



Computer Organization and Software Systems

CONTACT SESSION 2.....



BITS Pilani
Pilani Campus

Prof. C R Sarma
WILP.BITS-PILANI

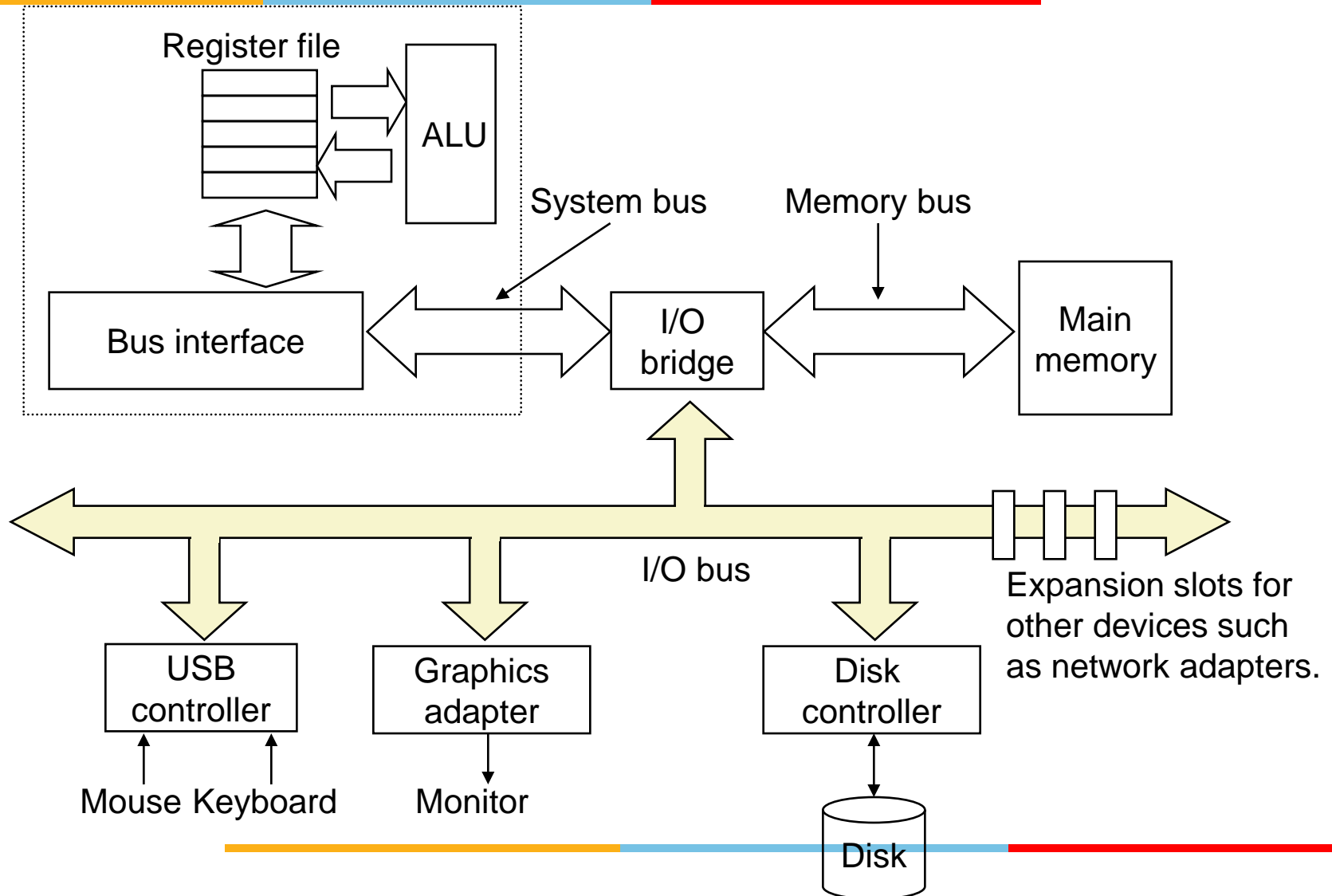
Logical Disk Blocks

- Modern disks present a simpler abstract view of the complex sector geometry:
 - The set of available sectors is modeled as a sequence of b-sized **logical blocks** (0, 1, 2, ...)
- Mapping between logical blocks and actual (physical) sectors
 - Maintained by hardware/firmware device called disk controller.
 - Converts requests for logical blocks into (surface, track, sector) triples.
- Allows controller to set aside spare cylinders for each zone.
 - Accounts for the difference in "formatted capacity" and "maximum capacity".

I/O Bus

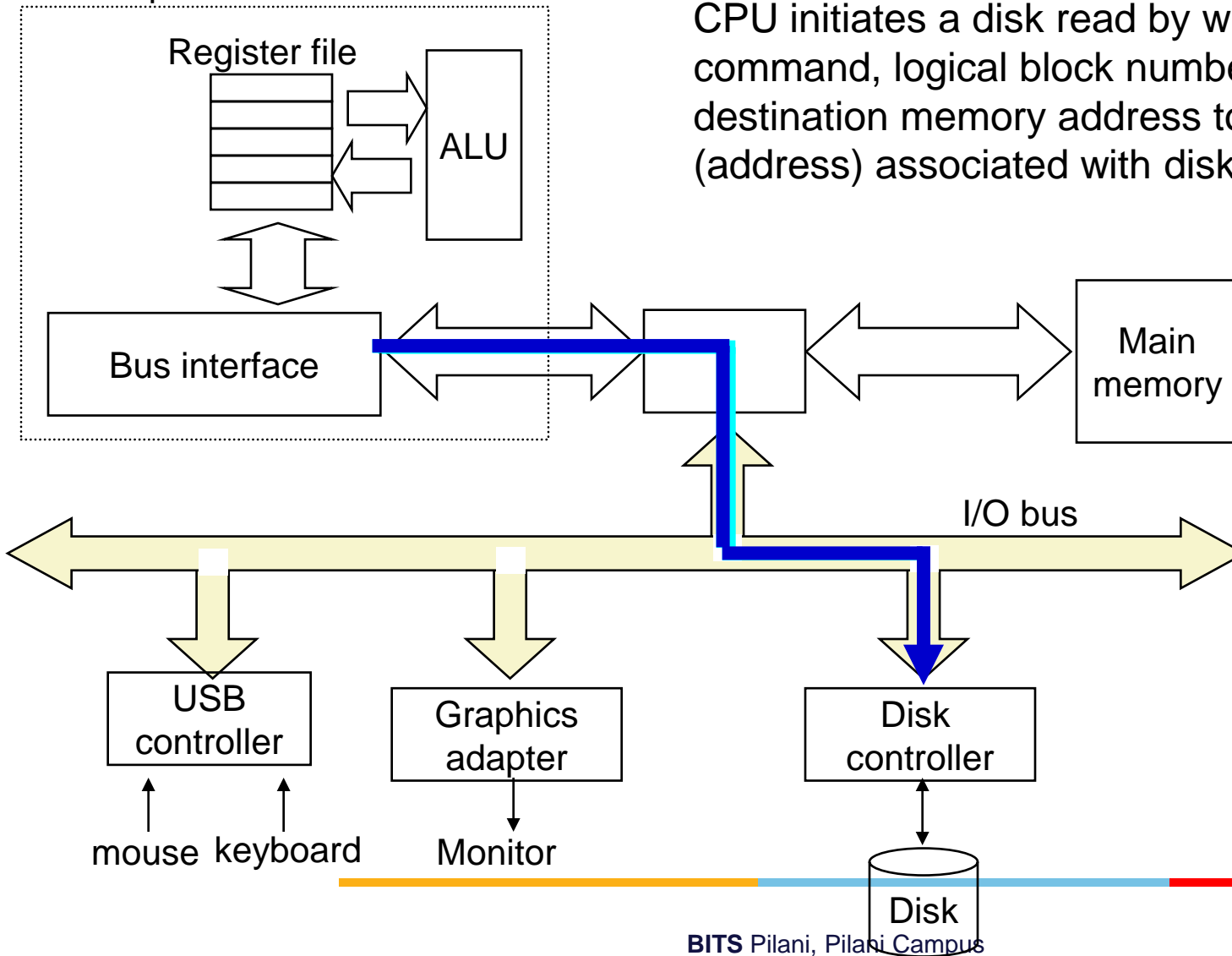


CPU chip



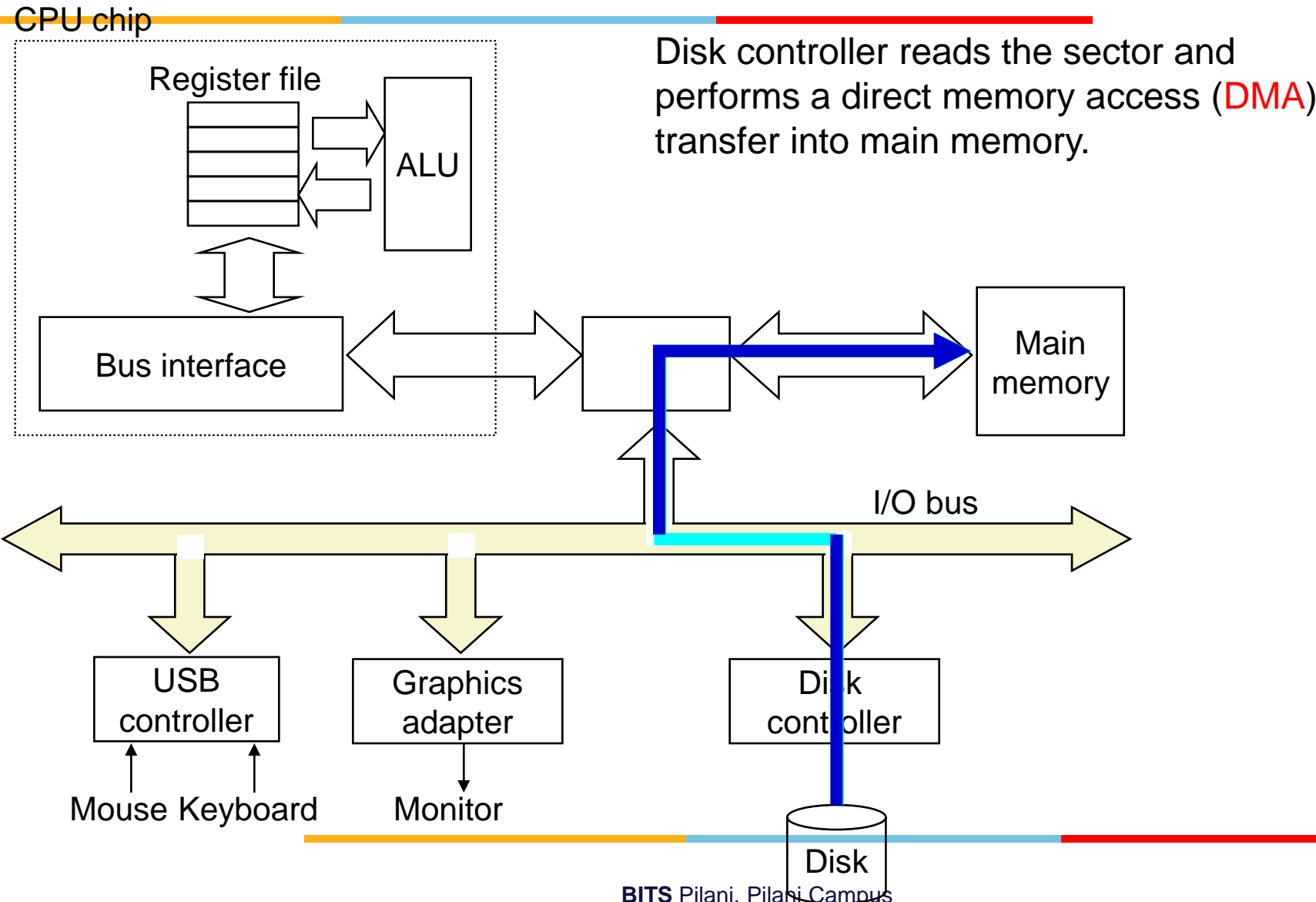
Reading a Disk Sector (1)

