

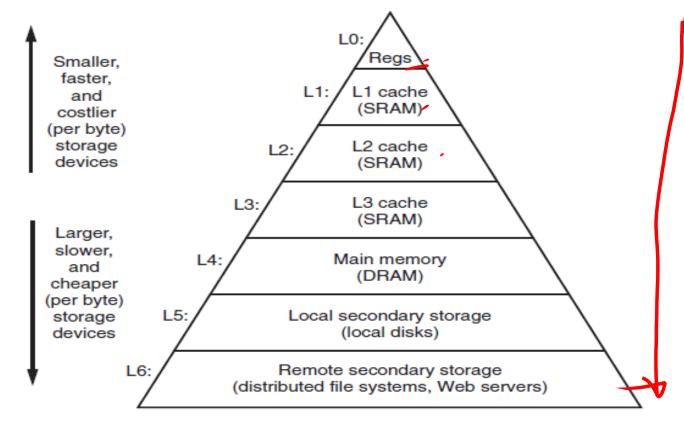
Computer Organization and Software Systems

Contact Session 4

Dr. Lucy J. Gudino

The Bottom Line

- How much?
 - Capacity
- / · How fast?
 - · Time is money
- å How expensive?



An example of a memory hierarchy.

- •Faster access time, ----- cost per bit
- •Greater capacity, ----access time



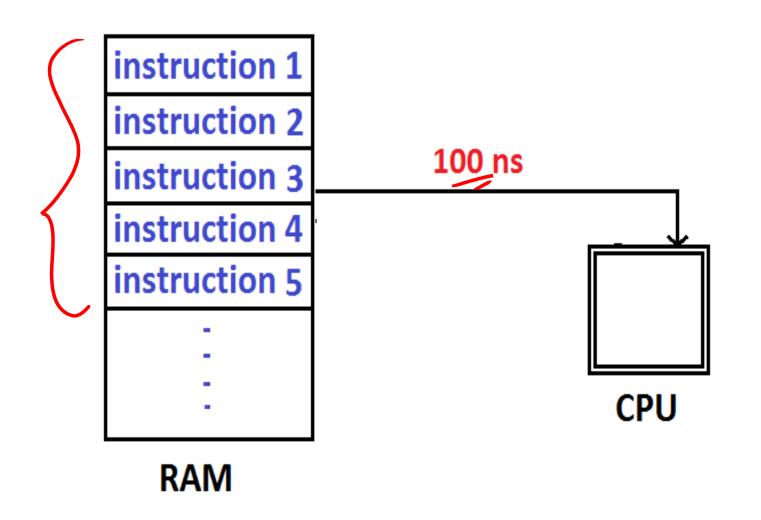
Memory Hierarchy

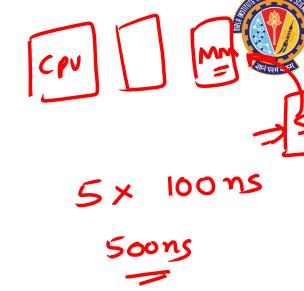


- Registers
 - In CPU
- Internal or Main memory
 - May include one or more levels of cache
 - "RAM"
- External memory
 - Backing store



Performance enhancement - Motivation







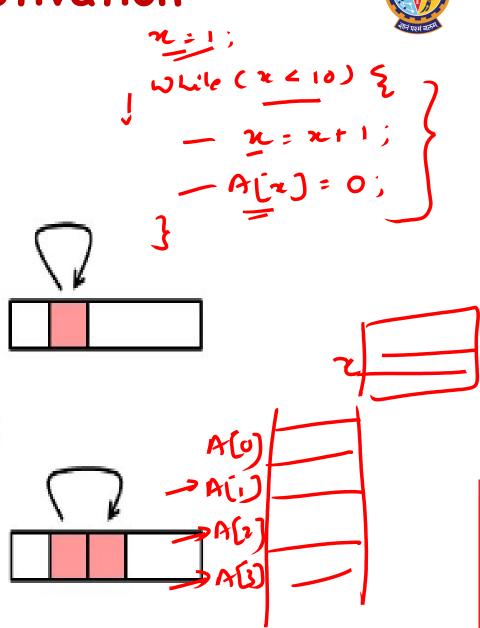
Performance enhancement - Motivation



Locality of Reference

During the course of the execution of a program, memory references tend to cluster

- Temporal locality: Locality in time
 - If an item is referenced, it will tend to be referenced again soon
- Spatial locality: Locality in space
 - If an item is referenced, items whose addresses are close by will tend to be referenced soon.





Example

```
product = 1;
for (i = 0; i < n-1; i++)
    product = product * a[i];</pre>
```

temporal: i,

product - temporal
i - temporal
n - temporal

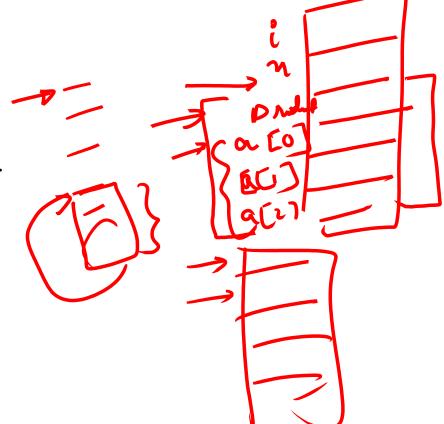


•Data:

- · Access array elements in succession spatial locality
- Reference to "product" "i" and "n" in each iteration -Temporal locality

•Instructions:

