D) 28 bits and 0 bits

~~(D) 28 bits and 0 bits~~



- Q.45 A certain processor deploys a single-level cache. The cache block size is 8 words and the word size is 4 bytes. The memory system uses a 60-MHz clock. To service a cache miss, the memory controller first takes 1 cycle to accept the starting address of the block, it then takes 3 cycles to fetch all the eight words of the block, and finally transmits the words of the requested block at the rate of 1 word per cycle. The maximum bandwidth for the memory system when the program running on the processor issues a series of read operations is $\times 10^6$ bytes/sec.

Total clock cycle to fetch = $3 \times 1 = 12$ CC. CPU Cache 8 words $\times 4$ Main Memory

32 byte

Time taken by 12 CC = $12 \times 16 \text{ } \mu\text{s}$.

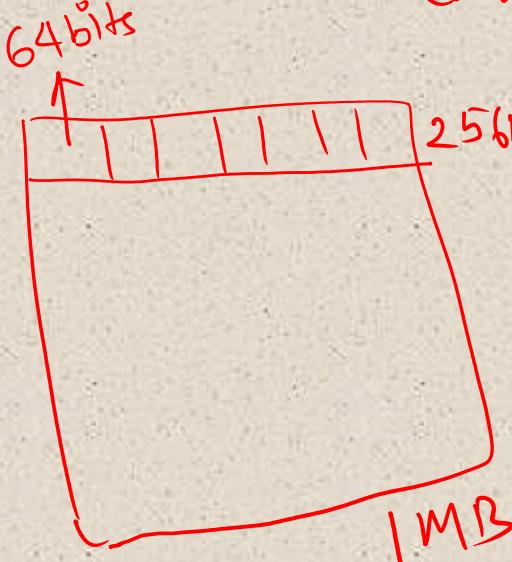
160

rate =

$$\frac{32 \times 16}{6 \times 12 \times 16 \text{ } \mu\text{s}} \approx \frac{0.16 \times 10^6}{160000} = 160000$$

$$\text{C.P.} = 60 \text{ MHz} \\ T = \frac{1}{f} = \frac{1}{60 \text{ MHz}} = 0.16 \times 10^{-6} = 16 \text{ } \mu\text{s.}$$

Q.No. 21 A direct mapped cache memory of 1 MB has a block size of 256 bytes. The cache has an access time of 3 ns and a hit rate of 94%. During a cache miss, it takes 20 ns to bring the first word of a block from the main memory, while each subsequent word takes 5 ns. The word size is 64 bits. The average memory access time in ns (round off to 1 decimal place) is _____.



Access

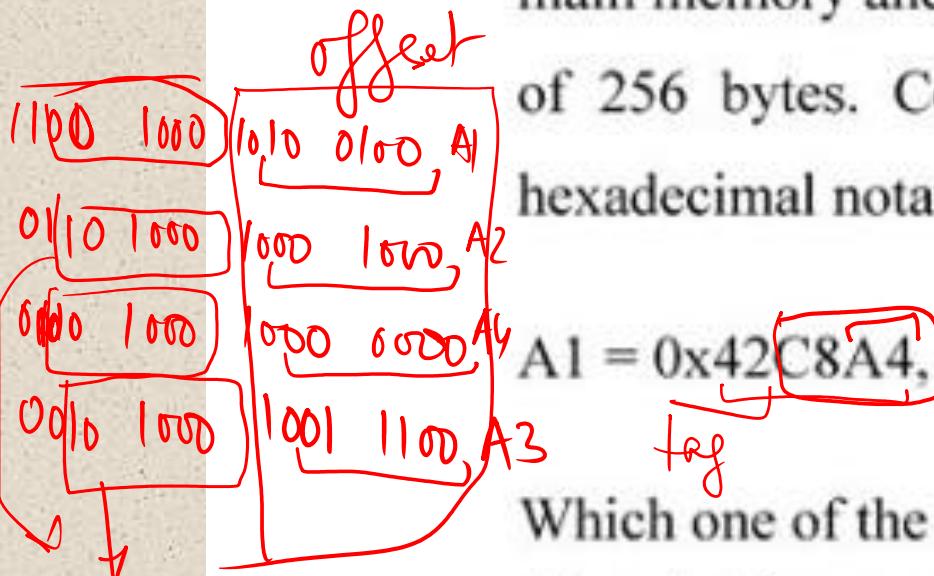
Cache Memory Time = Time taken during HIT + Time taken in Miss

$$94\% \times 3\text{ns} + 6\% \times \frac{(20\text{ns} + 3 \times 5) + 3}{\text{miss time}}$$

total no. of words/block = $\frac{256 \times 8}{64} = 32$

$$= .94 \times 3 + .06(20 + 155) + 3 = 13.5 \text{ ns.}$$

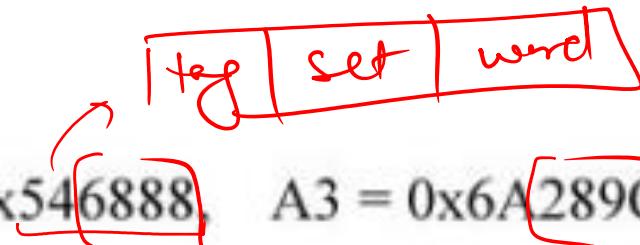
Q.No. 30 A computer system with a word length of 32 bits has a 16 MB byte-addressable main memory and a 64 KB, 4-way set associative cache memory with a block size of 256 bytes. Consider the following four physical addresses represented in hexadecimal notation.



Which one of the following is TRUE?

- (A) A1 and A4 are mapped to different cache sets.
- (B) A2 and A3 are mapped to the same cache set.
- (C) A3 and A4 are mapped to the same cache set.
- (D) A1 and A3 are mapped to the same cache set.

A₁, A₄
A₂, A₃



$$A1 = 0x42C8A4, \quad A2 = 0x546888, \quad A3 = 0x6A289C, \quad A4 = 0x5E4880$$

$$\text{lines} = \frac{64 \text{ KB}}{256 \text{ B}} = 64 \times 4$$

$$\text{Sets} = \frac{\text{lines}}{4} = 64 = 2^6$$

bits for set = 6 bits.

$$\text{offset} = \log_2 \text{Blocksize} = \log_2 256 = 8$$