CPU chip

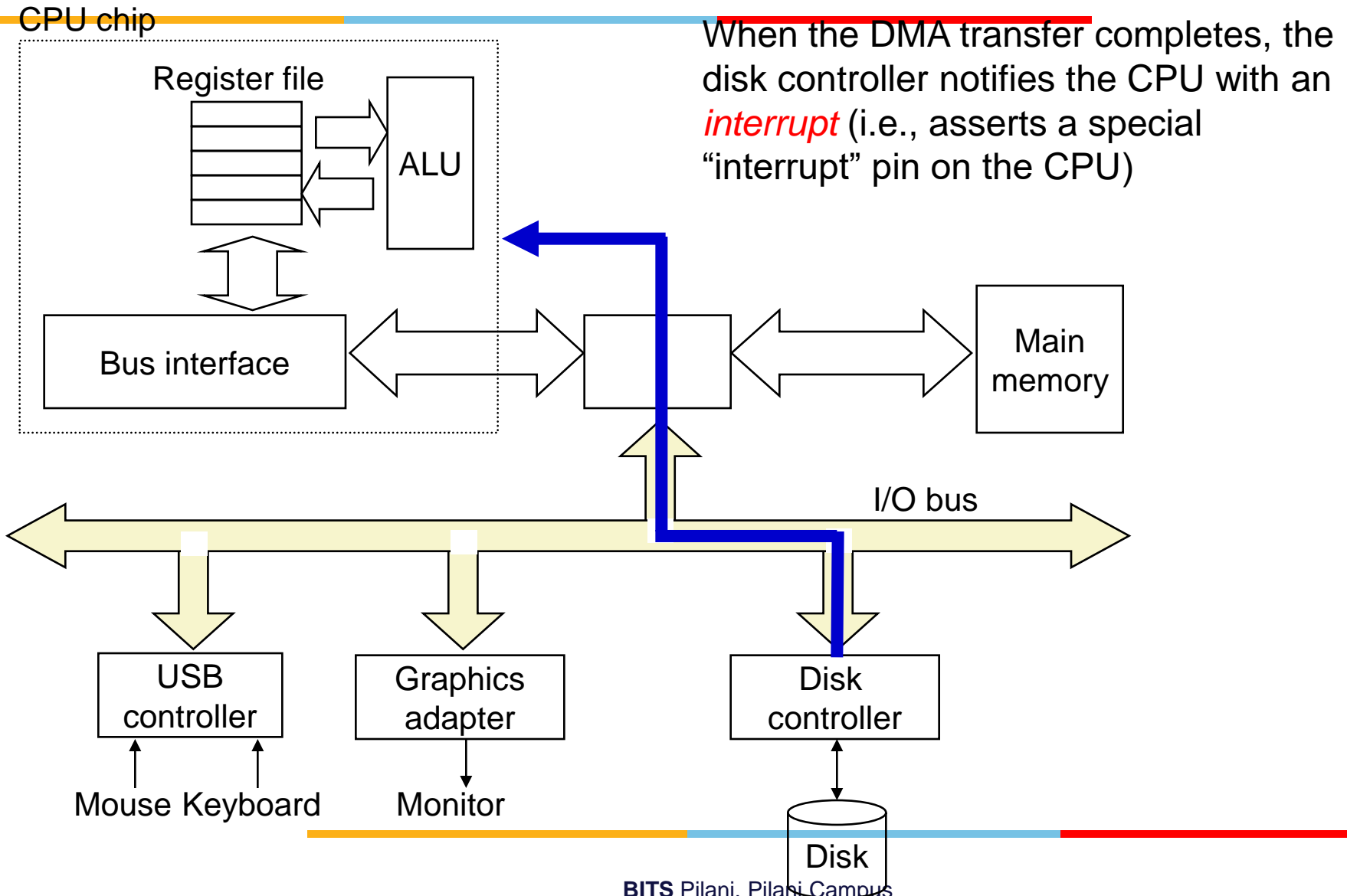


CPU initiates a disk read by writing a command, logical block number, and destination memory address to a **port** (address) associated with disk controller.

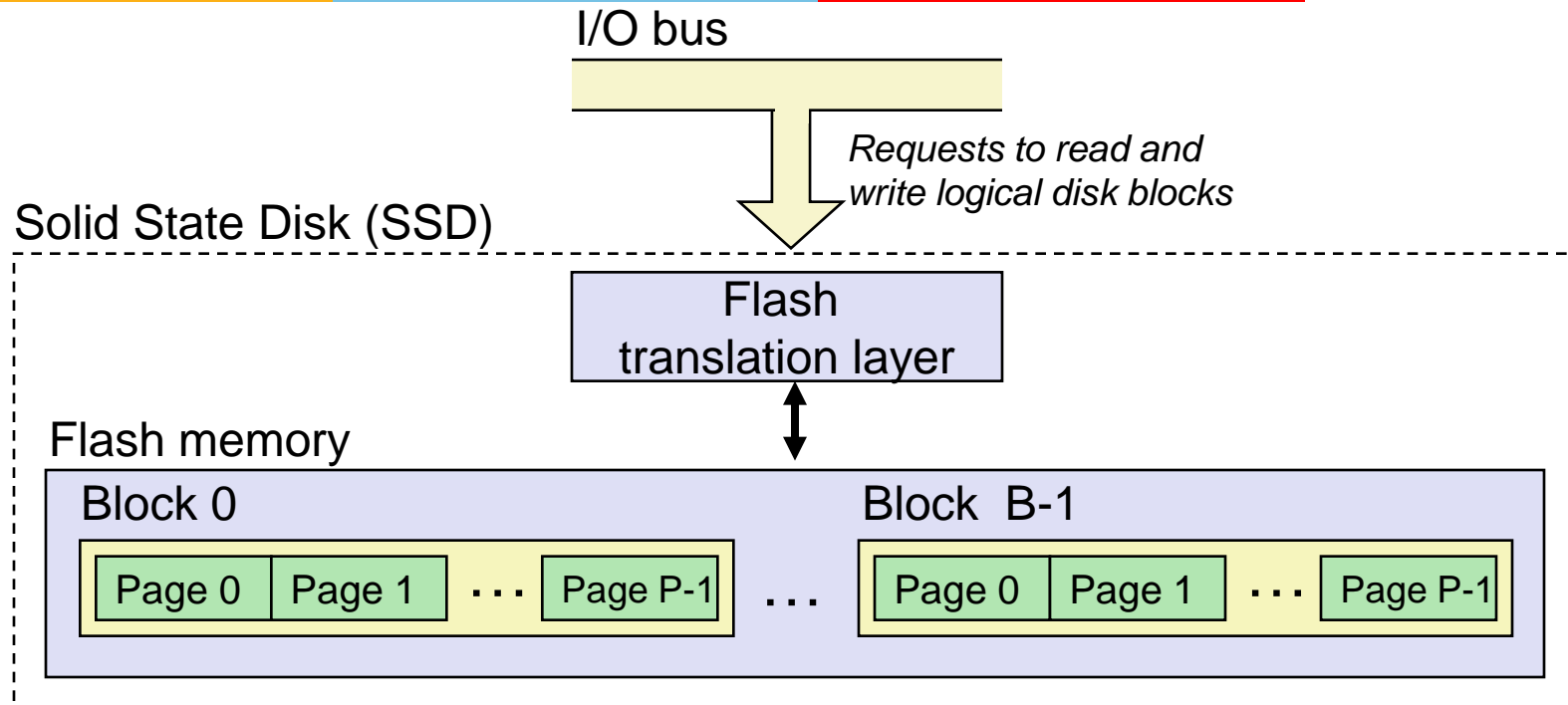
Reading a Disk Sector (2)



Reading a Disk Sector (3)



Solid State Disks (SSDs)



- Pages: 512KB to 4KB, Blocks: 32 to 128 pages
- Data read/written in units of pages.
- Page can be written only after its block has been erased
- A block wears out after about 100,000 repeated writes.



SSD Performance Characteristics

Sequential read input	550 MB/s	Sequential write output	470 MB/s
Random read input	365 MB/s	Random write output	303 MB/s
Avg seq read time	50 us	Avg seq write time	60 us

- Sequential access faster than random access
 - Common theme in the memory hierarchy
- Random writes are somewhat slower
 - Erasing a block takes a long time
 - Modifying a block page requires all other pages to be copied to new block
 - In earlier SSDs, the read/write gap was much larger.