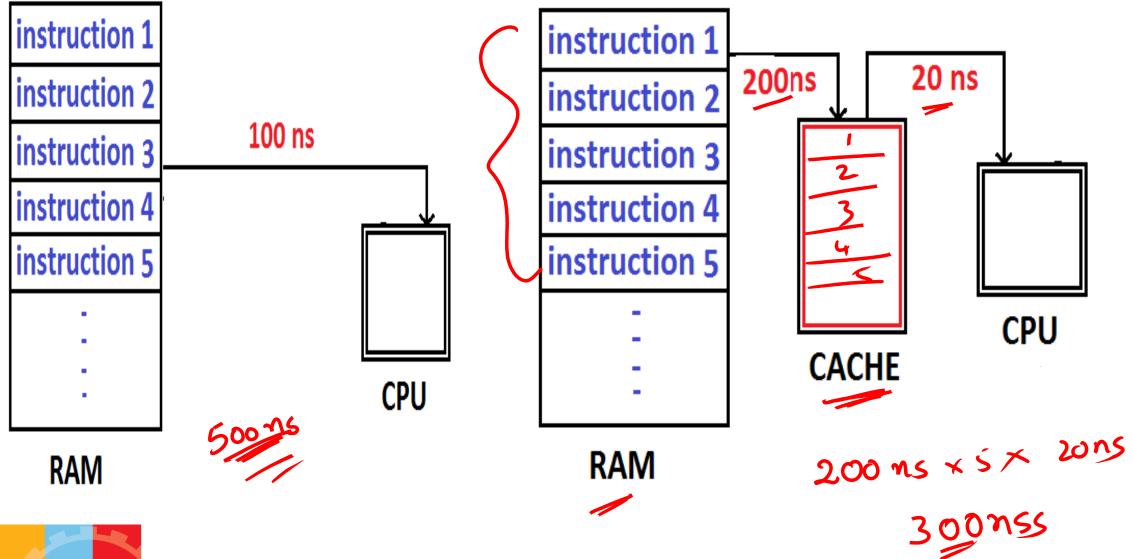
- Reference instructions in sequence: Spatial locality
- Looping through: Temporal locality





Performance enhancement - Motivation







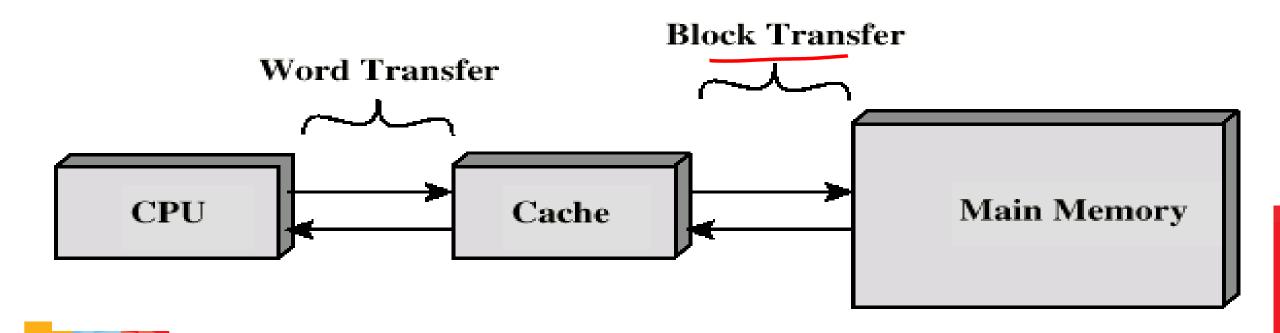


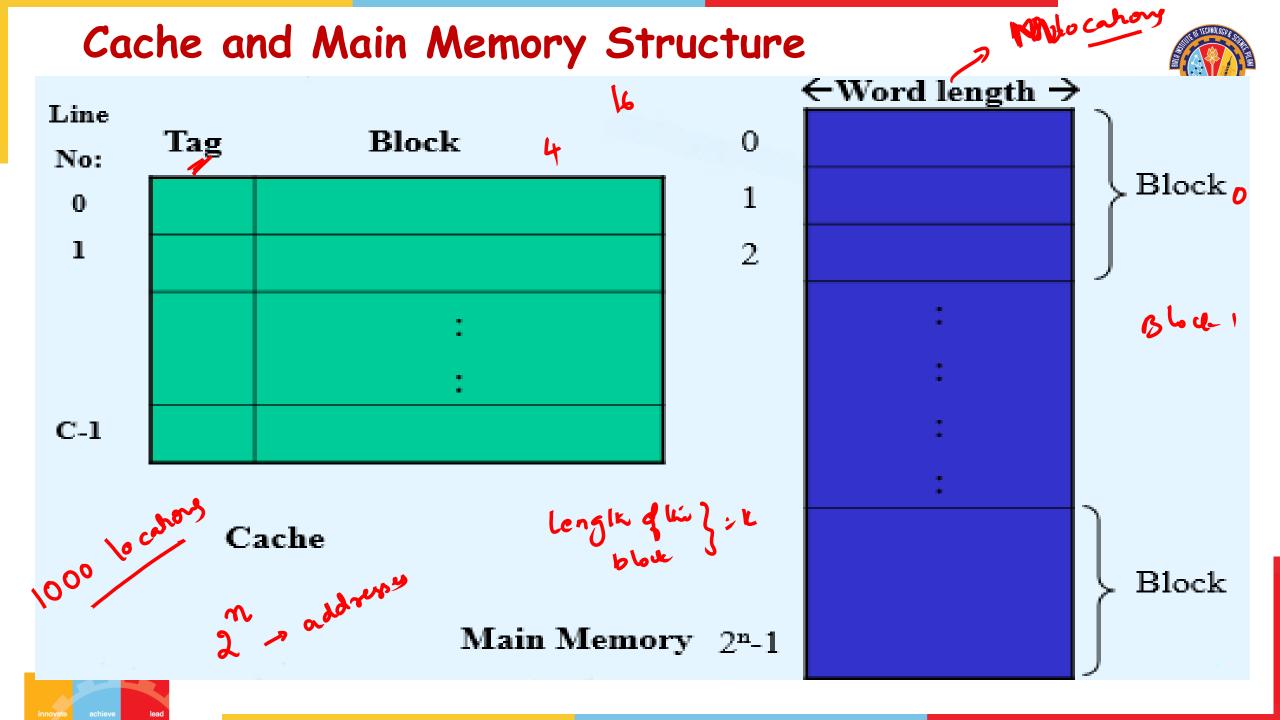
Cache

Cache



- Small, fast memory
- Sits between normal main memory and CPU
- · May be located on CPU chip or separate module





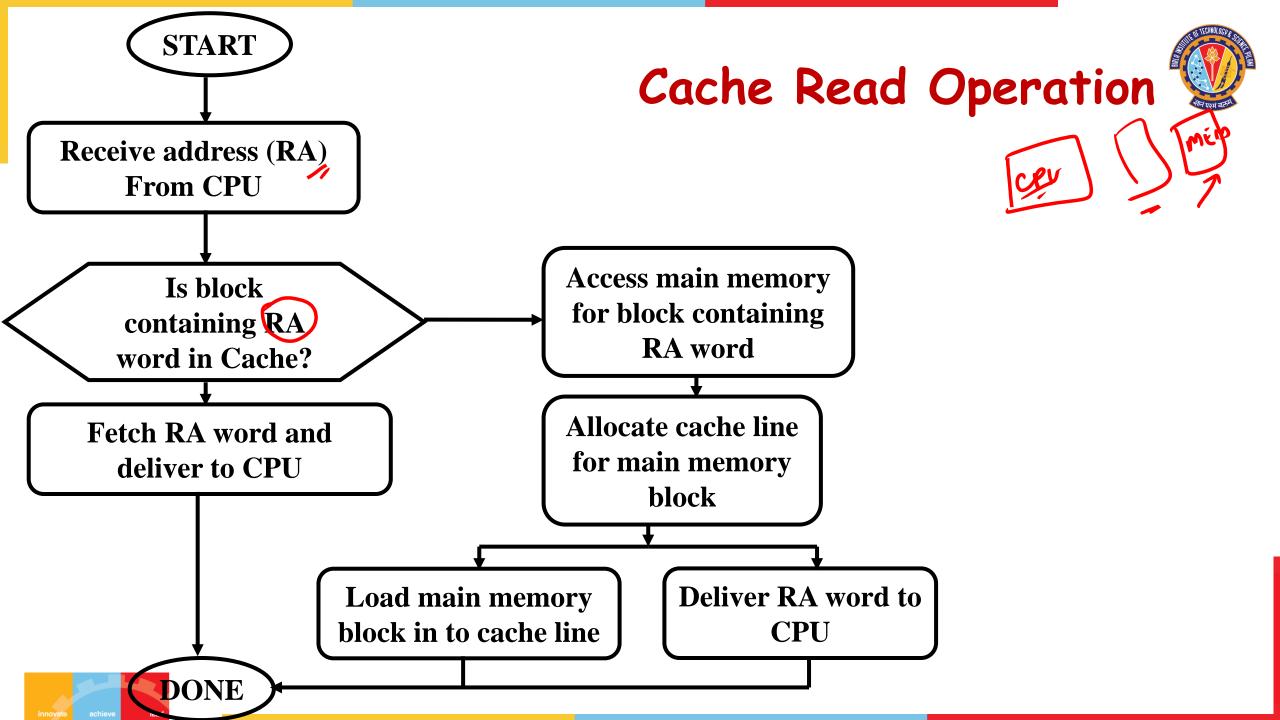
$$2^{\circ}$$
: 1
 2° : 2
 2° : 4
 2° : 8
 2° : 8
 2° : 32
 2° : 64
 2° : 128
 2° : 512
 2° : 612
 2° : 1024 = 112

1000

$$2^{11} = 2^{1} \cdot 2^{10} = 2 \cdot 10$$
 $2^{12} = 2^{2} \cdot 2^{10} = 4 \cdot 10$
 $2^{13} = 9 \cdot 10$
 $2^{14} = 16 \cdot 10$
 $2^{15} = 32 \cdot 10$
 $2^{11} = 128 \cdot 10$
 $2^{17} = 128 \cdot 10$
 $2^{19} = 256 \cdot 10$
 $2^{19} = 512 \cdot 10$
 $2^{19} = 1024 \cdot 10$