Source: Intel SSD 730 product specification.

SSD Tradeoffs vs Rotating Disks

- **Advantages**
 - No moving parts → faster, less power, more rugged
- **Disadvantages**
 - Have the potential to wear out
 - Mitigated by “wear leveling logic” in flash translation layer
 - E.g. Intel SSD 730 guarantees 128 petabyte (128×10^{15} bytes) of writes before they wear out
 - In 2015, about 30 times more expensive per byte
- **Applications**
 - MP3 players, smart phones, laptops
 - Beginning to appear in desktops and servers



Computer Organization and Software Systems

CONTACT SESSION 3

Prof. C R Sarma
WILP.BITS-PILANI



BITS Pilani
Pilani Campus

Today's Session



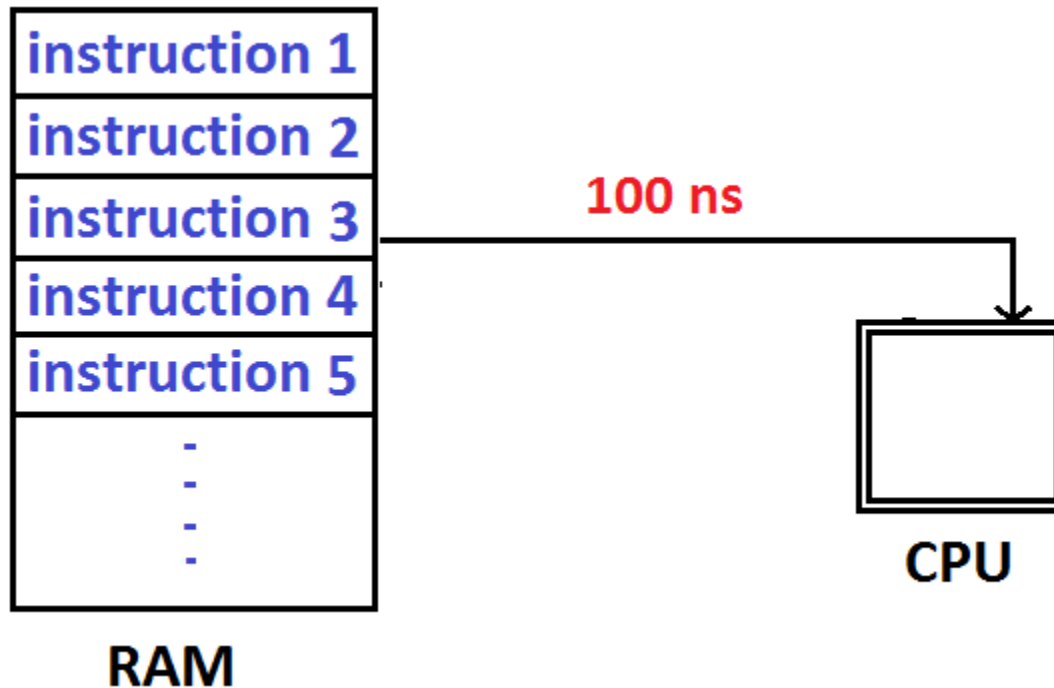
Contact Hour	List of Topic Title	Text/Ref Book/external resource
5-6	<ul style="list-style-type: none">• Memory Hierarchy• Locality<ul style="list-style-type: none">• Locality of Reference to Program Data• Locality of instruction fetches• Cache Memories<ul style="list-style-type: none">• Generic Cache Memory Organization• Direct-Mapped Caches• Fully Associative Caches	T1

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



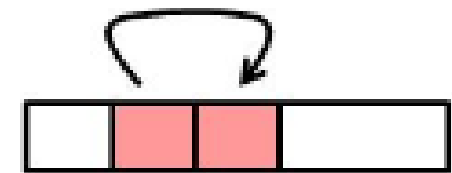
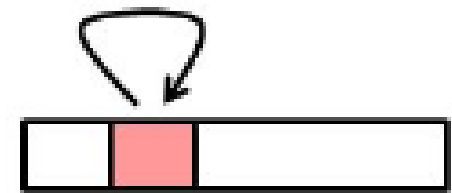
Performance enhancement - Motivation



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] ;
```

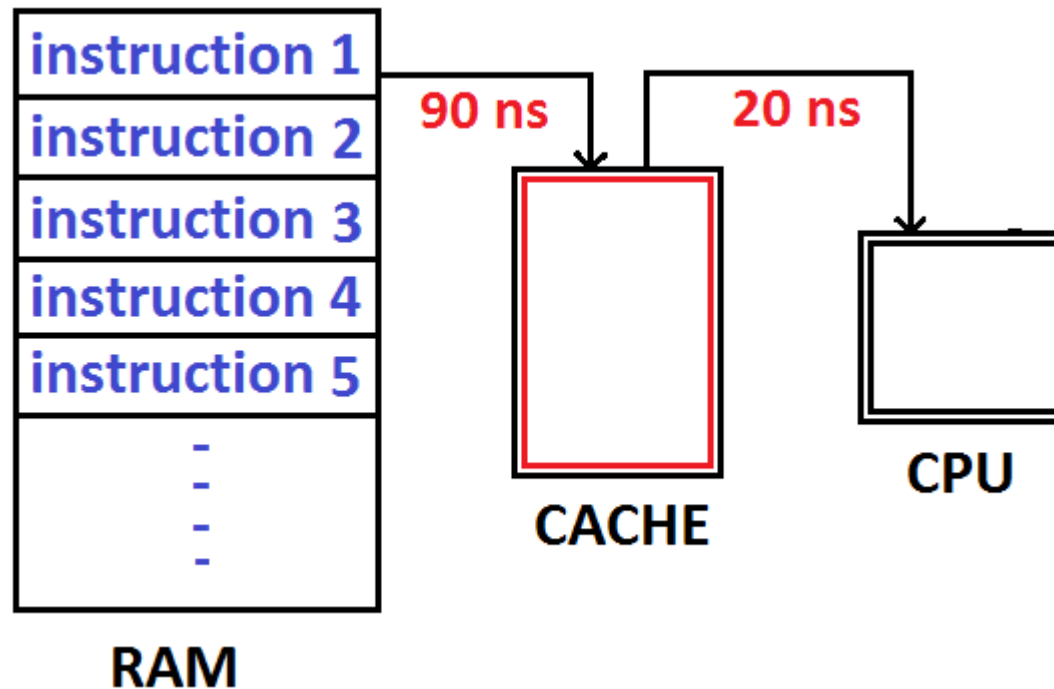
Data :

- Access array elements in succession - spatial locality
- Reference to "product" in each iteration - Temporal locality

Instructions :

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

Performance enhancement - Motivation





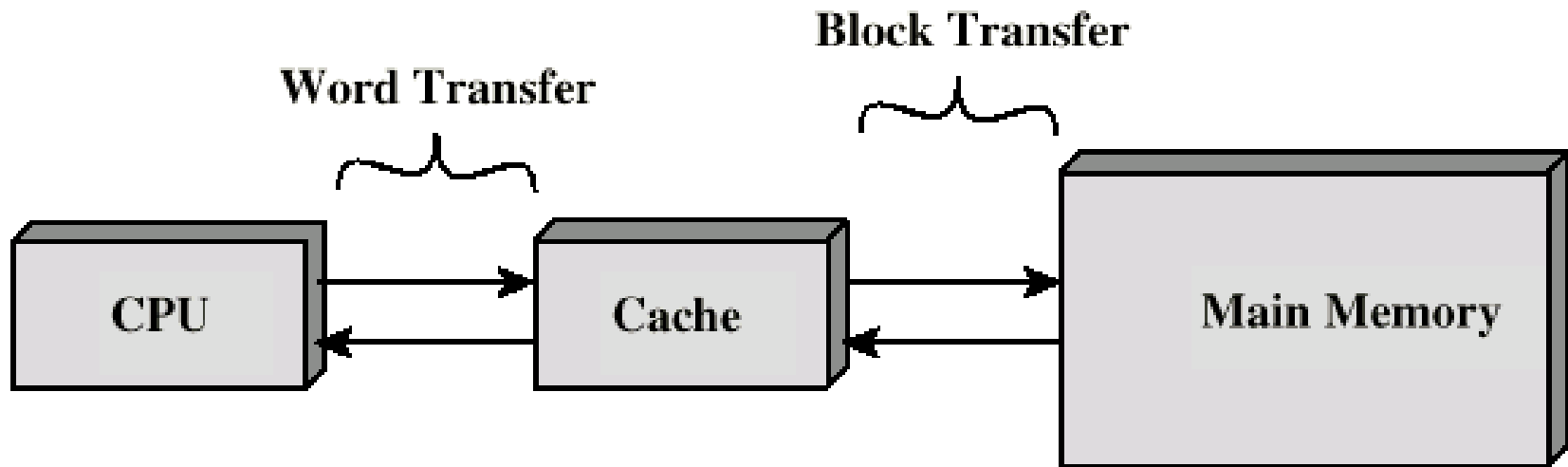
BITS Pilani
Pilani Campus

Cache

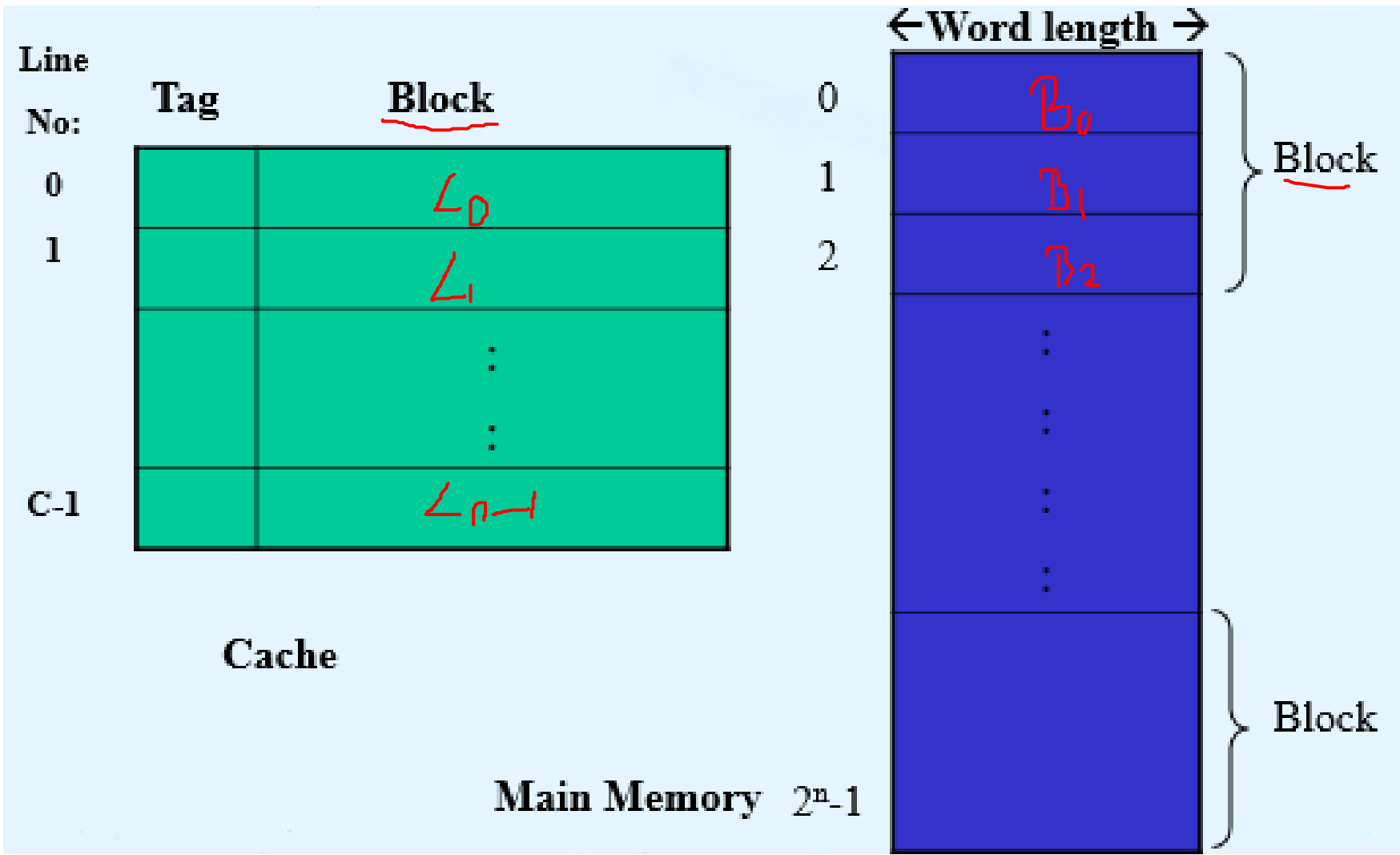
Cache



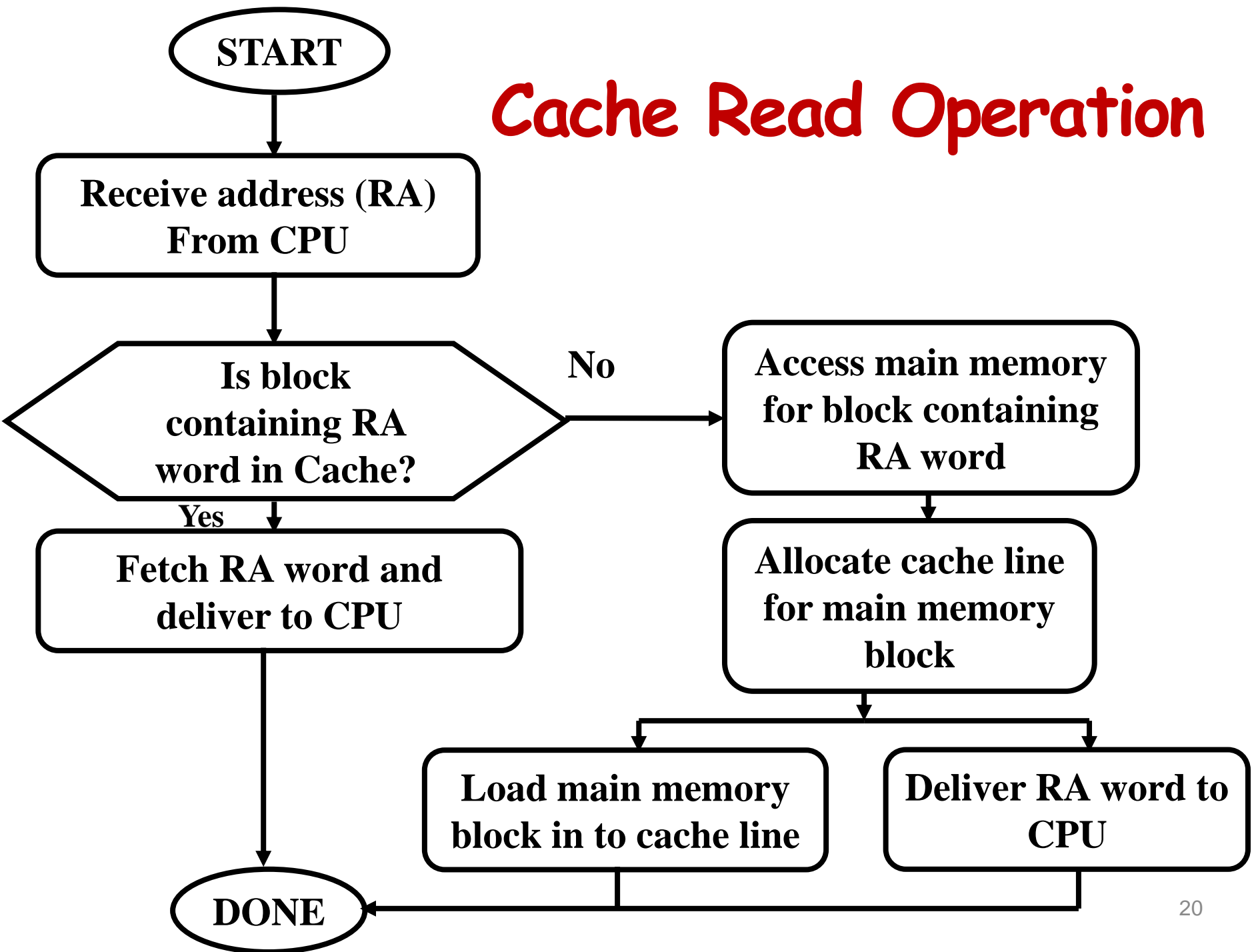
- Small, fast memory
- Sits between normal main memory and CPU
- May be located on CPU chip or separate module



Cache and Main Memory Structure



Cache Read Operation

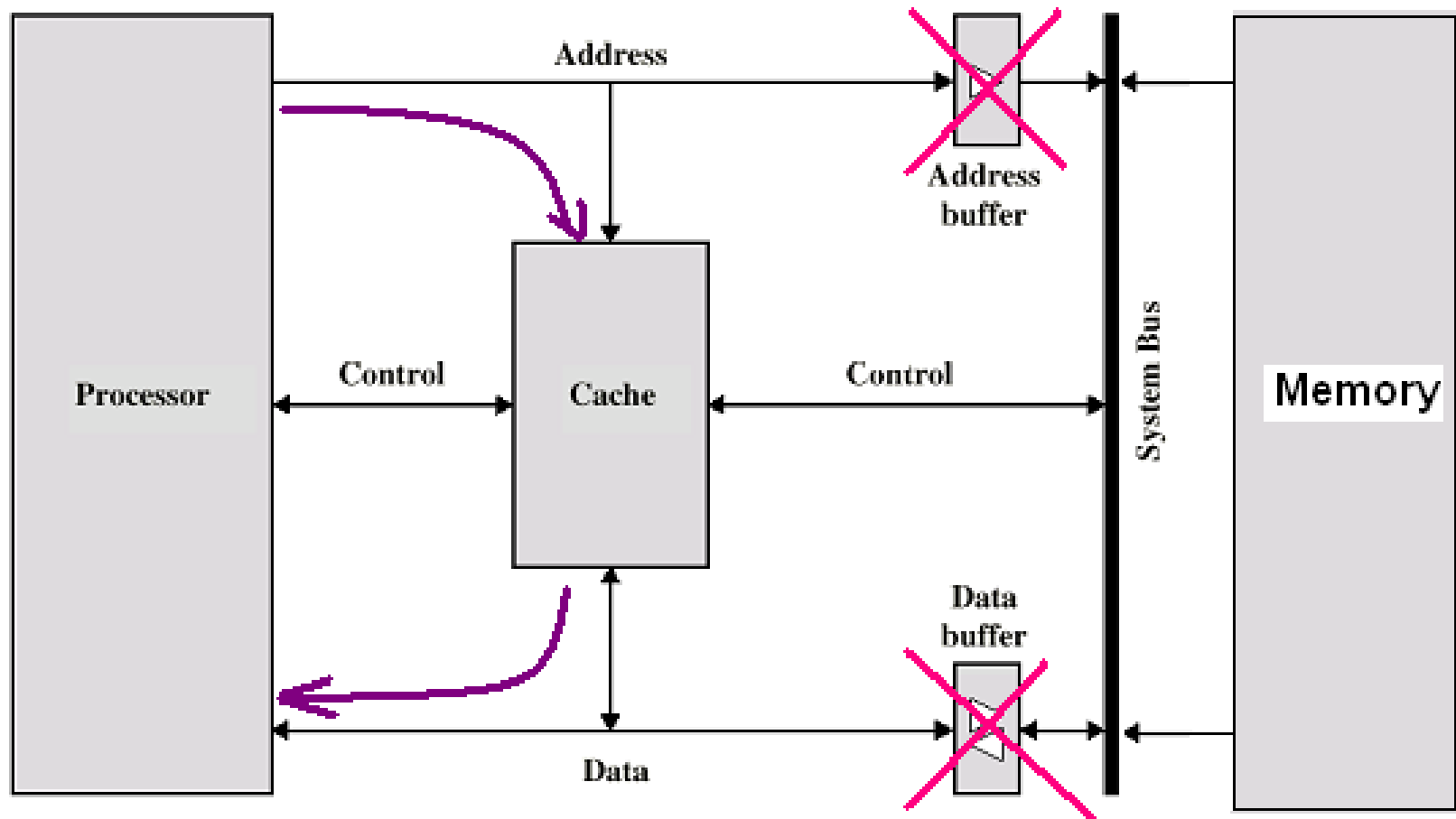


Performance of cache

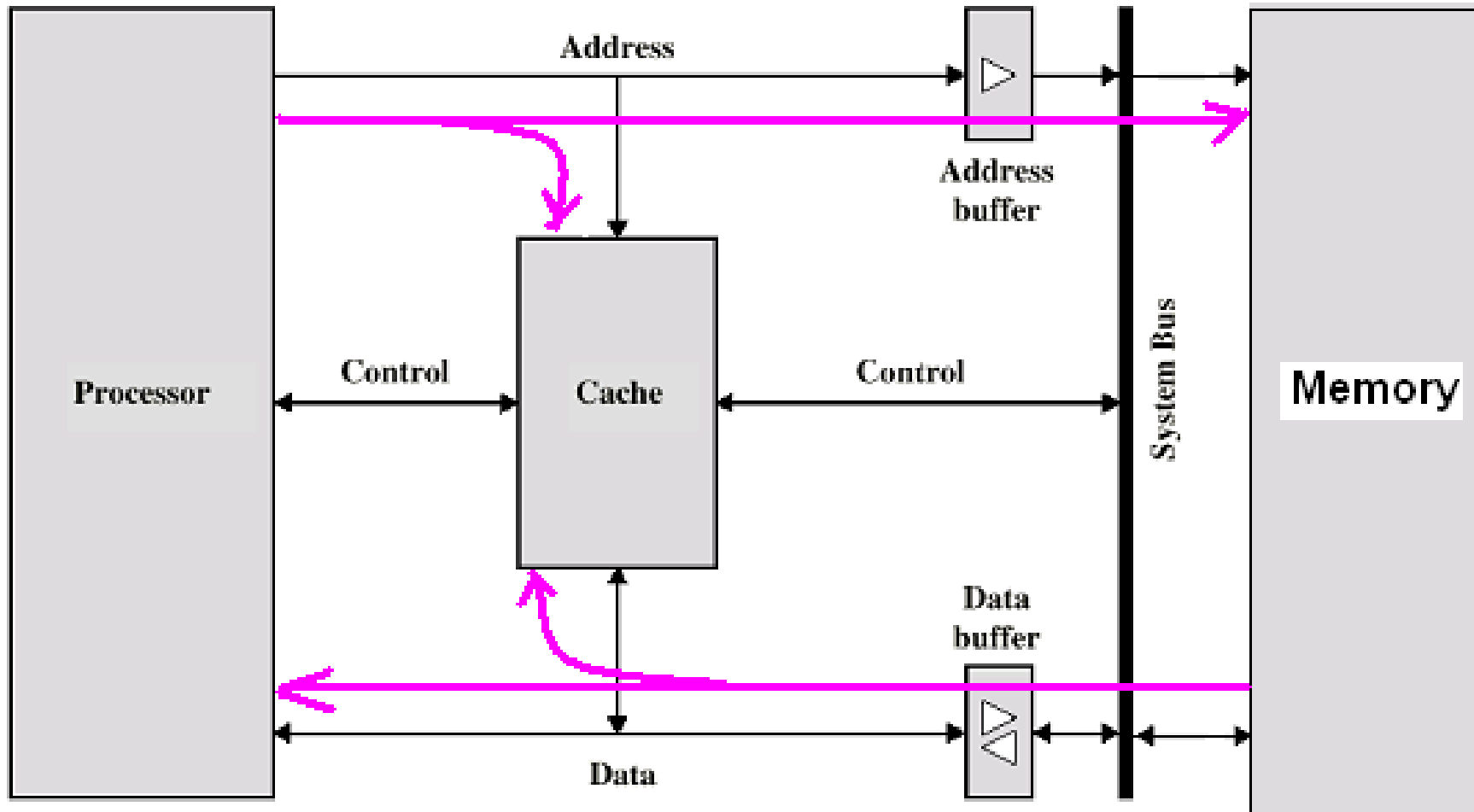


- Hit ratio : Number of Hits / total references to memory
- Hit
- Miss

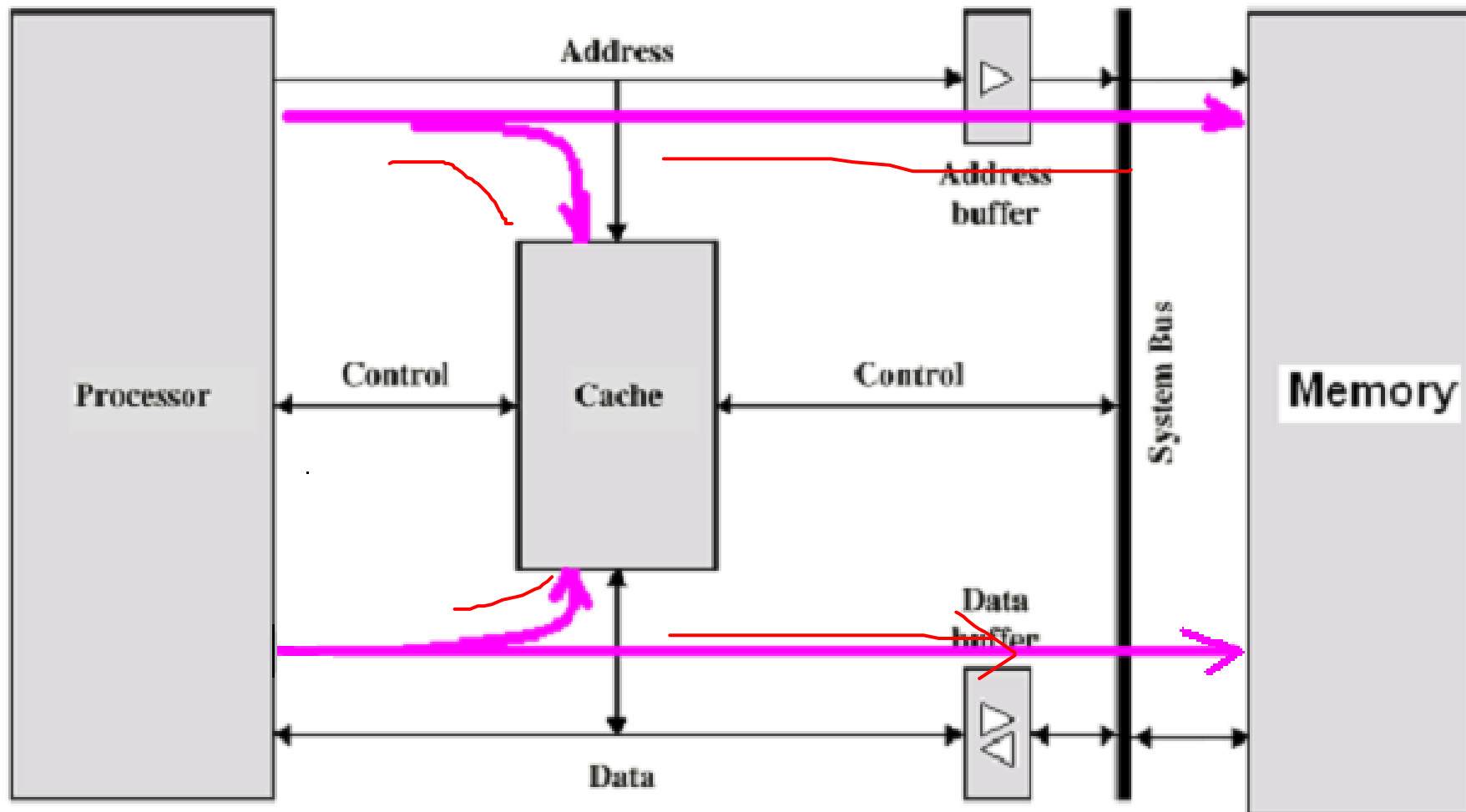
Read Hit



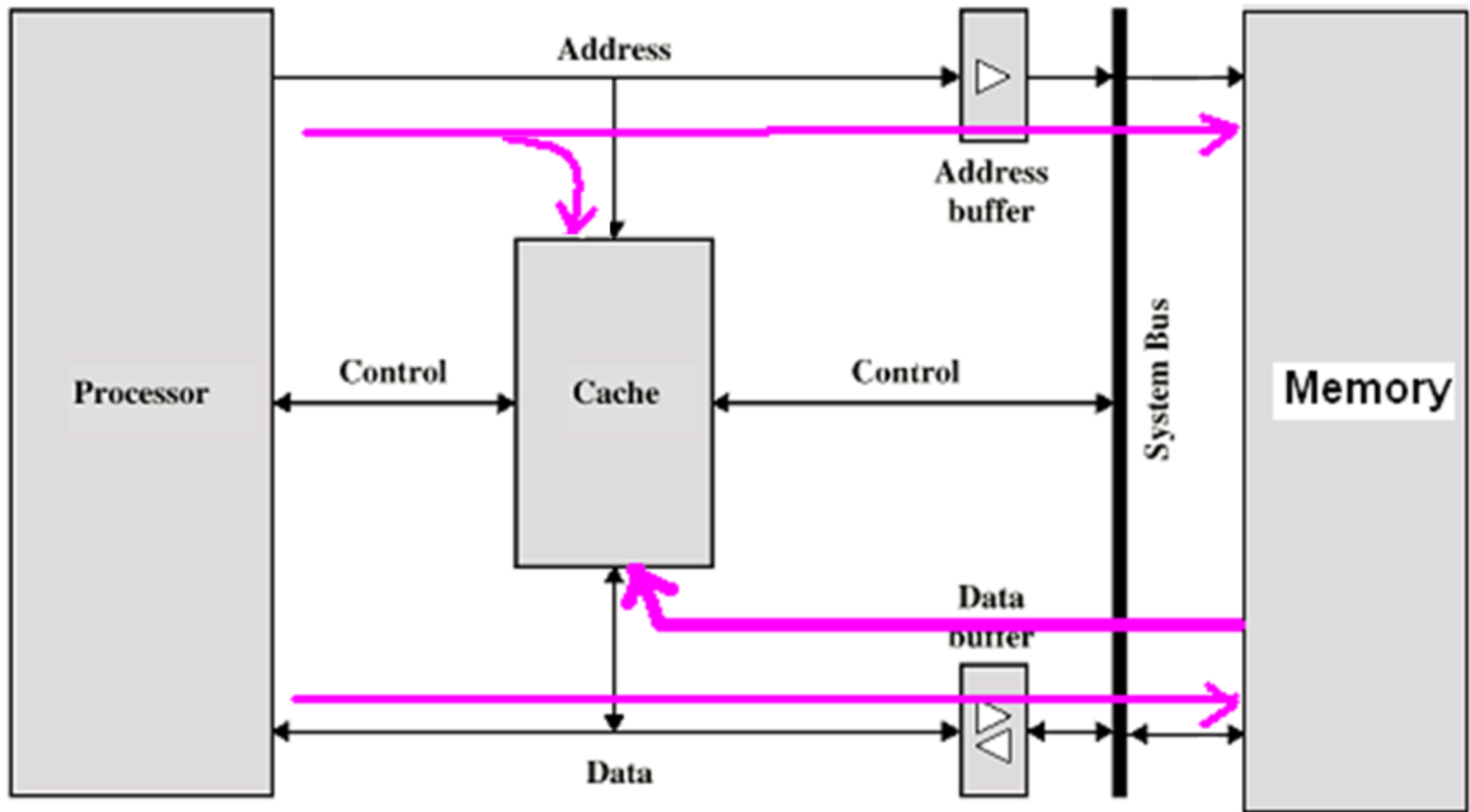
Read Miss



Write hit



Write miss



Mapping Function

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

Mapping Function

- 16 Bytes main memory
 - How many address bits are required?
- Memory block size is 4 bytes
- Cache of 8 Byte
 - How many cache lines?

Mapping Function

- 16 Bytes main memory
 - How many address bits are required?
- Memory block size is 4 bytes
- Cache of 8 Byte
 - How many cache lines?



4 bits

Mapping Function

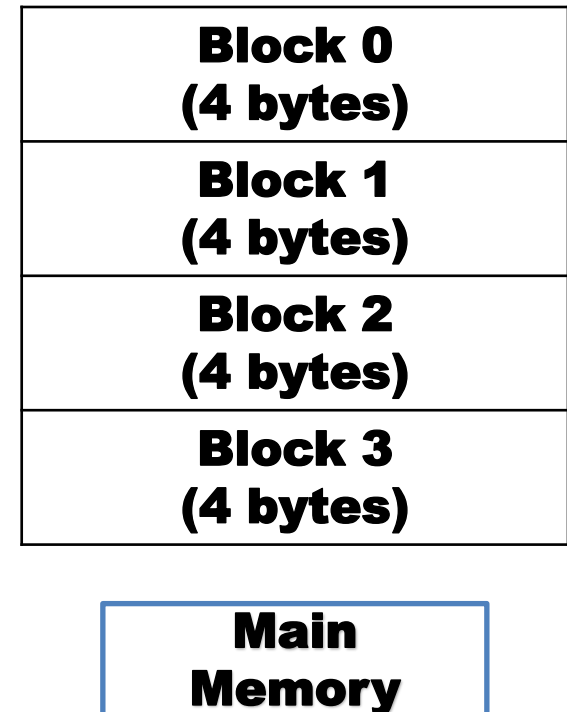
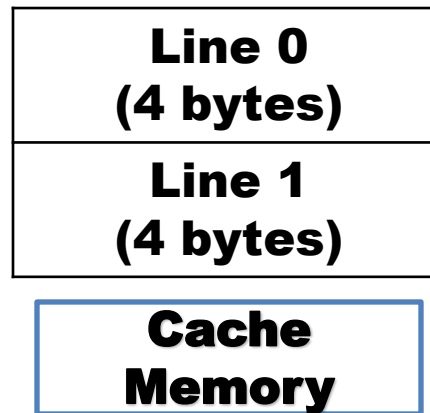
- 16 Bytes main memory
 - How many address bits are required?
- Memory block size is 4 bytes
- Cache of 8 Byte
 - How many cache lines?



**cache is 2 lines
(4 bytes per Line)**

Mapping Function

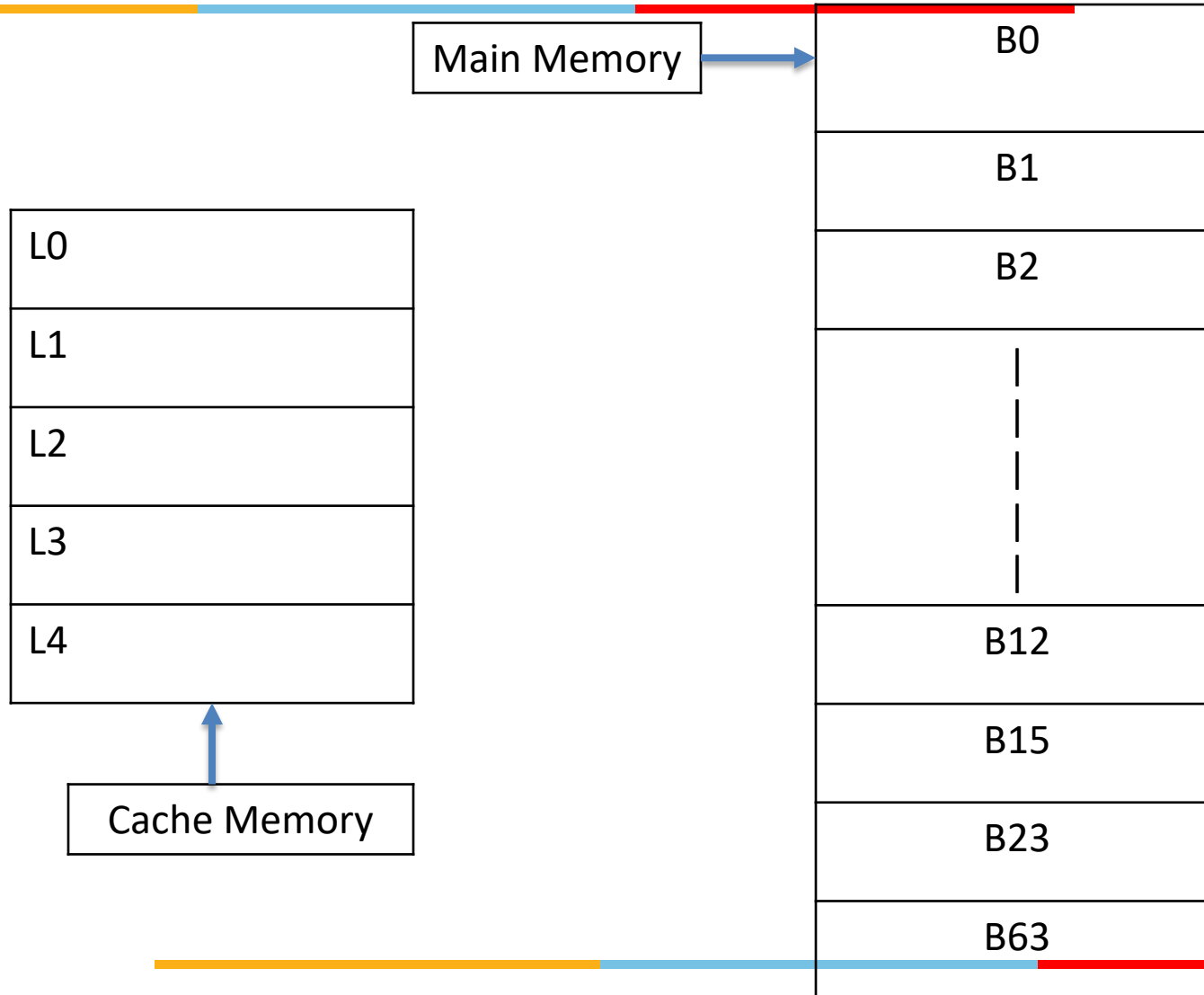
- 16 Bytes main memory
 - How many address bits are required?
- Memory block size is 4 bytes
- Cache of 8 Byte
 - How many cache lines?
 - cache is 2 lines (4 bytes per Line)



Direct Mapping

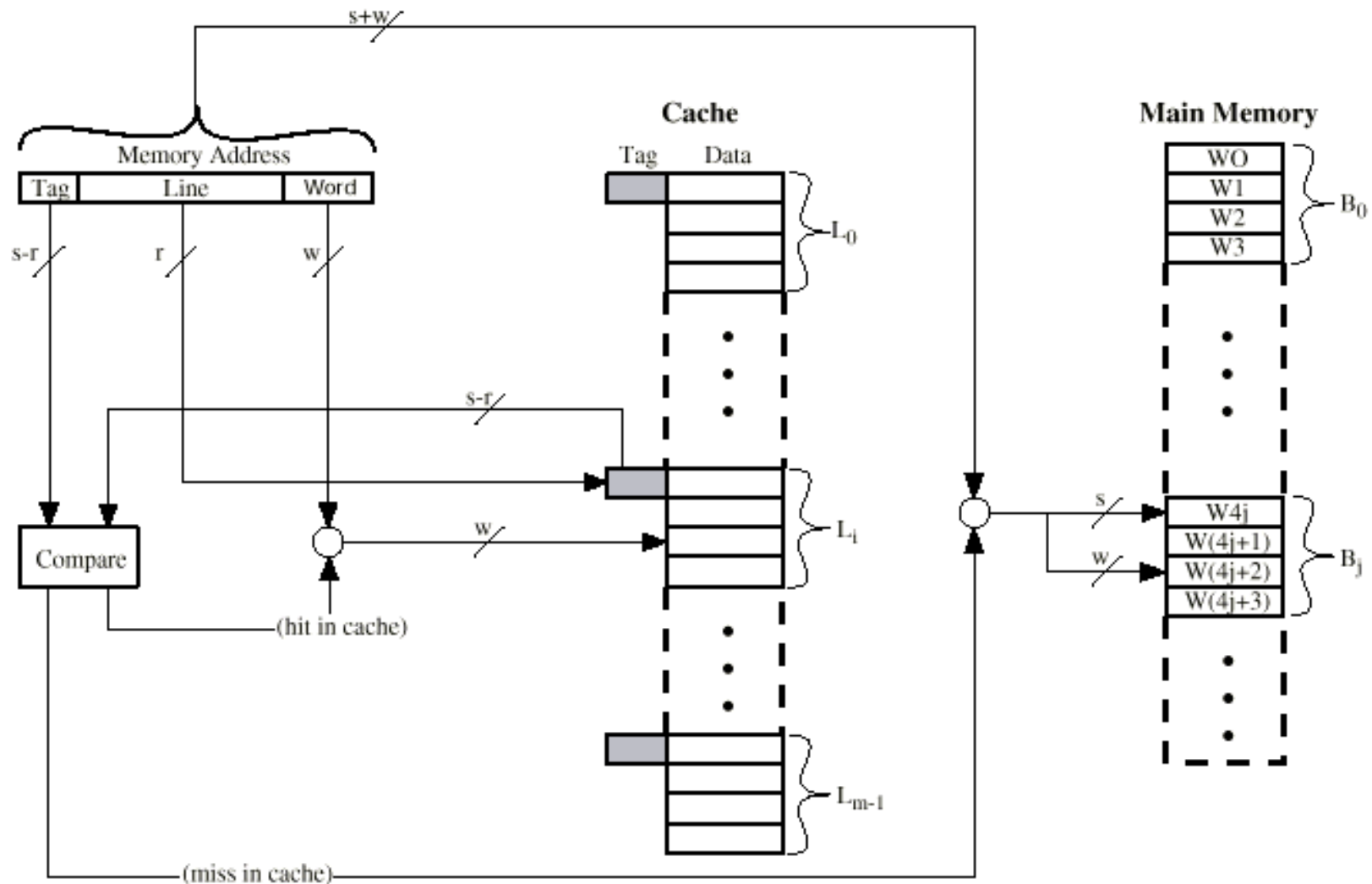
- Each block of main memory maps to only one cache line
 - i.e. if a block is in cache, it must be in one specific place