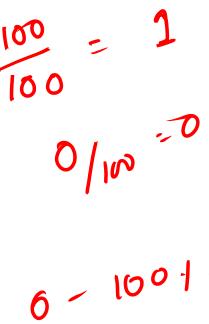
230: 1GB 240: 1TB



Performance of cache

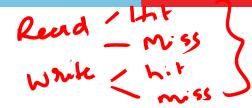


- Hit ratio: Number of Hits / total references
 to memory
- Hit
- Miss

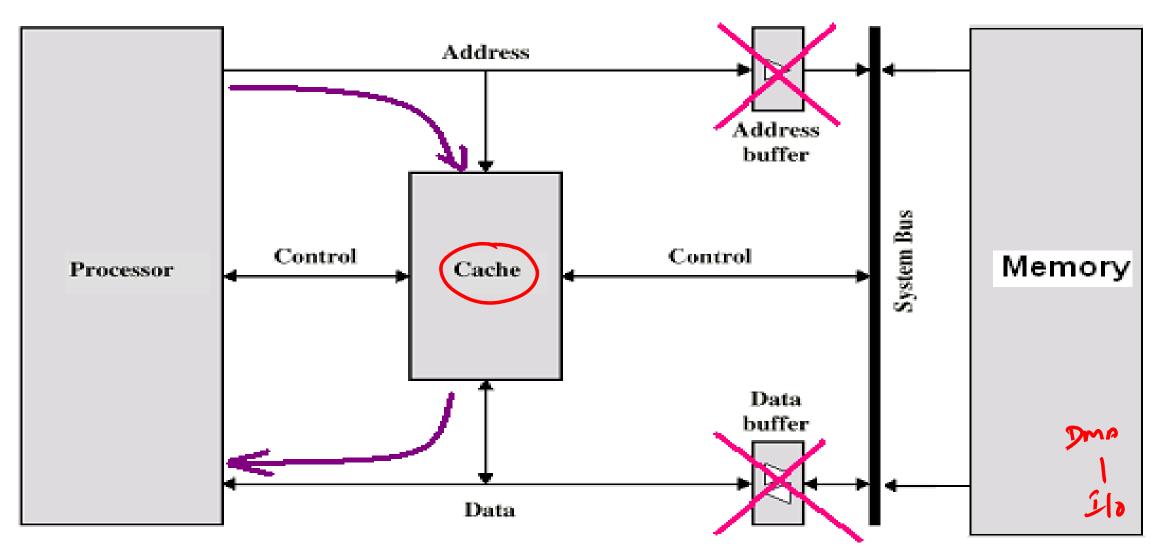




Read Hit

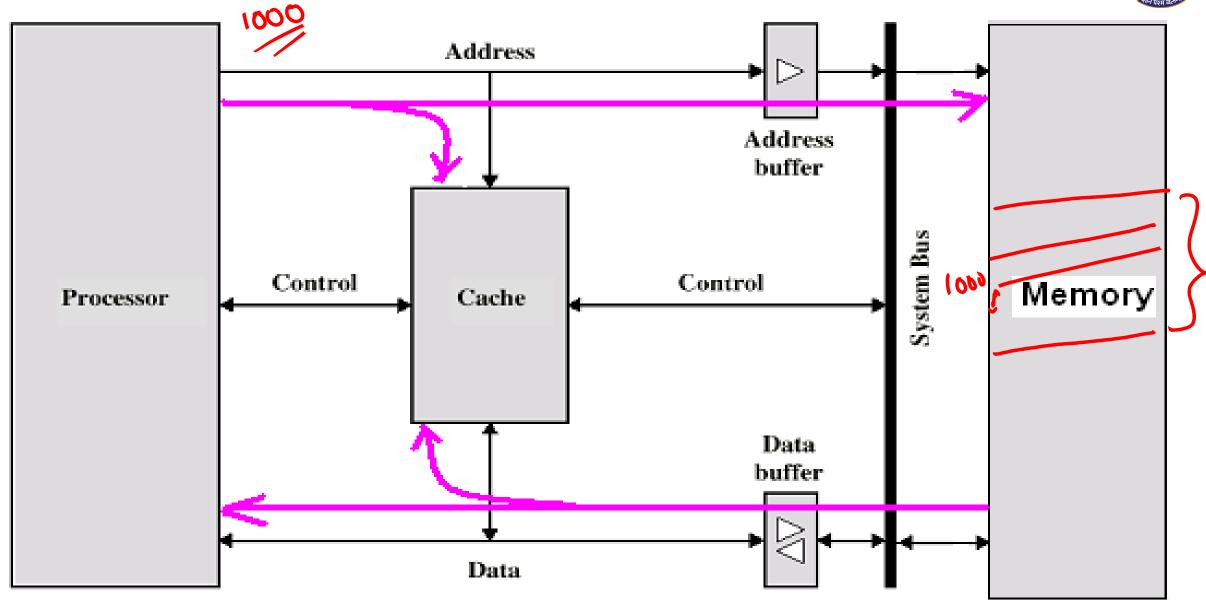


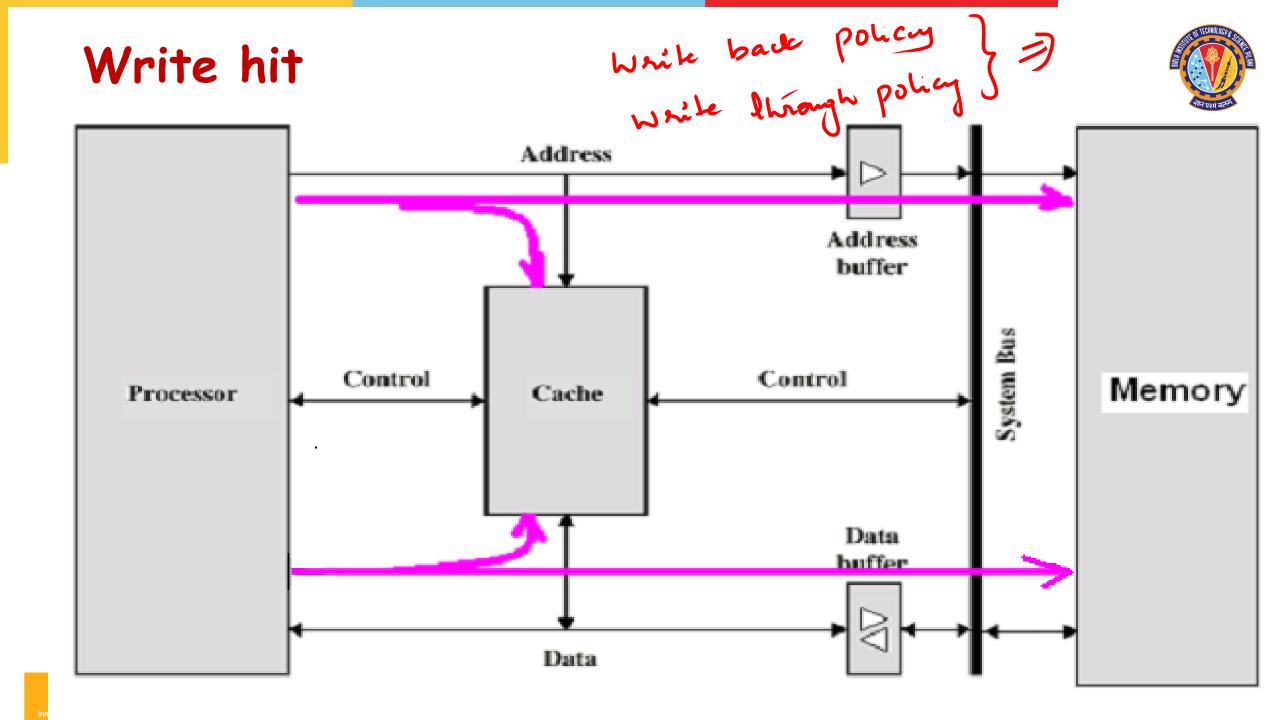




Read Miss

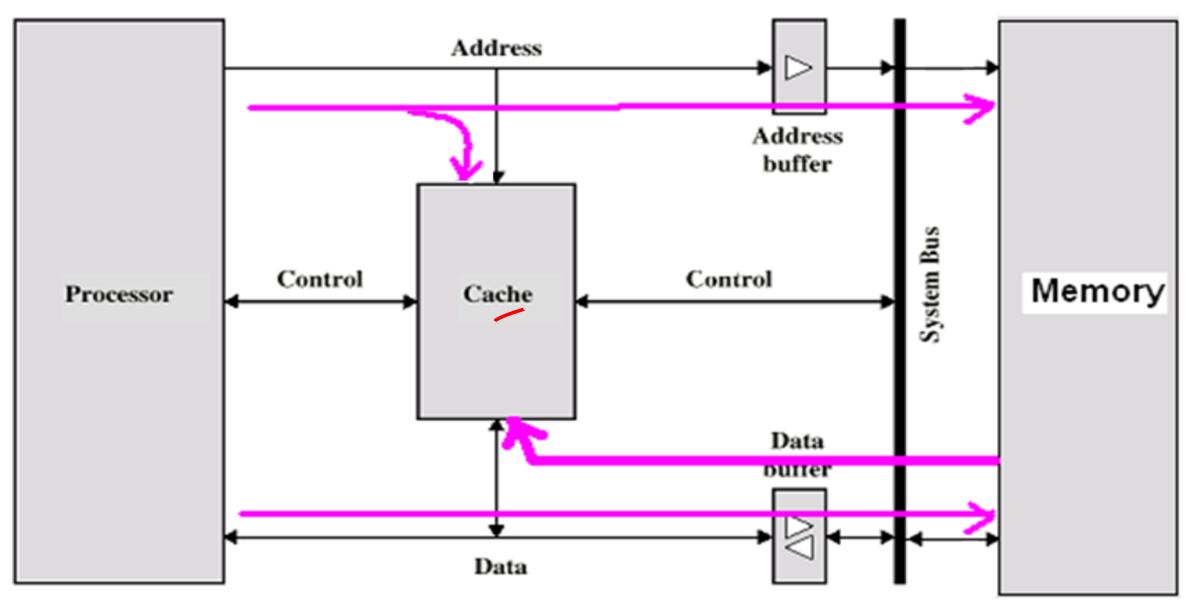






Write miss





Mapping Function



- · How memory blocks are mapped to cache lines
- Three types
 - -Direct mapping
 - -Associative mapping
 - -Set Associative mapping /

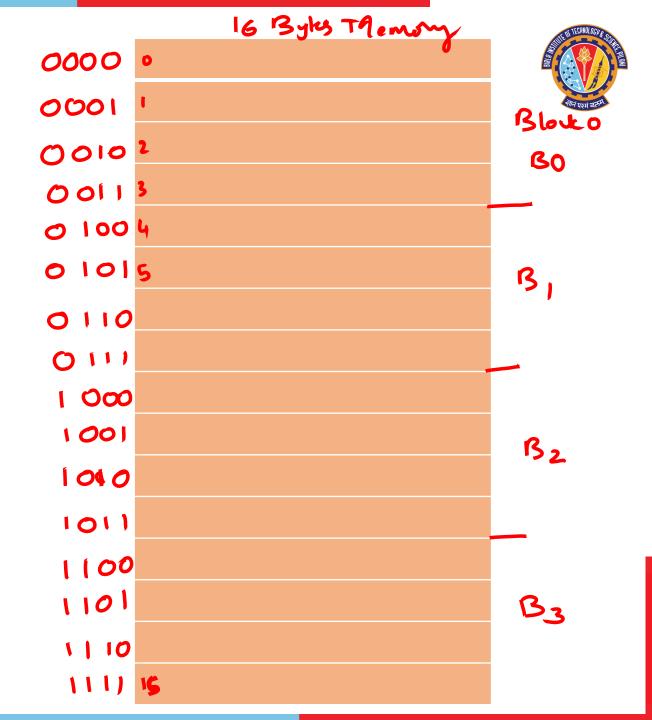


Direct Mapped Cache

- 16 Bytes main memory => # addm & 4 616
 - How many address bits are required?
- · Memory block size is 4 bytes
- Cache of 8 Byte

lo

- How many cache lines? >> 2 2
- cache contains 2 lines (4 bytes per Line)



Direct Mapped Cache



- · Each block of main memory maps to only one cache line
 - · i.e. if a block is in cache, it must be in one specific place
 - $i = j \mod u \log m$

```
where i = cache line number
j = main memory block no.
m = no.of lines in the
cache
```