 - $i = j \text{ modulo } m$where i = cache line number
 j = main memory block no.
 m = no.of lines in the cache
- Address is split in three parts:
 - Tag
 - Line
 - Word

Direct Mapping Cache Organization





Direct Mapping Cache Organization



Direct mapping- Summary



Address length = $(s+w)$ bits

Number of addressable units = 2^{s+w} words or bytes

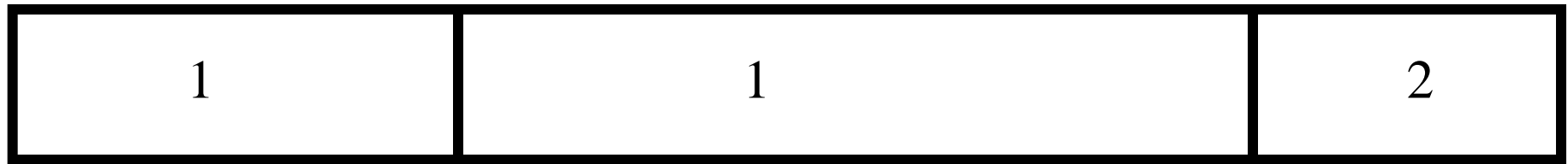
Block size = line size = 2^w words or bytes

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

Direct mapping Summary



Tag s-r

Line or Slot r

Word w

Direct Mapping 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

Given :

- Cache of 64kByte, Cache block of 4 bytes
- 16MBytes 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) Number of bits to identify a block
 - f) Tag, Line, Word
-

Solution

Given :

- Cache of 64kByte, Cache block of 4 bytes
- 16MBytes main memory

Find out

a) Number of bits required to address the memory

24 bits

b) Number of blocks in main memory

4M blocks

c) Number of cache lines

16K lines



Solution 1

Given :

- Cache of 64kByte, Cache block of 4 bytes
- 16MBytes main memory



Find out

d) Number of bits required to identify a word (byte) in a block?

2 bits

e) Number of bits required to identify block

22 bits

f) Tag, Line, Word

Tag $s-r$	Line r	Word w
8	14	2

Problem 2



Consider a machine with a byte addressable main memory of 2^{16} 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?

Solution 2



Consider a machine with a byte addressable main memory of 2^{16} 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?

TAG =8	LINE=5	WORD=3
--------	--------	--------

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

Solution 2



Consider a machine with a byte addressable main memory of 2^{16} 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?

Solution 2



Consider a machine with a byte addressable main memory of 2^{16} 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?

: Number of cache line 32

Block size : 8 bytes