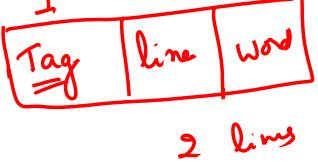
Direct Mapped Cache

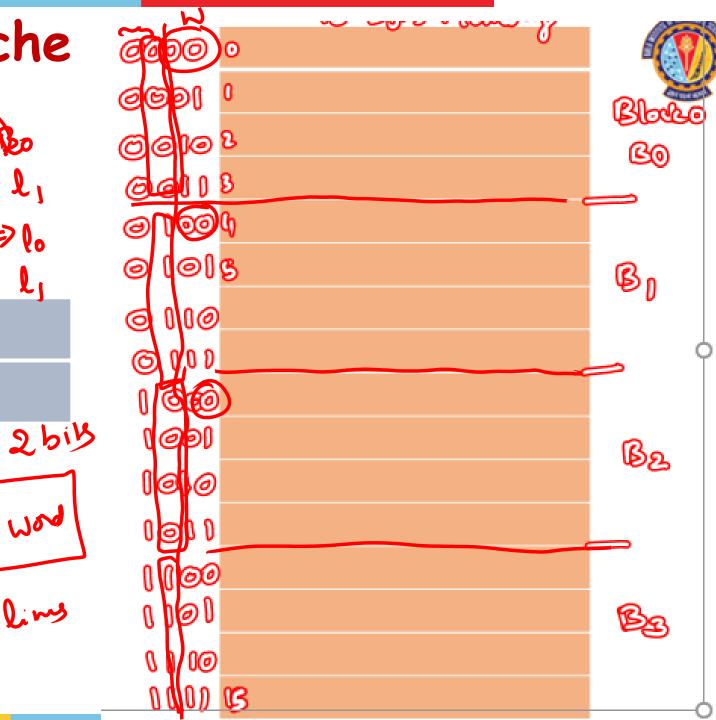
i = j modulo m $(a,bo): 0.7.2.90 \Rightarrow (3.0.3) & 6.0.0000$ $b_1 1.7.2.91 \Rightarrow (3.1.9) & 6.1.91$

b2 2/200 => B2 => lo b2 3/201 => B3 = l1

h3 3 1.2 今1 今 人の 3 Bo KB2 人の 3 Bo KB3

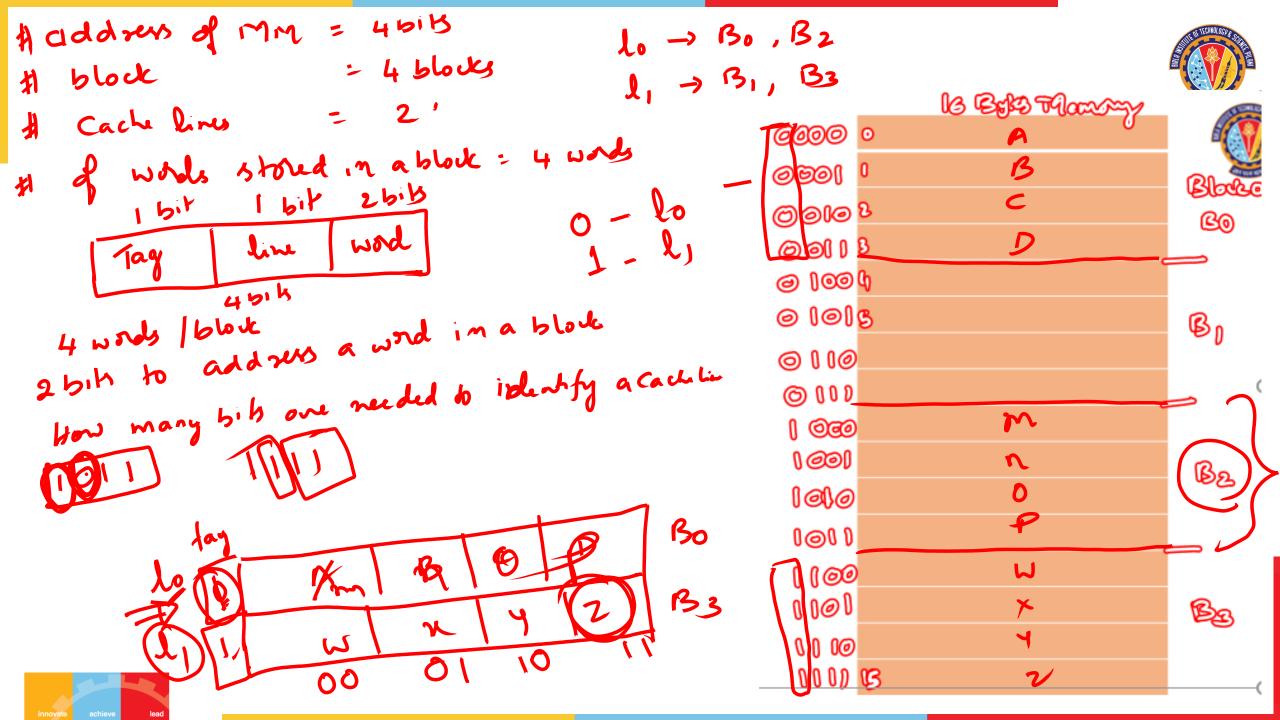
- Address is split in three parts:
 - Tag
 - Line
 - Word





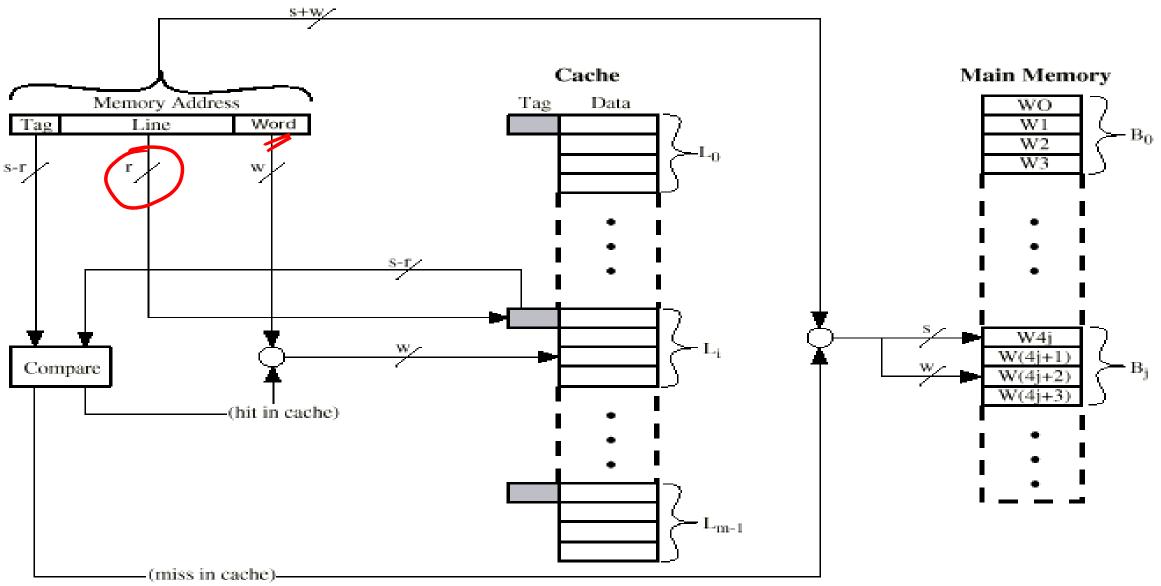






Direct Mapped Cache Organization





Direct mapped cache- 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





Direct mapped cache- Summary



1	1	2

Tag s-r

Line or Slot r

Word w



Direct mapped cache-pros & cons



- Simple
- Inexpensive
- Fixed location for given block
 - If a program accesses 2 blocks that map to the same line repeatedly, cache misses are very high



Problem 1: Direct Mapped Cache



•Given:

- Cache of 64KByte, Cache block of 4 bytes
- 16MBytes main memory
- Find out
- a) Number of bits required to address the main memory
- b) Number of blocks in main memory
- c) Number of cache lines
- d) Number of bits required to identify a word (byte) in a block
- e) Number of bits to identify a block
- f) Tag, Line, Word



Cache block 5% = MM block 51%
Slot





- ·Given:
 - · Cache of 64kByte, Cache block of 4 bytes
 - 16MBytes main memory
- Find out
- a) Number of bits required to address the main memory = 24 515

Capacity of the main memory:
$$16178$$

$$2^{n} = 1671$$

$$420$$

•b) Number of blocks in main memory

$$\frac{16M}{4} = 4M block$$

•c) Number of cache lines



block six = 4 bytes



•Given:

- Cache of 64kByte, Cache block of 4 bytes
- 16MBytes main memory
- Find out
- = 2 bits •d) Number of bits required to identify a word (byte) in a block?
- •e) Number of bits required to identify a block = 24 bith 25 ih :=
- Tag, Line, Word

Tag s-r Line r

A cache hous = 16 10 line

Word w

8



Consider a machine with a byte addressable main memory of 2¹⁶ bytes and block size of 8 bytes. Assume that a direct mapped cache consisting of 32 lines is used with this machine.

- a. How is a 16-bit memory address divided into tag, line number, and byte number?
- b. Into what line would bytes with each of the following addresses be stored?

0001	0001	0001	1011
1100	0011	0011	0100
1101	0000	0001	1101
1010	1010	1010	1010

- c. Suppose the byte with address 0001 1010 0001 1010 is stored in the cache. What are the addresses of the other bytes stored along with it?
- d. How many total bytes of memory can be stored in the cache?
- e. Why is the tag also stored in the cache?





Consider a machine with a byte addressable main memory of 2¹⁶ bytes and block size of 8 bytes. Assume that a direct mapped cache consisting of 32 lines is used with this machine.

a. How is a 16-bit memory address divided into tag, line number, and byte number?



b. Into what line would bytes with each of the following addresses be stored?

	ag_	المسللا	. his d	16.8421	. 0
0001	0001	0001	1011	00011	3
1100	0011	0011	0100	00110	うしょうしょ
1101	0000	0001	1101	00011	=> 13
1010	1010	1010	1010	. 10101	=> 121
			1		



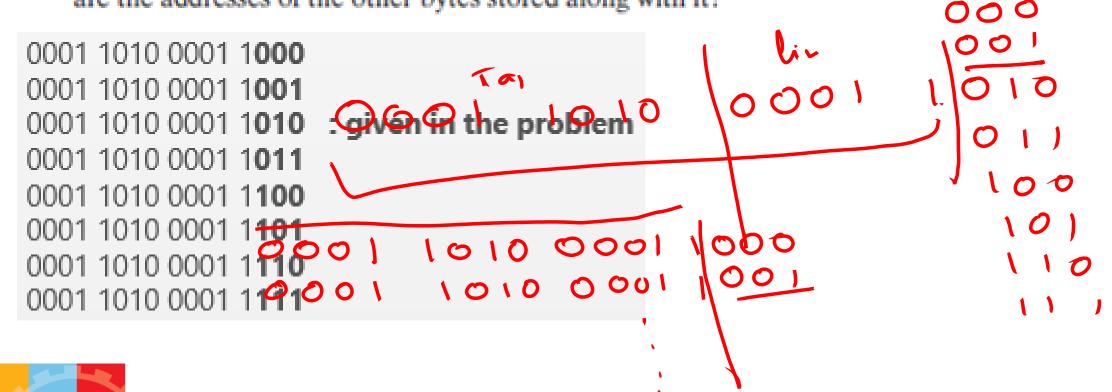
13





Consider a machine with a byte addressable main memory of 2¹⁶ bytes and block size of 8 bytes. Assume that a direct mapped cache consisting of 32 lines is used with this machine.

c. Suppose the byte with address 0001 1010 0001 1010 is stored in the cache. What are the addresses of the other bytes stored along with it?