Total bytes saved in cache = $32 \times 8 \text{ bytes} = 256 \text{ bytes}$ (excluding tag)

Tag bits saved : $32 \times 8 \text{ bits} = 256 \text{ bits} = 32 \text{ bytes}$

Total bytes saved in cache = $256 \text{ bytes} + 32 \text{ bytes} = 288 \text{ Bytes}$ (including tag)

Solution 2



Consider a machine with a byte addressable main memory of 2^{16} bytes and block size of 8 bytes. Assume that a direct mapped cache consisting of 32 lines is used with this machine.

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

Two or more blocks can be mapped to same cache line.

To distinguish between the blocks, tag bits are used.

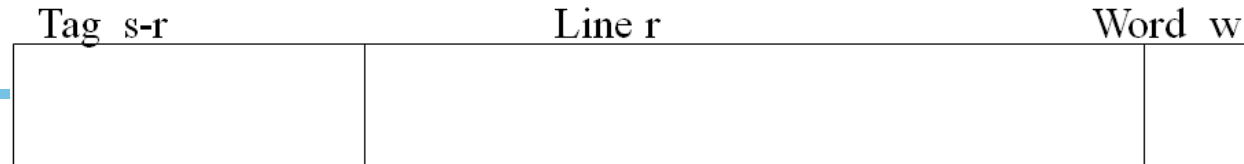
Problem 3



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?

32nd line.

Problem 4



The system uses a L1 cache with direct mapping and 32-bit address format is as follows:

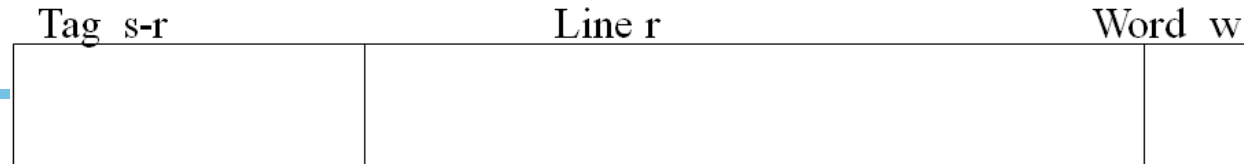
bits 0 - 3 = offset (word)

bits 4 - 14 = index bits (Line)

bits 15 - 31 = tag

- What is the size of cache line?
- How many Cache lines are there?
- How much space is required to store the tags for the L1 cache?
- What is the total Capacity of cache including tag storage?

Problem 4



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)

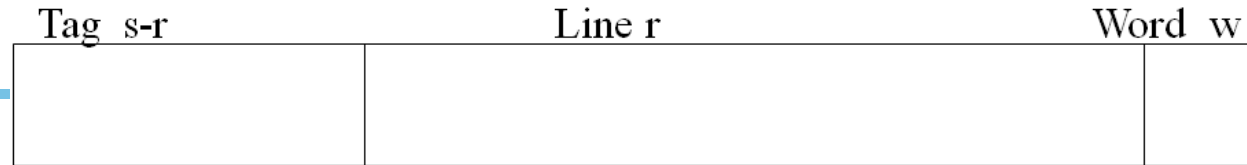
bits 15 - 31 = tag

- What is the size of cache line?
- How many Cache lines are there?
- How much space is required to store the tags for the L1 cache?
- What is the total Capacity of cache including tag storage?

a) Size of cache line

= size of the block = $2^4 = 16$ bytes

Problem 4



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?

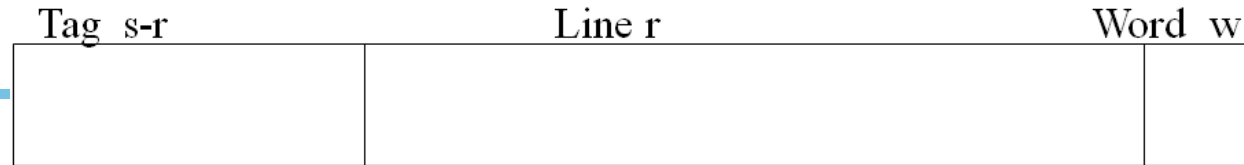
b) **How many Cache lines are there?**

b) No. of cache lines = $2^{11} = 2K$ cache lines

c) How much space is required to store the tags for the L1 cache?

d) What is the total Capacity of cache including tag storage?

Problem 4



The system uses a L1 cache with direct mapping and 32-bit address format is as follows:

bits 0 - 3 = offset (word)

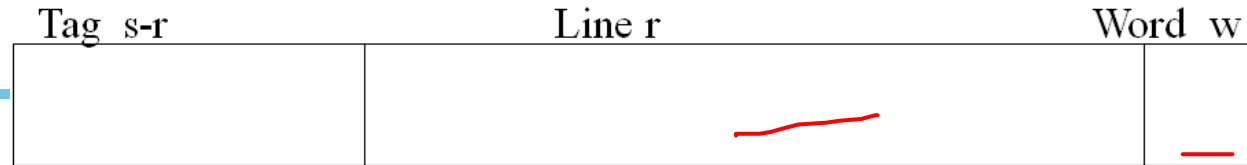
bits 4 - 14 = index bits (Line)