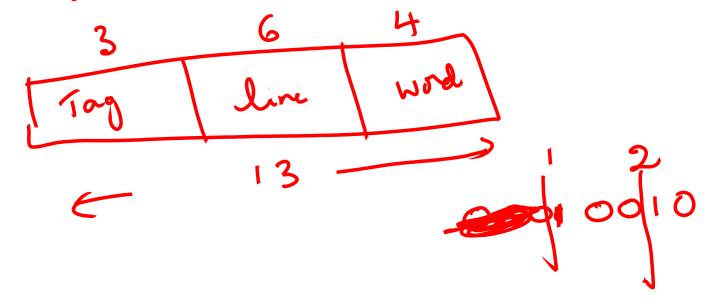
Consider a machine with a byte addressable main memory of 2¹⁶ bytes and block size of 8 bytes. Assume that a direct mapped cache consisting of 32 lines is used with this machine.

d. How many total bytes of memory can be stored in the cache?

e. Why is the tag also stored in the cache?



 Consider a direct-mapped cache with 64 cache lines and a block size of 16 bytes and main memory of 8K (Byte addressable memory). To what line number does byte address 1200H map?



21		The My Technology &
plexadecima	Binary	_Decimal_
0	0000	रुमं बर्लर
1 2 4	0001	1
2 3 7	0010	2
	6 Y ₀₁₁	3
42	0100	4
5 2 2	0101	5
6	0110	6
7	0111	7
8 2 ~ >	1000	10 8
9	1001	9
A	2010	10
B 2	1011	11
CO	1100	12
D	1101	10
500C	OOO	1432
F	1111	15
4	32	1842
	10	5000





- The system uses a L1 cache with direct mapping and 32-bit address format is as follows:
- bits 0 3 = offset (word)
- bits 4 14 = index bits (Line)
- bits 15 31 = tag
- a) What is the size of cache line?
- b) How many Cache lines are there?
- c) How much space is required to store the tags in the L1 cache?
- d) What is the total Capacity of cache including tag storage?

Tag s-r	Line r	Word w
•		





- 16 Bytes main memory, Memory block size is 4 bytes, Cache of 8 Byte (cache is 2 lines of 4 bytes each)
- Block access sequence:

•0 2 0 2 2 0 0 2 0 0 0 2 1

Find out hit ratio.



Problem 5 - Direct Mapped Cache



 16 Bytes main memory, Memory block size is 4 bytes, Cache of 8 Byte (cache is 2 lines of 4 bytes each)

Block access sequence:

0202200200021

Find out hit ratio.

0 2 0 2 2 0 0 2 0 0 0 2 1





- Suppose a 1024-byte cache has an access time of 0.1 microseconds and the main memory stores
 1 Mbytes with an access time of 1 microsecond. A referenced memory block that is not in cache must be loaded into cache.
- Answer the following questions:
- a) What is the number of bits needed to address the main memory?

a) If the cache hit ratio is 95%, what is the average access time for a memory reference?

Avg access time = hit ratio * cache access + (1- hit ratio) * (cache access + memory access)



Associative Mapping

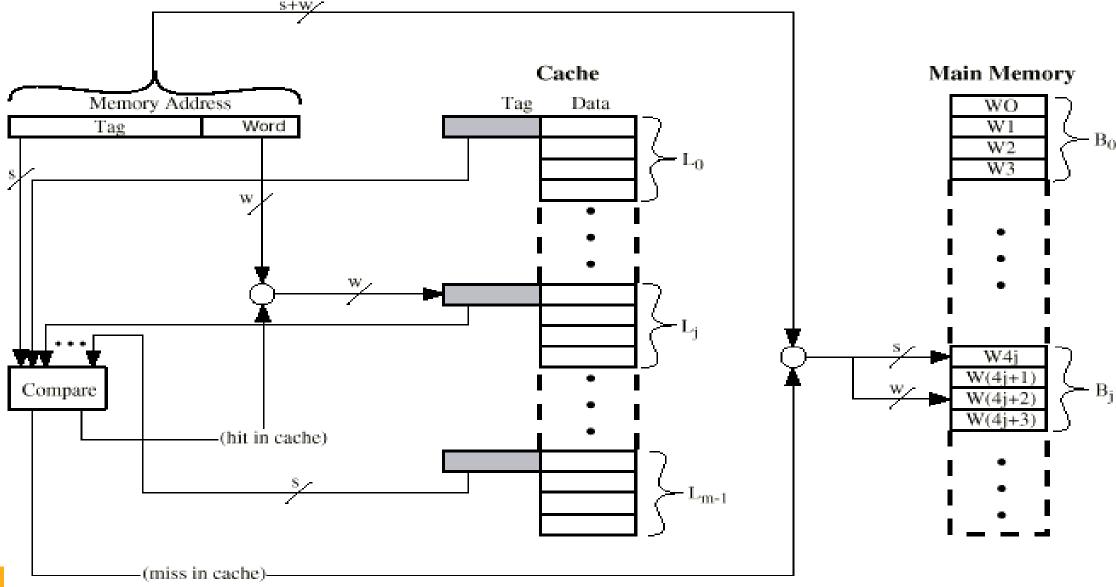


- · A main memory block can load into any line of cache
- Memory address is interpreted as tag and word
- Tag uniquely identifies block of memory
- · Every line's tag is examined for a match
- Cache searching gets expensive



Associative Cache Organization





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



•Given:

- Cache of 128KByte, Cache block of 8 bytes
- 32 MBytes main memory
- Find out
- a) Number of bits required to address the memory
- b) Number of blocks in main memory
- c) Number of cache lines
- d) Number of bits required to identify a word (byte) in a block?
- e) Tag, Word



•Cache of 64KByte, Cache block of 4 bytes, 16 M Bytes main memory and associative mapping.

Fill in the blanks:

Number of bits in main memory address = _____

Number of lines in the cache memory = _____

Word bits = _____

Tag bits = _____





16 Bytes main memory, Memory block size is 4 bytes, Cache of 8
Byte (cache is 2 lines of 4 bytes each) and associative mapping.
Block access sequence:

0 2 0 2 2 0 0 2 0 0 0 2 1

Find out hit ratio.

0 2 0 2 2 0 0 2 0 0 0 2 1