bits 15 - 31 = tag

- What is the size of cache line?
- How many Cache lines are there?
- How much space is required to store the tags for the L1 cache?
- What is the total Capacity of cache including tag storage?

Space for tag = No. cache lines * Tag length
= $2K * 17 \text{ bits} = 34 \text{ K bits}$

Problem 4



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 for the L1 cache?

d) What is the total Capacity of cache including tag storage?

Total capacity = 34Kbits + 32 Kbytes

Problem 5



- 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



- 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

Block 0 will be placed in line 0 $0\%2 \rightarrow 0$

Block 1 will be placed in line 1 $1\%2 \rightarrow 1$

Block 2 will be placed in line 2 $2\%2 \rightarrow 0$

- Find out hit ratio.

Problem 5



- 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



- 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.

$$4/13=30.76\%$$

Problem 6



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

Problem 6



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?

20 bits

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

Solution 6



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

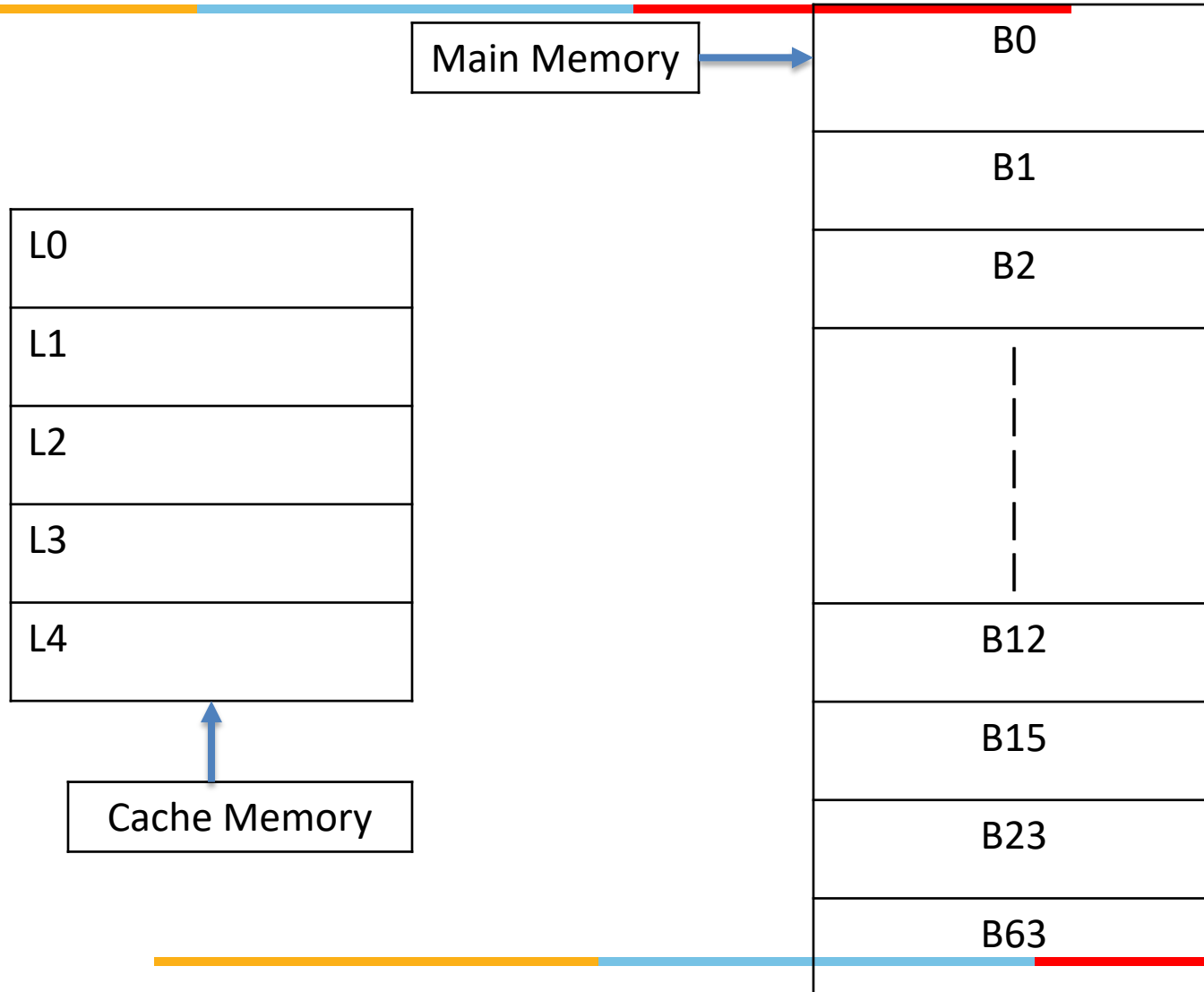
$$\begin{aligned}\text{Avg access time} &= \text{hit ratio} * \text{cache access} + \\ &\quad (1 - \text{hit ratio}) * (\text{cache access} + \text{memory access}) \\ &= .95 * 0.1 \text{ microsec} + .05 * (1 + 0.1) \text{ microsec} \\ &= .095 + .055 \text{ microsec} \\ &= 0.15 \text{ microseconds}\end{aligned}$$

Associative Mapping

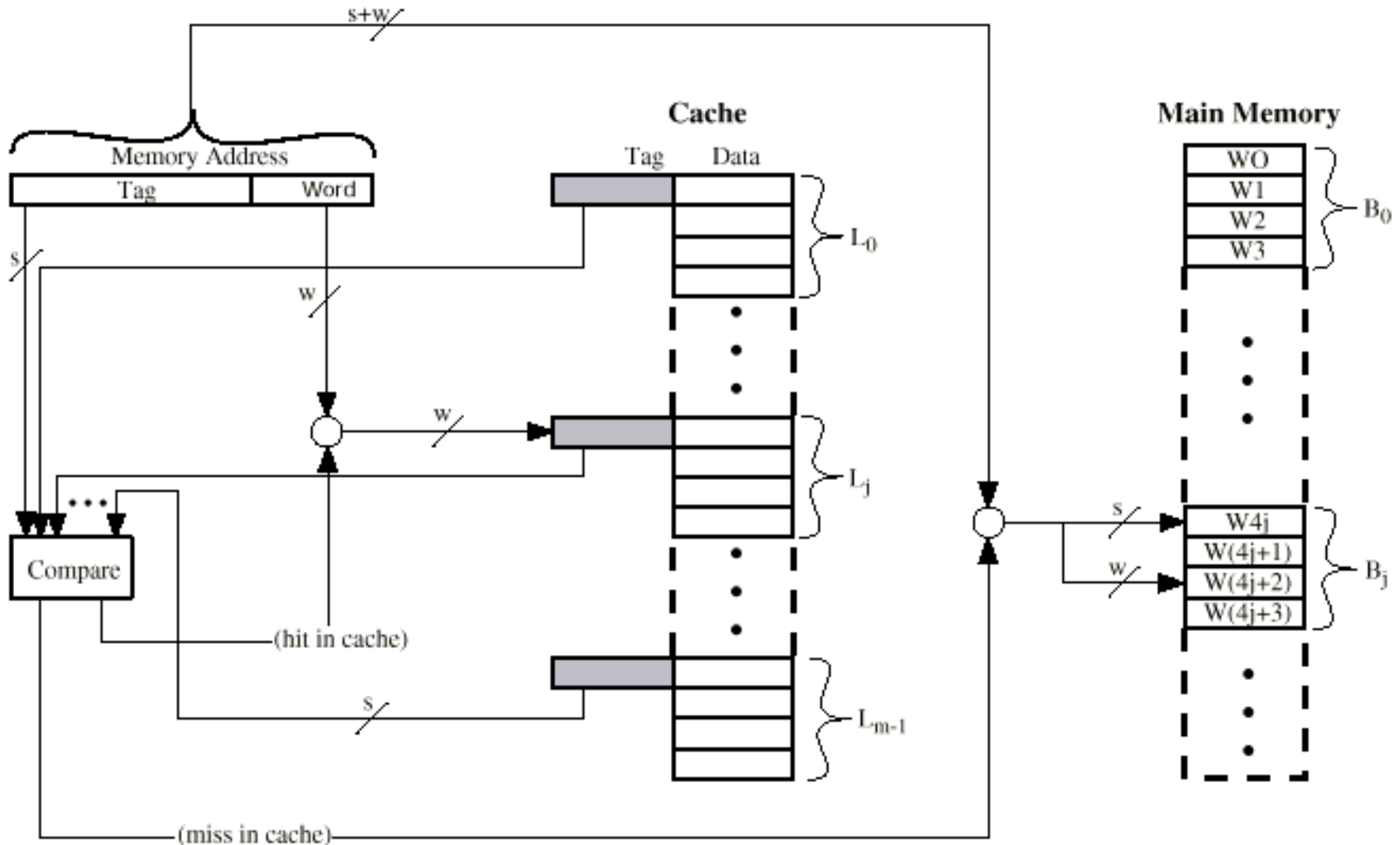


- A main memory block can load into any line of cache
- Memory address is interpreted as tag and word
- Tag uniquely identifies block of memory
- Every line's tag is examined for a match
- Cache searching gets expensive

Associative Mapping Cache Organization



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

Problem 7

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

Problem 7

Given :

- Cache of 128kByte, Cache block of 8 bytes
- 32 MBytes main memory

Find out

- a) Number of bits required to address the memory=25 bits
- 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

Problem 7

Given :

- Cache of 128kByte, Cache block of 8 bytes
- 32 MBytes main memory

Find out

- a) Number of bits required to address the memory=25 bits
- b) **Number of blocks in main memory**= 2^{22} blocks or
4M blocks
- c) Number of cache lines
- d) Number of bits required to identify a word (byte) in a block?
- e) Tag, Word

Problem 7

Given :

- Cache of 128kByte, Cache block of 8 bytes
- 32 MBytes main memory

Find out

- a) Number of bits required to address the memory=25 bits
- b) Number of blocks in main memory= 2^{22} blocks or
4M blocks
- c) Number of cache lines= 16k Lines
- d) Number of bits required to identify a word (byte) in a block?
- e) Tag, Word

Problem 7

Given :

- Cache of 128kByte, Cache block of 8 bytes
- 32 MBytes main memory

Find out

- a) Number of bits required to address the memory=25 bits
- b) Number of blocks in main memory= 2^{22} blocks or
4M blocks
- c) Number of cache lines= 16k Lines
- d) Number of bits required to identify a word (byte) in a
block?= 3bits
- e) Tag, Word

Problem 7

Given :

- Cache of 128kByte, Cache block of 8 bytes
- 32 MBytes main memory

Find out

- a) Number of bits required to address the memory=25 bits
- b) Number of blocks in main memory= 2^{22} blocks or
4M blocks
- c) Number of cache lines= 16k Lines
- d) Number of bits required to identify a word (byte) in a
block?= 3bits
- e) Tag, Word = 22 bits, 3bits

Problem 8

Cache of 64kByte, Cache block of 4 bytes and 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 = _____

Problem 8

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

Fill in the blanks:

Number of bits in main memory address = 24bits

Number of lines in the cache memory = 16K lines

Word bits = 2 bits

Tag bits = 22 bits

Problem 9



- 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.

Problem 9



- 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.
$$10/13 = 76.92\%$$